

Commissioned by:  
**NATIONAL RESEARCH FOUNDATION**  
PRIME MINISTER'S OFFICE  
SINGAPORE

A FOUNDATIONAL RESEARCH CAPABILITIES REPORT ON AI For Science 2025

A Foundational Research  
Capabilities Report on

# AI For Science

2025



## FRC TEAM LEADERS:



**Associate Professor  
Kedar Hippalgaonkar**  
Nanyang Technological University,  
Agency for Science Technology  
and Research



**Professor  
Yang Zhang**  
National University of Singapore

## TECHNICAL AND RESEARCH TEAM:



**Dr Eleonore Vissol-Gaudin**  
Nanyang Technological University



**Dr Beatrice Soh**  
Agency for Science  
Technology and Research



**Mr Andre Low**  
Nanyang Technological University



**Dr Ruiming Zhu**  
Nanyang Technological University



**Ms Cynthia Chew**  
Nanyang Technological University

## COPYRIGHT & DISCLAIMERS

This document was prepared for a study commissioned by the Foundational Research Capabilities Directorate of the National Research Foundation of Singapore (NRF).

The contents and opinions expressed here belong to the study team and may not represent the opinions of NRF. Please cite the report for non-commercial academic use; for other uses, please seek written permission from NRF. Citation: National Research Foundation, Singapore, 2025.

*A Foundational Research Capabilities Report on AI for Science 2025*

Neither NRF nor the study team shall be liable for any consequence resulting from the use of this report. © NRF 2025

## DESIGN & PRODUCTION

Redbean De Ptd Ltd

A Foundational Research  
Capabilities Report on

# AI For Science 2025

# TABLE OF CONTENTS

<b>EXECUTIVE SUMMARY</b>	<b>8</b>
<b>1. INTRODUCTION</b>	<b>10</b>
1.1. The Birth of AI for Science	13
1.2. Fast Forward to Modern Times – Where Do We Stand?	14
<b>2. SINGAPORE'S AI FOR SCIENCE INITIATIVE</b>	<b>15</b>
2.1. Global Efforts in AI for Science	16
2.2. Opportunity for Singapore: Scope and Goals	19
2.3. Grand Challenge –Developing an AI for Science Framework	20
2.4. Need of the Hour for the Singapore AI for Science Initiative	22
2.4.1. Data	22
2.4.2. Quantum (including hybrid compute)	22
2.4.3. High-performance computing (HPC)	23
2.4.4. AI algorithm development	23
2.4.5. Self-Driving Laboratories	23
2.4.6. AI agents and leveraging on large language models (LLMs)	24
2.4.7. Enhancing collaborations towards collective action	24
<b>3. STRATEGIC PLAN OF INITIATIVE DEVELOPMENT</b>	<b>25</b>
3.1. Foundation	26
3.1.1. Building peaks of basic AI research excellence	27
3.1.2. Integration with NAIS 2.0	27
3.1.3. Government support and policy frameworks	28
3.1.4. Potential benefits and global trends	29
3.2. Capacity Gaps	29
3.3. Potential for Industrial and Academic Alignment	30
3.4. Vision	30
3.5. Specific Policy Recommendations	31
<b>4. AI4SCI DOMAINS</b>	<b>32</b>
4.1. Biomedical and Health Sciences	33
4.1.1. AI-based protein and RNA structure modelling	33
4.1.2. RNA biology	34
4.1.3. Synthetic biology	35
4.1.4. Computer-aided drug design	36
4.1.5. Health and medicine	37
4.2. Advanced Materials and Sustainability	39
4.2.1. Future materials science and development	39
4.2.2. Quantum materials and quantum computing	40
4.2.3. Imaging and high-resolution microscopy	41
4.2.4. Advanced semiconductor technology development	41
4.2.4. Sustainability	42
4.3. Natural Sciences – Chemistry, Physics and Earth/Climate	43
4.3.1. Physics	43
4.3.2. Chemistry	44
4.3.3. Earth and climate	44
4.4. AI Frameworks – Mathematics, Theory and Model development	45
4.5. Finance	46
4.6. Education	47
4.7. Software and Security	48
4.8. Robotics	49
<b>5. RECOMMENDATIONS AND CONCLUSIONS</b>	<b>51</b>
5.1. Cross-Domain Learning and Collaboration	52
5.1.1. Knowledge integration and generalization	52
5.1.2. Data integration and standardization	53
5.1.3. Resource optimization	53
5.1.4. AI ethics and governance	53
5.2. Common Challenges for AI Algorithms Across Domains	54
5.2.1. Handling large and diverse datasets	54
5.2.2. Training AI models on limited data	55
5.2.3. Bias and data imbalance	55
5.2.4. Collaboration gap between domain experts and AI practitioners	55
5.2.5. AI model interpretability	56
5.2.6. Computational complexity and scalability	56
5.2.7. Ethical and privacy concerns	56
5.3. Conclusions	57
5.3.1. Importance of AI for Science	57
5.3.2. Mission of Singapore's AI4SCI Initiative	57
5.3.3. Opportunities and challenges of Singapore	58
5.3.4. Strong support is essential for the success of AI4SCI Initiative	58
<b>6. AI4SCI WORKSHOPS AND WHITEPAPERS</b>	<b>59</b>
<b>REFERENCES</b>	<b>61</b>

## APPENDICES

<b>Appendix i. Earth/Climate Sciences</b>	<b>66</b>
Executive Summary	67
Introduction	67
Background	68
Grand Challenges	68
AI Methods and Data - Challenges and Opportunities	72
Singapore's Role	72
Conclusions	72
References	73
<b>Appendix ii. Physics/Complexity</b>	<b>74</b>
Executive Summary	74
Introduction	74
Background	76
Grand Challenges	77
AI Methods and Data - Challenges and Opportunities	79
Singapore's Role	82
Conclusion	83
References	84
<b>Appendix iii. RNA Biology &amp; Therapeutics (AI4 RBT)</b>	<b>85</b>
Executive summary	85
Introduction	85
Background	86
Grand challenges	87
AI methods and data - challenges and opportunities	92
Concluding remarks	93
References	94
<b>Appendix iv. Biomedical Sciences</b>	<b>95</b>
Executive Summary	95
Introduction	95
Background	95
Grand Challenges	96
AI Methods and Data - Challenges and Opportunities	99
Singapore's Role	99
Conclusions	100

<b>Appendix v. Healthcare and Imaging</b>	<b>101</b>
Executive Summary	101
Introduction	101
Background	103
Grand Challenges	104
AI Methods and Data - Challenges and Opportunities	108
Singapore's Role	109
Conclusions	110
References	110
<b>Appendix vi. Genomics</b>	<b>111</b>
Executive Summary	111
Introduction	111
Background	112
Grand Challenges	112
AI Methods and Data - Challenges and Opportunities	116
Singapore's Role	117
Conclusion	117
References	118
<b>Appendix vii. Digital Phenotyping</b>	<b>119</b>
Executive Summary	119
Introduction	119
Background	121
Grand Challenges	121
Singapore's Role	126
Conclusion	127
References	127
<b>Appendix viii. Materials/Chemistry</b>	<b>129</b>
Executive Summary	128
Introduction	128
Background	129
Grand Challenges	130
Role of Singapore in AI for Materials Science and Chemistry	136
Conclusion	137
<b>Appendix ix. Chemical and Biomanufacturing</b>	<b>138</b>
Executive Summary	138
Background	138
Grand Challenges	139
Impact of AI Methods and Data - Challenges and Opportunities	141
Singapore's Role	143
Conclusion	143
References	143

<b>Appendix x. Financial Services</b>	<b>144</b>
Executive Summary	144
Introduction	144
Background	146
Grand Challenges	148
AI Methods and Data – Challenges and Opportunities	152
Singapore’s Role	155
Conclusions	156
References	157
<b>Appendix xi. Electronics/Semiconductors</b>	<b>158</b>
Executive Summary	158
Introduction	159
Background	160
Grand Challenges	161
Conclusion	164
Acknowledgments	165
References	165
<b>Appendix xii. Hybrid Quantum Computing</b>	<b>166</b>
Executive Summary	166
Introduction	166
Background	167
Grand Challenges	168
Singapore’s Role	173
Conclusions	173
References	174
<b>Appendix xiii. Sustainability</b>	<b>175</b>
Executive Summary	175
Introduction	175
Background	176
Grand Challenges	177
AI Methods and Data – Challenges and Opportunities	179
Singapore’s Role	181
Conclusions	182
References	183
<b>Appendix xiv. Education</b>	<b>184</b>
Executive Summary	184
Introduction	185
Background	186
Grand Challenges	187
AI Methods and Data - Challenges and Opportunities	190
Singapore’s Role	193
Conclusion	194
References	194
Annex	195

<b>Appendix xv. Science, Software and Security</b>	<b>197</b>
Executive Summary	197
Introduction	197
Background	199
Grand Challenges	200
AI Methods and Data – Challenges and Opportunities	205
Singapore’s Role	206
Conclusion	207
References	207
<b>Appendix xvi. Robotics</b>	<b>208</b>
Executive Summary	208
Introduction	208
Background	209
Grand Challenges	209
AI Methods and Data – Challenges and Opportunities	210
Singapore’s Role	212
Conclusions	213
References	213
<b>Appendix xvii. AI Methods and Mathematics</b>	<b>214</b>
Executive summary	214
Introduction	214
Grand Challenges	215
Conclusion	222
References	222
<b>Appendix xviii. Global R&amp;D efforts in AI for Science</b>	<b>223</b>
North America	223
United Kingdom	224
Asia-Pacific	224
European Union	224
South America	225
Africa	225
ANNEX A. Global R&D Efforts in AI for Science – North America	226
Annex B. Global R&D Efforts in AI for Science – UK	236
Annex C. Global R&D Efforts in AI for Science – Asia-Pacific	238
Annex D. Global R&D Efforts in AI for Science – European Union	242
Annex E. Global R&D Efforts in AI for Science – World infographics and Table of references	246





# EXECUTIVE SUMMARY

The AI for Science – AI4SCI - Initiative represents a strategic effort to harness the transformative potential of artificial intelligence (AI) in advancing scientific research within the context of Singapore's National AI strategy for Singapore, the public good and the world (1). The Initiative

aims to promote and position Singapore as a global leader in AI-driven scientific advancements, ensuring that the nation's scientific community is well-equipped to address contemporary challenges and drive innovation.

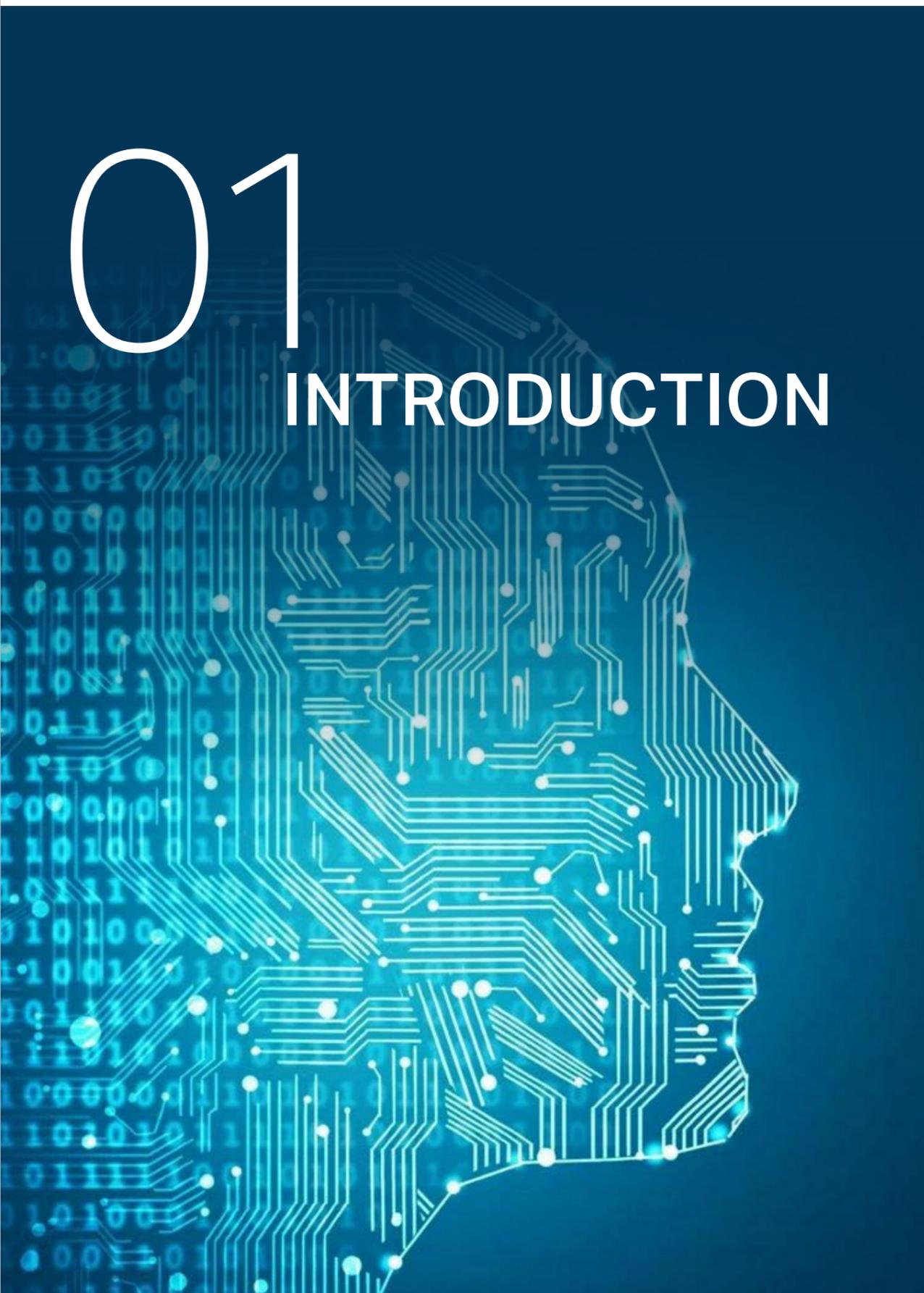
This Initiative is a collaborative endeavour led by National Research Foundation (NRF) and helmed by A/Prof Kedar Hippalgaonkar from Nanyang Technological University (NTU) and the Agency for Science Technology and Research (A\*STAR) and Prof Yang Zhang from the National University of Singapore (NUS). The primary objective is to explore AI's current impact and future potential in multiple research domains, including materials science and chemistry, life sciences, environmental science, physics, finance, software & security and computational sciences, especially<sup>1</sup>.

Through a series of workshops, comprehensive scoping studies, and active consultations with the scientific community; the Initiative aims to identify key opportunities and challenges in integrating AI into scientific research within Singapore.

The first workshop was held on February 29, 2024, which served as a lighthouse event bringing together researchers and relevant stakeholders from public and private sector institutions as well as government agencies to discuss and share insights on AI's integration into the scientific realm. Following this, the domain-specific workshops were conducted in a hybrid format (both in person and online) to maximize engagement and inclusivity. The community engagement ended with the flagship AI4Science and Nobel Turing Challenge Initiative Conference<sup>2</sup>, which was held on July 23-24, 2024. This included keynote speeches from renowned leaders in the fields of AI and scientific research, panel discussions on the integration of AI technologies in science, and breakout sessions focused on specific challenges such as knowledge extraction from scientific literature, representation of scientific knowledge, mathematical and methods' advancements in the pursuit of scientific discovery and the automation of experiments.

The findings and recommendations from these activities, together with domain experts' thoughts, are now compiled into this full Foundational Research Capabilities (FRC) study report on AI for Science, which will inform Singapore's Research, Innovation, and Enterprise (RIE) ecosystem<sup>3</sup>.

# 01 INTRODUCTION



Science provides a systematic and evidence-based understanding of the natural world, fostering technological advancements and medical breakthroughs, ultimately enhancing the well-being and progress of human civilization. Despite significant advancements in scientific research and developments in the past centuries, discovering new scientific knowledge persists as a time-consuming and resource-intensive undertaking. *Conventional avenues of scientific research and discovery often involve multiple steps including hypothesis generation, design of experiments, data collection, and result analysis, taking years or even decades to unveil meaningful insights. Furthermore, the escalating volume of data produced by modern experimental means has heightened the need for researchers to automate methods to discern patterns and extract insightful conclusions from their findings.*

The past decade has already witnessed exponential growth in the number of scientific breakthroughs enabled by artificial intelligence (AI)- assisted research, particularly in fields such as biomedicine<sup>4</sup>, materials<sup>5</sup> and climate science<sup>6</sup> relevant to Singapore's research plan (RIE2025<sup>7</sup> and beyond). This progress has been fueled by the availability of big data, the evolution of rapid and massively parallel computing powered by graphics processing units (GPUs), innovations in robotics for autonomous and high-throughput experimentation, as well as the development of cutting-edge deep learning protocols. *By enhancing and accelerating every stage of the scientific process, from hypothesis generation and experimental design to data analysis and model construction, AI holds the potential to completely redefine the landscape and paradigm of scientific discovery<sup>8</sup>.*

Nevertheless, full integration and optimal utilization of cutting-edge AI techniques in the scientific research process face critical challenges that demand thoughtful deliberation<sup>9</sup>. These challenges encompass the development of advanced AI algorithms capable of efficiently analyzing intricate and large datasets, constructing deep and accurate models for complex systems, and addressing ethical and privacy concerns. Special attention should be given in gathering and utilizing personal data, such as in clinical applications, which is typically more sensitive and distinct from other datasets obtained through natural experiments. There are also pressing needs for the establishment of common standards and digital platforms to facilitate the exchange and sharing of datasets and models among different systems and entities. Addressing these challenges will be pivotal in unlocking the full potential of AI as a transformative tool for advancing science.

Amidst the promising developments across various domains of AI applications, the Singapore AI for Science (AI4SCI) Initiative stands to benefit from sustained research investments in AI for Science. Here, "AI for Science" is defined as the general effort to utilize and extend the AI technology for solving scientific problems towards scientific discoveries and advancements. The term "AI for Science" is different from "Science of AI" (or "Research for AI"), which refers to the development and advancement of AI technology and algorithms themselves. It is also distinct from the more general term of "AI for X" which refers to the implementation of AI to solve practical problems and challenges in different domain sectors ('X'), known as "applied AI", as we want specifically to leverage AI advancements for enhancing basic science studies.

One ideal example of *AI for Science* is the 2024 Chemistry Nobel-prize winning algorithm, AlphaFold2<sup>4</sup>, where the Google DeepMind team utilized a deep neural network architecture to learn the evolutionary and structural patterns from the Protein Data Bank (PDB) library<sup>10</sup> (10) and help solve the problem of protein structure prediction. Although substantial problems in the field remain<sup>11</sup>, the advancement in AI-based protein structure predictions has brought about significant impact on structural biology and general life sciences. *Following this example, our goal of the Singapore AI4SCI Initiative is to unite "Science of AI" experts with scientific domain specialists to establish an "AI for Science" community, who will work collaboratively to identify, characterize, and ultimately solve the most critical and important problems in fundamental sciences, including mathematics, physics, chemistry, climate science, biomedicine, and advanced materials.*

## 1.1

# THE BIRTH OF AI FOR SCIENCE

The birth of AI for science is closely tied to the broader development of AI as a field of study, which can be traced back to the mid-20th century. The term "artificial intelligence" was coined during a seminal workshop at Dartmouth College in 1956, organized by John McCarthy, Marvin Minsky, Nathaniel Rochester, and Claude Shannon<sup>12</sup>. While the term was new, the idea of creating machines capable of intelligent behaviour had been explored in various forms earlier. Alan Turing, in his 1950 paper titled "Computing Machinery and Intelligence"<sup>13</sup>, proposed the Turing Test as a measure of a machine's ability to exhibit human-like intelligence.

Early AI research focused on symbolic reasoning and problem-solving, leading to the development of programs like the Logic Theorist by Allen Newell and Herbert Simon in 1956<sup>14</sup> to prove mathematical theorems. Over the years, AI research has evolved through different phases, including the enthusiasm of the 1950s and 1960s, the "AI winter" of reduced funding and interest in the 1970s and 1980s, and the resurgence and rapid progress seen from the 1990s onwards<sup>15</sup>. Key milestones in AI development include the creation of expert systems<sup>16</sup>, the advent of machine learning algorithms and breakthroughs in neural network research<sup>17</sup>, particularly with the development of deep learning techniques<sup>18</sup>. Today, AI is a multifaceted field with applications ranging from natural language processing and computer vision to robotics and autonomous systems<sup>8</sup>.

In the context of scientific discovery, the development of the field of "Automated Scientific Discovery" in the 1980s, with pioneering work by researchers like Pat

Langley and Herbert Simon<sup>19,20</sup>, brought the concept closer to reality. Their work involved creating programs that could generate hypotheses and theories based on empirical data, thus automating aspects of the scientific process.

The first computer simulation conducted by John von Neumann<sup>21</sup>, in collaboration with Enrico Fermi, John Pasta, Stanislaw Ulam, and Mary Tsingou, marked a significant milestone in computational physics and was instrumental in the way we should think about the role of AI in scientific discovery. Known as the Fermi-Pasta-Ulam-Tsingou (FPUT) problem<sup>22,23</sup>, this experiment in the early 1950s utilized MANIAC I, one of the earliest electronic computers. The team set out to study a highly nonlinear system of springs and masses, simulating a one-dimensional crystal lattice, expecting to observe thermalization, where energy would evenly distribute among the system's modes. However, the results defied expectations, revealing that the system did not reach thermal equilibrium. Instead, energy remained localized in a few modes, challenging the then-prevailing understanding in statistical mechanics, in a system that was not solvable by purely theory.

This surprising outcome underscored the potential of computational models in scientific inquiry to simulate complex systems, especially when contextualized for understanding the current integration of self-driving labs and AI in scientific research. The transition from the FPUT experiment to today's technology-driven research landscape reflects a significant evolution in the methods and speed of scientific investigation.

## 1.2

# FAST FORWARD TO MODERN TIMES – WHERE DO WE STAND?

In the era of the FPUT experiment, computational power was a novel resource, enabling scientists to simulate and study complex systems in ways that were previously impossible. In the 1980s, the development of expert systems marked AI's foray into mimicking human expert decision-making, significantly impacting fields such as medicine and geology. This laid the groundwork for AI's application in science, which was further propelled by IBM's Deep Blue defeating world chess champion Garry Kasparov in 1997<sup>24</sup>, showcasing AI's potential in complex problem-solving.

The advent of big data in the 2000s and the computational power to process it expanded AI's role in scientific research. The deep learning revolution<sup>18</sup>, ignited by the success of deep neural networks in the 2012 ImageNet competition, opened up new avenues for AI across various scientific disciplines. In 2016, DeepMind's AlphaGo's victory over Lee Sedol in Go game demonstrated AI's advanced problem-solving capabilities<sup>25</sup>, inspiring applications in science, such as in predicting protein structures with AlphaFold in 2018<sup>26</sup>, and following up with AlphaFold2 in 2020<sup>4</sup>. Throughout the 2020s, AI-driven discoveries have been successfully demonstrated in

various fields including materials science, drug discovery, climate science, astronomy, physics and complex systems analysis.

Today, the landscape of scientific research has been revolutionized by the advancements of High Performance Computing (HPC) combined with AI and self-driving labs, providing vast and rich datasets in various domains. Indeed, AI can process vast amounts of data rapidly, uncovering patterns and insights that would be impossible for humans to discern within a reasonable timeframe. This synergy between data-driven labs, simulations and AI mirrors the initial impact of the MANIAC I computer in the FPUT experiment, but on a much grander scale. This represents an opportunity to continue the trend initiated by pioneers like von Neumann, Turing, Simon and others: leveraging cutting-edge data-driven technology to delve deeper into nature's mysteries, thus defining the era of *AI for Science*. AI's capability to compose axioms, often in ways surprising to humans, would be grounded in the meticulous efforts of scientists and engineers who formulate and test these foundational principles rigorously by the scientific method. In this context, AI acts as a collaborative partner to human scientists and engineers.



## 2.1

# GLOBAL CONTEXT: AI FOR SCIENCE EFFORTS AROUND THE WORLD

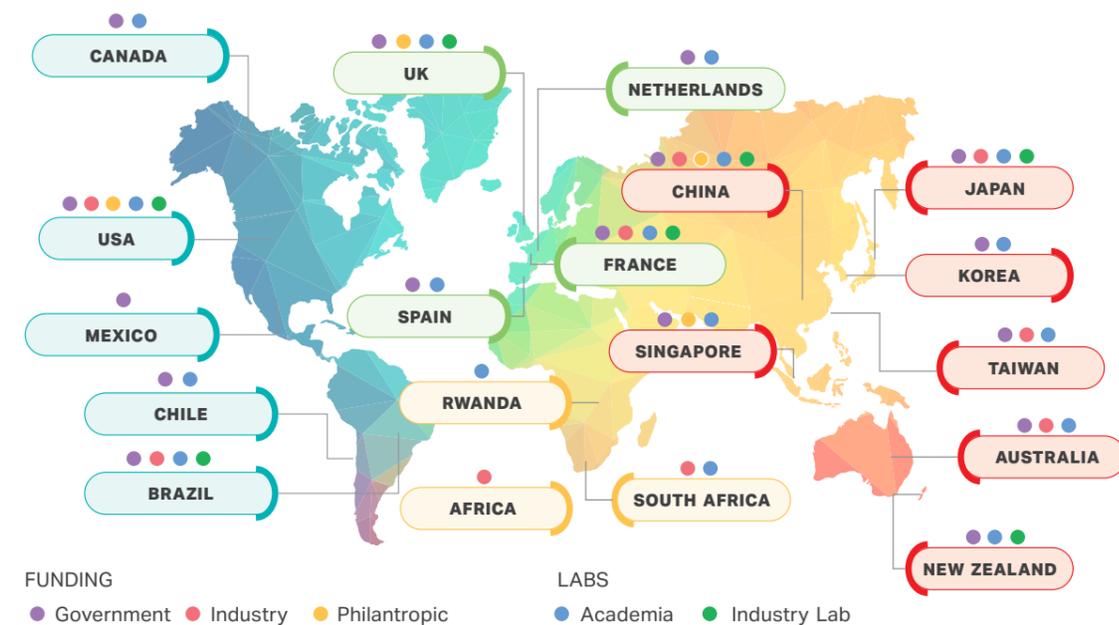
Global efforts in AI for Science encompass a multitude of initiatives, projects, and collaborations aimed at leveraging AI to enhance scientific research and discovery on an international scale. These efforts involve organizations, governments, research institutions, and industry players from around the globe, with examples including Horizon Europe (AIEOSC, EuroHPC), the United States National AI Initiatives, and global AI Challenges and Competitions (such as the Data Science Bowl and Kaggle competitions). Several large technology companies have also driven efforts in AI for Science, such as Microsoft Research with its AI4SCI Initiative aiming to develop AI tools for scientific discovery, and IBM Research focusing on AI-enabled scientific discovery through foundation models and multi-cloud computing.

Furthermore, the AI for Science community actively promotes collaboration among researchers by organizing workshops and investigating AI's transformative impact across disciplines such as Chemistry, Biology, and Physics. These efforts extend both within individual fields and through participation in prominent international AI research conferences, including NeurIPS, ICML, ICLR, AAAI and ECCV. Other global initiatives, including our AI4SCI Initiative in Singapore, are well-positioned to drive a transformative shift in scientific progress by developing emulators capable of rapidly iterating and predicting complex scientific phenomena.

A detailed, but non-exhaustive table of the international and national initiatives in AI for Science, along with references to the relevant documentation, can be found in Annexes A, B, C, D and E, Appendix xxviii – Global R&D Efforts In AI for Science.

### NORTH AND SOUTH AMERICA

### ASIA - PACIFIC



AI for Science Efforts infographics. Details in Annex E.

### NORTH AND SOUTH AMERICA

#### USA

- **Government Funding (Washington all)**  
DOE AI 4 Science Initiative  
NIH Strategic Plan for Data Science  
NIH Common Fund's Bridge to AI program(Bridge2AI)
- **Industry Funding**  
Toyota Research Institute (TRI)  
Microsoft Research  
Google  
IBM Research
- **Philanthropic Funding**  
Chang-Zuckerberg Foundation  
Bill and Melinda Gates Foundation  
Margot and Tom Pritzker Foundation  
Eric and Wendy Schmidt
- **Academic Labs / Main Funding Source / Year Of**  
MIT / NSF / 2020, 2024 (Cambridge, MA)Illinois-Urbana Champaign / NSF / 2023 (Illinois)  
Columbia / NSF / 2023 (New York City)  
UC Berkeley / DOE / 2023 (California)  
Uni. Chicago DSI / NSF / 2022 (Chicago)  
Caltech / AWS / 2018 (California)  
Cornell / TRI / 2021 (New York State)  
Northwestern – AMDD / TRI / 2023 (Illinois)
- **Industry Lab**  
Microsoft Research – AI4Science (Washington State)  
Google AI: Science AI, Quantum AI (California)  
Google Deepmind (California)  
InSitro (San Francisco)

#### CANADA

- **Government Funding**  
Quebec Research and Innovation Investment Strategy (SQRI2) - (Montreal)  
Ontario's Critical Technology Initiatives program (CTI) - (Ottawa)  
Canada First Research Excellence Fund (CFREF) - (Ottawa)
- **Academic Labs/Main Funding Source / Year Of**  
Vector Institute for AI / CTI / 2023 (Toronto)  
University of Toronto - Acceleration Consortium / CFREF / 2023 (Toronto)  
Mila Quebec / SQRI2 / 2023 (Montreal)

#### BRAZIL

- **Government Funding**  
São Paulo Research Foundation (Sao Paulo)
- **Industry Funding**  
IBM - C4AI (Sao Paulo)  
Petrobras - Litcomp-IA (Rio de Janeiro)
- **Academic Labs / Main Funding Source / Year Of**  
Uni. of Sao Paulo, C4AI / IBM + FAPESP / 2020 (Sao Paulo)  
BCRP – Litcomp-IA Lab / Petrobras, FACC / 2022 (Rio de Janeiro)  
UFRJ-Coppe / FAPESP, CGI.br, MCTI / 2024 (Rio de Janeiro)

#### INDUSTRY LAB

- **Industry Lab**  
IBM-Research (Sao Paulo, Rio de Janeiro)

#### MEXICO

- **Government Funding**  
Agencia Nacional de Inteligencia Artificial (ANIA) (Mexico City)

#### CHILE

- **Government Funding**  
Agencia Nacional de Investigacion y Desarrollo (ANID) (Santiago de Chile)
- **Academic Labs / Main Funding Source / Year Of**  
National Center for Artificial Intelligence Research (CENIA) / ANID ++ / 2022 (Santiago de Chile)

#### UK

- **Government Funding (London - all)**  
National AI Strategy  
UK Quantum Strategy  
UKRI AI and Data Science for Engineering, Health and Government  
AI Life Science Accelerator Mission  
Industry Funding  
Google Deepmind  
Microsoft Research
- **Philanthropic Funding**  
Eric and Wendy Schmidt
- **Academic Labs / Main Funding Source / Year Of**  
Oxford AI4Science Lab / UKRI, NASA, ESA, US Department of Energy, UK Space Agency / 2022 (Oxford)  
Cambridge Accelerate Science / Schmidt Future / 2020 (Cambridge)  
Imperial College London – I-X Centre for AI in Science / Schmidt Future / 2022 (London)  
Uni. of Southampton (Lead) / EPSRC / 2024 (Southampton)
- **Industry Lab**  
Google Deepmind (London)  
Microsoft Research - AI4science (Cambridge)

#### FRANCE

- **Government Funding**  
ANR - National AI Research Program (PNRIA) (Paris)
- **Industry Funding**  
Sanofi (Paris)  
Thales (Paris)
- **Academic Labs / Main Funding Source / Year Of**  
CNRS – AISSAI / ANR-PNRIA / 2024  
Uni. Côte d'Azur - 3IA Côte d'Azur / ANR-PNRIA / 2022 (Nice)  
Uni. de Toulouse - 3AI ANITI / ANR-PNRIA / 2022 (Toulouse)  
INRIA LaborIA / ANR-PNRIA / 2021 (Paris)

#### INDUSTRY LAB

- **Industry Lab**  
Data Innovation Lab at EDF R&D (Paris)  
Bioptimus (Paris)

#### NETHERLANDS

- **Government Funding (Amsterdam)**  
Dutch Research Council (DRC)
- **Academic Labs / Main Funding Source / Year Of**  
Uni. of Amsterdam AI4Science Lab / NA / 2020 (Amsterdam)  
Uni. of Leyden Vision4AI (Lead) / EU Horizon 2020+ / 2020 (Leyden)

#### SPAIN

- **Government Funding**  
Ministro de Ciencia e Innovacion MCI (Madrid)  
Tecnio Catalonia – ACCIO (Barcelona)
- **Academic Labs / Main Funding Source / Year Of**  
Agencia Estatal Consejo Superior de Investigaciones Científicas – AI4EOSC (Lead) / Horizon Europe / 2022 (Madrid)  
AI Research Institute IIIA-CSIC / ACCIO, EU, MCI / 2020 (Barcelona)

## AFRICA

- **International Industry Funding:**  
AI Africa Consortium - Cirrus Foundry fund (incubator – up to 7.5M grants)

## SOUTH AFRICA

- **Industry Funding:**  
Google Deepmind (4.5M USD) – AI4science Master's Program  
Cortex Logic
- **Academic Labs / Main Funding Source / Year Of**  
Cirrus Foundry / cortexlogic / 2019  
Wits University – Cirrusai (lead) / Cirrus Foundry fund / 2019

## RWANDA

- **Academic Labs / Main Funding Source / Year Of**  
CAIR / United Nations / 2023

## ASIA - PACIFIC

### JAPAN

- **Government Funding (Tokyo - all):**  
Ministry of Education, Culture, Sports, Science and Technology (MEXT) - KAKENHI  
Japan Society for the Promotion of Science (JSPS) - KAKENHI  
Ministry of Economy, Trade and Industry (METI)
- **Industry Funding:**  
Microsoft, Microsoft Research
- **Academic Labs / Main Funding Source / Year Of**  
OIST - Correspondence and Fusion of AI and Brain Science / KAKENHI / 2016-2023 (Okinawa)  
UTokyo - Next Generation AI Research Center / ? / 2023 (Tokyo)  
RIKEN - AIP / MEXT / 2016  
RIKEN TRIP - AI for Science Platform / MEXT / 2023  
GENIAC / METI / 2024
- **Industry Lab**  
Microsoft Research Asia lab / Microsoft / 2024

### CHINA

- **Government Funding (Beijing – all):**  
Ministry of Science and Technology (MOST)  
National Natural Science Foundation (NSFC)  
Chinese Academy of Science (CAS)
- **Industry Funding:**  
Microsoft Research  
Alibaba
- **Philanthropic Funding (Shanghai):**  
Tianqiao and Chrissy Chen Institute (TCCI)
- **Academic Labs / Main Funding Source / Year Of**  
AI for Science Institute / MOST,NSFC / 2021 (Beijing)  
Zhejiang lab / MOST,NSFC, Alibaba / 2017 (Hangzhou)  
Chen Frontier lab - AI + Brain Science / TCCI / 2021 (Shanghai)
- **Industry Lab**  
Microsoft Research - AI4Science (Beijing)  
Xtalpi (Shanghai)  
Insilico (Hong Kong)  
Wanhua Chemical (Yantai)  
CATL (Hong Kong)

## TAIWAN

- **Government Funding:**  
National Science and Technology Council (NSTC)  
Ministry of Science and Technology (MOST)
- **Industry Funding:**  
NVIDIA
- **Academic Labs / Main Funding Source / Year Of**  
AI Labs / MOST, private / 2017  
Taiwan AI Center of Excellence / NSTC / 2023  
AI Research and Development Center / NVIDIA / 2023

## KOREA

- **Government Funding:**  
Ministry of Science and ICT (MSIT)
- **Academic Labs / Main Funding Source / Year Of**  
AI Research Hub Project / MSIT, private / 2024  
CSRC KIST - Smart Lab / KAIST, UNIST / 2020  
Global AI Frontier lab (NYU-KAIST) / MSIT, NSF / 2024

## ASIA - PACIFIC

### AUSTRALIA

- **Government Funding (Adelaide – all):**  
Department of Industry, Innovation and Science (DIIS)  
Department of Health and Aged Care - MRFF
- **Industry Funding:**  
Google Australia
- **Academic Labs / Main Funding Source / Year Of (Melbourne - all)**  
CSIRO Data61 - Mixed Reality lab / DIIS / 2019  
CSIRO Data61 - AI for Missions Program  
CSIRO Data61 - Google Partnership / Google / 2021  
Australian Alliance for Secure Genomics and AI in Rare Disease / MRFF / 2024

### NEW ZEALAND

- **Government Funding (Auckland):**  
Ministry of Business, Innovation and Employment (MBIE)
- **Academic Labs / Main Funding Source / Year Of**  
AUT - Center for AI Research (CAIR) / MBIE / 2000 (Auckland)
- **Industry Lab**  
Litmaps

## SINGAPORE

- **Government Funding:**  
National Science Foundation (NSF)  
Agency for Science, Technology and Research (A\*STAR)
- **Philanthropic Funding:**  
Schmidt Future - Eric and Wendy Schmidt AI in Science Fellowships
- **Academic Labs / Main Funding Source / Year Of**  
Center for Frontier AI Research / A\*STAR / 2022

## 2.2

# OPPORTUNITY FOR SINGAPORE: SCOPE AND GOALS

Singapore benefits from unique research expertise and state-of-the-art facilities across its Institutes of Higher Learning - National University of Singapore (NUS); Nanyang Technological University (NTU); Singapore University of Technology and Design (SUTD); Singapore Management University (SMU) - as well as its national laboratories within the Agency for Science, Technology and Research (A\*STAR). The Singapore AI4SCI Initiative will provide a national platform to address key national challenges by leveraging on institutional research strengths, along with enabling collaborations with key industries and international partners. A measure of the Initiative's success is a significant AI-driven acceleration of fundamental breakthroughs in specific scientific domains, carefully selected based on the unique expertise of the current Singapore-based laboratories. This will put Singapore at the forefront of global research and development in interdisciplinary science and AI.

We posit that the time is ripe for a national investment into AI for Science, and Singapore is well-placed, with the right mix of resources and expertise to contribute to AI-driven discovery. The *AI4SCI* Initiative should therefore aim to:

- Unite AI experts and leading researchers across various basic science disciplines, fostering and supporting collaborative efforts between these communities within Singapore.
- Establish Singapore as a global hub for cutting-edge research in developing methods and algorithms that address fundamental challenges in the basic sciences, with a particular emphasis on interdisciplinary scientific fields.
- Offer infrastructure and hardware computing resources, including storage, CPUs, and GPUs, to support researchers participating in the Initiative, enabling them to conduct outstanding and innovative AI for Science research.
- Nurture a new generation of interdisciplinary specialists skilled in both AI technologies and scientific fields, driving innovation and advancement within their respective domains.
- Foster a community of AI-literate scientists who can apply advanced computational methods to complex research challenges and bridge the gap between AI and various scientific disciplines.

## 2.3

# GRAND CHALLENGE –DEVELOPING AN AI FOR SCIENCE FRAMEWORK

The hierarchy of AI model development for scientific discovery entails progressively advanced levels of learning and integration of scientific knowledge<sup>27</sup>. We envision a domain-agnostic framework that can be demonstrated in areas of strategic relevance to Singapore. This framework could be structured across five levels, with levels 1-2

primarily leveraging scientific research studies, levels 3-4 necessitating foundational work in the science of AI, and level 5 representing a convergence of efforts from both AI and scientific disciplines. To illustrate this, we present non-exhaustive examples across biomedical, physics, chemistry, and materials sciences:

## Examples Across Domains

### LEVEL 1

**Feature Importance and Identification** - Utilizing statistical analysis techniques to identify and assess the importance of features and recognize patterns within datasets.

- Healthcare: distinguishing benign from malignant tumours in radiology scans.
- Biomedical Sciences: Identification of significant gene expressions linked to diseases.
- Materials Science: Pattern recognition and process optimization in materials informatics to identify best performance characteristics of catalysts

### LEVEL 2

**Transfer Learning to Bridge Fields** - Reusing and adapting AI models across different domains to leverage existing knowledge and data.

- Biomedical Sciences: Predict disease susceptibility based on data from different populations in Genomics
- Chemistry: Using knowledge from known chemical reactions to predict outcomes in new, untested reactions

### LEVEL 3

**Integrating Physical Principles into Data-driven Models** - Embedding established physical laws and equations within AI models to align predictions with known scientific principles.

- Physics: Incorporating fluid dynamics equations into climate prediction models.
- Chemistry: Using chemical laws to forecast the results of reactions.
- Materials Science: Physics Inspired Neural Networks that build in partial differential equation form of solutions.

### LEVEL 4

**Learning Interpretations and Emergent Knowledge from Data** - Developing AI models that can infer the underlying physical laws from data without explicit programming.

- Biomedical Science: Deriving insights into biological processes and disease progression from patient data.
- Physics: Theorizing new thermodynamic coordinates from experimental data on dynamics and stochastic systems.
- Chemistry: Proposing new synthetic chemical pathways through reaction data analysis.
- Materials Science: Discovering materials with desired properties by learning from atomic to macroscopic structure relationships.

### LEVEL 5

**Closed loop integration and validation** - Incorporating these four levels via an 'emulator' into an active learning framework for few/zero shot learning and/or optimization and validation, towards interpretable discovery.

- Biomedical Science: Highly accurate protein folding forward model that directly predicts functionality, verified by ex-situ testing.
- Physics/Materials Science: Stimuli-response behaviour of materials following learnt physical principles, interacting with environment and time.

Progressing through these levels, the complexity and sophistication of the AI models increases, enabling them to perform tasks ranging from simple feature

identification to the discovery of new scientific principles, thus transforming the paradigm of scientific exploration across various domains.

## 2.4

# NEED OF THE HOUR FOR THE SINGAPORE AI FOR SCIENCE INITIATIVE

To enable a successful AI4SCI program in Singapore, strategic investments on developing enabling resources are essential. We summarise these enabling resources below.

### 2.4.1 Data

Leveraging real-world data in diverse domains such as healthcare, materials and environmental sustainability is crucial. In addition, a national framework that incorporates the ability to generate balanced, accessible and integrated datasets will need to be devised, including key players such as NSCC, A\*STAR and NRF. Collaboratory networks with national institutions from other countries such as Japan's RIKEN, and the USA's NIST can be established to devise an internationally aligned effort with peaks of excellence established in Singapore.

A centralized, multi-modal data generation and validation initiative that addresses the complexity of the domain's grand challenges can drive rapid advancements and unite the research community. A notable example is the CASP experiment in the protein and RNA structure prediction field, which features a well-defined dataset, a clear problem statement, and independent experimental validation, and therefore effectively monitors and stimulates the significant advancements of the field of AI-based structural bioinformatics<sup>28,29</sup>.

### 2.4.2 Quantum (including hybrid compute)

Enhancing computational capabilities is crucial for efficiently collecting, organizing, and analysing vast amounts of domain-specific data. One promising direction to advance national computational capacity involves integrating quantum computing with classical systems, which is essential for conducting complex simulations and analyses in scientific research. Quantum processors with increased qubit coherence and

error-correction abilities, high-performance classical computers with powerful processing capabilities, and sophisticated software for the smooth interoperability of quantum-classical computations are required. This would involve creating algorithms optimized for quantum acceleration and classical systems capable of handling the complex data integration necessary for such hybrid technology.

### 2.4.3 High-performance computing (HPC)

High-performance computing resources, like those provided by the NSCC (Aspire 1a and Aspire 2a), are vital for handling large-scale computations and data processing. While providing local practitioners access to reliable and fast HPC, a strong collaboration with the global HPC frameworks in Japan and Europe

can be envisioned. Computational science is a critical enabler in Singapore's RIE efforts. Singapore's AI4SCI Initiative aims to usher in the era where a purely empirical/experimental approach can robustly incorporate theoretical foundations to the observations, as well as provide independent insight into the science.

### 2.4.4 AI algorithm development

Interdisciplinary research is critical in AI algorithm development. Emergent algorithms such as geometric deep learning<sup>30,31</sup>, self-supervised learning<sup>32,33</sup>, and generative AI models<sup>9,34,35</sup> are being researched to directly extract knowledge and formulate hypotheses. Geometric deep learning leverages neural message-passing within graphs to create latent representations that encapsulate the complex structural information of geometric data, making it particularly valuable for analysing 3D structures. Self-supervised learning, exemplified by contrastive learning methods, excels in distinguishing and enhancing the similarities and differences within diverse datasets, such as those in molecular predictions, by refining the data representations for better downstream task performance. Generative AI models, on the other hand, learn the probability distributions of data to innovate in design tasks, like creating new drugs, materials

and proteins, by synthesizing information from various modalities, including visual and sequential data<sup>36</sup>.

These generative methods are crucial for transforming raw data into scientific insights and theories, enabling a deeper grasp of intricate scientific phenomena. Encompassing mathematical and physical principles directly into the machine learning architecture, as well as extracting symbolic equations from the most generalized predictions are also enticing possibilities. Furthermore, as hardware progress in large-scale data computations, like GPU architectures, nears saturation and becomes more energy-intensive and potentially scarce, the imperative for advancing novel AI algorithms tailored for specialized, domain-aware, and efficient processing becomes increasingly critical for the success of AI for Science studies.

### 2.4.5 Self-Driving Laboratories

Automating research processes using AI can accelerate experiments, data collection, and analysis, leading to faster scientific breakthroughs<sup>37</sup>. These labs harness advanced algorithms to design and plan experiments, make decisions, and execute processes with minimal (or zero) human intervention, thereby drastically accelerating the cycle of scientific discovery<sup>38</sup>. By employing machine learning to analyse experimental data in real-time, self-driving labs can iteratively refine hypotheses,

adapt experimental protocols, and navigate the vast experimental space more efficiently than traditional methods<sup>39</sup>. This continuous, closed-loop system facilitates a more dynamic and responsive research environment, allowing for high-throughput testing and the ability to uncover complex patterns and relationships within the data. In essence, self-driving labs operationalize the vision of AI-driven science in propelling forward the frontiers of scientific research and innovation.

## 2.4.6 AI agents and leveraging on large language models (LLMs)

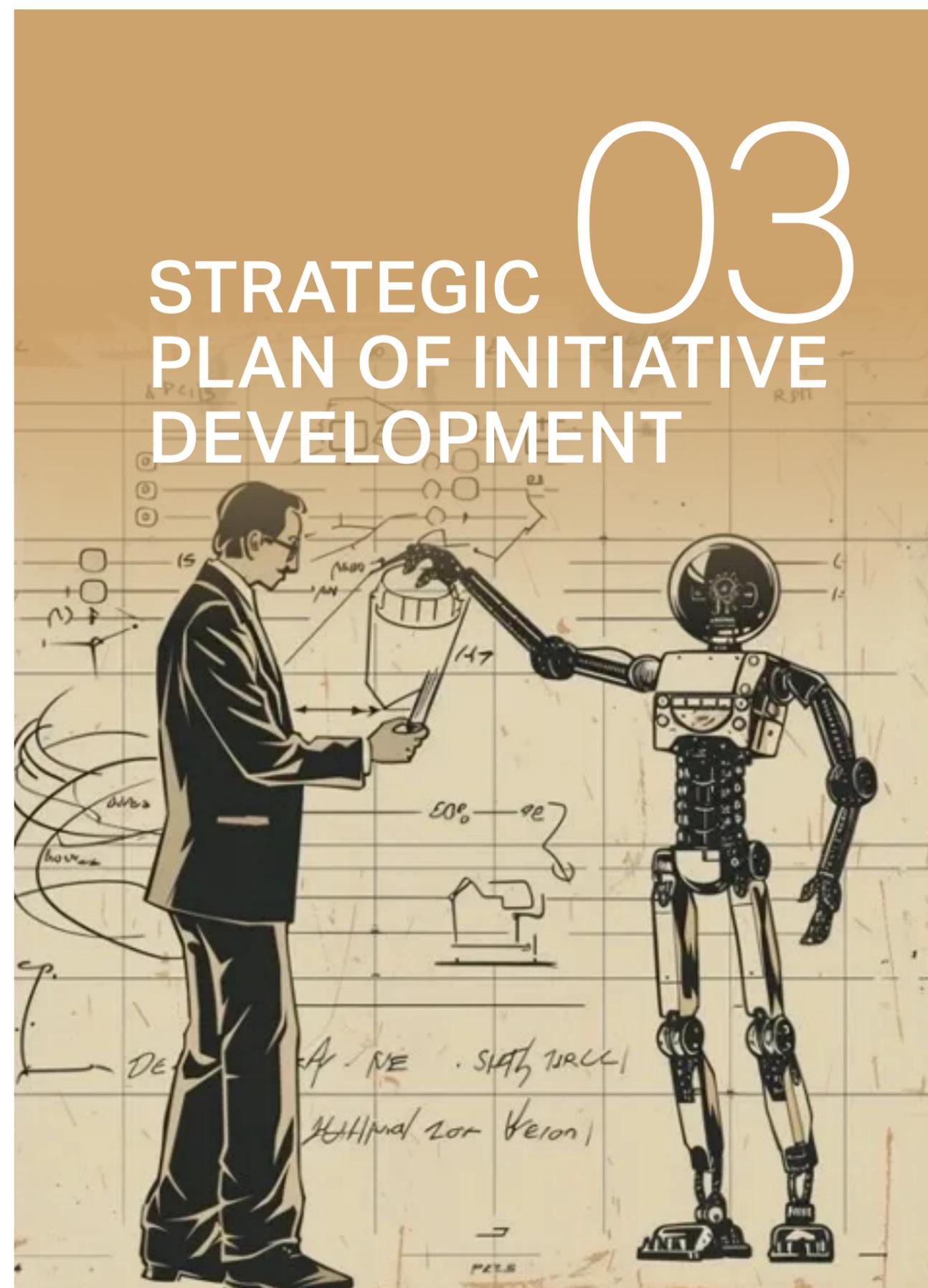
Developing skills in data mining, pattern recognition, and machine learning is necessary for harnessing the power of Large Language Models (LLMs) in scientific research. They streamline knowledge synthesis, digesting vast amounts of scientific literature to highlight key findings and research trends. Agentic AI can potentially transform AI for science by enabling autonomous systems that actively explore, learn, and make decisions. These AI agents adapt to real-time feedback, optimizing experiments and refining models to accelerate discovery while reducing resource use. A fully self-driven lab where recipes were suggested by LLMs was recently demonstrated<sup>40</sup>. Another exceptional example is Evolutionary Scale Modeling (ESM) which enables high-quality

protein structure prediction and functional design by training LLMs using a transformer network on large-scale protein sequence and structure data<sup>41,42</sup>. LLM agents also aid in hypothesis generation, leveraging expansive literature databases to propose novel research directions and connections. However, a concerted effort needs to be made in order to avoid the inherent bias of 'successful results only' in the repository of literature knowledge, managing the current issue of training LLMs on increasing body of LLM-generated literature<sup>43</sup>, as well as fine-tuning LLMs for scientific discovery. Coordinated AI agents, working together with humans and robots, can indeed accelerate innovation towards true scientific discovery.

## 2.4.7 Enhancing collaborations towards collective action

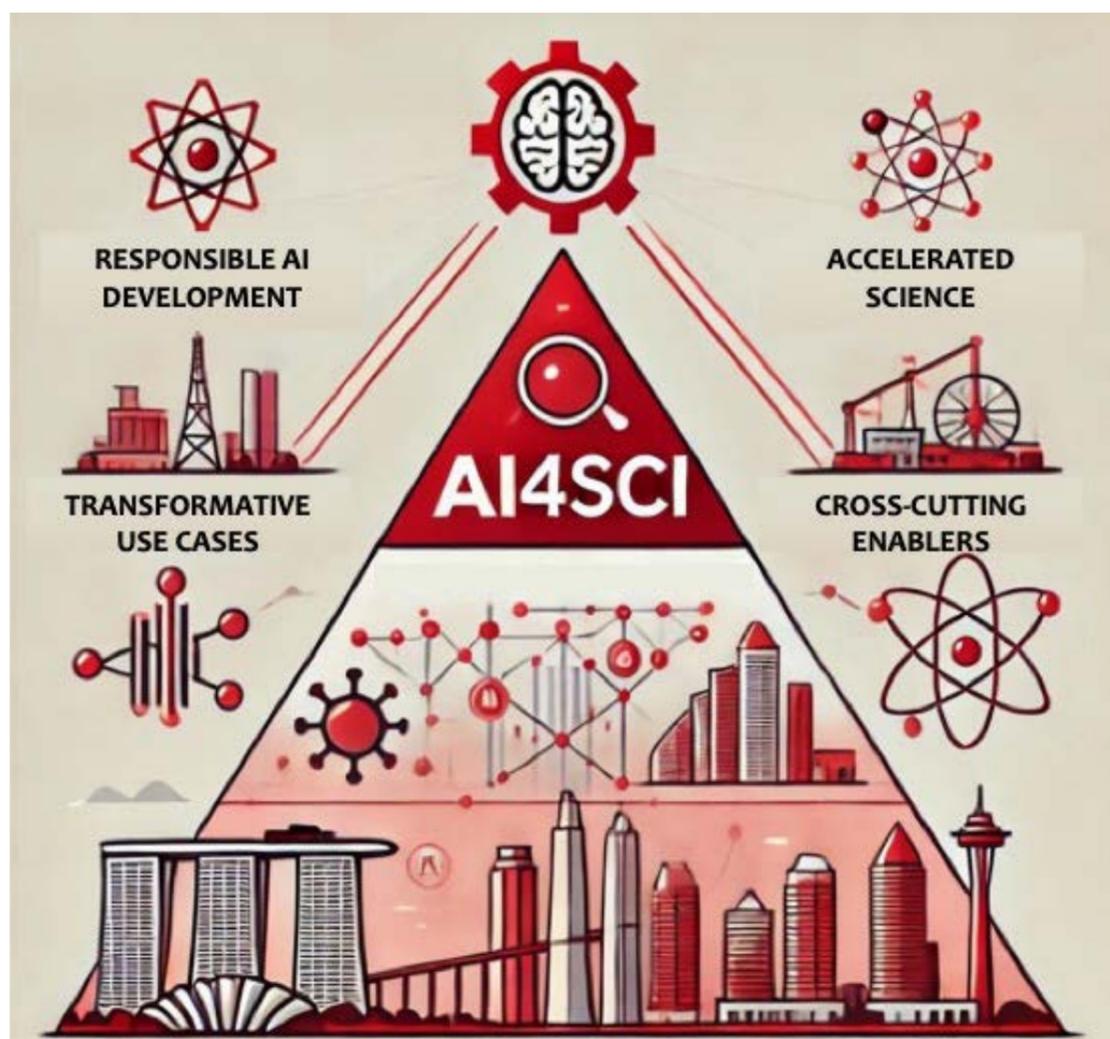
An initiative seeking to foster collaboration among siloed efforts would need to invest in a number of initiatives and sharing platforms, namely: interdisciplinary research grants, data sharing infrastructure, AI toolkits, collaborative workshops, cross-institutional partnerships, incentives for data sharing, training, community building, evaluation, and ethical frameworks. This can be operationalized

through a "gymnasium", where datasets and models are shared and commonly updated, with clearly defined problem statements and tracking of data and model development<sup>44</sup>. The presentation and discussion of research aims and data analysed can be supported by the research grants and funding of local and international workshops.



## 3.1 FOUNDATION

Singapore's AI4SCI Initiative is firmly anchored in the country's comprehensive National AI Strategy 2.0 (NAIS 2.0: <https://www.smartnation.gov.sg/nais/>), which underscores the transformative potential of AI across various scientific domains. The foundational aspects of the Singapore AI4SCI Initiative are designed to foster a robust AI ecosystem, capable of accelerating scientific research and driving innovative discoveries.



AI generated.

### 3.1.1 Building peaks of basic AI research excellence

The Singapore AI4SCI Initiative focuses on developing peaks of fundamental scientific research excellence and a conducive environment for advanced AI research in Singapore. This foundation is built on four key pillars:

- **Responsible AI Development:** Ensuring AI technologies are developed with a strong emphasis on explainability, trustworthiness, and resource efficiency. This involves investing in research that makes AI systems more transparent, reliable, and environmentally sustainable.
- **Transformative Use Cases:** Applying AI to drive transformative effects in selected sectors by leveraging AI's capabilities to solve complex scientific problems in various domains. The Initiative seeks to implement AI solutions that can significantly enhance research productivity and innovation, transforming fields and pushing boundaries in the specified domains.
- **Acceleration of the Scientific Process:** Utilizing AI to boost the research methodology itself, making scientific inquiry faster and more efficient. This includes automating routine tasks, enhancing hypothesis generation, and enabling new experimental methodologies.
- **Cross-cutting horizontals:** the development of AI models and methodologies should be applicable across various domains, ensuring their adaptability to a wide array of problems. Advancements in computational capabilities are critical, which include data (and the ability to generate data) and quantum computing, to support AI research.

Overall, the aim is to formulate ideas and impactful research challenges that are only solvable through innovative application of AI, augmented through data, compute and/or experiments.

### 3.1.2 Integration with NAIS 2.0

Singapore's National AI Strategy (NAIS) was officially launched in November 2019 as part of the country's broader Smart Nation initiative<sup>45</sup>. The strategy aims to position Singapore as a global hub for AI by developing its research, infrastructure and governance. The plan identified five key sectors for AI development: transport and logistics, smart cities, healthcare, education, and safety and security. The aim is to drive economic growth and improve citizens' quality of life by fostering collaboration between the public, private, and academic sectors, Singapore's AI strategy emphasizes responsible AI use, talent development, and establishing an

ecosystem for innovation. The government's proactive approach has positioned Singapore as a leader in the ethical development and deployment of AI technologies.

The development and implementation of a Singapore AI4SCI Initiative is an integral part of Singapore's NAIS 2.0, announced in December 2023<sup>1</sup>, which aims to harness AI for the public good, both within Singapore and globally. NAIS 2.0 builds on Singapore's 2019 strategy with key initiatives such as the establishment of AI Singapore (AISG) and a focus on infrastructures supporting research and product development in AI.

Within this context, the AI4SCI Initiative is expected to play a pivotal role by targeting the development of AI models and methodologies that accelerate scientific discovery, especially in areas of national and international interest. This approach will allow Singapore to contribute towards solving complex global challenges and ensuring that AI is utilized as a transformative tool for the public good, along the NAIS 2.0 strategic agenda<sup>1</sup>.

Furthermore, the AI4SCI Initiative will support Singapore's commitment, outlined in the

NAIS, to building a cohesive ecosystem that promotes collaboration among industry, government, and academia. Specifically, the Initiative aims to create intellectual property that engages global industry into investing in AI to improve productivity and innovation. This strategy aligns with Singapore's vision of economic growth and aims to address pressing societal needs, enhancing the quality of life for its citizens. We expect that through the Initiative, Singapore will become a foundational hub for AI research and development.

### 3.1.3 Government support and policy frameworks

Singapore's government has recently outlined its commitment to supporting the AI4SCI Initiative through a robust framework involving the NRF, the Economic Development Board (EDB), Ministry of Digital Development and Information (MDDI) and the Ministry of Trade and Industry (MTI). The NRF plays a critical role by providing funding and facilitating research collaborations, aiming to enhance the nation's AI capabilities in various sectors, through discussions with the ministries. This includes hosting research programs like the AI4SCI Initiative and decarbonization.

The EDB complements this effort by focusing on attracting investments and partnerships that leverage AI technology to drive economic growth. It has established various Centers of Excellence (CoEs) to develop tailored AI solutions for different industries such as manufacturing<sup>46</sup> and finance<sup>47</sup>, thereby addressing sector-specific challenges while promoting the adoption of AI. Meanwhile, MDDI and MTI work to create a conducive policy environment that encourages innovation and cross-sector collaboration whilst ensuring that AI Initiatives align with national priorities. Together, these agencies embody a comprehensive approach to harnessing AI's transformative potential, positioning Singapore as a leader in the global AI landscape.

Along with investments in and support of AI-related initiatives, the Singapore government recognises the importance of ensuring that AI is developed and use without compromising data security and research excellence. This has led to a recent national policy is the Model AI Governance Framework<sup>48</sup>, which provides guidelines for the responsible use of AI technologies. The framework was updated to address the emerging challenges posed by generative AI and has garnered international support, reflecting Singapore's proactive stance in ensuring that AI developments are aligned with public good.

In addition to CoEs and government agencies, technology-driven organisations such as the Defence Science Organisation (DSO) and Temasek play a significant role in Singapore's AI ecosystem. The DSO leverages its expertise in defence and security to develop advanced AI technologies. For instance, it has recently established a joint France – Singapore R&D Lab, based in Singapore, to develop AI Capabilities<sup>49</sup>. Temasek, as an investment company, actively supports AI Initiatives by investing in companies and technologies that support innovations in AI driving economic growth. By integrating efforts from government agencies and technology receptacles, Singapore is uniquely positioned as a synergistic environment to accelerate AI advancements and enhances its application in scientific research and other areas.

### 3.1.4 Potential benefits and global trends

AI has the capability of boosting scientific productivity, augmenting human cognitive capabilities, and accelerating breakthroughs, and can especially be impactful in a small nation like Singapore with a strong foundation in research, innovation and enterprise. This Initiative is aligned with global trends where AI applications in science are growing at a significant pace and providing a significant value add. Countries like China, the EU, and

the US are leading in AI scientific production. According to recent bibliometric analyses<sup>50</sup>, the number of AI publications is increasing rapidly, indicating a significant shift towards AI-driven scientific research. The AI4SCI Initiative aims to position Singapore in a leadership position within this global context, leveraging AI to drive innovation and maintain competitive advantage.

## 3.2

### CAPACITY GAPS

To fully realize the potential of the Singapore AI4SCI Initiative on significantly promoting scientific research and innovations in Singapore, addressing the following capacity gaps will be essential:

- **Multi-lingual Talent with AI and Domain Expertise:** There is a critical demand for bilingual scientists who possess both domain-specific expertise in basic sciences and proficiency in AI methodologies. Developing and attracting such talent is essential to closing the gap between conventional scientific research and AI-driven innovations. Some examples of this include academic groups (e.g., AI for Science at Caltech, Acceleration Consortium at University of Toronto) and industry (e.g., META FAIR, DeepMind, Microsoft Research AI4Science, IBM Research, Tencent AI Lab, Sony AI).
- **Algorithm and Software Efficiency:** Current AI algorithms often require excessive computational resources due to inefficiencies. Developing more efficient algorithms and management protocols is necessary to optimize resource usage

and reduce the computational burden. In addition, the ability to leverage on specific data uniquely generated in Singapore, relevant to various domains, and not necessarily in the data-rich regime could prove a key differentiator. Investment into efficient hardware will also be important.

- **Data Management:** The labour-intensive process of big data curation and acquisition presents considerable challenges. Optimizing data management workflows and improving data accessibility across computing infrastructures, including the National Supercomputing Center (NSCC) and institution-based HPC clusters, are vital for advancing AI research. Additionally, Singapore's AI4SCI Initiative should establish dedicated infrastructure to support its funded projects and facilitate coordination to meet the Initiative's specific requirements.

### 3.3

## POTENTIAL FOR INDUSTRIAL AND ACADEMIC ALIGNMENT

By fostering deep collaborations between Science of AI experts and AI for Science domain researchers, the Initiative aims to create synergies that enhance the application of AI across various sectors. Two key outcomes are expected:

- **Cross-Domain AI development:** Encouraging the generalization of AI models and solutions developed for specific problems to other scientific domains. This approach promotes learning and accelerates innovation across multiple fields. This can then provide foundational AI development that acts as the engine of innovation.
- **Industry engagement:** Strengthening collaborations between academic institutions and industry players to drive research and development in this exciting space. This is especially pertinent because AI for Science is of significant interest to various industry leaders due to the critical roles that big data and AI techniques play in transformative goals. With Singapore's reputation as a safe IP haven, with access to world markets, this could be an opportunity to capitalize on key domains through the Singapore AI for Science Initiative.

### 3.4

## VISION

The vision of the Singapore AI4SCI Initiative is to revolutionize the scientific discovery process of the national community of Singapore through the strategic application of cutting-edge AI techniques. By accelerating the pace of scientific discovery and improving research outputs, the Initiative aims to position Singapore as a global leader in AI-driven research. This vision encompasses:

- **Accelerated Discoveries:** Harnessing AI to automate the scientific discovery process across different disciplines, thereby driving faster development in scientific knowledge and solutions to previously intractable challenges. A focus on accelerated discoveries aligns with Singapore's transition to a knowledge-based economy, leveraging its robust foundation in digitization and Industry 4.0.
  - **Global Leadership:** Establishing Singapore as a hub for cutting-edge AI research and innovation, attracting leading researchers and fostering international collaborations.
- This vision builds of a dynamic and well-connected academic ecosystem, anchored by institutions such as NUS, NTU, and A\*STAR, which thrive on a solid foundation of excellence and growth.
- **Sustainable Development:** Ensuring that AI technologies are developed and deployed responsibly, with a focus on ethical considerations and environmental sustainability. This approach reflects Singapore's dedication to responsible innovation, striking a balance between economic progress, sustainability, and societal well-being.

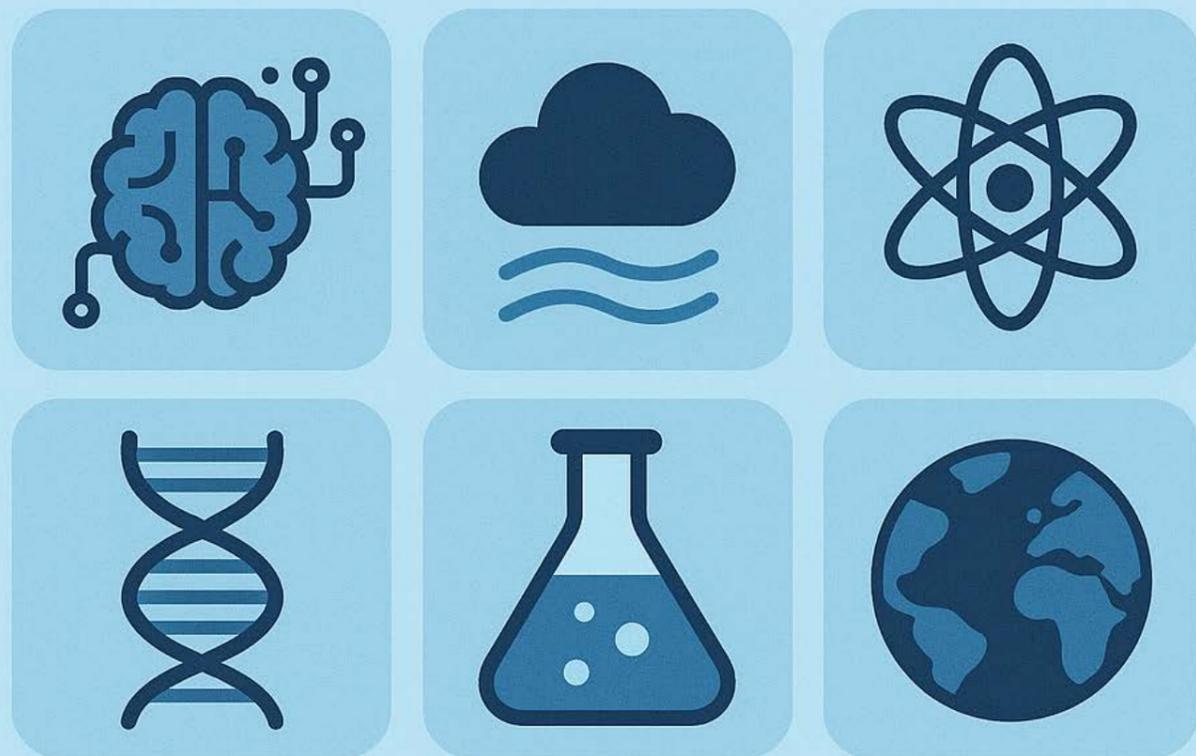
### 3.5

## SPECIFIC POLICY RECOMMENDATIONS

To achieve the vision and goals of the Singapore AI for Science Initiative, several targeted policy recommendations are essential:

- **Investment in Talent Development:** Establish programs to train and support multilingual scientists who can integrate domain expertise with AI capabilities. This includes scholarships, grants, and partnerships with educational institutions to develop specialized AI and domain-specific curricula.
  - **Enhancement of Computational Resources:** Increase funding for high-performance computing infrastructure and develop more efficient algorithms to optimize resource usage. Encourage the development of software that reduces computational requirements while maintaining research accuracy and efficiency.
  - **Streamlined Data Management:** Implement policies that facilitate the curation, procurement, and sharing of high-quality datasets. Develop standardized protocols for data management and incentivize the creation of open-access data repositories to support collaborative research efforts.
  - **Promotion of Industry-Academia Collaboration:** Foster partnerships between academic institutions and industry players through joint research programs, innovation hubs, and collaborative projects. Establish frameworks for intellectual property management and commercialization of AI-driven research outcomes.
  - **Support for Interdisciplinary Research:** Encourage interdisciplinary research initiatives that bring together AI experts and domain scientists to tackle complex scientific problems. Provide funding and resources for collaborative projects that span multiple scientific fields and leverage AI methodologies.
- By implementing these policy recommendations, the Singapore AI for Science Initiative can overcome capacity gaps, enhance industrial and academic alignment, and realize its vision of transforming the scientific discovery process through the power of AI.

# 04 AI4SCI DOMAINS



Given the relatively low scientist population and limited infrastructural resources in Singapore, the AI4SCI initiative can start with the development of research areas that represent some of the most urgent and important scientific domains within which AI could play a critical role. These are also the areas where many Singapore based research laboratories have already achieved (and/or approached) state of the art of the field and therefore build an important base to enhance the leadership position of Singapore relative to the global scientific community in *AI for Science*. This section provides an overview of the existing AI for Science ecosystem in Singapore, along with summaries of domain-specific white papers and references, which are included in the appendices. It is important to note that the list of principal investigators (PIs) cited here is intended to be illustrative and non-exhaustive; we acknowledge contributions from other PIs that are inadvertently not mentioned herein, and we apologize for any inadvertent omissions.

## 4.1 BIOMEDICAL AND HEALTH SCIENCES

### 4.1.1 AI-based protein and RNA structure modelling

Protein and RNA are two of the most important types of molecules in life. While proteins are 'workhorse' molecules carrying out most biological activities of living organisms<sup>51</sup>, RNAs, especially the non-coding RNAs, are recently discovered to perform critical cellular functions beyond their canonical roles in gene transcription<sup>52</sup>. Because the biological functions of proteins and RNAs are specified by their spatial shape, determination of the 3D structure is essential to annotate their biological functions and develop new drugs to regulate their functional roles<sup>53</sup>. There are by now close to 250 millions of proteins/RNAs with known sequences in the UniProt/RNACentral libraries; but less than 0.1% of them have structures solved experimentally in the PDB<sup>10</sup>. Therefore, computer-based structure prediction is the only means to alleviate the gap and high demanding of protein/RNA structures by the biomedical community.

The research of **Prof. Yang Zhang** (NUS) involves pioneering bioinformatics algorithms that harness AI and physical force fields for protein and RNA structure modeling

<sup>54,55,56,57</sup>. His team's work on D-I-TASSER led to its recognition as the top server/predictor surpassing AlphaFold2 in the 15th CASP experiment<sup>58</sup>, reflecting significant strides in the field of computational structural biology. Led by Executive Director **Dr. Sebastian Maurer-Stroh**, the Bioinformatics Institute (BII) specializes in computational protein sequence and structure analysis<sup>59</sup>. Their work includes predicting various aspects of molecular and cellular functions, identifying biologically important residues and disease-related mutations, and developing predictors for short functional motifs in protein sequences<sup>60</sup>. **Prof. Yansong Miao** (NTU) pioneers AI-based protein assembly and structure prediction, aiming at uncovering new therapeutic approaches and understanding complex biological systems<sup>61</sup>. **Prof. Lanyuan Lu** (NTU) focuses on intrinsically disordered proteins, shedding light on their roles in diseases and potential drug discovery avenues<sup>62</sup>. To capitalize on these AI approaches, Singapore provides an exceptional environment including multiple academic partners, national platforms and large or small companies focused on drug discovery and clinical development.

## 4.1.2 RNA biology

RNA is a type of fundamental nucleic acid molecule involved in multiple functions within living organisms, while RNA biology seeks to understand the structure, function, and regulation of these molecules. Due to the important roles found recently<sup>52</sup>, RNA biology and its applications in developing new therapeutics represent a critical emerging field, attracting tremendous academic and industrial investments worldwide, including Singapore. It is estimated that targeting RNAs with small molecules will expand the drug design landscape by more than an order of magnitude compared to traditional protein-targeted drug discovery<sup>63</sup>. Using AI to analyze vast RNA-related data holds transformative potential for treating infections, cancer, and heart-related diseases. With the significant strengths in AI research and RNA biology, Singapore has the opportunity to lead globally in this field.

Three critical questions will be addressed in this area of AI for Science: 1) how does the diversity in Singaporean genome influence RNA diseases? 2) Can we develop smarter AI algorithms to precisely model atomic structures of RNA and RNA-protein interactions? 3) How do we use AI to predict RNA and small molecule interactions and enhance success rate of drug discovery targeting these interactions? To address these questions, a multidisciplinary team, comprising researchers from different principles of computer, medicine, and biological sciences, will engage national-level capabilities and programmes such as the AI-Singapore consortium, National Super Computing Centre, National Precision Medicine & SG100K programmes, emerging National Initiative for RNA Biology & Its Applications, and recently formed Nucleic Acid Therapeutics Initiative.

The development of RNA Biology is led by Dr. Ashok Venkitaraman, the Distinguished Professor of Medicine at NUS, Director of CSI Singapore, and Director of NUS Centre for Cancer Research<sup>64,65</sup>, where the key players include Prof. Yang Zhang who pioneered AI-based RNA structure modeling and RNA design and developed DeepFoldRNA<sup>57</sup>; Prof. Polly Chen, who is the leading expert in RNA editing and RNA-based therapeutics<sup>66</sup>; **Prof. Yvonne Tay**, **Prof. Anthony Khong** and **Prof. Daniel G. Tenen** lead programmes focused on RNA alterations promoting cancer<sup>67</sup>; and **Dr. Jason Pitt**, a new PI of CSI with strong expertise in RNA and cancer genome data science<sup>68</sup>. In addition, **Prof. Melissa Fullwood** (SBS, NTU) has worked on AI applied to biology, particularly epigenomics, by developing new AI methods for predicting chromatin interactions and transcription, and is now combining biological methods for understanding epigenetics, together with the development of AI systems for predicting important RNAs and epigenetic components towards small molecule inhibitor drugs<sup>69</sup>. **Prof. Meng How Tan** (CCEB, NTU) is working on deep learning methods to detect RNA modifications in native transcripts sequencing using their nanopore platforms<sup>70</sup>. **Dr. Chuan-Sheng Foo** (I2R, A\*STAR) has developed AI-based RNA secondary structure prediction, including methods that can integrate constraints from multiple experimental structure-probing data sources (CONTRAFold-SE)<sup>71</sup>.

## 4.1.3 Synthetic biology

Synthetic Biology is an interdisciplinary field that combines principles from biology, engineering, and computer science to design and construct new biological systems or redesign existing ones for specific purposes. It involves the assembly of biological components, such as genes and pathways, to create synthetic organisms or modify existing systems to achieve desired functions. The significance of the studies lies in the potential to engineer biological entities with tailored functionalities, paving the way for innovative solutions to address challenges in medicine, energy, environment, and beyond. While the major challenge consists of the complexity of biological systems and unpredictability of their behavior when manipulated, AI represents the most effective and powerful tool to address these challenges by providing computational models that can analyze vast biological datasets, predict the outcomes of genetic modifications, and optimize the design of synthetic organisms. AI's ability to handle complex biological information and accurately predict system behaviors enables more efficient and precise synthetic biology designs, paving the way for groundbreaking advancements in biotechnology and medicine.

To address the challenges, the new National Centre for Engineering Biology (NCEB), as SynBio-AI Innovation Accelerator with NSCC/AI-SG, led by **Prof. Matthew Chang** and **Prof. Wen Shan Yew**, have outlined a strategic roadmap spanning both mid-term (5 years) and long-term (>5 years) objectives aimed at expediting advancements in Synthetic Biology within Singapore, called SynCTI<sup>72</sup>. This strategic approach is specifically geared towards facilitating the AI-based translation and commercialization of Synthetic Biology innovations, with emphasis on identifying pivotal areas, engaging relevant stakeholders, and implementing measures to bridge critical gaps, to harness the full potential of synthetic biology for the benefit of Singapore's scientific and economic landscape. The NRF-supported Singapore Consortium for Synthetic Biology (SINERGY) brings together Principal Investigators from the local ecosystem, linking them closely with industrial partners and could serve as a launchpad for new AI-driven discoveries<sup>73</sup>. In related work, **Prof. Kang Hao Cheong** (NTU) has applied AI in biomaterials and drug discovery to develop a bio-functional polymer for preventing retinal scarring and a nanomicelle system for delivering anti-VEGF agents to the retina, showcasing non-invasive treatments and enhancing drug delivery efficiency. In combinatorial optimization, they utilized evolutionary algorithms with divide-and-conquer strategies, significantly improving problem-solving efficiency in complex systems compared to traditional methods<sup>74</sup>.

#### 4.1.4 Computer-aided drug design

Computer-Aided Drug Design (CADD) is a computational approach that employs computer algorithms and simulations to facilitate the discovery and optimization of new candidates in drug discovery industry<sup>75</sup>. The significance of the CADD studies lies in their ability to predict the interactions between potential drugs and biological targets, assess their efficacy and safety, and simulate their behavior within the human body, all before entering extensive experimental phases. However, a critical challenge in CADD involves navigating the vast chemical space to identify optimal drug candidates and to also understand the underpinning drug-target interactions.

An integrated approach can be envisioned with **Prof. Richard Ming Wah Wong**, a NUS professor in Chemistry, **Prof. Saif Khan**, at NUS Department of Chemical and Biomolecular Engineering, and **Dr. Henry Yang**, at CSI NUS. Expertise in computational quantum chemistry, molecular dynamics simulations, omics data analysis, computational genomics and AI for biological chemistry/materials can be leveraged on, in addition to robotic flow-based manufacturing of pharmaceutical and advanced materials. **Prof. Kelin Xia** (NTU) employs mathematical AI, leveraging differential geometry, algebraic topology, and statistical learning for studying structures and dynamics in molecular sciences, with a focus on data analysis and deep learning models applied to drug and materials design<sup>76</sup>. The Experimental Drug Development Centre

(EDDC) led by **Prof. Damian O'Connell** has fully integrated capabilities with high-throughput screening and testing as well as discovery components<sup>77</sup>. **Prof. Peter Preiser** (NTU) is dedicated to identifying new chemical entities targeting the Plasmodium falciparum parasite, with the goal of developing more effective treatments for malaria and understanding the mechanisms of resistance<sup>78</sup>. **Prof. Yuguang Mu** (NTU) develops deep learning algorithms for modelling protein and ligand binding interactions<sup>79</sup>. **Dr. Xiaoli Li** from (I2R, A\*STAR) has developed network-based deep learning techniques in the realms of drug repositioning and the identification of novel drug targets<sup>80,81</sup>. Their research encompasses a wide array of applications, such as predicting synthetic lethality, forecasting disease-related genes, uncovering microbe-drug associations, and elucidating drug-target interactions. Generative AI approaches can also be used to leverage multi-omics data for the prediction of drug-patient interactions for personalized medicine.

Importantly, to guide and provide clinical line-of-sight for CADD efforts on cancer, clinician-scientists at the CSI (e.g., **Prof. Boon Cher Goh**, **Prof. Wee Joo Chng**, **Prof. Soo Chin Lee**, **Prof. David Tan**) lead research programmes on experimental therapeutics at the CSI closely linked to research efforts there<sup>82</sup>. Together, their work exemplifies the dynamic and interdisciplinary nature of modern biology, contributing to advances in drug design, medicine and beyond.

#### 4.1.5 Health and medicine

While the field of health and medicine is vast, one can envision the application of AI to solve complex problems and improve healthcare outcomes. AI systems can aid to analyse vast amounts of biological, genomics and medical multi-modal phenotype data for new insights, to personalize medicine by tailoring treatments to individual patient profiles, and to enhance diagnostic and prognostic tools. It could be synergized through collaborations with the automation and optimization of drug discovery and development processes, from initial screening of drug compounds to clinical trials and patient monitoring, including digital bio- and behavioural markers. The aim would be to create more efficient, effective, and precise medical interventions, potentially revolutionizing the fields of drug development, diagnostics, treatment protocols, and disease prevention.

**Prof. Dean Ho** (NUS) specializes in integrating AI with personalized medicine, developing the CURATE.AI platform which optimizes individual patient treatments, and leading clinical trials in areas such as cancer and COVID-19<sup>83</sup>. He has also innovated in nanodiamond-based drug delivery and imaging, enhancing chemotherapy and diagnostic safety. As Director at The N.1 Institute for Health and WisDM at the National University of Singapore, Prof. Ho drives advancements in digital medicine and cost-effective drug development. In addition, **Prof. Yueming Jin** (NUS) has developed an AI system that can accurately identify and predict different stages of a surgical procedure, which can provide support to surgeons and facilitate the development of intelligent robotic systems and automated surgeries. She has also worked on spatial-temporal representation learning, data efficient learning, multi-modality learning, and personalized federated learning for medical image segmentation<sup>84,85</sup>. **Prof. Ray Najjar** (NUS) is an AI and vision scientist who studies the eye-brain connection and the effects of light on health. He has invented methods to diagnose serious diseases

with eye measurements and deep learning, and to treat circadian disorders with light therapy<sup>86,87</sup>. At the Lee Kong Chian School of Medicine (LKCMedicine) in NTU, the Centre for Biomedical Informatics uses machine learning (ML) based methodologies as a means to extract meaningful combinations of features predictive of some endpoints from biomedical data which is usually of high dimensionality e.g. genetics, transcriptomics, proteomics, and metabolomics. **Prof. Bernett Lee** (NTU) generates predictive multivariate models for biomedical data with low sample sizes and applies such methods to electronic medical records as well as genomic sequence data<sup>88</sup>. **Prof. Wilson Goh** (NTU) works on a personalized AI for mental health prediction through development of Emotional Variance Analysis and on combining network modelling and transductive personalised modelling towards tackling heterogeneity of multi-omics profiling in elderly dementia<sup>89</sup>. Further, at the Dementia Research Centre at NTU, integration of AI technologies in its multidimensional approach to cognitive health and detection of dementia at the earliest stages is being performed by **Prof. Nagaendran Kandiah** (NTU). Other relevant work include a data platform to enable AI approaches to multimodal health information by **Prof. Balazs Gulyas** (NTU), data-driven approaches for precision medical imaging (MI) tasks and medical data interpretation by **Prof. Si Yong Yeo** (NTU), and **Prof. Wynne Hsu's** (NUS) work involving applying data analytics and machine learning techniques to various medical applications, such as retina image analysis, disease spread visualization, and primary care automation. **Prof. Joanne Ngeow** (NTU) works on cancer genomics and genetic mutations, where AI and machine learning techniques are increasingly being used to analyze large genomic datasets, identify patterns, and predict cancer risks more accurately<sup>90</sup>. **Dr. Mile Šikić** (GIS) is involved in development of foundation models for RNA and DNA nucleotides, de novo assembly of complex genomes and analytics<sup>91</sup>.

Singapore's healthcare clusters and agencies have been forward-looking in the use of AI for healthcare applications, with priorities in addressing the burden of clinical administration, coaching/nudging for health behaviour change, sensing and digital phenotyping, clinical decisions support, and application of LLMs for clinical and patient advice. The National University Health System (NUHS), through efforts led by Prof. Kee Yuan Ngiam, Group CTO, and **Prof. James Yip**, Head, Academic Informatics Office, has developed the Discovery AI Tribrid platform that provides strategic infrastructure for AI development and deployment<sup>92</sup>. At Singapore Health Services (SingHealth), efforts are led by **Prof. Say Beng Tan**, Group Chief Research Officer, Prof. **Khung Keong Yeo**, Deputy Group Chief Medical Informatics Officer, and **Prof. Daniel Ting**, Chief Data and Digital Officer. Prof. Daniel Ting has led successful healthcare innovation projects, such as SELENA+, an AI-powered image reader to analyze eye scans and diagnose diabetic eye diseases, which was co-developed with the Singapore Eye Research Institute and the NUS School of Computing<sup>93</sup>. Synapse, Singapore's national HealthTech agency, supports our healthcare sector in development and deployment of such healthcare innovation projects leveraging AI.

Trust is a critical determinant of successful adoption of AI in health and medicine. Hence, fundamental work to address the explainability and reliability of AI remains key, for instance to mitigate and overcome the problem of hallucinations by generative AI models. On the data front, high quality data made accessible on trusted platforms is crucial. The TRUST<sup>94</sup> data exchange, set up by the MOH Office for Healthcare Transformation (MOHT) helmed by **Dr. Koh Mingshi** (Head, TRUST Office), provides a trusted environment for data access for national clinical data and simplifies data access by providing a single point of governance and enforcing data interoperability. This enables and accelerates the use and development of data analytics and AI for the healthcare sector, anchored by **Prof. Robert Morris** (MOHT/NUS) and his team working at the intersection of technology and data. There exists a strong partnership between Institute of Mental Health (IMH) and MOHT, which has developed world-leading competence in digital phenotyping and interventions over the past 5 years.

## 4.2

# ADVANCED MATERIALS AND SUSTAINABILITY

### 4.2.1 Future materials science and development

Complexity and variability in materials science present a significant set of challenges that could be addressed through AI, from modeling complex systems and discoveries in new material design to improving synthesis and characterization processes. The large range of spatial and time scales makes modeling and reverse engineering even more complicated. Data-driven AI-enabled approach is required to address such variability to enhance predictive modeling and optimize material properties. AI can help bridge the gap between different scales of molecular modeling, such as Density Functional Theory (DFT) and larger scale continuum mechanics, by providing surrogate multiscale modeling techniques, including physics-based interpretability<sup>95</sup>. Additionally, AI can assist in generative and inverse design processes, allowing for rapid synthesis, testing, and characterization of a wide range of materials. This includes the use of AI in autonomous systems that can speed up the integration and characterization of novel materials<sup>96</sup>. Through the integration of new AI techniques in this area, Singapore is poised to fast-track materials science advancements, enhancing its competitive edge in technology and logistics on the international stage.

**Sir Prof. Kostya Novoselov**, the 2010 Nobel Laureate in Physics and NUS Professor in Materials Science and Engineering is co-directing the Institute of Functional Materials (I-FIM, NUS) along with **Prof. Antonio Castro-Neto** (I-FIM, NUS), focused on creating and controlling intelligent functional materials. **Prof. Kedar Hippalgaonkar**, (MSE, NTU/IMRE, A\*STAR), is a leader in materials-by-design through AI and high-performance computations combined with high-throughput, robotic experiments<sup>97</sup> it is natural to wonder what lessons can be learned from other

fields undergoing similar developments. In this Review, we comparatively assess the evolution of applied ML in materials research, gameplaying and robotics. We observe ML being integrated into each field in three phases: first into discrete hardware and software tools (toolset integration). He is also developing property-directed generative models and working on building a foundation model for materials science. **Prof. Qianxiao Li** (NUS) in Mathematics is an expert in the interplay of machine learning, AI for sciences and dynamical systems. NUS Visiting **Prof. Andrey Ustyuzhanin** is also the Director of AI/ML research at Acronis and a PI at I-FIM. In addition, **Dr. Tan Teck Long**, a computational material scientist who has developed AI surrogate models to bridge length-scales in materials simulations at the Institute of High Performance Computing (IHPC), along with **Dr. Ivor Tsang**, the Director of the Centre For AI Research (CFAR) at A\*STAR, are using generative models to link alloy structure to properties. **Dr. Tan Teck Leong** and **Dr. Benjamin Chen**, at A\*STAR's IHPC, lead a team of computational material scientists/chemists to build up an informatics platform for alloys and catalysts via a broad range of methods from high-throughput quantum simulations to data-driven approaches. Some of their areas of application include alloy and composite design, an area where **Prof. Yeong Wai Yee**, **Prof. Hejun Du**, **Prof. Xuan Liang** and **Prof. Kun Zhou** (MAE, NTU) are active in. **Prof. Yonggang Wen** (SCSE, NTU) working with **Prof. Alex Yan Qingyu** (MSE, NTU) are utilizing optimization algorithms for battery performance enhancements. Similarly, **Prof. Bo An** (SCSE, NTU), working with **Prof. Guan Cuntai** (SCSE, NTU) and **Prof. Liu Zheng** (MSE, NTU), are developing a generalist materials discovery

platform. Along similar lines, **Prof. Kee Woei Ng** (MSE, NTU) is developing a general AI-driven platform for urban farming. **Prof. Zhi Ning Chen** (NUS), Director of Advanced Research and Technology Innovation Centre and **Prof. Cheng-Wei Qiu** (NUS), **Prof. Xinchao Wang** (NUS) and **Prof. Zhaogang Dong** (SUTD/IMRE) have worked on prior-knowledge-guided deep learning-enabled methods for modeling, optimization, and synthesis of metaphotonics, including metacells, metasurfaces, frequency-selective surfaces, and metasurface antennas, pioneering this

## 4.2.2 Quantum materials and quantum computing

Quantum materials are a class of substances that exhibit unique and exotic properties arising from quantum mechanical phenomena, leading to groundbreaking applications in areas such as superconductivity and quantum computing<sup>102</sup>. In the field of quantum materials, AI can augment the way scientists generate hypotheses, discover new properties, and develop advanced materials. By using unsupervised learning, AI can reveal the structure of unlabelled datasets, potentially highlighting novel features or behaviours of materials that can form the basis for new hypotheses. AI accelerates the design and testing of these materials by predicting their behaviors and prioritizing the most promising ones for further study. Additionally, Bayesian optimization and other active learning methods allow for the dynamic improvement of candidate prioritization, leading to a more efficient exploration of the parameter space<sup>103</sup>. In the area of quantum simulations, we can tap onto analog and digital simulations. Analog simulations are designed to provide insights into specific scientific models, while digital simulations can be divided into noisy and fault-tolerant ones. Some hybrid solutions have portions of the computations performed on (noisy) quantum processors and others on a classical one. These hybrid solutions have seen significant interest from multi-national corporations, and there is significant progress towards the realization of fault tolerant quantum solutions<sup>104</sup>.

field by not just better performance, but by creating new functions in existing designs, and exploring new knowledge [ref]<sup>98,99,100,101</sup>.

With the strong and multidisciplinary expertise in the ecosystem, involving AI, material science, and device physics, such science-relevant methods are well placed to auger a paradigm shift in new material and device development, ensuring heightened innovation, enabling the discovery of new materials and new functionality, discovered through accelerated processes.

**Prof. Jiong Lu** at Chemistry, NUS, has built scanning tunneling probes, providing a solid base for generating data at the single-atom and quantum level, and developed effective methods and built their own datasets<sup>105,106,107</sup>. **Prof. Pinaki Sengupta** (NTU) works on complex quantum states, useful for topological quantum computation and has recently used group equivariant convolutional neural network (GCNN) ansatz and variational Monte Carlo simulations to reveal the ground state of the spin-1/2 kagome lattice antiferromagnet as a spinon pair density wave (PDW) state, marking a significant advancement in understanding these novel quantum phases and their potential applications<sup>108</sup>. **Prof. Lam Ping Koy** at the Center for Quantum Technologies (CQT) and IMRE, A\*STAR has worked on quantum information and metrology, while **Prof. Dimitris Angelakis** (CQT) has been a leading figure in the area of quantum and hybrid (classical-quantum) simulations<sup>109</sup>. **Prof. Jose Ignacio Latorre**<sup>110</sup> has been leading developments in hybrid quantum computing, and **Dr. Joe Fitzsimons** is an expert of compilers involved in the future use of fault tolerant quantum computers<sup>111</sup>. **Prof. Dacheng Tao** (NTU) and **Prof. Mile Gu** (NTU) are developing large scale quantum AI microprocessors and algorithms, as well as quantum-native AI that can solve complex problems with lower energy costs and deepen our understanding of the energetic limits of learning and decision-making in quantum environments<sup>112,113,114,115</sup>.

## 4.2.3 Imaging and high-resolution microscopy

The human eye is a powerful visual tool, but it does not have the resolution to bring microscopic images into focus. This is where microscopy can help us view and understand nanoscale objects from pandemic-inducing viruses to the increasingly miniaturized electronics. Currently, techniques like electron microscopy (with many variants, scanning transmission (3D and 4D), energy electron-loss spectroscopy, electron diffraction, etc.) are effective, yet can be expensive and require extensive preparation. Optical imaging is a more cost-effective and less invasive alternative, but it suffers from limitations in resolution/throughput and expensive costs in light sources/optics settings. Rapid characterization via a combination of such techniques, requiring both advances in hardware and software, will allow the generation of rich datasets capturing the complexity of materials.

**Dr. Xianwen Mao**, a NUS Presidential Young Professor in MSE, is an expert in optical imaging of energy materials<sup>116,117</sup>. His work focuses on AI-based preprocessing involving denoising, smoothing, and image

enhancement to improve accuracy and robustness. The development of AI models for rapid identification and classification of materials will utilize extraction of detailed visual, textural, and shape features, enhancing the screening process. By training AI to correlate data across imaging methods like SEM and OM, image resolution and quality can be enhanced for cost-effective high-throughput material analysis<sup>118</sup>. **Prof. Thorsten Wohland** in NUS, Biological Sciences and **Prof. Adrian Rollin** in Statistics and Data Science, have built on deep learning-based image analyses to extend single molecule spectroscopy via FCS to the next frontier<sup>119,120</sup>. **Prof. Yeng Ming Lam** at NTU MSE is leading the Facility for Analysis, Characterisation, Testing and Simulation (FACTS), which provides access to cutting art facilities in imaging and has started to explore AI-augmented image analysis<sup>121</sup>. **Prof. Tong Ling** (CCEB, NTU) has developed a novel subpixel motion correction method for Fourier-domain optical coherence tomography (OCT) and an unsupervised learning algorithm that extracts the spatiotemporal patterns of the OCT signals.

## 4.2.4 Advanced semiconductor technology development

As semiconductor technology advances deep into the Angstrom era of scaling, complexity and cost of chip technology are escalating rapidly. Ironically, the chip technology advancement is also the very hardware foundation that accelerates the AI and computational sciences. In recent years, the fault isolation and diagnostic of the complex chips have become daunting due to the nanometer component sizes, highly miniaturized circuitries, and the multitude of tightly integrated atomic-scale electronic materials. Working with various chip companies in Singapore (E.g. AMD, Qualcomm, Global Foundry and Micron), it is clear that fault diagnostic of advanced chips is a high-value problem that gates the yield and manufacturing of their next-generation products. To this end, there have

been research projects at NUS investigating Physics-guided-AI - based chip diagnostic methods<sup>122</sup>. Working with industrial partners, major approaches have been developed to address fundamentals related to scarcity of defects, non-destructive vs. destructive validation, AI-Physical Model integrated digital twinning, atomic variability, and statistical machine learning coupled to real-time testing. **Prof. Aaron Thean**, Director of the NRF's Singapore Hybrid Integrated Next-Generation Electronics (SHINE) centre at NUS has a team that investigates and develops such methods<sup>123</sup>, in collaboration with **Prof. Yeow Kheng Lim** (NUS) and **Dr. J. Senthilnath** (I2R, A\*STAR). They have reported success on advanced technologies 3nm and beyond, working in collaboration with AMD. One of the latest projects funded under the Singapore

Industry Alignment fund is a multi-institution (A\*STAR& NUS) and company project (AMD, Qualcomm, POET Technologies, Neocera Inc.) led by NUS, hosted at SHINE, looking into AI-Enabled Magnetic Tomography for online non-destructive fault testing of next-generation 3D Chips. They are also developing multi-level multi-physics models, and generative AI to serve as an inverse solver and for data augmentation, sustainable AI for data-efficient and resource-efficient processing, and optimization methods for semiconductor manufacturing process. In related work, **Prof. Mohamed M. Sabry** (SCSE, NTU) focuses on enhancing computing performance and

#### 4.2.4 Sustainability

Sustainability refers to the ability to meet our planet's growing needs without compromising the ability of future generations to live a high quality life, ensuring long-term ecological, social, and economic balance. This involves managing resources in a way that minimizes environmental impact and supports economic resilience for the future. AI plays a critical role in advancing sustainability, aligning with Singapore's Net Zero 2050 goal, especially on themes like decarbonization, energy management, and urban sustainability<sup>125</sup>. Four key areas emerge - Future Grid, Carbon Management, Electrification, and Urban Nexus. These are pillars of Singapore's low-carbon roadmap, including Singapore's 2030 green plan<sup>126</sup>, and exploring AI's contributions to grid resiliency, carbon capture, electrification, and urban management are critical technical challenges that need to be addressed.

**Prof. Madhavi Srinivasan** (NTU) is the Executive Director of Energy Research Institute (ERI@N) and NTU's Sustainability office. Her research focuses on circular economy with an emphasis on novel sustainable energy storage solutions and recycling of e-waste,

energy efficiency across various scales through system-level design, optimization, and resource management, using deep learning and hardware-software co-optimization. **Dr. Xiaoli Li** (I2R, A\*STAR) and his team have advanced the semiconductor industry by integrating deep learning with 3D X-ray microscopy for defect detection and metrology, and creating a deep learning-aided circuit design process to reduce costs. They also apply deep learning in co-designing heterogeneous integrated packages, optimizing chiplet placements and package configurations to improve design efficiency<sup>124</sup>.

which could be potentially be augmented by the application of data analytics and AI<sup>127,128</sup>. As Chief Sustainability Officer of A\*STAR, **Prof. Yeoh Lean Weng** sets A\*STAR's strategy to address Singapore's demands and opportunities in sustainability, including regulation. He has discussed the role of AI in not only providing solutions, but also challenges such as energy needs for data centres. **Assoc. Prof. Yan Xu** (NTU) works on optimization of renewable energy power systems, data-analytics for smart grid and the potential for electric vehicles, all done through a data-driven AI approach<sup>129</sup>. **Prof. Vish Vishwanathan** (NTU) has worked on leveraging AI for carbon management and decarbonization. **Dr. Zhiquan Yeo** (SIMTech, A\*STAR) focuses on industrial sustainability, ranging from Life Cycle Assessment (LCA) analyses to data science applications for enabling industrial symbiosis and circular economy. **Prof. Koh Lian Pin** (NUS) and **Prof. Lee Poh Seng** (NUS), working in the AI+X institute in NUS, are addressing urban transportation and mobility in the context of safety and sustainability, in collaboration with **Prof. Anthony Tung** (NUS).

### 4.3

## NATURAL SCIENCES – CHEMISTRY, PHYSICS AND EARTH/CLIMATE

The Natural Sciences encompass multiple disciplines like chemistry, physics, earth, and climate sciences, each dedicated to exploring the fundamental principles and processes governing the natural world. While chemistry focuses on the composition, structure, and properties of matter, physics investigates the laws of nature and the universe, earth sciences study the physical constitution of the Earth and its atmosphere, and climate sciences examine the climate system and its variations over time. While they are each immensely broad fields, we illustrate below some key examples of research in some areas relevant to Singapore that can benefit from a sustained effort in data-driven approaches and the development of an *AI for Science* framework.

#### 4.3.1 Physics

Physics is a branch of natural science that focuses on the study of non-living systems, encompassing the fundamental principles and laws governing the behaviour of matter and energy in the universe. AI can significantly aid physics studies by processing vast datasets, optimizing experiments, and enhancing simulations<sup>130,131</sup>. In data analysis, AI excels at identifying patterns and extracting insights from complex experimental or observational data. It accelerates the optimization of experimental parameters, reducing resource-intensive trial and error. In simulations, AI improves modelling accuracy, enabling researchers to explore intricate physical phenomena more efficiently<sup>132,133</sup>. Furthermore, AI-driven robotics automates experimental tasks, while smart sensors enhance the precision of scientific instruments. Collectively, these applications empower physicists with

tools to streamline research processes, uncover hidden patterns, and accelerate scientific discoveries<sup>88</sup>.

As examples of AI-assisted physical studies, **Prof. Juan-Pablo Ortega Lahuerta** at NTU Physics has developed structure-preserving machine learning algorithms for applications in control theory, mechanical engineering, and mathematical physics. **Prof. Baile Zhang** (NTU) is using AI algorithms for topological classifications of matter phases<sup>134</sup>. **Prof. Wai Lee Chan** (MAE, NTU) and **Prof. Adams Wai Kin Kong** (SCSE, NTU) work on physics-informed artificial intelligence (PiAI), focusing on the topic of computational fluid dynamics (CFD)<sup>135</sup>. Prof. Duane Loh (NUS) has developed machine learning frameworks that elucidate structural motifs in disordered materials, pioneering open-source computational lenses for lens-less imaging with X-ray lasers.

### 4.3.2 Chemistry

AI is already revolutionizing many aspects of chemical science by automating laboratory tasks, optimizing chemical reactions, facilitating inverse material design, and generating new chemical structures. By handling repetitive processes and analysing vast datasets, AI not only accelerates drug and innovative materials discovery, but also contributes to more sustainable chemical processes. The transition from molecule/material discovery to viable, sustainable manufacturing is a scientific challenge that demands a multi-physics and multi-scale modelling approach, which couples the ordinary differential equations (ODEs) describing complex reaction chemistry at the molecular scale to the partial differential equations (PDEs) describing, for example, turbulent fluid flows and system-level energy transport at macroscopic length scales. AI is revolutionizing scale-cognizant HPC simulations of such problems, thus enabling the prospect of a priori generative design of novel and sustainable manufacturing technologies. This transformation in chemical and engineering research, powered by AI,

promises to enhance prediction accuracy, positioning AI as a critical partner in driving future breakthroughs in the field, alongside significant commercial potential.

**Prof. Sergey Kozlov** (NUS ChBE) and **Prof. Tej Choksi** (NTU ChBE) have worked on machine learning accelerated ab-initio simulations for atom-level modelling, while **Prof. Jianwen Jiang** (NUS ChBE) uses Molecular Dynamics. **Prof. Gianmarco Mengaldo** and **Prof. Karl Erik Birgersson** (NUS MechE) is an expert at CFD and systems, with scalable manufacturing being done by **Prof. Saif Khan** (NUS ChBE). **Prof. Wei Chen** (NUS) and **Prof. Lu Jiong** (NUS) are leveraging AI in chemistry to enhance molecular and materials research, including in-situ spectroscopy through a data-driven approach. Prof. Chen specializes in interface engineering for nanocatalysis and electronics<sup>136</sup>, while Prof. Lu develops AI-driven tools like the Chemist-Intuited Atomic Robotic Probe (CARP)<sup>137</sup> for precise material design. Their work advances innovations in catalysis, quantum materials, and electronic devices.

### 4.3.3 Earth and climate

Various groups in NUS and NTU conduct systems-level research in climate change and sustainability. Expertise lies in downscaling global climate models to regional and local models to understand microclimate, including extremes in temperature and rainfall, sea level rise, as well as urban heat island effect. Satellite images from the region, along with observed datasets, are valuable for validating climate models used for downscaling. Other related research activities at NUS' Environmental Research Institute include studies on the effects of land use change on climate, atmospheric monitoring of aerosols, and a wide range of hydrological studies.

AI algorithms enable the processing and analysis of vast amounts of satellite imagery, aerial data, and ground-level observations. Significant advancements and challenges are being addressed at NTU. Research at NTU employing AI include but are not limited to (1) forecasting and downscaling climate and sea-level rise projections; (2) improving humanitarian assistance and disaster relief; (3) predicting the impact of climate-related threats such as invasive species, habitat destruction, and the spread of disease; (4) analyzing patterns in vegetation health, deforestation rates, urban expansion, water resource availability, and other key indicators of environmental degradation or unsustainable practices; (5) projecting the economic impacts of climate events; and their projected economic impacts; (6) creating smarter cities with reduced emissions and improved resilience to extreme weather events.

## 4.4

# AI FRAMEWORKS – MATHEMATICS, THEORY AND MODEL DEVELOPMENT

Fundamental development at the intersection of mathematics, hypothesizing and codifying via AI is paramount for advancement of the AI for Science paradigm<sup>8,138</sup>. For example, different from the static setting in classical machine learning contexts where data was passively generated/collected (e.g. nature images in ImageNet), in scientific applications the data used to train models have to be judiciously generated for the training process – often requiring an interplay of the learning algorithm and high-throughput experimentation or scientific computing. Further, the translation of empirical correlations into interpretable models is an active area of exploration. Researchers are also exploring the use of mathematical principles and physical models to generate synthetic data and enhance the credibility and understanding of machine learning models. Addressing the fundamental limits of hybrid models trained with adaptively generated samples, improving the reduction of correlations into phenomenological models, and refining the integration of physical principles into AI architectures are crucial challenges<sup>139</sup>. Finally, understanding machine learning processes as guided by underlying physical principles remains an ongoing endeavour. To fully realize the potential of AI for Science, collaboration and innovation across these domains are essential in shaping the future of scientific discovery and knowledge generation.

NUS' Centre for Data Science and Machine Learning (CDSML) has researchers (e.g. **Adrian Roellin**, **Alex Thiery**, **Zhigang Yao**, **Ying Chen**, **Hui Ji**, **Jonathan Scarlett**, **Thompson Tong**, **Wanjie Wang**) working on the exploitation of machine learning in data analytics, feature identification, data assimilation, uncertainty quantification, and high-dimension systems<sup>140</sup>. Also at NUS, **Prof. Qianxiao Li**, **Prof. Zuwei Shen**, and **Prof. Yong Sheng Soh** are interested in the mathematics of machine/deep learning. **Prof. Dario Poletti** (SUTD), **Prof. Leong Chuan Kwek** (NIE), **Prof. Dimitris Angelakis**, **Prof.**

**Jiangbin Gong**, all affiliated with the Centre for Quantum Technologies (CQT)<sup>141</sup>, work on quantum complex systems, quantum many-body systems, quantum computing algorithms and quantum machine learning. **Prof. Alvin Chua** (NUS) is using machine learning in his research on gravitational wave astronomy. **Prof. N. Duane Loh** (NUS) has been developing machine learning to coarse-grain complex non-reciprocal many-body dynamics, identify transient quasi-stable states in disorder-order transitions, and reverse sample damage in high-resolution microscopy. **Prof. Xavier Bresson** from School of Computing at NUS is a leading researcher in the field of graph neural networks (GNNs)<sup>142,143</sup>, focusing on developing advanced computational methods for analysing complex data structures across various applications, including social networks, neuroscience, and computer vision and can potentially be expanded to symmetry-aware systems in materials and physics. **Prof. Kenji Kawaguchi** (NUS) has worked on physics-inspired neural networks and his research interests bridge both theoretical and practical aspects of deep learning and machine learning [ref]<sup>5,144</sup>.

**Prof. Luke Ong's** (NTU) research spans areas such as computation semantics, programming languages, and verification, with significant contributions to game semantics, and higher-order model checking. He is recognized for integrating semantics with automated verification, and his current research focuses on computer security, higher-order logic, and probabilistic programming. **Prof. Ong Yew Soon's** (CCDS, NTU & A\*STAR) work, noted for its innovative use of evolutionary algorithms, seeks to solve complex optimization problems, thereby impacting engineering and environmental sciences. **Prof. Bryan Low** (CSE, NUS) focuses on the development of active learning and experimental design machine learning models that can predict and optimize outcomes in resource-limited settings, a critical aspect of

sustainability science. Similarly, **Prof. Ng See Kiong's** (CSE, NUS) efforts in network science and big data analytics aim at unravelling the complexities of social and biological networks, contributing to advancements in epidemiology and public health. **Prof. Bo An's** (CCDS, NTU) contribution lies in his pioneering work on multi-agent systems, enhancing collaborative AI capabilities to address strategic decision-making challenges in complex environments. **Prof. Guan Cuntai** (CCDS, NTU) has made significant strides in the field of neural engineering, leveraging AI to develop brain-computer interfaces that promise revolutionary impact in healthcare and neurotechnology. Together, these scholars, along with the strong Computer Science school/departments at NUS, NTU and SUTD exemplify the vital role that AI plays to advance fundamental scientific knowledge, showcasing the exciting paradigm of "AI for Science" in leading to novel discoveries and solutions that span disciplinary boundaries.

## 4.5 FINANCE

Financial services refer to a wide range of economic services related to capitals, which are delivered by the financial industry. It encompasses a wide variety of business areas including hedge funds, banks, insurance companies and accountancy firms. With billions of users and investors worldwide, the financial industry has become a paramount pillar for Singapore as an international financial hub. In the last decade, we have witnessed the significant development of AI-powered financial applications due to advanced AI techniques' capability of analysing and processing large volume of high dimensional financial data<sup>146</sup>. There are fruitful AI-based Fintech applications in quantitative investment, mobile banking, blockchain, credit rating and finance regulations, which leverage machine learning approaches. Complementarily, the field of economics examines how theories and models influence

In addition, the strong foundation established by AI Singapore<sup>145</sup> is expected to complement the AI4SCI initiative by building upon and fostering further collaborations between academia, industry, and government to solve complex, real-world problems through innovative AI research and applications. Their focus on developing AI talent and cutting-edge technologies has not only advanced the science of AI itself but is also providing the tools and expertise needed to fuel scientific discoveries across various domains, thereby bridging the gap between theoretical AI advancements and practical scientific applications.

A collaborative effort would, therefore, be well-placed to uncover new insights into the fundamental structures and dynamics of various scientific domains, paving the way for innovative solutions and breakthroughs in mathematics and the natural sciences.

social structures, businesses, policies, and individual behaviours, focusing on the allocation of resources, wealth distribution, and the impact of economic activities on societal and economic well-being.

The Asian Institute of Digital Finance (AIDF) focuses on leveraging artificial intelligence (AI) to propel advancements in financial technologies, aiming to address complex financial challenges through innovative AI applications. This encompasses developing AI-driven solutions for financial inclusion, risk management, smart banking, and regulatory technologies. The institute's work extends to using AI to analyse financial data, predict market trends, and optimize financial services and products and could serve as a platform for integrated data to be used for the challenges listed above. The FinTech Lab led by **Prof. Hahn Jungpil** (NUS) and **Prof. Ke-**

**Wei Huang** (SCE, NUS & Director, AIDF) are examples of such leadership in this sector. **Prof. Bo An** (CCDS, NTU) has led a project to explore AI for financial services with a focus on financial trading. The team has released a reinforcement learning based trading platforms called TradeMaster with over 1k users<sup>147</sup>. Many advanced deep learning, reinforcement learning and ensemble learning algorithms are proposed for various tasks including high frequency trading, portfolio management and intraday trading. Recently, there is ongoing research trying to use powerful large language models (LLMs) in the finance domain as well. **Prof. Xiaohui Bei** (SPMS, NTU) integrates AI with economic theory and computational economics to enhance auction design and the fair division of resources, developing frameworks that improve model efficiency while ensuring fairness. **Prof. Ariel Neufeld** (SPMS, NTU) also at SPMS in NTU has

developed deep learning based algorithms together with their theoretical justifications which help the financial and insurance industry to price and hedge financial derivatives, to optimize portfolios, and to reduce risk. **Dr. Rajaraman Kanagasabai** and **Dr. Adam Westerski** (A\*STAR I2R), lead a team of data scientists to develop an award-winning platform for procurement data analytics using a broad range of data-driven AI techniques, that is capable of detecting patterns at multiple levels, and also across multiple procurement systems.

By applying advanced AI methodologies to analyse financial data, predict market movements, and create equitable economic models, one can envision pioneering a data-driven, computational approach that significantly enhances decision-making and policy development in financial ecosystems.

## 4.6 EDUCATION

AI is transforming education by offering personalized learning experiences, automating administrative tasks, and providing intelligent tutoring systems. A shining example of that is from Sal Khan, founder of Khan Academy<sup>148</sup> and Khanmingo, which pairs generative AI with a graphical user interface as a learning aid. By analysing individual learning patterns, AI enables customized instruction that adapts to each student's needs, enhancing engagement and comprehension. Additionally, AI streamlines administrative processes, allowing educators to focus more on teaching and student interaction. However, the integration of AI in education also presents challenges, including ethical considerations, data privacy concerns, and the need for teacher training to effectively utilize these technologies<sup>149</sup>. As AI continues to evolve, its role in education is expected to expand, offering new opportunities to improve learning outcomes and educational equity to foster knowledge creation and critical thinking in educational settings.

AI for Science can enhance education by equipping students and researchers with essential AI skills tailored to scientific inquiry, promoting interdisciplinary learning, and helping learn the skills necessary to accelerating scientific discoveries. Workshops and training programs could foster AI literacy in scientific contexts, allowing learners to apply AI techniques to solve complex problems across domains like chemistry, biology, and physics. By integrating AI tools into science education, students gain hands-on experience with real-world applications, preparing them for future scientific challenges and innovation<sup>150</sup>. AI technologies, such as deep learning models, can also serve as tools to automatically detect and categorize pedagogical features in lecture recordings, providing real-time feedback to enhance student learning. AI can act as a means to personalize learning experiences and maintain student interest.

Assoc. Prof. Seng Chee Tan (NTU) was among the pioneers who helped to set up the Centre for Research and Development in Learning (CRADLE) at Nanyang Technological University, Singapore has studied the use of Generative AI towards enabling sustainable student discourse and knowledge creation in education. Assoc. Prof. Ben Leong (NUS) acts as the director of the Centre for Computing for Social Good and Philanthropy (CCSGP) as well as director of the AI Centre for Educational Technologies (AICET) and has focused recently on LLM-based AI agents for providing feedback to students<sup>151</sup>. Dr. Alwyn Lee (NIE) has worked on learning analytics and machine learning for understanding knowledge building discourse within computer-supported collaborative learning (CCSL). Dr. Lim Fun Siong (NTU), as the head of the Centre for Applications of Teaching & Learning Analytics for Students (ATLAS) has been exploring the use of generative AI in classroom teaching and evaluation and developing an end-to-

end data and AI platform to facilitate research and application development. In the area of data analytics, Assoc. Prof. Carol Hargreaves (NUS) has been studying the variation of student learning outcomes influenced by different demographic factors and the benefits of synthetic data. Dr. Nancy Chen (CFAR & I2R, A\*STAR) is active on developing multimodal, multilingual large language models in education.

Overall, integration of AI in education requires continued exploration and integration of teaching pedagogy and the scientific analysis of various aspects. By addressing challenges such as data privacy, ethical considerations, and the need for institutional support, AI has the potential to significantly enhance educational outcomes and support lifelong learning. An AI4Sci program can foster collaboration between educators, researchers, and institutions to effectively leverage AI for the benefit of all learners.

## 4.7

# SOFTWARE AND SECURITY

AI4SCI is expected to play a critical role in software and security by leveraging fundamental scientific principles to enhance functionality and resilience<sup>152, 153</sup>. In software, AI-driven techniques optimize code development, automate testing, and enable predictive maintenance, improving quality and reducing time to deployment. For security, AI enhances threat detection, anomaly identification, and vulnerability assessment by analyzing vast amounts of data and adapting to evolving cyber threats in real time. At the intersection of these fields, scientific methodologies guide the development of robust AI algorithms, ensuring transparency and scalability while addressing ethical concerns. This synergy between AI, software, and security fosters innovative solutions to complex problems in both domains.

Key challenges include the integration of computational tools, data protection, and scientific methodologies to advance

research and innovation. This field addresses challenges such as safeguarding intellectual property, improving software reliability, and enhancing information system security, fostering interdisciplinary approaches to tackle emerging technological and digital threats; take for instance, cryptography, which leverages scientific principles to address complex security challenges. **Prof. David Lo** (SMU) works at the intersection of software engineering, cybersecurity, and data science<sup>154</sup>. His research aims to enhance software quality, security, and developer productivity by analyzing various software artifacts. Home Team X (HTX) is a statutory board within the Ministry of Home Affairs dedicated to advancing science and technology for homeland security. HTX operates in five key areas: solving crimes, enhancing public safety and security, saving lives, securing borders, and safeguarding data and systems. **Dr. Yidi Yuan** (HTX) and **Dr. Wong Swee Liang** (HTX)<sup>155</sup> amongst

others in HTX have led the exploration of generative AI, vision algorithms and Natural Language Processing (NLP) for key security issues including spam detection, CCTV face recognition and speech recognition. **Prof. Lwin Khin Shar** (SMU) has been addressing cyberbiosecurity, defined as vulnerabilities to unwanted surveillance, intrusions, and malicious activities within or at the interfaces of combined life and medical sciences, cyber systems, cyber-physical systems, and infrastructure. **Prof. Jun Sun** (SMU) works on AI for software science, especially in the area of formal modeling, analysis and synthesis, while **Prof. Christoph Truede** (SMU) works on code quality and development efficiency, focused on AI-driven research into software metrics and development processes. **Prof. Reza Shokri** (NUS) is exploring the intricate intersection of AI, privacy, and scientific endeavors. **Prof. Prateek Saxena** (NUS)

works in the space of computer security and privacy, especially interested in principled and algorithmic approaches to practical security problems such as automatic program generation. **Prof. Yang Liu** (NTU) is advancing AI applications in security and software engineering, emphasizing the significance of systematic approaches and leveraging Language Model (LLM) capabilities. **Prof. Sudipta Chattopadhyay** (SUTD) is an expert on secure computing, software design as well as network and automated systems.

Overall, there are diverse applications of AI across various domains, particularly in software and security, critical for Singapore's needs in both the near-term and future. From enhancing software development processes to bolstering cybersecurity measures, the potential for AI-driven fundamental innovations is vast.

## 4.8

# ROBOTICS

Robotics and AI have co-evolved since their birth. In the early days, robot intelligence and artificial intelligence were synonymous. Recent advances in Foundation Models and Generative AI opened explosive opportunities for new-generation robot systems, with capabilities unimaginable before: open-world perception, natural-language communication with humans, common-sense reasoning, dexterous manipulation, etc. Such capabilities are critical for tasks in sectors such as manufacturing, healthcare, facilities management and logistics. AI allows robots to learn from data (collected from both physical and virtual worlds) to improve performance in dynamic and uncertain environments, working for and with humans as assistants and collaborators, driving innovation and efficiency in scientific research and commercial adoption.

Key challenges for the application of AI in Robotics (AI x Robotics) include (i) new robot architectures for data-driven robotics, (ii) data-scarcity, (iii) ability to generalize over broad

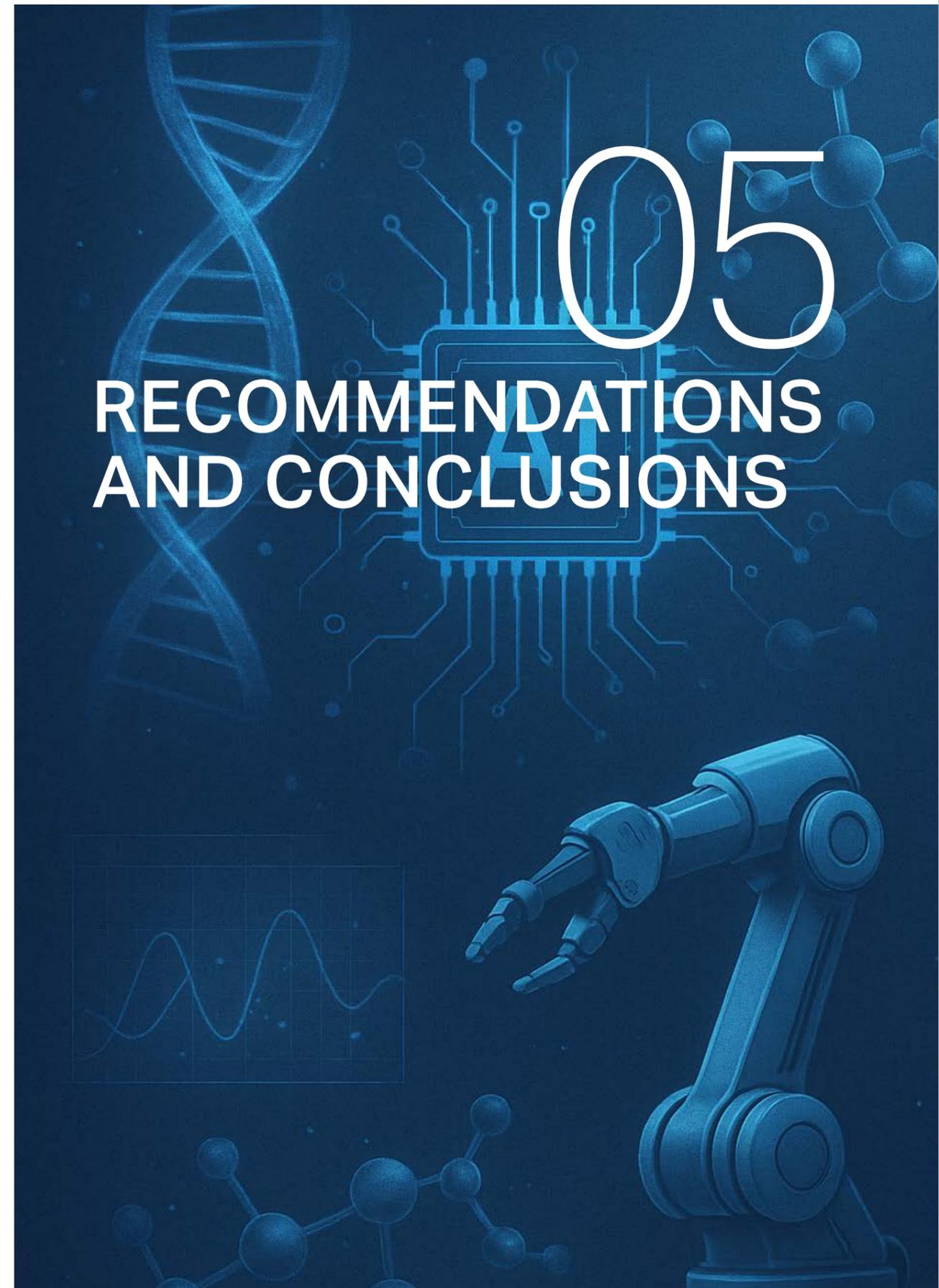
task domains, (iv) task model specification and (v) robustness & safety. The classic robot architecture of sense-reason-act was designed more than half a century ago and cannot benefit effectively from the internet-scale data available today. The recent end-to-end architecture of mapping sensor input directly to motor commands through a monolithic neural network is impossibly data-hungry and generalizes poorly over even slightly novel task domain. New robot architectures whose performance scale up gracefully with increasing data are critically needed. While natural language processing and computer vision have benefited enormously from the internet-scale data availability, robotics faces a fundamental difficulty: the high cost of gathering large amounts of data from the physical world. With sufficient data, robots today can already succeed in many narrowly specified tasks that are difficult or dangerous for humans. Generalization is, however, critical as it enables robots to apply learned knowledge to new, unseen tasks

and environments, ensuring adaptability and robustness in diverse real-world scenarios. Specifying models for robot tasks is difficult because of the inherent complexity and variability of real-world environments, which present robots with a wide range of unpredictable situations. Robustness ensures reliable performance and resilience in the face of unpredictable and dynamic real-world conditions, thereby enhancing safety and operational effectiveness. AI foundation models are prone to hallucination. Their negative, and potentially dangerous consequences in robotics must be contained for adoption of any robot systems in natural human environments.

**Prof. David Hsu** (NUS) works in the intersection of robotics and AI. He is recently working on augmenting foundation models with classic AI reasonability for scalable and robust robot reasoning in natural human environments. **Prof. Lee Wee Sun** (NUS) works in machine learning and planning with applications in embodied AI and robotics. **Prof. Harold Soh** (NUS) works on machine learning / AI for human-robot interaction. He is interested in developing trustworthy robots that are both performant and robust in interactive settings. **Prof. Lin Shao** (NUS) works on robot learning, robotic manipulation, and world models, focusing on developing machine learning methods and systems to train robots to intelligently perform diverse manipulation tasks, enabling them to understand and interact effectively with complex physical environments. **Prof. Ang Wei Tech** (NTU)'s research focuses on robotics technology in biomedical applications, and assistive and rehabilitation technologies.

His team holds the world's largest human ability database that will be used for RFM for assistive applications. **Prof. Wang Ziwei** (NTU) works in embodied AI with robotic foundation models. He is currently working on designing generalizable mobile and bimanual manipulation policy across diverse tasks, environments and embodiments. **Dr. Zhang Jingbing** (A\*STAR/ARTC) leads the robotics R&D division at ARTC and SIMTech. He is currently working on domain-specific robotics foundation models and embodied AI designed to generalize across common tasks in smart and sustainable manufacturing. **Dr. Cheston Tan** (A\*STAR/CFAR) works on neuro-symbolic visual reasoning and intuitive psychology in embodied simulations, as well as data- and parameter-efficient training of foundation models. **Dr. Kenneth Kwok** (A\*STAR/IHPC) who leads the Cognitive and Human-like AI group working on human behaviour understanding and modelling for Human-Robot Interaction. **Dr. Wu Yan** (A\*STAR/I2R) works in service and assistive robotics, dexterous manipulation, robot learning and human-robot-interaction.

Robotics is important to Singapore considering the challenge of a reducing workforce from our rapidly aging population and declining birthrates. For robots to be widely deployed in Singapore and contribute to our economy, there is a need for SG to keep pace with the new AI x Robotics space and potentially lead in selected application domains. The good news is: we do have the required ingredients to meet this challenge, but we will need coherent and focused strategies to develop AI x Robotics capabilities to deliver economic and societal impact to Singapore.



# 05 RECOMMENDATIONS AND CONCLUSIONS

This FRC report outlines Singapore's focus on developing and integrating advanced AI techniques into basic science research and their applications across various scientific domains. Led by NRF and supported by key institutions such as NUS, NTU, A\*STAR, SUTD and SMU, the AI4SCI Initiative aims to launch a nationwide research endeavour harnessing AI technologies to accelerate scientific discovery and innovation in Singapore. Through a series of workshops, consultations, surveys and investigations, the Initiative has identified numerous opportunities and challenges in incorporating AI into scientific workflows.

The Initiative's strategy focuses on building a strong AI4SCI ecosystem in Singapore, promoting interdisciplinary collaboration between AI experts and research scientists, and addressing talent and capacity gaps, particularly where AI can enhance traditional scientific methods. It places special emphasis on using AI to tackle complex and transformative challenges across large datasets, leveraging high-throughput experiments and advanced computational models. This will create a cohesive framework that effectively integrates policy support, computational resources, and collaboration among academia, industry, and government agencies. The ultimate goal is to bring together top AI experts and talented research scientists to position Singapore as a global leader in AI-driven science and innovation.

As the Initiative unites researchers from diverse institutions with varying expertise, two critical components are vital to the success of Singapore's AI4SCI initiative. First, **cross-domain learning**, which involves sharing insights, methodologies, and problem-solving approaches that can be applied across multiple scientific disciplines. Second, **identifying common challenges for AI algorithms across different fields**. Below, we highlight the key cross-domain learnings and common challenges that Singapore AI4SCI Initiative should focus on in its future development and implementation.

## 5.1

# CROSS-DOMAIN LEARNING AND COLLABORATION

### 5.1.1 Knowledge integration and generalization

An essential part of cross-domain learning is the sharing of scientific insights between researchers from different fields. Given the multi-disciplinary nature of Singapore's AI4SCI Initiative, such cross-domain learning is crucial for effective knowledge transfer. The initiative therefore needs to not only enable, but actively promote and support interdisciplinary collaborations between researchers from academia, across AI and domains, as well as industry. These collaborations should focus

on how models and techniques developed in one domain can be adapted to accelerate research in others. For example, can AI models used in healthcare be repurposed for environmental science? By sharing methodologies and insights across domain boundaries, researchers may be able to address complex challenges more effectively, thereby enhancing the overall impact of AI on scientific discovery and innovation.

### 5.1.2 Data integration and standardization

A key trend across domains is the growing volume of data being generated. To address this, it is crucial to create a centralized, open-access platform that integrates diverse scientific datasets from multiple disciplines, encouraging collaboration among researchers, institutions, and industry. The platform should utilize AI-driven tools to automate data collection, cleaning, and standardization, ensuring seamless interoperability. Standardizing data across

domains would facilitate comparisons and benchmarking, whilst potentially reducing inconsistencies. By leveraging cross-domain techniques, researchers can develop robust AI systems capable of managing complex, multi-modal data, leading to more effective knowledge transfer between disciplines and enhancing the potential for innovation and discovery across scientific fields.

### 5.1.3 Resource optimization

Given the limitations in computing and data resources, optimizing their use is crucial for the success of AI developments in science. Within the Initiative's cross-domain learning framework, resource optimization can be achieved by applying AI tools and techniques across various fields to maximize efficiency. By sharing computational models and methods between disciplines such as biology, physics, and social sciences, researchers can avoid duplicating efforts and streamline processes.

For example, AI algorithms developed for high-throughput data analysis in genomics can be adapted to process large datasets in materials science, conserving both time and computational resources. This approach enables the more efficient use of shared infrastructure and expertise, reducing costs while enhancing research output. Cross-domain learning thus promotes smarter allocation of resources, driving progress across multiple scientific domains.

### 5.1.4 AI ethics and governance

The Initiative will endeavour to focus on AI ethics and governance integrating ethical frameworks from areas like healthcare and social sciences, and apply lessons learned to other domains such as environmental science or chemistry. Implementing standard practices will ensure that AI models are not only effective but also ethically sound across disciplines. This collaborative approach helps establish unified standards for transparency, accountability, and governance, promoting trust in AI-driven research. By sharing best

practices, researchers can create AI systems that respect individual rights while accelerating innovation in multiple fields.

Within the proposed AI4SCI initiative, cross-theme intersection of the ideas above will be encouraged; the initiative will emphasize horizontal domain-agnostic learnings through establishment of datasets, protocols and data/code management via a 'gymnasium' that executes upon many of the ideas described in the table below:

Themes	Cross-Domain Learning	Data Integration	Knowledge Transfer and generalization	Resource Optimization	AI Ethics and Governance
Collab-oration	Develop adaptable AI models for multiple fields.	Build shared, standardized data platforms.	Facilitate cross-disciplinary training programs.	Streamline shared resources and workflows.	Establish shared ethical standards.
Data Integration	Enhance multi-modal data handling capabilities	Automate data cleaning and standardization.	Share methods for cross-domain data utilization.	Optimize AI tools for multi-domain data use.	Implement ethical data privacy practices.
Knowledge Sharing	Document and share best AI practices.	Enable open access to integrated datasets.	Transfer AI methods to solve new challenges.	Provide training on shared tools and resources.	Promote ethical use of shared knowledge.
Scalability	Repurpose models to scale across disciplines.	Scale analysis with centralized datasets.	Apply scalable methods across domains.	Share infrastructure for scalable solutions.	Ensure scalable systems uphold fairness.
Ethics	Integrate ethics into AI collaboration frameworks.	Standardize data governance protocols.	Address biases during knowledge transfer.	Align resource use with ethical guidelines.	Develop unified AI governance frameworks.

## 5.2

# COMMON CHALLENGES FOR AI ALGORITHMS ACROSS DOMAINS

### 5.2.1 Handling large and diverse datasets

High throughput and high-quality data are essential for AI model development, making one of the biggest challenges in modern science the management of the vast volumes of data generated across various fields. Large datasets demand substantial computational resources and infrastructure for storage and efficient processing. The diversity of data adds complexity to integration efforts, as it requires handling different formats,

inconsistencies, and missing values. Ensuring data quality and accuracy while mitigating biases and noise is crucial, especially with those domains where experimental validation is part of the scientific process. Additionally, integrating diverse data types into AI models is challenging, necessitating advanced techniques to process multi-modal data and maintain consistency across datasets.

### 5.2.2 Training AI models on limited data

While managing large-scale and high-throughput data is a significant challenge, many research areas face the opposite issue: insufficient high-quality data for reliable AI model development. With too little data, models may learn patterns specific to the small dataset rather than generalizable trends, leading to poor performance on new, unseen data. Furthermore, limited data often fails to capture the full diversity of real-world scenarios, hindering the model's ability to

be generate realistic results. To address this challenge, techniques such as data augmentation, transfer learning, and synthetic data generation can improve learning from limited data while reducing overfitting. Additionally, the development of foundation models trained on large-scale datasets, which can subsequently be fine-tuned on smaller, specialized datasets, offers a promising solution.

### 5.2.3 Bias and data imbalance

AI algorithms are highly sensitive to data biases, especially when datasets are imbalanced with over- or under-represented classes or outcomes. The core challenge in bias and data imbalance stems from uneven representation across different groups within datasets, leading to biased models. The issue of dataset bias is relatively prevalent in experimental datasets where poor results or failure may not be reported. Therefore any new experiment or data collection process undertaken within the AI4SCI Initiative should

to be designed with the aim of generating unbiased, AI-compatible datasets from the onset. For historical data collected before the beginning of the initiative or from experiments where it is impossible to achieve balance and zero bias, techniques such as data augmentation, re-sampling and algorithmic adjustments can be implemented. This can be enabled through an AI4SCI 'gymnasium' where data, best practices, models etc. can be shared.

### 5.2.4 Collaboration gap between domain experts and AI practitioners

The primary goal of Singapore's AI4SCI initiative is to unite science domain experts and AI experts. A key challenge is to address the communication gap between these two groups. Domain experts have deep knowledge in sciences but may not be familiar with AI techniques, while AI practitioners may lack an understanding of the complexities within specific scientific or industrial domains. This disconnect can lead to misalignment in defining objectives, interpreting outcomes, and setting goals. Successful integration of domain knowledge with AI approaches

requires a common language and framework to develop models that are both accurate and applicable. Overcoming these challenges necessitates effective communication through continual engagements such as AI4SCI workshops, interdisciplinary learning, and strong collaboration throughout the project. This can be facilitated via a common physical location, perhaps hosted by NRF, for funded research programs where people can intermingle, cross-pollinate ideas and resources can be shared.

## 5.2.5 AI model interpretability

Advanced science demands that scientists not only know the “what” but also understand the “why”. While AI models can often deliver high-quality predictions, the sheer number of parameters in deep network architectures poses a significant challenge to their interpretability. These models often function as “black boxes,” with little transparency into their decision-making processes. This lack of clarity undermines trust, especially in

critical sectors like healthcare and finance, where understanding the model’s reasoning is crucial for accountability and compliance. Enhancing interpretability often involves trade-offs with model complexity and performance. Techniques such as explainable AI (XAI), feature importance analysis, and physics-based hybrid models seek to improve transparency, but achieving a balance between accuracy and interpretability will be a key challenge that will need to be addressed under the AI4SCI Initiative.

## 5.2.6 Computational complexity and scalability

As the complexity of models grows (e.g., in climate simulations, chemical reactions, metagenome sequences, or drug interactions), scaling becomes a significant challenge. AI techniques such as surrogate modelling and distributed computing allow scientists to scale their models and conduct large-scale simulations or computations that were previously unfeasible. In addition, AI has the ability to handle large datasets, enabling researchers to perform simulations that

consider a broader range of variables and scenarios. However, developing scalable AI models with high computational complexity will need to account for rising energy consumption and the need for model interpretability. Targeting this for the AI4SCI initiative, within the available HPC resources in Singapore, such as NSCC and others will require efficiency, as well as complementarity and collaboration with other global approaches.

## 5.2.7 Ethical and privacy concerns

AI4SCI model development frequently involves the collection of sensitive data across fields such as medicine, healthcare, social science, and finance. A key challenge in ethical and privacy concerns stems from the extensive collection and use of personal data in AI systems. Safeguarding user privacy while maintaining transparency and accountability is difficult, as AI models often rely on large datasets containing sensitive information.

Ensuring data security and sovereignty, preventing unauthorized access, and adhering to regulations like GDPR<sup>156</sup> is crucial. Moreover, in social sciences, data biases can result in unfair or discriminatory outcomes, raising ethical concerns. Work undertaken under the Initiative’s umbrella will need to balancing innovation with respect for individual rights and public trust.

## 5.3

# CONCLUSIONS

In summary, cross-domain learnings through the series of workshops and the scoping study reveal that many scientific disciplines encounter similar challenges in data management, model development, and problem-solving, where Singapore’s AI4SCI Initiative is well-positioned to address these issues. By promoting the transfer of knowledge and techniques and data sharing and standardization across different fields, the AI4SCI Initiative can greatly accelerate AI-driven scientific discovery in Singapore and contribute to the AI for Science efforts globally. The subsections below summarize the first part of this document, which is followed by domain specific white papers and an extensive list of national and international AI for Science initiatives in the *APPENDICES*.

### 5.3.1 Importance of AI for Science

This FRC report has highlighted the need for a national and global effort to further develop the field of AI for Science. Countries such as the USA, China and the UK have already started to invest heavily in this field of research, along with large technological companies such as IBM, Google, META, Sony and Microsoft. The importance given to AI for Science research is justified by its potential to revolutionize scientific research methodologies and accelerate discoveries and knowledge generation. Innovation across the globe is already driven by advances in computational hardware, software and

modelling. At the same time, AI has proven to be very useful in practical applications such as automating decision processes and computer vision. Combining AI development and domain-specific scientific research under the Singapore AI4SCI Initiative is therefore a timely endeavour that aligns with current research and investments efforts worldwide. It is also a necessary endeavour, as accelerating the research process across different fields of basic science has the potential to accelerate the understanding of the world around us, and to tackle global challenges including the climate emergency.

### 5.3.2 Mission of Singapore’s AI4SCI Initiative

The AI4SCI Initiative aims to bring together AI experts and talented research scientists to foster creative collaboration across disciplines. Its mission is to leverage AI technologies to both identify and solve critical and urgent scientific problems. To achieve this, a domain agnostic implementation framework has been proposed. It is structured across five levels, with levels 1-2 leveraging scientific research studies, levels 3-4 necessitating foundational work in the science of AI, and level 5

representing a convergence of efforts from both AI and scientific disciplines. In addition, the Initiative aims to create a platform for cross-domain learning, data sharing and the identification of hardware resource needs for Singapore. Ultimately, AI4SCI aims to enhance Singapore’s digital ecosystem and help position the country as a global leader in AI-driven scientific research, putting Singapore at the forefront of the international research community.

### 5.3.3 Opportunities and challenges of Singapore

In this dynamic era of AI, Singapore faces both significant opportunities and challenges in developing AI for Science. On the opportunity side, the timing is ideal as AI technologies are rapidly advancing, allowing Singapore to catch up and innovate in this space. The country boasts world's top talents in both AI technology and scientific research, providing it with a strong potential to lead in AI-driven discovery. Additionally, Singapore's strong economic and research environment help attract global scientists, fostering a vibrant ecosystem for talent recruitment and collaborations.

Nevertheless, the research community and the number of research institutions in

Singapore are relatively small, limiting the scope and opportunity for broader and deep collaborations. There is still a shortage of AI talent and international-level AI4SCI leaders, given their importance to lead the multiple domain research areas. Access to and management of computing resources, particularly high-end GPUs required for training large AI models, remain challenging for the majority of AI4SCI laboratories. Furthermore, the lack of a comprehensive platform like this AI4SCI Initiative to facilitate collaboration across disciplines presents a barrier. Addressing these challenges while leveraging existing opportunities will be critical to advancing Singapore's leadership in AI-driven science.

### 5.3.4 Strong support is essential for the success of AI4SCI Initiative

For the AI4SCI Initiative to succeed and achieve its mission, strong and sustaining support from the national government and the NRF is crucial. This includes adequate budget allocation to fund key AI4SCI research projects, AI-related infrastructure, and technological innovations. Significant investments in computing resources, such as high-performance computing facilities (large-scale, high-end GPU clusters) and data storage, dedicated to the AI4SCI laboratories, are needed to manage the complexity of large AI models and applications.

Talent recruitment is equally vital, requiring efforts to attract world-level leading AI experts

and research scientists from around the globe. The involvement of these leading experts in high-level research and mentorship is also crucial for nurturing the next generation of AI researchers and leaders. Additionally, effective manpower management and coordination of the Initiative are necessary to ensure smooth collaboration across institutions and disciplines.

Overall, these combined efforts will bring significant benefits to the nation and help propel AI4SCI toward achieving its goal of positioning Singapore as a global leader in the predominant field of AI for science.





To identify the key scientific domains where AI can make significant contributions, along with domain leaders working on implementing and developing the AI for Science landscape in Singapore, we organized multiple AI4SCI workshops from February to July 2024. The process involved bringing together leading researchers, AI practitioners, and domain experts to curate topics that align with current challenges and opportunities within these fields. The workshop sessions were structured to facilitate interdisciplinary collaboration, with a combination of keynote talks, panel discussions, and breakout sessions to encourage deep dives into specific problems. Each workshop was designed not only to disseminate cutting-edge knowledge but also to foster open dialogue and shared insights among participants.

During these workshops, experts from diverse scientific backgrounds presented AI-driven solutions that address specific challenges within their domains. Common themes include the necessity of improving data accessibility and fostering collaboration between AI experts and domain scientists. Discussions lead to the identification of emerging challenges, such as the need for developing a gymnasium for AI4SCI efforts, explainable AI models, science-based metrics and the importance of cross-disciplinary skill sets in future research teams.

The insights gained from these workshops not only deepen the understanding of how AI can advance science but also inform the AI4SCI Initiative's future strategies, ensuring the continued alignment of AI research with real-world scientific applications. These insights, along with actionable recommendations, are captured by this report, with detailed domain-specific reports by the domain leaders and workshop organizers provided in the APPENDICES section at the end of this FRC report.

## REFERENCES

- National AI Strategy [Internet]. Available from: <https://www.smartnation.gov.sg/nais/>
- AI4SCIENCE AND NOBEL TURING Challenge INITIATIVE CONFERENCE [Internet]. Available from: <https://ai4science.sg/ai4sci%2Fntci-conf>
- RIE Ecosystem [Internet]. Available from: <https://www.nrf.gov.sg/rie-ecosystem/ecosystem/>
- Jumper J, Evans R, Pritzel A, Green T, Figurnov M, Ronneberger O, et al. Highly accurate protein structure prediction with AlphaFold. *Nature*. 2021;596(7873):583–9.
- Merchant A, Batzner S, Schoenholz SS, Aykol M, Cheon G, Cubuk ED. Scaling deep learning for materials discovery. *Nature*. 2023;624(7990):80–5.
- Rolnick D, Donti PL, Kaack LH, Kochanski K, Lacoste A, Sankaran K, et al. Tackling climate change with machine learning. *ACM Comput Surv*. 2022;55(2):1–96.
- RIE2025 Handbook [Internet]. Available from: <https://www.nrf.gov.sg/rie-ecosystem/rie2025handbook/>
- Wang H, Fu T, Du Y, Gao W, Huang K, Liu Z, et al. Scientific discovery in the age of artificial intelligence. *Nature*. 2023;620(7972):47–60.
- Yohsua B, Daniel P, Tamay B, Rishi B, Stephen C, Yejin C, et al. International Scientific Report on the Safety of Advanced AI. Department for Science, Innovation and Technology; 2024.
- Berman HM, Westbrook J, Feng Z, Gilliland G, Bhat TN, Weissig H, et al. The protein data bank. *Nucleic Acids Res*. 2000;28(1):235–42.
- Pearce R, Zhang Y. Toward the solution of the protein structure prediction problem. *J Biol Chem*. 2021;297(1).
- McCarthy J, Minsky ML, Rochester N, Shannon CE. A proposal for the dartmouth summer research project on artificial intelligence, august 31, 1955. *AI Mag*. 2006;27(4):12.
- Turing AM. *Computing machinery and intelligence*. Springer; 2009.
- Newell A, Simon H. The logic theory machine--A complex information processing system. *IRE Trans Inf theory*. 1956;2(3):61–79.
- Toosi A, Bottino AG, Saboury B, Siegel E, Rahmim A. A brief history of AI: how to prevent another winter (a critical review). *PET Clin*. 2021;16(4):449–69.
- Waterman DA. *A guide to expert systems*. Addison-Wesley Longman Publishing Co., Inc.; 1985.
- Hopfield JJ. Neural networks and physical systems with emergent collective computational abilities. *Proc Natl Acad Sci*. 1982;79(8):2554–8.
- LeCun Y, Bengio Y, Hinton G. Deep learning. *Nature*. 2015;521(7553):436–44.
- Blum AL, Langley P. Selection of relevant features and examples in machine learning. *Artif Intell*. 1997;97(1–2):245–71.
- Simon HA. *The sciences of the artificial*. 1969;
- Von Neumann J, others. Various techniques used in connection with random digits. *John von Neumann, Collect Work*. 1963;5:768–70.
- Benettin G, Christodoulidi H, Ponno A. The Fermi-Pasta-Ulam problem and its underlying integrable dynamics. *J Stat Phys*. 2013;152:195–212.
- Dauxois T. Fermi, Pasta, Ulam, and a mysterious lady. *Phys Today*. 2008;61(1):55–7.
- Hsu F. IBM's deep blue chess grandmaster chips. *IEEE micro*. 1999;19(2):70–81.
- Silver D, Schrittwieser J, Simonyan K, Antonoglou I, Huang A, Guez A, et al. Mastering the game of go without human knowledge. *Nature*. 2017;550(7676):354–9.
- Senior AW, Evans R, Jumper J, Kirkpatrick J, Sifre L, Green T, et al. Improved protein structure prediction using potentials from deep learning. *Nature*. 2020;577(7792):706–10.
- Kitano H. Nobel Turing Challenge: creating the engine for scientific discovery. *NPJ Syst Biol Appl*. 2021;7(1):29.
- Kryshtafovych A, Schwede T, Topf M, Fidelis K, Moulton J. Critical assessment of methods of protein structure prediction (CASP)—Round XV. *Proteins Struct Funct Bioinforma*. 2023;91(12):1539–49.
- Zhang Y. Progress and challenges in protein structure prediction. *Curr Opin Struct Biol*. 2008;18(3):342–8.
- Atz K, Grisoni F, Schneider G. Geometric deep learning on molecular representations. *Nat Mach Intell*. 2021;3(12):1023–32.
- Bronstein MM, Bruna J, Lecun Y, Szlam A, Vandergheynst P, May C V. Geometric deep learning: going beyond Euclidean data. :1–22.
- Huang S-C, Pareek A, Jensen M, Lungren MP, Yeung S, Chaudhari AS. Self-supervised learning for medical image classification: a systematic review and implementation guidelines. *NPJ Digit Med*. 2023;6(1):74.
- Chen T, Kornblith S, Norouzi M, Hinton G. A simple framework for contrastive learning of visual representations. In: *International conference on machine learning*. 2020. p. 1597–607.
- Goodfellow I, Pouget-Abadie J, Mirza M, Xu B, Warde-Farley D, Ozair S, et al. Generative adversarial networks. *Commun ACM*. 2020;63(11):139–44.
- Vaswani A. Attention is all you need. *Adv Neural Inf Process Syst*. 2017;
- Sengar SS, Hasan A Bin, Kumar S, Carroll F. Generative artificial intelligence: a systematic review and applications. *Multimed Tools Appl*. 2024;1–40.
- MacLeod BP, Parlange FGL, Morrissey TD, Häse F, Roch LM, Dettelbach KE, et al. Self-driving laboratory for accelerated discovery of thin-film materials. *Sci Adv*. 2020;6(20):eaaz8867.
- Abolhasani M, Kumacheva E. The rise of self-driving labs in chemical and materials sciences. *Nat Synth*. 2023;2(6):483–92.

- 39 Low AKY, Mekki-Berrada F, Gupta A, Ostudin A, Xie J, Vissol-Gaudin E, et al. Evolution-guided Bayesian optimization for constrained multi-objective optimization in self-driving labs. *npj Comput Mater*. 2024;10(1).
- 40 Boiko DA, MacKnight R, Kline B, Gomes G. Autonomous chemical research with large language models. *Nature*. 2023;624(7992):570–8.
- 41 Rives A, Meier J, Sercu T, Goyal S, Lin Z, Liu J, et al. Biological structure and function emerge from scaling unsupervised learning to 250 million protein sequences. *Proc Natl Acad Sci*. 2021;118(15):e2016239118.
- 42 Lin Z, Akin H, Rao R, Hie B, Zhu Z, Lu W, et al. Evolutionary-scale prediction of atomic-level protein structure with a language model. *Science* (80- ). 2023;379(6637):1123–30.
- 43 Shumailov I, Shumaylov Z, Zhao Y, Papernot N, Anderson R, Gal Y. AI models collapse when trained on recursively generated data. *Nature*. 2024;631(8022):755–9.
- 44 Towers M, Kwiatkowski A, Terry J, Balis JU, De Cola G, Deleu T, et al. Gymnasium: A standard interface for reinforcement learning environments. *arXiv Prepr arXiv240717032*. 2024;
- 45 Smart Nation Singapore [Internet]. Available from: <https://www.smartnation.gov.sg>
- 46 New AI Centre of Excellence to drive innovation in manufacturing [Internet]. Available from: <https://www.edb.gov.sg/en/about-edb/media-releases-publications/new-ai-centre-of-excellence-to-drive-innovation-in-manufacturing.html>
- 47 American Express expands Singapore Decision Science Center of Excellence [Internet]. Available from: <https://www.edb.gov.sg/en/about-edb/media-releases-publications/american-express-expands-singapore-decision-science-center-of-excellence.html>
- 48 Model AI Governance Framework for Generative AI [Internet]. Available from: <https://aiverifyfoundation.sg/wp-content/uploads/2024/05/Model-AI-Governance-Framework-for-Generative-AI-May-2024-1-1.pdf>
- 49 Singapore and France Defence Ministries to Establish First Joint R&D Lab in Singapore to Develop Artificial Intelligence Capabilities [Internet]. Available from: [https://www.mindef.gov.sg/news-and-events/latest-releases/20apr23\\_nr](https://www.mindef.gov.sg/news-and-events/latest-releases/20apr23_nr)
- 50 Arranz D, Bianchini S, Di Girolamo V, Ravet J. Trends in the use of AI in science: a bibliometric analysis. 2023. 34 p.
- 51 Eisenberg D, Marcotte EM, Xenarios I, Yeates TO. Protein function in the post-genomic era. *Nature*. 2000;405(6788):823–6.
- 52 Statello L, Guo C-J, Chen L-L, Huarte M. Gene regulation by long non-coding RNAs and its biological functions. *Nat Rev Mol cell Biol*. 2021;22(2):96–118.
- 53 Zhang Y. Protein structure prediction: when is it useful? *Curr Opin Struct Biol*. 2009;19(2):145–55.
- 54 Zhou X, Zheng W, Li Y, Pearce R, Zhang C, Bell EW, et al. I-TASSER-MTD: a deep-learning-based platform for multi-domain protein structure and function prediction. *Nat Protoc*. 2022;17(10):2326–53.
- 55 Zheng W, Wuyun Q, Li Y, Zhang C, Freddolino PL, Zhang Y. Improving deep learning protein monomer and complex structure prediction using DeepMSA2 with huge metagenomics data. *Nat Methods*. 2024;21(2):279–89.
- 56 Li Y, Zhang C, Feng C, Pearce R, Lydia Freddolino P, Zhang Y. Integrating end-to-end learning with deep geometrical potentials for ab initio RNA structure prediction. *Nat Commun*. 2023;14(1):5745.
- 57 Pearce R, Omenn GS, Zhang Y. De novo RNA tertiary structure prediction at atomic resolution using geometric potentials from deep learning. *BioRxiv*. 2022;2005–22.
- 58 Zheng W, Wuyun Q, Freddolino PL, Zhang Y. Integrating deep learning, threading alignments, and a multi-MSA strategy for high-quality protein monomer and complex structure prediction in CASP15. *Proteins Struct Funct Bioinforma*. 2023;91(12):1684–703.
- 59 Maurer-Stroh S, Debulpaep M, Kuemmerer N, De La Paz ML, Martins IC, Reumers J, et al. Exploring the sequence determinants of amyloid structure using position-specific scoring matrices. *Nat Methods*. 2010;7(3):237–42.
- 60 Selvarajoo K, Maurer-Stroh S. Towards multi-omics synthetic data integration. *Brief Bioinform*. 2024;25(3):bbae213.
- 61 Miao Y, Guo X, Zhu K, Zhao W. Biomolecular condensates tunes immune signaling at the host-pathogen interface. *Curr Opin Plant Biol*. 2023;74:102374.
- 62 Zhang H, Lim EJK, Lu L. Investigating the stability of dengue virus envelope protein dimer using well-tempered metadynamics simulations. *Proteins Struct Funct Bioinforma*. 2020;88(5):643–53.
- 63 Warner KD, Hajdin CE, Weeks KM. Principles for targeting RNA with drug-like small molecules. *Nat Rev Drug Discov*. 2018;17(8):547–58.
- 64 Venkitaraman AR. Cancer susceptibility and the functions of BRCA1 and BRCA2. *Cell*. 2002;108(2):171–82.
- 65 Kong LR, Gupta K, Wu AJ, Perera D, Ivanyi-Nagy R, Ahmed SM, et al. A glycolytic metabolite bypasses “two-hit” tumor suppression by BRCA2. *Cell*. 2024;187(9):2269–87.
- 66 Chen L, Li Y, Lin CH, Chan THM, Chow RKK, Song Y, et al. Recoding RNA editing of AZIN1 predisposes to hepatocellular carcinoma. *Nat Med*. 2013;19(2):209–16.
- 67 Chan JJ, Zhang B, Chew XH, Salhi A, Kwok ZH, Lim CY, et al. Pan-cancer pervasive upregulation of 3 UTR splicing drives tumorigenesis. *Nat Cell Biol*. 2022;24(6):928–39.
- 68 Perera AR, Warriar V, Sundararaman S, Hsiao Y, Ghosh S, Kularatnarajah L, et al. Melvin is a conversational voice interface for cancer genomics data. *Commun Biol*. 2024;7(1):30.
- 69 Tang SC, Vijayakumar U, Zhang Y, Fullwood MJ. Super-enhancers, phase-separated condensates, and 3D genome organization in cancer. *Cancers (Basel)*. 2022;14(12):2866.
- 70 Tan MH, Li Q, Shanmugam R, Piskol R, Kohler J, Young AN, et al. Dynamic landscape and regulation of RNA editing in mammals. *Nature*. 2017;550(7675):249–54.
- 71 Foo C-S, Pop C. Learning RNA secondary structure (only) from structure probing data. *bioRxiv*. 2017;152629.
- 72 Synthetic Biology for Clinical and Technological Innovation [Internet]. Available from: <https://syncti.org>
- 73 Singapore Consortium for Synthetic Biology [Internet]. Available from: <https://www.nrf.gov.sg/tech-consortia/sinergy/>
- 74 Zhao J, Cheong KH. Obfuscating community structure in complex network with evolutionary divide-and-conquer strategy. *IEEE Trans Evol Comput*. 2023;27(6):1926–40.
- 75 Sadybekov A V, Katritch V. Computational approaches streamlining drug discovery. *Nature*. 2023;616(7958):673–85.
- 76 Meng Z, Xia K. Persistent spectral--based machine learning (PerSpect ML) for protein-ligand binding affinity prediction. *Sci Adv*. 2021;7(19):eabc5329.
- 77 Experimental Drug Development Center [Internet]. Available from: <https://www.eddc.sg/>
- 78 Siau A, Ang JW, Sheriff O, Hoo R, Loh HP, Tay D, et al. Comparative spatial proteomics of Plasmodium-infected erythrocytes. *Cell Rep*. 2023;42(11).
- 79 Lam HYI, Pincket R, Han H, Ong XE, Wang Z, Hinks J, et al. Application of variational graph encoders as an effective generalist algorithm in computer-aided drug design. *Nat Mach Intell*. 2023;5(7):754–64.
- 80 Mei J-P, Kwok C-K, Yang P, Li X-L, Zheng J. Drug-target interaction prediction by learning from local information and neighbors. *Bioinformatics*. 2013;29(2):238–45.
- 81 Eldele E, Ragab M, Chen Z, Wu M, Li X. Tslanet: Rethinking transformers for time series representation learning. *arXiv Prepr arXiv240408472*. 2024;
- 82 Cancer Science Institute of Singapore [Internet]. Available from: <https://csi.nus.edu.sg/our-research/>
- 83 Truong ATL, Tan S-B, Wang GZ, Yip AWJ, Egermark M, Yeung W, et al. CURATE. AI-assisted dose titration for anti-hypertensive personalized therapy: study protocol for a multi-arm, randomized, pilot feasibility trial using CURATE. AI (CURATE. AI ADAPT trial). *Eur Hear Journal-Digital Heal*. 2024;5(1):41–9.
- 84 Wang H, Yang G, Zhang S, Qin J, Guo Y, Xu B, et al. Video-instrument synergistic network for referring video instrument segmentation in robotic surgery. *IEEE Trans Med Imaging*. 2024;
- 85 Zhou J, Jiang M, Wu J, Zhu J, Wang Z, Jin Y. MGI: Multimodal Contrastive pre-training of Genomic and Medical Imaging. *arXiv Prepr arXiv240600631*. 2024;
- 86 Milea D, Najjar RP, Jiang Z, Ting D, Vasseneix C, Xu X, et al. Artificial intelligence to detect papilledema from ocular fundus photographs. *N Engl J Med*. 2020;382(18):1687–95.
- 87 Leong Y-Y, Vasseneix C, Finkelstein MT, Milea D, Najjar RP. Artificial intelligence meets neuro-ophthalmology. *Asia-Pacific J Ophthalmol*. 2022;11(2):111–25.
- 88 Young BE, Fong S-W, Chan Y-H, Mak T-M, Ang LW, Anderson DE, et al. Effects of a major deletion in the SARS-CoV-2 genome on the severity of infection and the inflammatory response: an observational cohort study. *Lancet*. 2020;396(10251):603–11.
- 89 The AI will see you now: diagnosing mental disorders intelligently [Internet]. Available from: <https://www.ntu.edu.sg/research/research-hub/the-ai-will-see-you-now-diagnosing-mental-disorders-intelligently>
- 90 Tasnim S, Lim PXH, Griva K, Ngeow J. Identifying the psychosocial barriers and facilitators associated with the uptake of genetic services for hereditary cancer syndromes: a systematic review of qualitative studies. *Health Psychol Rev*. 2024;1–28.
- 91 Vaser R, Šikić M. Time-and memory-efficient genome assembly with Raven. *Nat Comput Sci*. 2021;1(5):332–6.
- 92 Dalakoti M, Wong S, Lee W, Lee J, Yang H, Loong S, et al. Incorporating AI into cardiovascular diseases prevention--insights from Singapore. *Lancet Reg Heal Pacific*. 2024;48.
- 93 Singapore Eye Lesion Analyser [Internet]. Available from: <https://www.synapxe.sg/healthtech/health-ai/selena>
- 94 MOH TRUST [Internet]. Available from: <https://trustplatform.sg/>
- 95 Chen X, Soh BW, Ooi Z-E, Vissol-Gaudin E, Yu H, Novoselov KS, et al. Constructing custom thermodynamics using deep learning. *Nat Comput Sci*. 2024;4(1):66–85.
- 96 Tom G, Schmid SP, Baird SG, Cao Y, Darvish K, Hao H, et al. Self-driving laboratories for chemistry and materials science. *Chem Rev*. 2024;124(16):9633–732.
- 97 Hippalgaonkar K, Li Q, Wang X, Fisher JW, Kirkpatrick J, Buonassisi T. Knowledge-integrated machine learning for materials: lessons from gameplaying and robotics. *Nat Rev Mater*. 2023;8(4):241–60.
- 98 Liu P, Chen ZN. Full-range amplitude--phase metacells for sidelobe suppression of metalens antenna using prior-knowledge-guided deep-learning-enabled synthesis. *IEEE Trans Antennas Propag*. 2023;71(6):5036–45.
- 99 Jin P, Xu L, Xu G, Li J, Qiu C-W, Huang J. Deep Learning-Assisted Active Metamaterials with Heat-Enhanced Thermal Transport. *Adv Mater*. 2024;36(5):2305791.
- 100 Fang G, Ma X, Song M, Mi MB, Wang X. Depgraph: Towards any structural pruning. In: Proceedings of the IEEE/CVF conference on computer vision and pattern recognition. 2023. p. 16091–101.

- 101 Dong Z, Jin L, Rezaei SD, Wang H, Chen Y, Tjiptoharsono F, et al. Schrödinger's red pixel by quasi-bound-states-in-the-continuum. *Sci Adv*. 2022;8(8):eabm4512.
- 102 Giustino F, Lee JH, Trier F, Bibes M, Winter SM, Valentini R, et al. The 2021 quantum materials roadmap. *J Phys Mater*. 2021;3(4):42006.
- 103 Wang Z, Chen A, Tao K, Han Y, Li J. MatGPT: A Vane of Materials Informatics from Past, Present, to Future. *Adv Mater*. 2024;36(6):1-44.
- 104 Fauseweh B. Quantum many-body simulations on digital quantum computers: State-of-the-art and future challenges. *Nat Commun*. 2024;15(1):2123.
- 105 Li X, Mitchell S, Fang Y, Li J, Perez-Ramirez J, Lu J. Advances in heterogeneous single-cluster catalysis. *Nat Rev Chem*. 2023;7(11):754-67.
- 106 Hai X, Zheng Y, Yu Q, Guo N, Xi S, Zhao X, et al. Geminal-atom catalysis for cross-coupling. *Nature*. 2023;622(7984):754-60.
- 107 Fang H, Mahalingam H, Li X, Han X, Qiu Z, Han Y, et al. Atomically precise vacancy-assembled quantum antidots. *Nat Nanotechnol*. 2023;18(12):1401-8.
- 108 Duric T, Chung JH, Yang B, Sengupta P. Spin-1/2 kagome Heisenberg antiferromagnet: Machine learning discovery of the spinon pair density wave ground state. *arXiv Prepr arXiv240102866*. 2024;
- 109 Chee CH, Leykam D, Mak AM, Bharti K, Angelakis DG. Resource-Efficient Hybrid Quantum-Classical Simulation Algorithm. *arXiv Prepr arXiv240510528*. 2024;
- 110 Efthymiou S, Orgaz-Fuertes A, Carobene R, Cereijo J, Pasquale A, Ramos-Calderer S, et al. Qibolab: an open-source hybrid quantum operating system. *Quantum*. 2024;8:1247.
- 111 Zhao L, Zhao Z, Rebentrost P, Fitzsimons J. Compiling basic linear algebra subroutines for quantum computers. *Quantum Mach Intell*. 2021;3(2):21.
- 112 Qian Y, Du Y, He Z, Hsieh M-H, Tao D. Multimodal deep representation learning for quantum cross-platform verification. *Phys Rev Lett*. 2024;133(13):130601.
- 113 Wang X, Du Y, Tu Z, Luo Y, Yuan X, Tao D. Transition role of entangled data in quantum machine learning. *Nat Commun*. 2024;15(1):3716.
- 114 Song M, Narasimhachar V, Regula B, Elliott TJ, Gu M. Causal classification of spatiotemporal quantum correlations. *Phys Rev Lett*. 2024;133(11):110202.
- 115 Takagi R, Tajima H, Gu M. Universal sampling lower bounds for quantum error mitigation. *Phys Rev Lett*. 2023;131(21):210602.
- 116 Mao X, Chen P. Inter-facet junction effects on particulate photoelectrodes. *Nat Mater*. 2022;21(3):331-7.
- 117 Fu B, Mao X, Park Y, Zhao Z, Yan T, Jung W, et al. Single-cell multimodal imaging uncovers energy conversion pathways in biohybrids. *Nat Chem*. 2023;15(10):1400-7.
- 118 Zhang L, Shao S. Image-based machine learning for materials science. *J Appl Phys*. 2022;132(10).
- 119 Sankaran J, Wohland T. Current capabilities and future perspectives of FCS: super-resolution microscopy, machine learning, and in vivo applications. *Commun Biol*. 2023;6(1):699.
- 120 Tang WH, Sim SR, Aik DYK, Nelanuthala AVS, Athilingam T, Röhlén A, et al. Deep learning reduces data requirements and allows real-time measurements in imaging FCS. *Biophys J*. 2024;123(6):655-66.
- 121 Zhong X, Qin Y, Liang C, Liang Z, Nong Y, Luo S, et al. Smartphone-Assisted Nanozyme Colorimetric Sensor Array Combined "Image Segmentation-Feature Extraction" Deep Learning for Detecting Unsaturated Fatty Acids. *ACS sensors*. 2024;9(10):5167-78.
- 122 Pan J, Low KL, Ghosh J, Jayavelu S, Ferdaus MM, Lim SY, et al. Transfer learning-based artificial intelligence-integrated physical modeling to enable failure analysis for 3 nanometer and smaller silicon-based CMOS transistors. *ACS Appl Nano Mater*. 2021;4(7):6903-15.
- 123 Zhou B, Jieming P, Sivan M, Thean AV-Y, Senthilnath J. Quantile Online Learning for Semiconductor Failure Analysis. In: ICASSP 2023-2023 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP). 2023. p. 1-5.
- 124 Ferdaus MM, Zhou B, Yoon JW, Low KL, Pan J, Ghosh J, et al. Significance of activation functions in developing an online classifier for semiconductor defect detection. *Knowledge-Based Syst*. 2022;248:108818.
- 125 Wang Q, Li Y, Li R. Ecological footprints, carbon emissions, and energy transitions: the impact of artificial intelligence (AI). *Humanit Soc Sci Commun*. 2024;11(1):1-18.
- 126 Singapore Green Plan [Internet]. Available from: <https://www.greenplan.gov.sg/>
- 127 Kim HJ, Madhavi S. A Reinforcement Learning Model for Quantum Network Data Aggregation and Analysis". *J Syst Manag Sci*. 2022;12(1):283-93.
- 128 Roy JJ, Phuong DM, Verma V, Chaudhary R, Carboni M, Meyer D, et al. Direct recycling of Li-ion batteries from cell to pack level: Challenges and prospects on technology, scalability, sustainability, and economics. *Carbon Energy*. 2024;e492.
- 129 Yan Z, Xu Y. Real-time optimal power flow with linguistic stipulations: integrating GPT-Agent and deep reinforcement learning. *IEEE Trans Power Syst*. 2023;
- 130 Krenn M, Pollice R, Guo SY, Aldeghi M, Cervera-Lierta A, Friederich P, et al. On scientific understanding with artificial intelligence. *Nat Rev Phys*. 2022;4(12):761-9.
- 131 Thais S. Physics and the empirical gap of trustworthy AI. *Nat Rev Phys*. 2024;1-2.
- 132 Huang B, von Rudorff GF, von Lilienfeld OA. The central role of density functional theory in the AI age. *Science (80- )*. 2023;381(6654):170-5.
- 133 Karagiorgi G, Kasieczka G, Kravitz S, Nachman B, Shih D. Machine learning in the search for new fundamental physics. *Nat Rev Phys*. 2022;4(6):399-412.
- 134 Long Y, Xue H, Zhang B. Unsupervised learning of topological non-abelian braiding in non-hermitian bands. *Nat Mach Intell*. 2024;6(8):904-10.
- 135 Jessica LSE, Arafat NA, Lim WX, Chan WL, Kong AWK. Finite Volume Features, Global Geometry Representations, and Residual Training for Deep Learning-based CFD Simulation. *arXiv Prepr arXiv231114464*. 2023;
- 136 Mao HY, Lu YH, Lin JD, Zhong S, Wee ATS, Chen W. Manipulating the electronic and chemical properties of graphene via molecular functionalization. *Prog Surf Sci*. 2013;88(2):132-59.
- 137 Su J, Li J, Guo N, Peng X, Yin J, Wang J, et al. Intelligent synthesis of magnetic nanographenes via chemist-intuited atomic robotic probe. *Nat Synth*. 2024;3(4):466-76.
- 138 Weinan E, others. The dawning of a new era in applied mathematics. *Not Am Math Soc*. 2021;68(4):565-71.
- 139 Karniadakis GE, Kevrekidis IG, Lu L, Perdikaris P, Wang S, Yang L. Physics-informed machine learning. *Nat Rev Phys*. 2021;3(6):422-40.
- 140 Centre for Data Science and Machine Learning [Internet]. Available from: <https://www.math.nus.edu.sg/cdsml/>
- 141 Centre for Quantum Technologies [Internet]. Available from: <https://nusgs.nus.edu.sg/booth-pages/cqt/>
- 142 Dwivedi VP, Luu AT, Laurent T, Bengio Y, Bresson X. Graph neural networks with learnable structural and positional representations. *arXiv Prepr arXiv211007875*. 2021;
- 143 Dwivedi VP, Joshi CK, Luu AT, Laurent T, Bengio Y, Bresson X. Benchmarking graph neural networks. *J Mach Learn Res*. 2023;24(43):1-48.
- 144 Kawaguchi K. Deep learning without poor local minima. *Adv Neural Inf Process Syst*. 2016;29.
- 145 AI Singapore [Internet]. Available from: <https://aisingapore.org/>
- 146 Cao L. AI in finance: challenges, techniques, and opportunities. *ACM Comput Surv*. 2022;55(3):1-38.
- 147 Sun S, Qin M, Zhang W, Xia H, Zong C, Ying J, et al. Trademaster: A holistic quantitative trading platform empowered by reinforcement learning. *Adv Neural Inf Process Syst*. 2023;36:59047-61.
- 148 Thompson C. How Khan Academy is changing the rules of education. *Wired Mag*. 2011;126:1-5.
- 149 Akgun S, Greenhow C. Artificial intelligence in education: Addressing ethical challenges in K-12 settings. *AI Ethics*. 2022;2(3):431-40.
- 150 Holmes W, Tuomi I. State of the art and practice in AI in education. *Eur J Educ*. 2022;57(4):542-70.
- 151 Ahmed UZ, Sahai S, Leong B, Karkare A. Feasibility Study of Augmenting Teaching Assistants with AI for CS1 Programming Feedback. 2025;
- 152 Borchert H, Schütz T, Verbosvzky J. The Very Long Game: 25 Case Studies on the Global State of Defense AI. Springer Nature; 2024.
- 153 Soremekun E, Udeshi S, Chattopadhyay S. Towards backdoor attacks and defense in robust machine learning models. *Comput & Secur*. 2023;127:103101.
- 154 Lo D. Trustworthy and synergistic artificial intelligence for software engineering: Vision and roadmaps. In: 2023 IEEE/ACM International Conference on Software Engineering: Future of Software Engineering (ICSE-FoSE). 2023. p. 69-85.
- 155 If you don't see the gun, can you pinpoint where it was fired? [Internet]. Available from: <https://www.htx.gov.sg/news/featured-news-if-you-don-t-see-the-gun-can-you-pinpoint-where-it-was-fired>
- 156 Voigt P, dem Bussche A. The eu general data protection regulation (gdpr). *A Pract Guid 1st Ed*, Cham Springer Int Publ. 2017;10(3152676):10-5555.



# APPENDICES

- Appendix i. Earth/Climate Sciences
- Appendix ii. Physics/Complexity
- Appendix iii. RNA Biology & Therapeutics (AI4 RBT)
- Appendix iv. Biomedical Sciences
- Appendix v. Healthcare and Imaging
- Appendix vi. Genomics
- Appendix vii. Digital Phenotyping
- Appendix viii. Materials/Chemistry
- Appendix ix. Chemical and Biomanufacturing
- Appendix x. Financial Services
- Appendix xi. Electronics/Semiconductors
- Appendix xii. Hybrid Quantum Computing
- Appendix xiii. Sustainability
- Appendix xiv. Education
- Appendix xv. Science, Software and Security
- Appendix xvi. Robotics
- Appendix xvii. AI Methods and Mathematics
- Appendix xviii. Global R&D efforts in AI for Science

## APPENDIX I. EARTH/CLIMATE SCIENCES

### AUTHORS:

Prof. Sang-Ho Yun  
(NTU, Earth Observatory of Singapore - Remote Sensing Lab),  
Dr. Dale Barker  
(Centre for Climate Research Singapore)

### Executive Summary

The workshop on Advancing Earth and Climate/Weather Science through Artificial Intelligence (AI) identified critical grand challenges and opportunities where AI can significantly contribute to understanding and mitigating global environmental and climate issues. Key themes included the need for AI in processing large datasets, the importance of ground-truth information, and the potential of AI in:

- identifying  **tipping points**  in Climate and Earth system,
- **automating data classification**  for comprehensive cataloging,
- **understanding complex system linkages**  across interconnected components of the Earth system.
- **reducing uncertainties in forecasts**  in climate and weather forecasts, and
- **improving risk communication**  to better inform policy and public action.

### Introduction

As global environmental and climate challenges grow in complexity and urgency, traditional scientific approaches increasingly benefit from the support of advanced technologies. AI presents transformative opportunities to enhance the analysis of large datasets, model complex interactions within the Earth system, and improve predictive capabilities crucial for decision-making and policy development.

This workshop convened experts in Earth sciences, climate and weather research, health science, and AI to discuss these pressing challenges and explore innovative AI-driven solutions. The details of each grand challenge are described below, and the overall design and implementation of the workshop, along with the post-workshop information compilation work, are provided in the appendices attached at the end of this white paper.

## Background

The workshop on Earth and Climate Science organized in Singapore brought together domain experts from academia and industry to address the pressing need for better exploitation of available multi-modal weather, health, and climate data. A core inquiry emerged regarding the appropriate hardware resources necessary to enable AI to effectively utilize this data, striking a balance between the demands for powerful computational machines and the imperative of environmental sustainability.

Central to the discussions was the role of AI in identifying not only correlations but also causal relationships between the collected data and target climate, weather, or health events. The workshop emphasized the necessity for AI to be scalable and capable of deciphering the complex interrelationships among various elements

within the studied systems. This capability is vital for advancing our understanding of these multifaceted dynamics.

A comprehensive vision was articulated for developing a robust, scalable, and explainable AI model for the Earth, designed to provide both short- and long-term predictions of regular and outlier events. This model would enhance our understanding of the mechanisms underlying weather and climate, while also facilitating early action in the face of disasters.

This white paper consolidates on the workshop, and the consensus among AI and domain experts indicating that a crucial first step toward this ambitious goal involves creating AI models capable of discerning causal relationships within the continually evolving data landscape.

## Grand Challenges

One of the outcomes from the workshop was the identification of five grand challenges, namely:

- Identification of Tipping Points in the Earth System
- Automated Identification and Classification of Large Volumes of Data
- Identification of Complex Links Between Earth System Elements

- Quantification and Reduction of Uncertainties in Forecasts
- Forecasting and Communication of Risk
- The motivations for choosing these challenges are detailed below, along with the requirements identified to address them.

---

### GRAND CHALLENGE 1: IDENTIFICATION OF TIPPING POINTS IN THE EARTH SYSTEM

---

AI can revolutionize our ability to identify tipping points in ecological systems, climate change<sup>1</sup>, earthquakes, and volcanic activity by analyzing long-term data and detecting early signs of instability.

Detecting critical thresholds in the Earth system is essential for preventing irreversible damage and mitigate the ever-increasing extreme events (e.g. *Figure 1*). Early identification of tipping points can enable timely interventions and preparation.

Singapore can lead by investing in the development of AI models tailored to regional data (Southeast Asia), where large data gaps exist. Current approaches rely heavily on historical data analysis, but AI can enhance this by integrating real-time data and improving predictive accuracy.

#### OBJECTIVE

Identifying specific tipping points in the Earth system requires long time series data and advanced AI models capable of recognizing subtle changes that precede significant shifts. Early identification of tipping points, enabling preemptive actions to mitigate potential disasters.

#### DATA REQUIREMENTS

Long-term, spatially dense datasets across ecological, climatic, and geological domains are required.

#### AI METHOD REQUIREMENTS

Machine learning algorithms for time series analysis, anomaly detection, and causal discovery.

#### RESOURCE REQUIREMENTS

High-performance computing, access to large datasets, and collaboration between AI experts and domain scientists.

---

### GRAND CHALLENGE 2: AUTOMATED IDENTIFICATION AND CLASSIFICATION OF LARGE VOLUMES OF DATA

---

Automation of data classification using AI can significantly speed up the analysis process, particularly in regions like Southeast Asia with large data gaps.

Many Earth science applications require large, labeled datasets, which are currently labor-intensive to produce. By developing AI tools that address regional data challenges and can be scaled globally. Existing AI tools are limited in their ability to classify new, unlabeled data, especially in diverse ecosystems<sup>2</sup>.

#### OBJECTIVE

Creating AI systems that can automatically identify and classify data, including previously unknown species or environmental samples. Faster, more accurate data classification, leading to better environmental monitoring and decision-making.

#### DATA REQUIREMENTS

Large, diverse datasets from ecological surveys, remote sensing, and field observations.

#### AI METHOD REQUIREMENTS

Deep learning for image recognition, natural language processing for species identification.

#### RESOURCE REQUIREMENTS

Data collection initiatives, computational resources, and AI model training.

**GRAND CHALLENGE 3:  
IDENTIFICATION OF COMPLEX LINKS  
BETWEEN EARTH SYSTEM ELEMENTS**

Understanding the complex interactions between different Earth system components can significantly enhance the accuracy of predictions related to ecological, climatic, and geohazard impacts. The Earth system is highly interconnected, and AI can help reveal hidden linkages and processes that traditional methods might miss<sup>3,4</sup>.

Singapore can become a hub for AI-driven Earth system science by fostering collaboration between AI and environmental researchers. Current models often focus on correlations, but AI can help establish causal relationships.

**OBJECTIVE**

Quantifying interactions between Earth system components by identifying interconnecting variables across different time and space scales. Improved understanding of Earth system dynamics, leading to better climate and environmental policies.

**DATA REQUIREMENTS**

Multivariate datasets encompassing variables of atmosphere, hydrosphere, biosphere, and geosphere.

**AI METHOD REQUIREMENTS**

Neural networks, Bayesian networks, and other AI techniques for causal inference.

**RESOURCE REQUIREMENTS**

Access to integrated datasets, interdisciplinary collaboration, and advanced AI tools.

**GRAND CHALLENGE 4:  
QUANTIFICATION AND  
REDUCTION OF UNCERTAINTIES  
IN FORECASTS**

AI can enhance the accuracy and reduce the uncertainty of forecasts related to climate, weather, and natural hazards whilst significantly reducing the time-to-solution thus enabling faster response times.

Reducing uncertainties in environmental predictions is crucial for effective risk management and policymaking. Whilst AI techniques have been applied for decades in climate and environmental science, in recent years a revolution in speed and skill of predictions has taken place due to rapid advances in development of AI deep-learning techniques and availability of novel supercomputing platforms<sup>5</sup>.

By developing AI-enhanced models that integrate physical and environmental data, Singapore can build on its existing significant capabilities in weather, climate and environmental modelling and prediction capabilities to set new standards in predictive accuracy. AI is currently used to speed up model computations, but integrating physical processes remains a challenge.

**OBJECTIVE**

Improving model accuracy and reliability while reducing computational costs and processing time. More accurate, reliable and timely predictions, reducing the impact of natural disasters and aiding in climate adaptation efforts.

**DATA REQUIREMENTS**

High-resolution climate, weather, and geohazard datasets, including real-time data.

**AI METHOD REQUIREMENTS**

Hybrid models combining AI with physics-based simulations<sup>6</sup>, foundation models developed for particular applications, etc.

**RESOURCE REQUIREMENTS**

Computational infrastructure, access to real-time data, and interdisciplinary expertise.

**GRAND CHALLENGE 5:  
FORECASTING AND COMMUNICATION  
OF RISK**

AI can improve the forecasting and communication of climate and environmental risks, making information more accessible and actionable. Effective communication of risks is essential for public preparedness and policy responses, and filtering disinformation is also crucial<sup>7</sup>.

By pioneering AI-driven risk communication tools that provide localized, personalized forecasts. AI is increasingly used for risk forecasting, but its application in personalized communication is still emerging. Developing AI systems that can assess the quality of risk

forecast and communication and personalize risk information and explore “what if” scenarios.

**OBJECTIVE**

Enhanced public understanding and response to climate risks, leading to better preparedness and mitigation strategies.

**DATA REQUIREMENTS**

Data from climate models, hazard simulations, and social demographics.

**AI METHOD REQUIREMENTS**

Large Language Models, geospatial analysis, and risk modeling algorithms.

**RESOURCE REQUIREMENTS**

Diverse data sources, AI development expertise, and user-centered design.

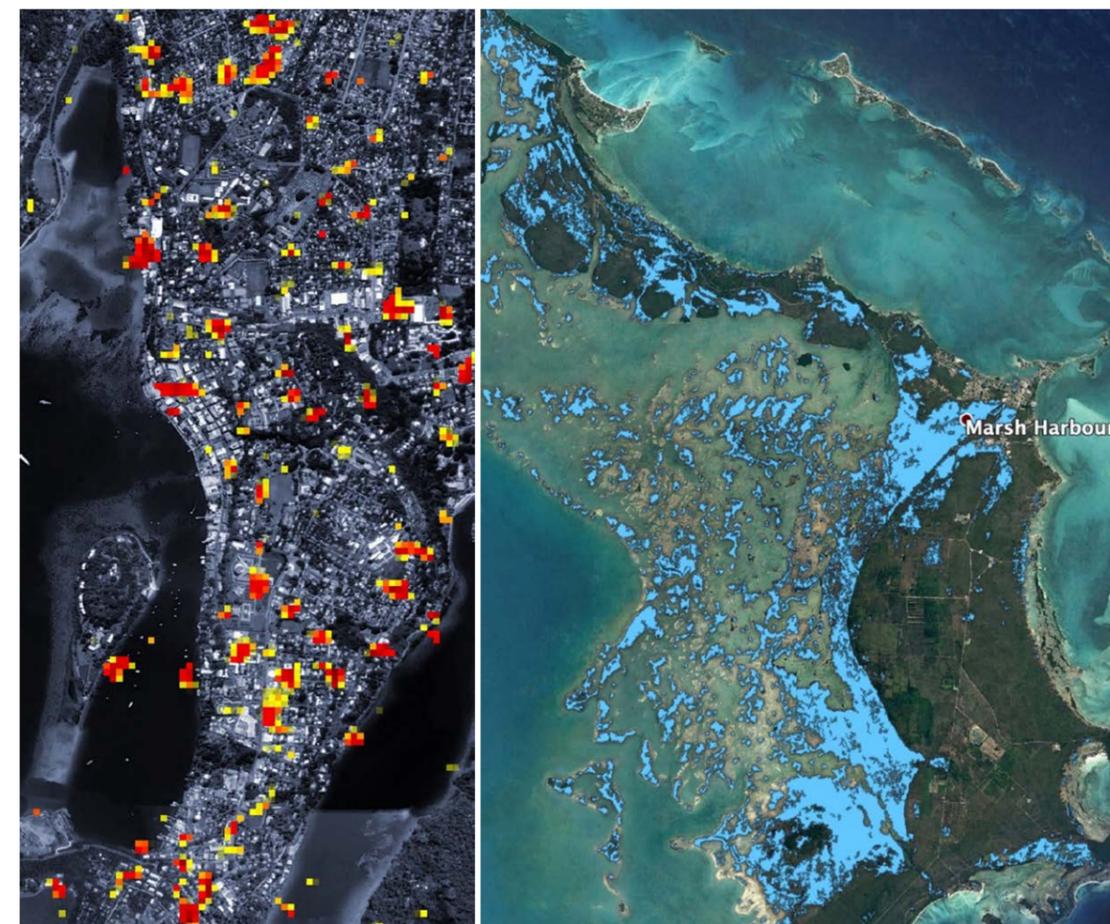


Figure 1: Climate change is increasing the frequency and severity of extreme events. A damage map of Vanuatu (left) to rapidly reveal the locations of likely damage due to Cyclone Judy and Cyclone Kevin in March 2023<sup>8</sup>, and a flood map (right) of the Bahamas to delineate the extent of inundation caused by Hurricane Dorian in September 2019<sup>9</sup>. Both maps were created by the Earth Observatory of Singapore - Remote Sensing Lab, derived from the radar data acquired by the Copernicus Sentinel-1 satellite operated by the European Space Agency.

## AI Methods and Data - Challenges and Opportunities

Constructing an AI capable of understanding the complex interconnections within Earth's systems, including causal relationships versus mere correlations, is key in constructing accurate global ecological and geological models. It is also necessary to develop AI methods capable of auto-identification of samples, such as fossils and insects, essential for reconstructing historical climate events.

In terms of data, there is a need for substantial investment in both AI-led data collection and storing infrastructure, and classic domain-led expert data collection. Data such as information about bird populations in specific regions,

must be captured in formats compatible with AI models to enhance predictive accuracy in long-term climate forecasting. For the monitoring of carbon stocks, it is necessary to implement robust global monitoring systems that incorporate wireless technology and land-based nodes for data analysis.

By harnessing AI methods effectively, researchers can enhance their understanding of the mechanisms driving climate and health, ultimately improving predictive modeling and informing targeted interventions in response to climate-related challenges.

## Singapore's Role

Singapore needs to dedicate research funding not only for specialized and potentially novel hardware, but also funding for traditional, domain-specific research efforts.

This is necessary to ensure collection of good quality data to train AI models, and ensure their predictions are robust against domain expertise.

## Conclusions

The workshop highlighted the transformative potential of AI in Earth and Climate Sciences, with the ability to address grand challenges such as identifying tipping points, automating data classification, understanding complex

linkages, reducing forecast uncertainties, and improving risk communication. For Singapore, leading in these areas presents an opportunity to contribute significantly to global environmental and climate solutions.

## REFERENCES

- 1 "How Close Are the Planet's Climate Tipping Points?", New York Times, 11 August 2024. Available online at: <https://www.nytimes.com/interactive/2024/08/11/climate/earth-warming-climate-tipping-points.html?searchResultPosition=2>
- 2 Han BA, Varshney KR, LaDeau S, Subramaniam A, Weathers KC, Zwart JA, "Synergistic future for AI and ecology", Proc Natl Acad Sci (PNAS), USA, 2023. <https://doi.org/10.1073/pnas.2220283120>
- 3 Hickmon, N., Varadharajan, C., & Hoffman, F. (2022). Artificial Intelligence for Earth System Predictability (AI4ESP) Workshop Report. Argonne National Laboratory. <https://doi.org/10.2172/1888810>
- 4 See, S., Adie, J. Challenges and opportunities for a hybrid modelling approach to earth system science. CCF Trans. HPC 3, 320–329 (2021). <https://doi.org/10.1007/s42514-021-00071-y>
- 5 "How AI models are transforming weather forecasting: a showcase of data-driven systems", ECMWF, 6 September 2023. Available online at: <https://www.ecmwf.int/en/about/media-centre/news/2023/how-ai-models-are-transforming-weather-forecasting-showcase-data>
- 6 Marković, D., Mizrahi, A., Querlioz, D., & Grollier, J. (2020). Physics solves a training problem for artificial neural networks. Nature Reviews Physics, 2(10), 499–510. <https://www.nature.com/articles/d41586-024-02392-8>
- 7 Karinshak, E., & Jin, Y. (2023). AI-driven disinformation: A framework for organizational preparation and response. Journal of Communication Management, 27(4), 539–562. <https://doi.org/10.1108/JCOM-09-2022-0113>
- 8 EOS-RS Damage Proxy Map: Vanuatu Cyclone Judy and Cyclone Kevin, 10 March 2023, v0.8. Available online at: [https://eos-rs-products.earthobservatory.sg/EOS-RS\\_202303\\_Vanuatu\\_Cyclone\\_Judy\\_Cyclone\\_Kevin/](https://eos-rs-products.earthobservatory.sg/EOS-RS_202303_Vanuatu_Cyclone_Judy_Cyclone_Kevin/)
- 9 Measuring global change, an interview with Science on the use of spaceborne Interferometric Synthetic Aperture Radar (InSAR) observations (19 March, 2021). <https://www.science.org/content/blog-post/measuring-global-change>

## APPENDIX II. PHYSICS/COMPLEXITY



### AUTHORS:

Assoc. Prof. Duane Loh  
(NUS, Centre for Bioimaging Sciences, Data Science Institute)

Prof. Lock Yue Chew  
(NTU, School of Physical & Mathematical Sciences)

### Executive Summary

Modeling the physical world across multiple length and time scales presents significant challenges due to the limitations of existing theories and models, which are effective only within narrow ranges. These models often fail when extrapolated beyond their specific domains, highlighting a fundamental issue: the absence of a unified approach that can seamlessly handle the complexities of different scales. Current scientific approaches rely on a patchwork of specialized models that address different scales independently, but this fragmented strategy exposes gaps, particularly in complex or real-world applications.

Artificial Intelligence (AI) offers a promising avenue to overcome these limitations through inspired coarse-graining, which effectively simplifies complex systems into manageable motifs and patterns. Unsupervised machine

### Introduction

Measuring, understanding, and modeling our complex physical world over multiple length/time scales is extremely challenging<sup>1,2</sup>. So far, we invented models and theories that are separately effective only over narrow time and length scales. There is usually a lower bound where the model's fundamental units lose rigidity or persistence, are no longer known, or all of the above. For example, while mathematical models of fluid dynamics describe the dynamics and interactions

learning may identify these motifs even in non-linear systems, which form the building blocks for more complex hierarchical and scalable models that can bridge the gaps across different physical scales.

Generative AI models can also help bridge the "theory-experiment gap" by helping create bridging phenomenological models from both experimental and theoretical data. These models can analyze vast amounts of data, recognize patterns, and generate models that may transcend traditional boundaries of scale, specialization, theoretical models, and experimental constraints. By integrating AI into scientific research, we can potentially develop more cohesive and adaptable models, bridging the gaps left by conventional theories and enabling a deeper understanding of complex systems across all scales.

between small volumes in the continuum, these models fail when the size of their volumes approaches several nanometers (i.e., the size of their molecular constituents). There is also an upper bound, where we can no longer compute our models' predictions within usable precision. For instance, whereas we know the classical kinematics and dynamics of a few rigid bodies remarkably well, it is impossible to predict their final states to arbitrarily long lengths and times.

Despite these challenges, science has collectively engineered a patchwork of specialized models and theories that span most of the observable length and time scales. This patchwork within physics comprises particle physics, atomic and molecular physics, condensed matter physics, classical mechanics, continuum mechanics of media, geophysics, planetary science, photonics, astrophysics, and cosmology. Each of these specializations provides a unique perspective, coarse-graining our physical world into rigid units that are effective at different length and time scales: fundamental particles, atoms, molecules, collective excitations (e.g., quasi-particles), massive aggregates, semi-rigid bodies, rigid bodies, and so on.

Although this patchwork has carried science far, we are constantly reminded of its large gaps in extrapolating into real, complex realizations. Thermal expansion of materials exemplifies this gap. While we can compute the linear thermal expansion coefficient of perfect crystals *ab initio*, doing the same for polycrystalline materials is practically impossible. As a result, this coefficient is still experimentally measured rather than calculated from first principles. More fundamentally, we still fail at re-formulating the coupling parameters of one model to those of another.

These gaps beg the obvious question: Why is it so difficult for a single model to describe our world across all observable lengths and times? The answer often concerns two critical aspects of complex systems: uncertainties and interactions.

First, inevitable uncertainties in measuring the properties of fundamental units can destabilize our predictions. This is especially true when these units harbor internal degrees of freedom that we coarse-grain away or are impossible to measure. Unfortunately, these uncertainties accumulate when we compute the interactions of many units over sufficiently long distances and times. Taken far enough, these cumulative errors can make these computations too imprecise to be useful. Compounded with the fact that humans are famously terrible at statistical reasoning (e.g.,

gambler's fallacy, base-rate fallacy, etc.), these uncertainties can make it difficult to hand-craft statistical models of complex systems.

Second, many interacting bodies (even with their internal details coarse-grained away) can become exponentially complex with their numbers. With strong and conditional interactions, dozens of bodies can already have behavioral "state spaces" that are complex and nuanced, far more than the familiar simple solids, liquids, and gases of weakly interacting bodies. In such circumstances, the boundary between meaningless disorder and persistent, complex but seemingly random states is quickly blurred.

### CONCEPTUAL CHALLENGES IN PHYSICS AND COMPLEXITY THAT AI CAN TACKLE

#### COARSE-GRAINING MULTISCALE COMPLEXITY WITH AI

A key to dealing with multiscale complexity is coarse-graining<sup>2,3,4,5</sup>. In fact, there are indications that Deep Learning is related to the variational renormalization group<sup>6</sup>. By grouping sufficiently many interacting bodies, the uncertainties at the intermediate scales are simply folded into phenomenological coupling parameters, response to impulses/stress, or decay and stability. For example, although the general three-body problem is analytically intractable (and sensitive to initial conditions), we can still describe the phenomenology of hundreds of millions of stars (and their accompanying planets) in a galaxy. The characteristics of each of the stars are effectively coarse-grained into a far smaller set of summarizing variables that can be used to predict the dynamics, stability, and properties of the entire galaxy. Similarly, we do not worry about the quantum uncertainty principle when we describe the friction between two surfaces. These uncertainties become coarse-grained away when such systems are sufficiently scaled up.

How do we coarse-grain complex systems? One strategy first identifies persistent, prevalent motifs and patterns in a complex system. Thereafter, the description of these complex systems is simplified to only model

the behaviors and properties of these motifs. Here, unsupervised machine learning can demonstrably help us define these persistent and prevalent motifs even in non-linear complex systems. More generally, AI can demonstrably learn the hierarchy of motifs (i.e., iteratively coarse-grain these motifs), not unlike learning how words compose sentences, and sentences make paragraphs.

#### **GENERATIVE AI MODELS TO LEARN MOTIF HIERARCHIES AND GRAPHS**

The era of powerful generative models is upon us. This power comes from the models' ability to learn and combine contextual graphs and hierarchies. As a result, these models can stochastically crawl through their hierarchies to generate endless, novel streams of words, sentences, and paragraphs steeped in context. The awesome power of these correlative models emerges when we start hybridizing them. For example, hybrid models learn the hierarchical correlations between the hierarchies between different modalities<sup>7</sup>.

A long-standing vision is to develop AI models that connect first-principles modeling with experimental observations across large gaps between their length and time scales<sup>8</sup>. Such

## **Background**

The rapid advancements in artificial intelligence (AI) and computational methods have unveiled significant potential for transforming research in physics and complexity science. Recognizing this opportunity, we organized a workshop bringing together leading experts in physics, data science, AI and complexity theory. This diverse group provided valuable perspectives that enriched the discussions and contributed to a comprehensive understanding of the current landscape and future directions for AI in scientific inquiry and specifically within the fields of Physics and Complexity.

A core inquiry arising from the workshop was: "To whom does explainable AI explain?" This prompted discussions about the subjective nature of explainability, which depends heavily

data must be sufficiently complex to embody meaningful hierarchies of correlations, and spanning multiple length and time scales. First-principles models, however, mostly guide us over very limited regions on these scales. So, we strive to create explainable phenomenology to leap over the chasm between first principles and observations.

#### **FINDING EMERGENT FIELDS FROM COMPLEXITY**

New fields of study often emerge from complex problems/systems<sup>9,10</sup>. These intricate interactions are usually too difficult for any one person to fully grasp. Quantum mechanics, for example, is seen as a great achievement in human discovery. Despite being counterintuitive and hard to interpret (like the Copenhagen interpretation), it still produces highly accurate calculations. When we have rules that can make strong predictions, we tend to believe we understand the underlying concepts. Can AI help us make sense of such complex scientific ideas in a more objective and falsifiable way<sup>11,12</sup>?

We need to create AI models that can be interrogated, validated, and extended in similar ways to our scientific models.

on prior knowledge and the interplay between extrapolation and reductionism. Emphasis was placed on AI's role in elucidating scientific principles, seen as axiomatic yet fundamentally interpretable frameworks that help explain the world. It was highlighted that AI has a huge potential as a companion in scientific inquiry, assisting in the discovery of emergent, interpretable, and explainable principles, with the ambition to uncover novel first principles that could surpass human intellectual capabilities.

However, the integration of AI into scientific disciplines comes with its own set of challenges. These include the need for robust data infrastructures, the importance of explainability in AI models, and the potential biases that may arise from algorithmic processes. Addressing

these challenges is crucial for ensuring that AI can effectively contribute to scientific discovery while maintaining transparency and trustworthiness. As the scientific community grapples with these issues, the focus shifts towards identifying grand challenges that will direct the development and application of AI in uncovering emergent principles and advancing our understanding of the natural world.

## **Grand Challenges**

Based on the outcome of the organised workshop, we decided one grand challenge spanning the different facets of research within the fields of Physics and Complexity:

---

#### **GRAND CHALLENGE: CREATING AN AI SCIENTIST TO DETECT EMERGENT PHENOMENA IN OBSERVATIONS OF COMPLEX DYNAMICAL SYSTEMS THAT SPAN MANY LENGTH AND TIME SCALES<sup>13</sup>**

---

Currently, this capability rests on the informed yet sporadic creativity of a community of researchers. Can this capability be learned by an AI? We are optimistic that the necessary tasks can be decomposed as a series of learnable intuitions, attention to complex details, and awareness to a broad spectrum of prior knowledge and literature<sup>14,15,16,17</sup>. The scope of this AI scientist spans atomic to macro scales, which is pivotal for enhancing our understanding of emergent behaviors. Fundamentally, the discovery of these emergent behaviors will seed new disciplines. Practically, it would help us rationally design new materials with potentially transformative applications across various industries. Together, this AI scientist's impact extends to improved material efficiencies, innovative product developments, and potentially solving some of the most pressing environmental and technological challenges of our time.

This white paper consolidates the workshop's insights and recommendations, aiming to guide future research and funding priorities in AI for physics and complexity within the Singapore ecosystem. By identifying these grand challenges, we hope to catalyze significant advancements in our understanding of the universe and enhance the role of AI in scientific discovery.

#### **CRITICAL REQUIREMENTS:**

To achieve the goal of an efficient AI scientist that can be used as a co-pilot to accelerate the generation of scientific hypotheses<sup>14</sup>, we would require the:

- Self-supervised, meaningful collection and labeling of massive and complex data<sup>14</sup>;
- Design and training of AI models with explicitly interpretable latent spaces that are explainable and extensible to scientists<sup>18</sup>;
- Efficient incorporation of prior knowledge and principles in physics into AI models<sup>19,20,21</sup> that do not involve showing explicit training data generated painfully from simulated forward models;
- Design and undertaking of multi-scale experiments that can be used to train, test and validate the AI models, in order to bridge the large gaps in the length and time scales of the phenomenon of interest;
- Learning of the hierarchy of correlations in multiscale systems<sup>22</sup>;

For an efficient and reliable AI scientist for physics, the following challenges need to be addressed.

## BRIDGING THE SCALES:

To accelerate the synthesis of functional quantum motifs in low-dimensional materials, we will need an atlas of multimodal maps (Figure 1).

**> 1000-fold acceleration in rational synthesis using AI guidance:** The actual bottleneck in discovering functional quantum defects is in their synthesis. While high-throughput simulations of periodic structures can identify many potential targets, synthesizing these targets in the laboratory is far from straightforward.

**High-content experimental atlas (10,000x more):** Amass the only *multi-scale structural map of synthesizable functional, complex quantum defects in non-periodic low-dimensional materials*. This atlas of maps will contain the largest and most diverse range of complex quantum defects that can be synthesized in metal-chalcogen-based low-dimensional materials.

• **Each structural map has 100x larger field-of-view.** To capture sufficiently many such defects, each structural map will span a range of  $10^5$  in spatial resolution. This

can be done with the help of an AI-assisted semi-autonomous microscope (Figure 2).

- **First-ever atlas (with > 200 maps).** We will need to collect 200 sufficiently distinct synthesis conditions.
- **10x faster route to 3D resolution.** Using physics-informed neural networks, we can infer the 3D structure with single 2D images<sup>24</sup> with computationally retrieved true phase contrast<sup>25</sup>.
- **Faster, reliable AI-based annotation of complex, non-periodic quantum defects with self-supervised learning [5].** This use of AI will reduce structural identification and reconstruction from hours to seconds (100x acceleration).
- **Unprecedented multi-modal co-annotations on each map.** Systematically co-characterizing each map with spatially-resolved photo-luminescence, Raman, and electron-energy loss spectroscopy (Figure 1).
- **Annotated maps with energetics and rates from first-principles models.** Create foundation AI-MD models for fast molecular dynamics<sup>26,27</sup>.

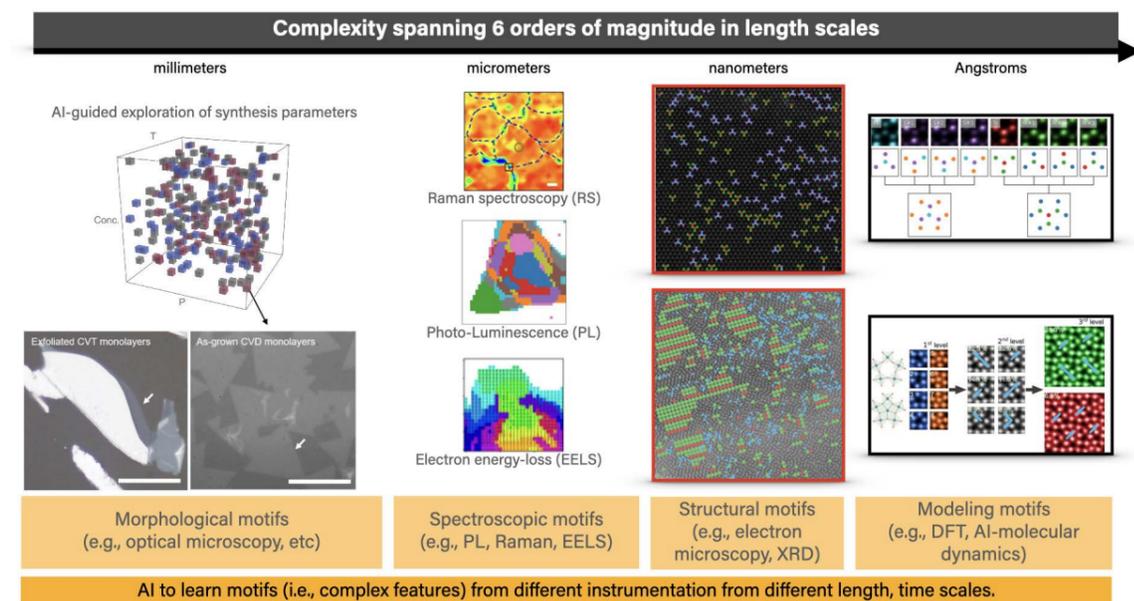


Figure 1: Overcome multi-scale complexity by combining characterization at different length scales with different measurement modalities and simulation capacities. These are almost never combined for a single sample in a field of view. Image credits<sup>5,23</sup>.

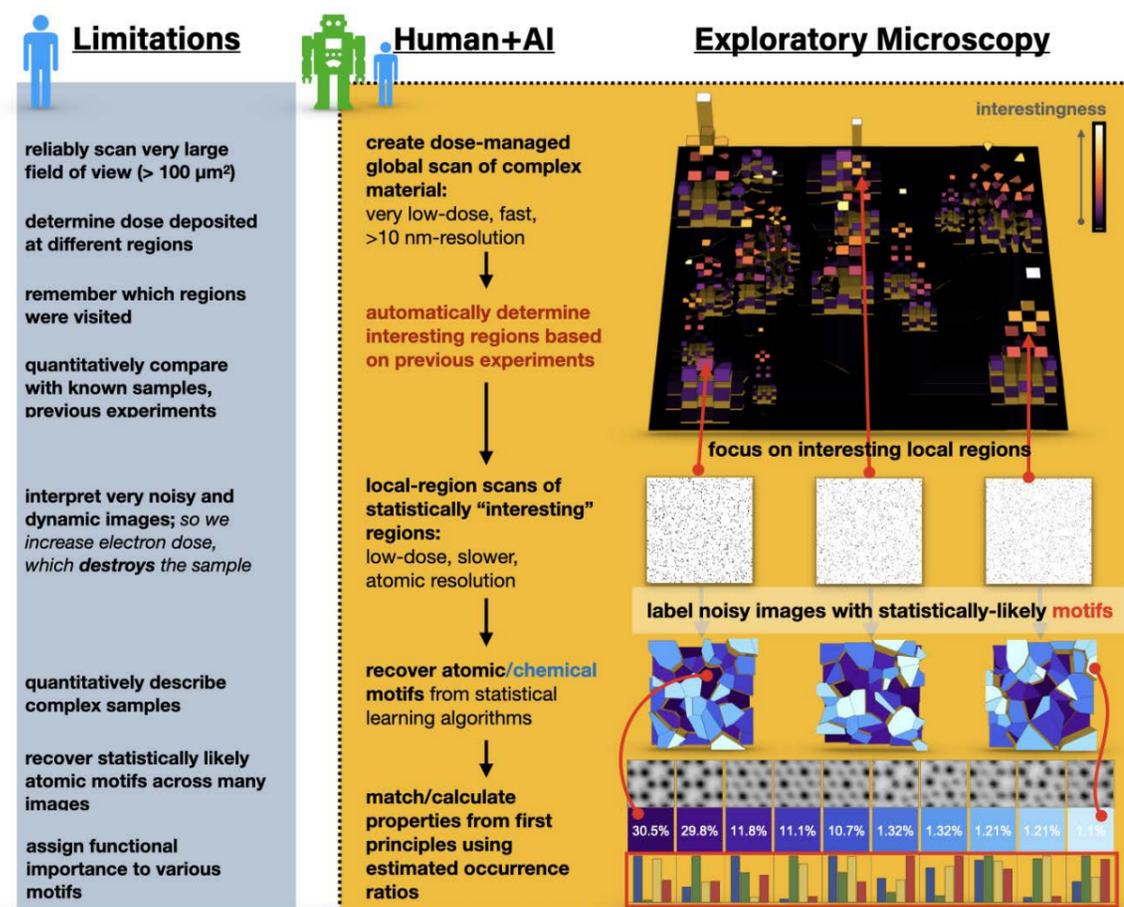


Figure 2: Using AI to create a semi-autonomous microscope that make statistical decisions about where to spend most time collecting high-resolution images.

## AI Methods and Data - Challenges and Opportunities

### DATA-DRIVEN DISCOVERY

If data are the currency of AI, then experiments are anchors of science. The data fed these AI models had to be measured, curated, and labeled by humans. When these AI models surpass conventional methods in suggesting novel materials or biologics, someone still has to synthesize them and evaluate their properties in a physical laboratory.

Experimental observations is also the primary source of inspiration for us to develop first-principles models that connect to experimental observations across large gaps in length and

time scales. Our theories and models reveal their blindspots only when tested against experiments. In many ways, they guide us across these gaps by providing a means to falsify our models.

Machine learning has become indispensable for understanding high-throughput, high-content, multi-modal experiments, which have become increasingly commonplace in the last decade. Such massive and rich measurements, rather than controlled experiments, are now crucial means of studying dynamically changing complex systems.

Here, however, a matching problem rears its ugly head: how can we find meaningful patterns in these massive experimental data to connect with limited insights from expensive first principles models? In machine learning jargon, data-driven discovery faces a double bottleneck in labeling and interpretability.

This matching problem can be decomposed as two complexity challenges: the complexity of matching multiple time and length scales spanned by experiment and simulations; the complexity of latent space induced in trained AI models. Existing unsupervised machine learning methods (e.g., manifold-learning or dimensionality reduction) will almost certainly find patterns in complex data that are not immediately interpretable in the context of prior knowledge. It is a constant battle to find complementary anchors in experimental data and our prior knowledge (e.g., using first-principles forward simulations, constraints, or invariants). An AI that can identify these anchors with limited human supervision will revolutionize our approach to data modelling and analysis, particularly in fields where scale and complexity are significant barriers. The ability to accurately model phenomena at multiple scales could lead to breakthroughs in understanding systems that are too complex for humans alone.

The fact that unsupervised learning methods should be this effective here is fundamentally enigmatic. This efficacy is popularly attributed to the poorly studied manifold hypothesis<sup>28,29,30</sup>: that real, experimental data exists in a smaller, more manageable subspace of all possible data spaces. Why should this hypothesis hold? For many “black-box” AI models in science that exploit dimensionality reduction, the hidden physical principles behind why each reduction was possible in that instance is opaque to its practitioners. Generalizing the principles that allow dimensionality reduction across applications could help us examine the unreasonable power of the manifold hypothesis and give crucial insights about the latent spaces embedded within our observations and, most importantly, the physical principles, constraints, and invariants that shape them.

---

## AI FOR COMPLEXITY SCIENCE

---

Complex systems can be defined by their ability to compute, which is broadly defined as a capacity to process information.

Here are three classes of highly anticipated applications of AI for complexity. First, digital twins or generative AI models (e.g., generative adversarial networks) based on forward models can augment and enrich incomplete or imbalanced training data. Such augmentation is especially critical in complex systems, which often show behaviors that are too diverse to be captured confidently. This diversity leads to impoverished datasets that often have biases and can be blind to critical outlier events. Second, generative AI models can play an important role in optimizing complex networks. For example, adaptive reinforcement learning can be used to optimize the scheduling of transportation systems, a classic NP-hard problem involving many agents whose behaviors change with time. Third, semi-supervised statistical learning methods can learn patterns of correlations in networks (e.g., the spread of infectious diseases like Dengue in Singapore, the neural activity of the visual cortex in response to visual stimuli, or the “information cascades” in social and media networks) and determine if these correlations are causal because of asymmetric time-ordering.

---

## CRITICAL LEARNING IN AI MODELS

---

Neural networks, like our brains, can be primed for optimal learning. Priming is a trade-off between adaptiveness (i.e., plasticity) and efficiency (i.e., incorporating existing priors), echoing the classic tension between order and disorder in statistical physics. There is mounting and intriguing evidence that certain neural network architectures behave optimally at the so-called “edge-of-chaos.” Put differently, these networks exhibit criticality that is reported in physical systems. Consequently, by re-framing neural network training as a critical phenomenon, we can rationally design and tune future AI or digital twins always to be critical.

Criticality may also lead to more explainable AI models. Criticality is often accompanied by spontaneous self-organization. The signatures of such critical self-organization in neural networks could form the hierarchical backbone within the learning graphs of such models.

---

## MANY-BODY PHYSICS

---

There have been several key applications of AI in many-body physics. In tensor neural networks the information propagation mirrors quantum physical principles. Hence, the action of these networks can be interpreted and made explainable through a physical analogy, where wave functions are decomposed into tensors<sup>[31, WANGYANG<sup>23</sup>]</sup>. Carleo & Troyer introduced the neural network quantum states [CARLEOTROYER<sup>17</sup>] where an input configuration (e.g., spin up & down) is processed to output complex numbers representing probability amplitudes. This technique supports both supervised and unsupervised learning with a physics-based loss function, allowing for the simulation of time evolution and the evaluation of ground state or thermal properties.

---

## METAPHOTONICS AND DEVICE DESIGN

---

It is very challenging, if not impossible, to design an ultrathin, highly multiplexed, on-chip and flat-profile advanced nanophotonic system using classical optimization approaches, in a “one-to-many” fashion (i.e., advanced multi-functions from a shared design using geometric perturbation, symmetry-broken physics and knowledge-expansion machine-learning approach). A knowledge-expansion foundation is needed to create a holistic advanced metaphotonic platform. Enhanced AI models can demonstrate various new concepts, including passive light control with disordered metaphotonics, ultracompact generation and manipulation of quantum light, on-chip photo-detection of multi-dimensional optical information, miniaturized spectroscopy combining AI and metasurfaces. Covering all significant aspects in optical science and nanotechnology, these research directions lay four pillars to support the edifice of AI-

empowered metaphotonics. Smaller, thinner, denser, and more advanced, AI-empowered metaphotonics can be envisioned.

From traditional machine learning models to neural networks and transformers, quantum-inspired techniques have evolved to describe rules and optimize loss functions through unsupervised learning. These models are adept at predicting complex evolutions and behaviors, including deterministic predictions of cellular automata via sequence-to-sequence learning.

Overall, we have identified several challenges and opportunities in AI of many-body physics. First, transfer learning between related systems and different scales is a key area for development. Second, the precision in the predictions by current many-body-physics AI models’ is inconsistent across applications, posing a challenge to the adaptability and transferability of models. Third, the actual dimensionality of the wave functions is sufficiently high that they must be extensively sampled as training examples for AI models by example. There is an opportunity to accelerate this process by including physics principles in networks and model architecture. Fourth, while tensor networks offer a degree of explainability, the overarching goal is to enhance the interpretability of quantum-inspired ML models, ensuring that insights and predictions can be easily understood and utilized.

---

## QUANTUM MACHINE LEARNING

---

Quantum computing introduces a paradigm shift in processing information, offering the potential to solve specific tasks with unprecedented speed and efficiency. Unlike classical computing bits, entangled quantum bits (i.e., qubits), can encode and elaborate information extremely efficiently. This capability is particularly useful for classes of computation in which one can obtain polynomial or even exponential speed ups compared to classical computations. Hybrid solutions in which part of the computations are performed on a quantum processor and other on classical ones are already used in industry, although a clear quantum advantage has not been proven.

This vision also involves an important hunt: to find novel machine learning applications that are only feasible and verifiable on quantum computers but utterly impractical on classical computers. Quantum algorithms, analog devices for computing, and quantum neural networks (QNNs) are at the forefront. Many problems at the atomic level and the nanoscale are fundamentally quantum in nature, positioning quantum computing as a tailor-made solution for these challenges. Machine learning tasks with a large amount of classical data is instead still an unsolved problem as encoding large amount of data into a quantum state still remains a very challenging task.

## Singapore's Role

The AI for science community has reached a tipping point. The ability to generate, interpret, and model complex data by the complexity and physics community in Singapore is unprecedented. We have a unique composition of researchers from various domains spread across reputable institutions on the island, each contributing a critical perspective on different scales of the problem. These individuals are characterized by their experience in AI for Science (AI4SCI), possessing the necessary data, experience, and computing resources to create ambitious AI models for physics, complexity science, and beyond.

Yet these exciting opportunities also raise serious concerns. Many of AI's potential dangers are widely known, especially deceptive AIs that persuade using asymmetric access to validated information and AIs that galvanize with populism. Any theory, right or wrong, will have purchase when enough people believe it. Despite negativism about this growing threat, science can become a trustworthy beacon in uncertain times.

In this regard, trustworthy AI in science will play an important role. This trustworthiness can manifest in the form of institutionally endorsed AI models. If created explainably, these models can also increase the scientific literacies of a much broader demographic.

Quantum-inspired solution, typically tensor-network based, are also a path to potentially accelerate machine learning and improve interpretability. Recently they have been used to compress neural networks and significantly reduce the number of parameters.

The development of machine learning algorithms for quantum computers is still at early stages. For example, fault-tolerant QML algorithms and processors are needed because practical implementations of qubits inevitably contain errors. Furthermore, there is a consensus that a theory for quantum learning needs to be developed. This theoretical framework will underpin the capabilities and limitations of quantum algorithms in learning and generalization.

This broader role, however, is not guaranteed, especially the models' embedded scientific principles are impenetrable to the general public. With social distrust mounting, the window to effect this role is rapidly shrinking.

However, this trust must involve critical thinking. Since scientific theories approximate our current understanding of the world, many will inevitably be supplanted by more accurate and robust ones. Our collective memories allow us to forget our falsified theories, but how will an AI forget them? Or will AI models forever circulate their traces, like the seemingly indestructible falsehoods and "urban legends" that will not dissipate in social media?

Retraining our models in science is also extraordinarily important for science to be self-correcting. This retraining requires an open availability of models and data.

There are even broader societal and political questions that Singapore must consider. Is it wise for Singapore to rely on data warehousing and models that are driven by commercial interests, or funded by philanthropy? If so, what restrictions come with such models? Further, who will be responsible for validating these models, checkpointing their learned biases, and regulating their use in different demographics and domains?

## Conclusion

In conclusion, the integration of AI in many-body physics, quantum machine learning, complexity and data-driven discovery offers transformative potential while also presenting significant challenges. The interplay between AI and experimental observations underscores the importance of data curation and labeling, as human expertise remains essential for synthesizing new materials and evaluating new theories and hypothesis suggested by AI models. As machine learning becomes more prevalent in the analysis of high-throughput and multi-modal experiments, the challenge of matching vast experimental data with limited insights from first-principles models creates a need for effective pattern recognition and interpretability methods. We want to stress the importance of addressing the double bottleneck of labeling and interpretability for leveraging AI's capabilities to model phenomena across multiple time and length scales, ultimately leading to breakthroughs in understanding complex systems.

This white paper has also highlighted the need for the creation of an atlas of multimodal maps that can accelerate the synthesis of functional quantum motifs in low-dimensional materials. By employing AI guidance, we envision achieving over a 1000-fold acceleration in rational synthesis, tackling the current bottleneck of synthesizing functional quantum defects identified through high-throughput simulations. This would be the first comprehensive multi-scale structural map of synthesizable complex quantum defects, encompassing diverse synthesis conditions and could be realised by incorporating advanced techniques such AI-assisted microscopy and physics-informed neural networks.

As Singapore navigates the challenges and opportunities created by AI, fostering an environment of critical thinking and open collaboration will be essential for harnessing the full potential of AI in scientific discovery. This would ensure that the country serves as a beacon of knowledge and innovation in uncertain times while addressing pressing environmental and technological challenges.

## REFERENCES

- 1 E. van der Giessen et al., *Roadmap on Multiscale Materials Modeling*, *Model. Simul. Mat. Sci. Eng.* **28**, 043001 (2020).
- 2 T. Schilling, *Coarse-Grained Modelling out of Equilibrium*, *Phys. Rep.* **972**, 1 (2022).
- 3 B. E. Husic et al., *Coarse Graining Molecular Dynamics with Graph Neural Networks*, *J. Chem. Phys.* **153**, 194101 (2020).
- 4 A. Tokuhisa, R. Kanada, S. Chiba, K. Terayama, Y. Isaka, B. Ma, N. Kamiya, and Y. Okuno, *Coarse-Grained Diffraction Template Matching Model to Retrieve Multiconformational Models for Biomolecule Structures from Noisy Diffraction Patterns*, *J. Chem. Inf. Model.* **60**, 2803 (2020).
- 5 J. Dan, X. Zhao, S. Ning, J. Lu, K. P. Loh, N. Duane Loh, and S. J. Pennycook, *Learning Motifs and Their Hierarchies in Atomic Resolution Microscopy*, *Science Advances* (2020).
- 6 P. Mehta and D. J. Schwab, *An Exact Mapping between the Variational Renormalization Group and Deep Learning*, <http://arxiv.org/abs/1410.3831>.
- 7 W. Chen, W. Wang, L. Liu, and M. S. Lew, *New Ideas and Trends in Deep Multimodal Content Understanding: A Review*, *Neurocomputing* **426**, 195 (2021).
- 8 R. K. Vasudevan, K. Choudhary, A. Mehta, R. Smith, G. Kusne, F. Tavazza, L. Vlcek, M. Ziatdinov, S. V. Kalinin, and J. Hattrick-Simpers, *Materials Science in the Artificial Intelligence Age: High-Throughput Library Generation, Machine Learning, and a Pathway from Correlations to the Underpinning Physics*, *MRS Commun.* **9**, 821 (2019).
- 9 P. Anderson, *More Is Different*, *Science* **177**, 393 (1972).
- 10 S. Strogatz, S. Walker, J. M. Yeomans, C. Tarnita, E. Arcaute, M. De Domenico, O. Artime, and K.-I. Goh, *Fifty Years of 'More Is Different,'* *Nature Reviews Physics* **4**, 508 (2022)
- 11 M. L. Leavitt and A. Morcos, *Towards Falsifiable Interpretability Research*, <http://arxiv.org/abs/2010.12016>.
- 12 A. Lavin et al., *Simulation Intelligence: Towards a New Generation of Scientific Methods*, <http://arxiv.org/abs/2112.03235>.
- 13 J. P. Crutchfield, *The Calculi of Emergence: Computation, Dynamics and Induction*, *Physica D* **75**, 11 (1994).
- 14 H. Wang et al., *Scientific Discovery in the Age of Artificial Intelligence*, *Nature* **621**, E33 (2023).
- 15 M. Krenn et al., *On Scientific Understanding with Artificial Intelligence*, *Nat. Rev. Phys.* **4**, 761 (2022).
- 16 J. Sourati and J. A. Evans, *Accelerating Science with Human-Aware Artificial Intelligence*, *Nat. Hum. Behav.* **7**, 1682 (2023).
- 17 Giuseppe Carleo, Matthias Troyer, "Solving the Quantum Many-Body Problem with Artificial Neural Networks", *Science* 355, 602 (2017)
- 18 P. Y. Lu, S. Kim, and M. Soljačić, *Extracting Interpretable Physical Parameters from Spatiotemporal Systems Using Unsupervised Learning*, *Phys. Rev. X.* **10**, (2020).
- 19 R. Yu and R. Wang, *Learning Dynamical Systems from Data: An Introduction to Physics-Guided Deep Learning*, *Proc. Natl. Acad. Sci. U. S. A.* **121**, e2311808121 (2024).
- 20 G. E. Karniadakis, I. G. Kevrekidis, L. Lu, P. Perdikaris, S. Wang, and L. Yang, *Physics-Informed Machine Learning*, *Nat Rev Phys* **3**, 422 (2021).
- 21 Z. Chen, Y. Liu, and H. Sun, *Physics-Informed Learning of Governing Equations from Scarce Data*, *Nat. Commun.* **12**, 6136 (2021).
- 22 J. Dan, M. Waqar, I. Erofeev, K. Yao, J. Wang, S. J. Pennycook, and N. D. Loh, *A Multiscale Generative Model to Understand Disorder in Domain Boundaries*, *Science Advances* **9**, eadj0904 (2023).
- 23 Maolin Wang, Yu Pan, Zenglin Xu, Guangxi Li, Xiangli Yang, and Andrzej Cichocki, "Tensor Networks Meet Neural Networks: A Survey and Future Perspectives", [arXiv:2302.09019](https://arxiv.org/abs/2302.09019) (2023).
- 24 D. Balakrishnan, S. W. Chee, Z. Baraissov, M. Bosman, U. Mirsaidov, and N. D. Loh, *Single-Shot, Coherent, Pop-out 3D Metrology*, *Communications Physics* **6**, 1 (2023).
- 25 K. X. Nguyen, Y. Jiang, C.-H. Lee, P. Kharel, Y. Zhang, A. M. van der Zande, and P. Y. Huang, *Achieving Sub-0.5-Angstrom-Resolution Ptychography in an Uncorrected Electron Microscope*, *Science* **383**, 865 (2024).
- 26 L. Zhang, J. Han, H. Wang, R. Car, and Weinan, *Deep Potential Molecular Dynamics: A Scalable Model with the Accuracy of Quantum Mechanics*, *Phys. Rev. Lett.* **120**, (2018).
- 27 S. Batzner, A. Musaelian, L. Sun, M. Geiger, J. P. Mailoa, M. Kornbluth, N. Molinari, T. E. Smidt, and B. Kozinsky, *E(3)-Equivariant Graph Neural Networks for Data-Efficient and Accurate Interatomic Potentials*, *Nat. Commun.* **13**, 2453 (2022).
- 28 Y. Ma, *Manifold Learning Theory and Applications* (CRC Press, Boca Roca, 2011).
- 29 B. T. Kiani, J. Wang, and M. Weber, *Hardness of Learning Neural Networks under the Manifold Hypothesis*, <http://arxiv.org/abs/2406.01461>.
- 30 G. Loaiza-Ganem, B. L. Ross, R. Hosseinzadeh, A. L. Caterini, and J. C. Cresswell, *Deep Generative Models through the Lens of the Manifold Hypothesis: A Survey and New Connections*, <http://arxiv.org/abs/2404.02954>.
- 31 T. Ayril, T. Louvet, Y. Zhou, C. Lambert, E. M. Stoudenmire, and X. Waintal, *Density-Matrix Renormalization Group Algorithm for Simulating Quantum Circuits with a Finite Fidelity*, *PRX Quantum* **4**, (2023) [WANGYANG23] Maolin Wang, Yu Pan, Zenglin Xu, Guangxi Li, Xiangli Yang, and Andrzej Cichocki, "Tensor Networks Meet Neural Networks: A Survey and Future Perspectives", [arXiv:2302.09019](https://arxiv.org/abs/2302.09019) (2023) [CARLEOTROYER17] Giuseppe Carleo, Matthias Troyer, "Solving the Quantum Many-Body Problem with Artificial Neural Networks", *Science* 355, 602 (2017)

## APPENDIX III. RNA BIOLOGY & THERAPEUTICS (AI4 RBT)

### AUTHORS:

Dr. Jason Pitt

(Principal Investigator, Cancer Science Institute of Singapore at National University of Singapore),

Prof. Ashok Venkitaraman

(Director, Cancer Science Institute at National University of Singapore and Director, Disease Intervention Technology Lab at Institute of Molecular and Cell Biology, A\*STAR)

### Executive summary

This white paper outlines Singapore's strategic plan to leverage Artificial Intelligence (AI) for advancements in RNA biology and therapeutics (AI4 RBT). This initiative builds upon recent government investments in RNA research and complements ongoing investigator-driven research across different Singaporean institutions. Key drivers of this initiative include: (i) the growing importance of RNA biology in human disease research; (ii) Singapore's ongoing commitment to AI

research; (iii) the current and future availability of rich RNA-related datasets through national initiatives. Following a multidisciplinary workshop attended by Singapore's leading researchers, clinicians and policy makers, we have formulated four key challenges in AI4 RBT. Addressing these challenges will help to position Singapore as a global leader in AI4 RBT, accelerating advancements in RNA biology and precision RNA medicine.

### Introduction

Recent advances in RNA biology & its applications for human disease have fueled a global scientific revolution. Our government has acknowledged that Singapore's future health and economic opportunities depend vitally on this field, and approved new investments into discovery science (National Initiative for RNA Biology & Its Applications, NIRBA) and clinical translation (Nucleic Acids Therapeutics Initiative, NATi). Besides these initiatives, a great deal of response-mode, investigator-initiated research funded in Singapore by agencies like the NMRC, MoE and other sponsors continues to focus in the area of RNA research.

In particular, NIRBA will generate unique large datasets whose exploitation promises to (a) provide powerful insights into the RNA biology of important diseases including cancer and cardiometabolic syndromes, specific to Asian genetic diversity in Singapore; (b) identify molecular targets for disease therapy & approaches for their modulation; and (c) potentiate the development of RNA-based tools for disease diagnosis, detection, stratification and prognosis. Collectively, these approaches will accelerate advances leading to precision RNA medicine. Complementary local and global datasets will also be generated by other funded grants and initiatives. These goals create a unique opportunity for synergy between RNA research and the AI4Science initiative.

## Background

In April 2024, AI4 Science along with NIRBA organized an AI4 RNA Biology and Therapeutics (AI4 RBT) workshop. Here, biologists, clinicians, computationalists, and policy experts gathered to discuss emerging opportunities at the intersection of RNA and AI. Particular emphasis was given to areas where existing local resources, infrastructure, and human capital offer global competitiveness. This was conducted as a series of short talks from domain experts followed by panel discussions. Major scientific themes included multimodal data analyses to elucidate transcriptional patterns, methods for predicting RNA structure – with particular emphasis on their interactions with proteins, how genetic variation may induce disease via disruption of RNA regulation, and how

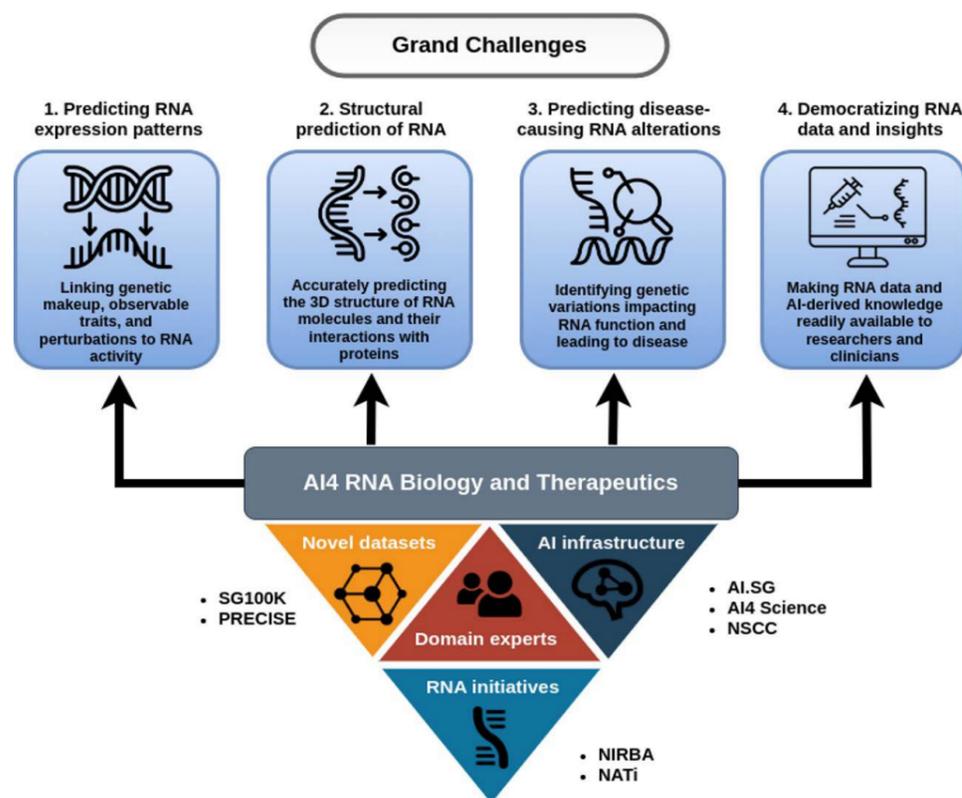
AI-driven RNA insights from vast collections of data can be leveraged by experimental biologists and clinicians.

Through these talks, panel discussions, and breakout sessions, these multidisciplinary experts provided thoughts on where Singapore is uniquely positioned to lead AI4 RBT at the global stage. How these efforts could be augmented via forward-thinking policy and state-of-the-art infrastructure was also discussed. Within this White Paper, we have consolidated these insights and opportunities into key challenges at the forefront of AI4 RBT. These key challenges aim to guide scientific efforts, resource allocation, and funding priorities for AI4 RBT in Singapore.

## Grand challenges

We identified *four AI4 RBT grand challenges* where Singapore is poised to lead globally impactful science:

Below we provide a foundation for these grand challenges, how cutting-edge AI approaches can address them, and the anticipated outcomes.



## Grand Challenges

### GRAND CHALLENGE 1: PREDICTING RNA EXPRESSION PATTERNS

For decades, geneticists have worked to decipher the genome to understand how patterns and variation within DNA sequence affect RNA transcripts and their regulation. While this can be explored via functional genomics, high-throughput technologies have largely rendered this a computational endeavor. Probabilistic approaches such as hidden Markov models (HMM) have been applied to reference DNA to use sequence context to predict functional genomic loci – including those likely to be transcribed<sup>1</sup>. However, the presence and activity of functional genomic elements within is far from deterministic. Across the 3.3 billion base pairs of the human genome, each individual contains millions of genetic variants – many of which can affect transcription<sup>2</sup>. The effect of DNA variants on the expression of RNA entities (expression quantitative trait loci [eQTLs]) has most commonly been inferred through linear models<sup>3</sup>. Notably, these methods study gene-transcript combinations in isolation and do not consider the influence of neighboring DNA regulatory elements such as transcription factor binding sites (TFBS). Recently, deep learning has been leveraged to combine genetic variation and sequence context to predict gene expression patterns within individuals. However, two key limitations have emerged: first, these models are trained to predict local gene expression (i.e. small windows of the genome) rather than global patterns; and second, state-of-the-art models are still underperforming at this task<sup>4</sup> – partially due to limited understanding of DNA's regulatory grammar<sup>5</sup>.

Modern generative AI models offer potential solutions to these limitations. Conditional variational autoencoders (CVAEs) and conditional generative adversarial networks (CGANs)<sup>6</sup> can be trained on large collections of RNA data to predict global transcriptional patterns based on user-defined features. This includes zero-shot learning scenarios where a

set of conditions (e.g. a combination of gene deficiencies) is not present in the training set. Conditions are not limited to genomic information and could be expanded by linking RNA to electronic medical health records (EMRs). Additionally, these approaches provide a non-linear latent space that can be used to infer broader RNA patterns across diseases, tissues, genotypes, and individuals<sup>7</sup>. Importantly, conditions of interest must be known a priori and provided as labels during model training. Datasets containing both DNA and RNA sequencing data will be critical. Capitalizing on other initiatives within the Singapore ecosystem – such as SG100k and PRECISE – would act as a force multiplier to understand the functional impact of genetic variants that are enriched within the local population.

The application of transformer-based large language models (LLMs) provide an opportunity to better understand DNA's regulatory grammar. Foundational LLMs have been created for DNA and have been fine-tuned to predict the impact of genetic variation on RNA expression within isolated segments of the genome (i.e. in *cis*)<sup>8,9</sup>. Such models can pinpoint transcripts likely affected by variants of unknown significance (VUS), which are substantially overrepresented in non-western populations<sup>10</sup>. The creation and integration of national datasets provides the opportunity to experimentally test hypotheses generated by task-specific LLMs to functionally characterize genomic variation relevant to the Singaporean population. Concerted efforts must be made to extract and tokenize biologically-relevant features from locally generated datasets in order to fine-tune LLM performance and maximize utility.

### OBJECTIVE

We propose to develop and train generative models to predict RNA expression – at the transcript, pathway, and transcriptome levels – based on biologically-relevant conditions. These conditions include gene inactivation, tissue type, chemical perturbations, therapeutic interventions, and disease.

Furthermore, foundational LLMs will be developed to understand DNA's regulatory grammar to better predict the functional effects of individual genetic variants. These foundational models will be fine-tuned to determine how variation within the Singaporean population can affect *cis* transcript regulation at disease-relevant genomic loci.

#### DATA REQUIREMENTS

This grand challenge requires thousands of RNA sequencing samples coupled with genotypic and phenotypic information. This can include large collections of biologically diverse data from The Cancer Genome Atlas (TCGA), the Gene Tissue Expression project (GTEx), the Encyclopedia of DNA Elements (ENCODE), etc. As mentioned above, locally generated DNA, RNA, and phenotype data from ongoing national initiatives such as SG100K, PRECISE, and NIRBA will be vital for model training and projecting findings onto the Singaporean population.

#### AI METHODS REQUIREMENTS

Conditional variational autoencoders, generalized adversarial networks, & diffusion models; multimodal data fusion; large language models

#### RESOURCE REQUIREMENTS

*Sequencing data processing and feature extraction*

- 12 million CPU hours (AMD EPYC 7713 or equivalent)
- 4 NVIDIA H100 GPUs (18 months)
- 300 TB of scratch disk space (6 months)

*Model training*

- 4 NVIDIA H100 GPUs (12 months)

#### GRAND CHALLENGE 2:

### STRUCTURAL PREDICTION OF RNA

Protein function is intimately linked to its tertiary (i.e 3D) structure. As such, decades of computational and experimental research have been dedicated to elucidate how protein structure dictates function and how this can be exploited therapeutically. In recent years, deep learning approaches have

vastly enhanced the accuracy of protein structure prediction<sup>11</sup>. RNA species have wide-spread functional roles within cells – including interactions with other RNAs, DNA, and proteins. Critically, cellular RNAs do not exist as linear structures as the biophysical properties of RNAs enable them to take a wide-variety of folding conformations<sup>12</sup>. Like with proteins, the fidelity of these structures are essential for their functional interactions. However, RNA structure prediction faces multiple unique challenges. **(a)** RNA has high tertiary structure complexity due to pseudoknot formation and promiscuous base pairing<sup>13</sup>. **(b)** Multiple stable intermediate, secondary, and tertiary structures can exist for a single RNA transcript<sup>14</sup>. **(c)** Chemical modifications to RNA transcripts may affect their structural conformations<sup>15</sup>. **(d)** The length and conformational flexibility of RNA transcripts creates a vast structural search space that is difficult to explore computationally<sup>16</sup>. In concert, these challenges have prevented the scientific community from properly characterizing, understanding, and manipulating RNA structures.

The diversity of RNA transcripts and the cellular conditions in which they exist makes systematic experimental resolution of 3D RNA structures a herculean, if not impossible, task. Nonetheless, even small sets of structures can help train sophisticated deep learning models to predict RNA structure from sequence alone. Capitalizing on the protein folding success of AlphaFold2, multiple research groups have developed end-to-end deep learning models for 3D RNA structure prediction<sup>16</sup>. These models are still in their nascent stages and have yet to experience the drastic improvements availed by protein folding.

RNA language models provide an opportunity to encode RNA folding grammar in order to more accurately predict structure<sup>17</sup>. Additionally, the generative nature of LLMs could be leveraged to design an RNA sequence that can mimic the tertiary structure of another. These generative approaches may enable RNA designs that can avoid degradation machinery or immune sensing<sup>18</sup>.

They may also provide a means of generating RNA aptamers that bind to proteins and ablate their function<sup>19</sup>. Importantly, there is increasing evidence that RNA modifications may affect RNA tertiary structures – particularly when those transcripts interact with proteins<sup>15</sup>. While RNA language models have been used to predict both structure and modifications<sup>20,21</sup>, RNA modifications themselves have not been comprehensively encoded within RNA language models to improve structural predictions. Sequencing via Oxford Nanopore Technologies (ONT) could provide long reads with in-phase modifications that could serve as auxiliary training data for more comprehensive RNA language models suitable for tertiary structure prediction.

#### OBJECTIVE

We propose to leverage unique RNA datasets that will be created through Singapore's investments in RNA research – which aim to incorporate overlooked features such as RNA chemical modifications and structure mapping – to build foundational AI models for accurate prediction of 3D RNA structures. Using SG100K to map genetic variation onto these structures, we can gain insight on how local diversity may affect RNA structure and subsequent function. Additionally, these foundational models will also be fine-tuned for therapeutically-relevant tasks such as the generation of stable and functional RNA aptamers.

#### DATA REQUIREMENTS

Known RNA transcripts from humans as well as other organisms will be collated from trusted sources such as Ensembl. Known and putative 3D protein and RNA structures will be taken from the Protein Data Bank (PDB) and FRABASE. Long-read RNA sequencing data from NIRBA and PRECISE will provide a rich source of transcript-specific RNA modifications that can augment structural prediction models.

#### AI METHODS REQUIREMENTS

Deep neural networks; large language models

#### RESOURCE REQUIREMENTS

*Model training*

- 6 NVIDIA H100 GPUs (24 months)
- 100 TB of scratch disk space (24 months)

#### GRAND CHALLENGE 3:

### PREDICTION OF DISEASE-CAUSING RNA ALTERATIONS

While DNA mutations are often linked to diseases by altering protein function, their impact extends beyond the proteome. Mutations can also disrupt RNA function leading to phenotypic consequences. Non-coding RNAs, which commonly regulate gene expression in a sequence-specific manner, are particularly susceptible to the effects of these mutations<sup>22</sup>. Their regulatory dysfunction can lead to diseases like cancer and neurodegenerative disorders<sup>23</sup>. Additionally, transcribed mutations can affect RNA processing, such as splicing, capping, and polyadenylation<sup>24</sup>. These aberrant transcripts may be unstable, poorly translated, or have altered functions. Furthermore, mutations can cause RNA misfolding<sup>25</sup>, leading to loss of function or the formation of toxic RNA aggregates. Some RNA molecules possess catalytic activity (e.g. ribozymes) or form intricate structures to interact with proteins. Mutations in these RNAs can disrupt their catalytic activity or protein-binding ability, which has been implicated in disease phenotypes<sup>26</sup>. Thus, considering the impact of mutations on the diverse, functional roles of RNA may be crucial for unraveling the molecular basis of many diseases.

Evolutionary scoring methods have been developed to predict pathogenicity of DNA mutations, including those transcribed into RNA. However, these approaches rely heavily on conservation, are not data-driven, and do not aim to predict specific functional effects<sup>27</sup>. Tools aiming to predict specific functional effects of mutations such as splicing, RNA structure, and RNA-protein interactions have been developed – with some even leveraging more modern AI methodologies<sup>28</sup>. Nonetheless, these tools remain limited in their predictive capabilities and consequently difficult to interpret. Their disparate nature is also not congruent with known biological phenomena. While distinct models have been constructed to predict mutational effects on RNA binding site affinity and RNA structure, these characteristics are biophysically coupled. A mutation within a transcript could

alter its 3D conformation such that an RNA binding site is no longer accessible. Similarly, DeepMind has recently demonstrated substantial benefits to pathogenicity predictions of amino acid substitutions when overlaid on protein structure<sup>29</sup>. Comprehensive AI models that simultaneously consider RNA's intertwined biophysical properties will be necessary to accurately predict relationships between mutation and transcript dysfunction – especially for the plethora of understudied variants within Singaporean genomes.

#### OBJECTIVE

We propose to construct multi-task convolutional deep neural networks to predict the phenotypic and disease consequences of mutations within known human RNA transcripts. Learned tasks will include transcript constitution (e.g. splicing, intron retention, etc.), stability, binding site disruption & affinity, and 3D structure. The shared representation will capture interdependencies amongst these tasks, which more accurately represents known biophysical properties of RNA. Once trained, the aim of these models is to take a mutation of interest, predict (in)direct functional consequences for RNA transcripts, and estimate pathogenicity with respect to disease.

#### DATA REQUIREMENTS

Large local (e.g. NIRBA & PRECISE) and public (e.g. GTEx) datasets that contain short- or long-read RNA sequencing and germline/somatic WGS/WES from the same individuals. Curated datasets of known RNA structures, splice isoforms, RNA binding motifs, regulatory RNAs, transcript disrupting mutations, and risk variants will also be utilized.

#### AI METHODS REQUIREMENTS

Multi-task convolutional deep neural networks; Large language models (e.g. RNA language models from Section 4.2)

#### RESOURCE REQUIREMENTS

Sequencing data processing and feature extraction

See Section 4.1.4

Model training

- 2 NVIDIA H100 GPUs (24 months)

---

#### GRAND CHALLENGE 4: DEMOCRATIZING RNA DATA AND INSIGHTS

---

The exponential increase of genomic data has served as a key substrate for scientific discovery over the past decade. However, the benefits reaped from this data have not been distributed uniformly. Those with the necessary technical prowess, namely computation scientists, have been able to wrangle, structure, and analyze these data. Experimental scientists and clinicians experience greater difficulty converting information into insights. To mitigate this disparity – scientists, institutions, and governments have developed user-friendly platforms in an attempt to level the playing field. For example, The Genomic Data Commons (GDC) is an initiative funded by the US National Institutes of Health to store, harmonize, and distribute cancer genomics data from major sequencing initiatives<sup>30</sup>. GDC also includes an interactive web portal where users can perform analyses over structured genomic data such as mutations or gene expression. Despite these efforts, users still require some technical knowledge of web portal usage and capabilities. Chat-like approaches via AI agents provide one avenue to circumvent this requirement<sup>31</sup>. Recently, such agents have been used for genomics to allow basic analyses and visualizations over high-value data using semi-structured text (e.g. DrBioRight<sup>32</sup>) and voice (e.g. Melvin<sup>33</sup>) queries. However, these existing agents are throttled by their: **(a)** basic natural language understanding (NLU) capabilities; **(b)** limited data and analysis methods; and **(c)** lack of data models to enable ongoing ingestion of data. The next generation of AI agents must circumvent these limitations in order to truly democratize and fully utilize genomic data in research and clinical settings.

The multi-headed attention mechanisms employed over large text corpuses have substantially improved NLU from free-text queries. This has enabled generalized LLMs such as ChatGPT and Gemini to understand, contextualize, and appropriately respond to users throughout multi-turn conversations. Existing foundational LLMs

can and should be fine-tuned in order to understand genomic queries related to research and medicine<sup>34</sup>. Notably, these fine-tuned LLMs must be connected to AI-powered analytical frameworks deployed over structured databases – similar to Google's Query understanding or Amazon's Kedra. Such frameworks can be designed or trained to perform domain-relevant queries in a statistically sound manner. Singapore's RNA efforts such as NIRBA and NATi – supplemented by auxiliary data from PRECISE and SG100K – offer a rich multimodal data source for such queries. The ability to seamlessly execute such queries with trustworthy results will hinge on a rigorously vetted data model. These data models will need to handle a litany of genomic technologies, clinical information, and data modalities (e.g. numerical, categorical, text, image, etc.). These data models will not only be critical for the AI-powered analytical framework but also will ensure standardization during the continuous ingestion of data. Overall, such a framework would ensure Singapore's investment in RNA will be maximized to drive scientific discovery and clinical insight across the island.

#### OBJECTIVE

To develop an AI-powered analytics platform that enables researchers and physicians to gain biological and clinical insight to Singapore's high-value RNA data through multi-turn conversations. This platform will consist of four modules: **(a)** a structured, multimodal database adhering to a domain-specific data model; **(b)** an AI engine to query this database and perform appropriate analytics and visualizations in real-time; **(c)** an AI agent powered by fine-tuned LLMs to understand, contextualize, and respond to user input; **(d)** an intuitive interface by which researchers can converse with this agent and view results.

#### DATA REQUIREMENTS

Hundreds of thousands of quantitative values representing RNA features (e.g. expression, splicing, modifications, intro retention, etc.) from thousands of locally generated RNA-sequencing samples with rich phenotypic annotations. Each sample can be supplemented with other data modalities (e.g. DNA sequencing, metabolomics, pathological/radiological imaging) where available and appropriate. In addition to locally generated data, features derived from publicly available RNA sequencing datasets (see Section 4.1) can be harmonized and ingested.

#### AI METHODS REQUIREMENTS

Large language models; domain-adapted text-to-query tools; automated data model generation

#### RESOURCE REQUIREMENTS

Sequencing data processing and feature extraction

See Section 4.1.4

Model training

- 2 NVIDIA H100 GPUs (36 months)
- 50 TB of persistent disk space (36 months)

Platform deployment (via commercial cloud providers)

- 1 NVIDIA T4 GPUs (24 months)
- 25 TB data storage within a low-latency database e.g. AWS RDS or Aurora (24 months)

## AI methods and data - challenges and opportunities

### DATA

#### DATA DIVERSITY AND QUALITY

While the grand challenges outlined by this whitepaper emphasize RNA-based discoveries, auxiliary data are critical to contextualize AI model predictions with respect to biology and disease. These additional modalities include genetic, genomic, phenotypic, health record, and medical imaging data. It is important to select samples/patients with the greatest coverage across each of the modalities of interest. This will bolster model training, improve predictions, and increase model interpretability. Nonetheless, specific model architectures and learning strategies can be employed to handle heterogeneous datasets with only partially overlapping modalities and feature sets (see Section 5.2). This is particularly important when leveraging large collections of public data in addition to locally generated datasets.

#### DATA INTEGRATION

The aforementioned grand challenges come with four key considerations with respect to data integration: **(a)** Data must be harmonized to minimize cross-dataset biases due to sample preparation, data processing, or feature definitions. These considerations are relevant to all RNA, genomic, and clinical data. **(b)** To leverage all available data, AI model architectures must be able to handle heterogeneity in modalities and features. Some of these strategies are discussed in Section 5.2. **(c)** The features underlying each of our modalities are not guaranteed to share the same statistical properties (e.g. mean & variance) or variable types (e.g. quantitative & categorical). Naively concatenating these variables may lose information critical for accurate inference. Employing sophisticated techniques to independently model and subsequently fuse modalities on a shared latent space will likely be necessary<sup>35</sup>. **(d)** As discussed in Section 4.4, well-defined data models will be necessary in order to develop database schemas that will facilitate AI-powered insights that combine RNA, genomic, and clinical modalities.

#### DATA UTILIZATION

The critical multimodal data sources fueling these grand challenges in RNA will come from Singapore national initiatives such as NIRBA, NATi, PRECISE, and SG100K. However, large collections of publicly available data will be harmonized to local datasets to enhance AI model training, generalizability, and inference. This large compendium of harmonized datasets will be made available to the broader Singapore clinical and research community (Section 4.4).

### AI METHODS

While the specific AI methods to be employed have been discussed in Section 4, one key technical consideration persists across each grand challenge. The heterogeneous nature of genomics and clinical datasets will influence model architectures and training strategies. Disparate sources of clinical and genomic data very rarely share all data modalities or features within a given modality. Properly handling datasets with non-overlapping modalities/features is crucial to maximize AI model learning and ensure generalizability. As such, AI4 RBT will need to consider multi-task or hierarchical learning strategies to mitigate this issue. Such approaches can derive task or dataset specific layers that can be fused into shared representation for downstream prediction tasks.

Importantly, generative models (VAEs, GANs, LLMs, etc.) can require substantial computational resources for training. The amount of resources required depends on the number of features and dimensionality of the training data as well as the complexity of the model (e.g. number of parameters, hidden layers, latent space size, etc.) Training the most advanced general purpose LLMs employed by OpenAI and Google have required \$78 and \$191 million USD in computational resources, respectively<sup>36</sup>. The United State's Oak Ridge National Laboratory has recently deployed Frontier, the world's fastest and first exascale supercomputer. The GPU-heavy resource

enables academic researchers to employ powerful LLMs to scoped yet critical scientific problems. Similar specialized resources, albeit of smaller scale, will be necessary to exploit the unique and high-value RNA data being generated within Singapore. Sufficient

computational infrastructure is necessary to enable domain experts, AI theorists, and HPC scientists to coalesce within Singapore to disentangle unknown relationships between RNA, DNA, and disease.

## Concluding remarks

We foresee that potentially critical advantages in the AI4RBT arena will be realized through (a) national investment in AI & the AI system offered by the National Supercomputing Centre's ASPIRE2A and other clusters, (b) symbiotic interactions between the AI4RBT initiative and national initiatives in RNA

research and its therapeutic application, to generate high-value multimodal datasets, and (c) the coalescence through these activities of domain experts in RNA biology, medicine, AI, & computing, creating a powerful environment for value creation in Singapore through health impact and commercialization opportunities.

## REFERENCES

- 1 Yoon, B.-J. Hidden Markov Models and their Applications in Biological Sequence Analysis. *Current Genomics* **10**, 402–415.
- 2 Karczewski, K. J. *et al.* The mutational constraint spectrum quantified from variation in 141,456 humans. *Nature* **581**, 434–443 (2020).
- 3 Gilad, Y., Rifkin, S. A. & Pritchard, J. K. Revealing the architecture of gene regulation: the promise of eQTL studies. *Trends in Genetics* **24**, 408–415 (2008).
- 4 Huang, C. *et al.* Personal transcriptome variation is poorly explained by current genomic deep learning models. *Nat Genet* **55**, 2056–2059 (2023).
- 5 Sasse, A. *et al.* Benchmarking of deep neural networks for predicting personal gene expression from DNA sequence highlights shortcomings. *Nat Genet* **55**, 2060–2064 (2023).
- 6 Alqahtani, H., Kavakli-Thorne, M. & Kumar, G. Applications of Generative Adversarial Networks (GANs): An Updated Review. *Archives of Computational Methods in Engineering* **28**, 525–552 (2019).
- 7 Kopf, A. & Claassen, M. Latent representation learning in biology and translational medicine. *Patterns (N Y)* **2**, 100198 (2021).
- 8 Avsec, Ž. *et al.* Effective gene expression prediction from sequence by integrating long-range interactions. *Nat Methods* **18**, 1196–1203 (2021).
- 9 Jaganathan, K. *et al.* Predicting Splicing from Primary Sequence with Deep Learning. *Cell* **176**, 535–548. e24 (2019).
- 10 Slavin, T. P. *et al.* Prospective Study of Cancer Genetic Variants: Variation in Rate of Reclassification by Ancestry. *J Natl Cancer Inst* **110**, 1059–1066 (2018).
- 11 Jumper, J. *et al.* Highly accurate protein structure prediction with AlphaFold. *Nature* **596**, 583–589 (2021).
- 12 Higgs, P. G. RNA secondary structure: physical and computational aspects. *Q Rev Biophys* **33**, 199–253 (2000).
- 13 Pleij, C. W. Pseudoknots: a new motif in the RNA game. *Trends Biochem Sci* **15**, 143–147 (1990).
- 14 Revolutions in RNA Secondary Structure Prediction. *Journal of Molecular Biology* **359**, 526–532 (2006).
- 15 Lewis, C. J. T., Pan, T. & Kalsotra, A. RNA modifications and structures cooperate to guide RNA-protein interactions. *Nat Rev Mol Cell Biol* **18**, 202–210 (2017).
- 16 Zhang, S., Li, J. & Chen, S.-J. Machine learning in RNA structure prediction: Advances and challenges. *Biophys J* **123**, 2647–2657 (2024).
- 17 Li, H.-L., Pang, Y.-H. & Liu, B. BioSeq-BLM: a platform for analyzing DNA, RNA and protein sequences based on biological language models. *Nucleic Acids Res* **49**, e129 (2021).
- 18 Bartok, E. & Hartmann, G. Immune Sensing Mechanisms that Discriminate Self from Altered Self and Foreign Nucleic Acids. *Immunity* **53**, 54–77 (2020).
- 19 Di Ruscio, A. & de Franciscis, V. Minding the gap: Unlocking the therapeutic potential of aptamers and making up for lost time. *Mol Ther Nucleic Acids* **29**, 384–386 (2022).
- 20 Yin, W. *et al.* ERNIE-RNA: An RNA Language Model with Structure-enhanced Representations. *bioRxiv* 2024.03.17.585376 (2024) doi:10.1101/2024.03.17.585376.
- 21 Liang, S. *et al.* Rm-LR: A long-range-based deep learning model for predicting multiple types of RNA modifications. *Comput Biol Med* **164**, 107238 (2023).
- 22 Diederichs, S. *et al.* The dark matter of the cancer genome: aberrations in regulatory elements, untranslated regions, splice sites, non-coding RNA and synonymous mutations. *EMBO Molecular Medicine* (2016) doi:10.15252/emmm.201506055.
- 23 Esteller, M. Non-coding RNAs in human disease. *Nature Reviews Genetics* **12**, 861–874 (2011).
- 24 Cooper, D. N. Human Gene Mutations Affecting RNA Processing and Translation. *Annals of Medicine* (1993) doi:10.3109/07853899309147851.
- 25 Russell, R. RNA misfolding and the action of chaperones. *Front Biosci* **13**, 1–20 (2008).
- 26 Cech, T. R. Ribozymes and their medical implications. *JAMA* **260**, 3030–3034 (1988).
- 27 Garcia, F. A. de O., Andrade, E. S. de & Palmero, E. I. Insights on variant analysis in silico tools for pathogenicity prediction. *Front. Genet.* **13**, 1010327 (2022).
- 28 Wagner, N. *et al.* Aberrant splicing prediction across human tissues. *Nature Genetics* **55**, 861–870 (2023).
- 29 Cheng, J. *et al.* Accurate proteome-wide missense variant effect prediction with AlphaMissense. *Science* **381**, eadg7492 (2023).
- 30 Heath, A. P. *et al.* The NCI Genomic Data Commons. *Nature Genetics* **53**, 257–262 (2021).
- 31 Nowogrodzki, J. ChatGPT for science: how to talk to your data. *Nature* **631**, 924–925 (2024).
- 32 Li, J., Chen, H., Wang, Y., Chen, M.-J. M. & Liang, H. Next-Generation Analytics for Omics Data. *Cancer Cell* **39**, 3–6 (2021).
- 33 Perera, A. R. *et al.* Melvin is a conversational voice interface for cancer genomics data. *Communications Biology* **7**, 1–6 (2024).
- 34 Thirunavukarasu, A. J. *et al.* Large language models in medicine. *Nature Medicine* **29**, 1930–1940 (2023).
- 35 Lahat, D., Adali, T. & Jutten, C. Multimodal data fusion: An overview of methods, challenges, and prospects. *Proc. IEEE Inst. Electr. Electron. Eng.* **103**, 1449–1477 (2015).
- 36 Chen, S. A day in the life of the world's fastest supercomputer. *Nature* **633**, 22–25 (2024).

## APPENDIX IV. BIOMEDICAL SCIENCES

### AUTHORS:

Dr. Hongye Ye

Dr. Sebastian Maurer-Stroh  
(Bioinformatics Institute BII, A\*STAR)

### Executive Summary

The workshop on “AI for Biomedical Sciences” brought together leading experts to identify and articulate the grand challenges in leveraging artificial intelligence (AI) for significant advancements in biomedical sciences. The workshop explored AI’s transformative potential in drug discovery,

synthetic biology, and disease modeling, with a specific focus on Singapore’s unique capabilities and opportunities. This white paper synthesizes the workshop’s findings and outlines a strategic vision to harness AI for groundbreaking scientific discoveries and improved healthcare outcomes.

### Introduction

The AI4SCI initiative aims to leverage AI technologies to address scientific problems and achieve significant breakthroughs. Distinguishing itself from AI applications in other domains, this initiative focuses on applying AI to challenges in biomedical sciences, with the ultimate goal of exploring

AI’s potential to achieve clinical, societal and economic impact. The workshop’s purpose was to define the grand challenges within AI for biomedical sciences, develop a vision statement, and outline the potential for AI to bridge foundational and applied research.

### Background

The burgeoning landscape of AI for Biomedical Sciences has over 4,000 companies globally engaged in research and development. Despite significant advances, successfully navigating drugs through clinical trials remains difficult, with a notable proportion of both AI-designed and human-designed drugs failing during this crucial phase. Singapore aims to capitalize on its strengths in data quality, expertise, and crucially the diversity of its gene pool, to carve out a competitive niche in AI-driven drug discovery.

A number of challenges were discussed, including limitations in AI models to leverage multi-modal/ multi-omic datasets and to

generate valid hypothesis for experimental planning. Improving computational resources and addressing the need for increased model explainability were identified as crucial areas for future development. It was mentioned that efforts are also underway to optimize AI models and datasets, with a focus on improving prediction accuracy and exploring AI’s potential in identifying novel drug targets from genomic data. Overall, the workshop underscored the promise of AI in biomedical sciences, the challenges it faces, and the opportunities for Singapore to lead the way in this critical domain, as will be delineated within this white paper.

## Grand Challenges

### GRAND CHALLENGE 1: AI IN DRUG DISCOVERY

AI has the potential to revolutionize drug discovery by accelerating the identification of viable drug candidates, reducing costs, and minimizing late-phase failures in clinical trials. This transformation could lead to more effective and affordable treatments for various diseases in over a shorter development period, significantly impacting global health.

The traditional drug discovery process is lengthy, costly, and prone to high failure rates, particularly in the later stages of clinical trials. It can take up to 10 to 15 years, and costs between US\$1-3B per drug developed, on average<sup>1,2,3</sup>. This results in significant financial losses and delayed patient access to new treatments. Over 90% of drugs fail during clinical trials, often due to unforeseen toxicity or lack of efficacy. An estimated 86% of drug candidates developed between 2000 and 2015 did not meet their stated endpoints due to major reasons including the lack of clinical efficacy (40%–50%), unmanageable toxicity (30%) and poor drug-like properties (10%–15%)<sup>4</sup>. Within the current Singapore ecosystem, many drug discovery and development projects tend to get stuck in the pre-clinical development phase, with limited avenues for further acceleration. Moreover, the gap in late-stage capabilities<sup>5,6</sup> in the local life sciences sector could impede our ability to tackle big unmet medical needs.

AI can address these challenges by enhancing the identification of drug targets, optimizing lead compounds, and predicting clinical outcomes more accurately in a clearly identified sub-population of patients. AI-driven drug discovery is also key in our advancement towards precision medicine.

The potential to apply AI to accelerate and improve drug discovery has garnered growing interest globally, with major breakthroughs in recent years such as the advent of AlphaFold<sup>7</sup> – which has been shown to be able to predict the 3D structure of proteins rapidly and accurately and its recent upgrade touted to have expanded in terms of scope to cover other biological molecules like RNA and DNA<sup>8</sup>. With the emergence in ChatGPT in recent years, and Insilico Medicine's generative AI drug design platform amassing a pipeline of 17 preclinical candidates with several in clinical trials, generative AI has rapidly gained prominence for its promise in drug discovery and development<sup>9</sup>. Large pharmaceutical companies are also rapidly adopting AI technologies to transform the drug discovery and development process.

According to GlobalData, the pharmaceutical industry is forecast to spend \$3.3B in AI by 2025 with a compound annual growth rate (CAGR) of 24.4% from 2019. Despite emerging value proofs, the potential of AI has yet to be demonstrated at scale, across population and diseases. This could be attributed to the key barriers identified in BCG's 2023 report<sup>10</sup>, which include shortcomings in the access, maturity and standardisation of data, tools and capabilities. There is need for **high-quality data, reproducible experiments, and robust validation** to ensure model performance<sup>11</sup>.

Over the last year, a number of AI-designed drugs by first generation biotechs such as Exscientia and BenevolentAI have fell short in clinical trials or have been deprioritised<sup>12</sup>, painting a more cautionary tale to the limits of AI capabilities in drug discovery. The foundation of effective drug discovery is ultimately underpinned by robust biological understanding and there is increasing recognition on the critical importance of generating and possessing own unique biological data at scale<sup>13</sup>.

### OBJECTIVE

- **Reduced Time and Cost:** Accelerated drug discovery processes can reduce the average development time by up to 50% and cut costs significantly.
- **Increased Success Rates:** Improved prediction models can enhance the success rates of clinical trials, reducing late-phase failures.
- **Faster Delivery of Effective Therapies:** Quicker development timelines will enable faster patient access to new treatments, improving overall health outcomes.

### DATA REQUIREMENTS

Access to high-quality clinical, genomic, and pharmaceutical data is essential. Singapore's existing databases and healthcare records can be leveraged for AI training and validation. Initiatives like the National Precision Medicine Programme provide a wealth of genomics and other omics data that can be harnessed for AI applications. In particular, Singapore's data as a unique representation of a mixture of different ethnic groups of Asian populations. This differentiates itself from other international efforts, resulting in more effective and safer drugs for the APAC region.

### AI METHOD REQUIREMENTS

There is a lot of research activities in Singapore on techniques such as machine learning for predicting drug-target interactions, generative models for novel compound design, and deep learning for clinical outcome prediction. AI is also being developed to facilitate virtual screening, molecular docking, and optimization of pharmacokinetic and pharmacodynamic properties.

### RESOURCE REQUIREMENTS

Advanced computational resources, interdisciplinary research teams, and substantial funding for AI and biomedical research integration. Investments in high-performance computing and cloud infrastructure are critical to support large-scale AI applications. A total of 600,000 hours of GPU usage, 1,650,000 hours of CPU usage and a total storage of 3,000 TB is estimated over a 4 year period. This could amount up to more than S\$7 million for the computation.

### GRAND CHALLENGE 2:

### AI IN SYNTHETIC BIOLOGY

Synthetic biology is a multidisciplinary field of science that involves redesigning organisms (including biological parts, devices, systems) for useful purposes by engineering them to have new abilities. AI can significantly enhance synthetic biology by automating and optimizing the design and development of new biological systems. This includes creating novel enzymes, pathways, and organisms for various applications, including healthcare, agriculture, and environmental sustainability. The integration of AI can transform synthetic biology into a more predictable and scalable engineering discipline.

Synthetic biology faces challenges such as the unpredictability of biological systems and the complexity of designing functional biological components. This results in high cost and long development times associated with synthetic biology projects. The current cycle times for developing new biological systems can span several years and cost millions of dollars.

AI can help reduce these timelines and costs by enhancing design accuracy and enabling high-throughput experimentation. AI can be integrated into predictive models and automation tools that accelerate and streamline the engineering process. AI-driven synthetic biology can lead to innovations such as biosynthetic production of pharmaceuticals, bio-based materials, and sustainable bioenergy solutions.

### OBJECTIVE

- **Accelerated Development:** AI-driven approaches can significantly reduce the development time for new biological products.
- **Reduced Costs:** Optimized workflows and predictive models can lower the costs associated with synthetic biology projects.
- **Enhanced Scalability:** AI can improve the scalability and reliability of synthetic biology applications, making them more viable for industrial and commercial use.

#### DATA REQUIREMENTS

High-quality, standardized datasets on genetic sequences, metabolic pathways, and bioprocess parameters. Data integration from various sources, including genomic, transcriptomic, proteomics and chemistry (e.g. mass-spectrometry), is crucial.

#### AI METHOD REQUIREMENTS

The advancement in machine learning for predictive modeling, deep learning for protein engineering, and AI-driven automation for high-throughput screening are all ways AI can help boost synthetic biology in Singapore. Especially with the recent development in large-language models and generative AI, the design of synthetic products and processes can be further optimized.

#### RESOURCE REQUIREMENTS

Investment in curated large datasets, linked data between genomic, chemistry and bioactivity data, enhanced data sharing protocols, and collaboration between AI researchers and synthetic biologists are required.

---

#### GRAND CHALLENGE 3: DATA HARMONIZATION FOR AI USE

---

A large factor that affects the accuracy of an AI model is the quality of its training data. Data harmonization is crucial for the effective application of AI in biomedical sciences. Harmonizing diverse datasets from different sources can improve the accuracy and generalizability of AI models, leading to better healthcare outcomes and more robust scientific discoveries.

The success of AI in biomedical sciences heavily depends on the quality and integration of data. Data fragmentation and lack of standardization result in inconsistencies that hinder the training and validation of AI models. This fragmentation can lead to

biased results, reduced model performance, and slower scientific progress. Harmonizing data can improve model accuracy and reproducibility, which is critical for translating AI research into more accurate and reliable practical applications.

#### OBJECTIVE

- Improved Model Accuracy: Harmonized data will improve data quality for AI model training and enhance the accuracy and reliability of these models, leading to better healthcare outcomes.
- Enhanced Collaboration: Standardized data protocols will facilitate data sharing and collaboration across institutions and even borders.
- Accelerated Scientific Discoveries: Comprehensive and high-quality datasets will enable more robust and reproducible research, accelerating scientific progress.

#### DATA REQUIREMENTS

Comprehensive and standardized datasets from clinical, genomic, socioeconomic and environmental data types are necessary. Singapore's healthcare systems and research institutions can provide high-quality data for AI applications.

#### AI METHOD REQUIREMENTS

Techniques such as data normalization, integration frameworks, and federated learning can facilitate the harmonization of diverse datasets. AI can also be used to identify and correct inconsistencies, ensuring high data quality. AI can also be used to impute data to reduce wastage due to missing data.

#### RESOURCE REQUIREMENTS

Investments in secured data infrastructure, standardization protocols, and collaborative frameworks are essential. Support for interdisciplinary teams to develop and implement harmonization strategies is also critical.

## AI Methods and Data - Challenges and Opportunities

Several challenges to AI methods development and data collection remain to be addressed before biomedical sciences can fully benefit from AI technology. Current AI methods, such as deep learning for medical imaging and natural language processing for electronic health reports (EHRs), face issues like interpretability, data bias, and limited generalizability. Furthermore, the collection of diverse, high-quality, and standardized data is hindered by privacy concerns, inconsistent data formats, and imbalances across populations. While AI shows promise in improving disease prediction and treatment, its widespread adoption is slowed by the need for clinical validation, data security, and

the resolution of these data and algorithmic challenges. However, AI in biomedical sciences presents several exciting opportunities. First, federated learning offers a way to train AI models across decentralized datasets while preserving patient privacy, addressing security concerns in data-sharing. Second, the use of wearable technologies and continuous health monitoring devices provides more consistent, high-resolution longitudinal data, improving the accuracy of AI models for disease prediction and management, especially in chronic conditions. These opportunities have the potential to accelerate innovation while safeguarding patient privacy and enhancing personalized care.

## Singapore's Role

Singapore's robust biomedical research infrastructure, combined with its strong AI capabilities, positions it uniquely to lead in AI-driven drug discovery. Singapore has world-class research institutions and a collaborative ecosystem that bridges academia, industry, and government. There is also close coupling of AI development with biological validation which is an important feature as AI models need to be validated and integrated<sup>14,15,16</sup> into the drug discovery process in a reliable, scalable, and efficient way. Tight coupling of biological validation to the AI models will be the cornerstone of the programme in order to ensure that results can be reproduced and reduce the chance of failure when entering clinical trials.

Singapore's leadership in biofoundries and integrated biomanufacturing facilities provides a solid foundation for advancing synthetic biology. Consortium like the Singapore Consortium for Synthetic Biology (SINERGY) and facilities in Synthetic Biology for Clinical and Technological Innovation (SynCTI) foster innovation and collaboration.

Singapore's commitment to data integration and its advanced healthcare infrastructure makes it an ideal leader in data harmonization efforts. The nation's initiatives in health data management and participation in international data-sharing collaborations strengthen its position. There has been ongoing efforts driven by the Singapore Ministry of Health to harmonize clinical data using the Observational Medical Outcomes Partnership (OMOP) data standard.

## Conclusions

The AI for Biomedical Sciences initiative presents a significant opportunity for Singapore to lead globally in leveraging AI for transformative advances in healthcare and biomedical research. By addressing the grand challenges in drug discovery, synthetic biology, and data harmonization, Singapore can drive innovation, improve public health outcomes, and establish itself as a leader in the AI-driven biomedical sciences domain. This white paper outlines the vision, challenges, and strategic approaches necessary to achieve these ambitious goals, emphasizing the need for interdisciplinary collaboration, robust data infrastructure, and sustained investment in AI and biomedical research.

To realize these opportunities, the following steps are recommended:

- **Strengthen Interdisciplinary Collaboration:** Foster partnerships between AI experts, biologists, clinicians, and industry stakeholders to address complex biomedical challenges.

- **Enhance Data Integration and Sharing:** Develop standardized protocols for data collection, sharing, and integration to support AI applications. Ensure data privacy and security are prioritized.
- **Invest in Advanced Infrastructure:** Expand computational resources, including high-performance computing and cloud infrastructure, to support large-scale AI research and applications.
- **Support Education and Training:** Develop interdisciplinary education and training programs to build a skilled workforce capable of integrating AI into biomedical research and healthcare.
- **Promote Public-Private Partnerships:** Encourage collaboration between public research institutions and private industry to accelerate the translation of AI innovations into practical applications.

By implementing these strategies, Singapore can harness the full potential of AI to drive transformative advances in biomedical sciences, ultimately improving health outcomes and contributing to global scientific progress.

## REFERENCES

- 1 Singh et al. (2023) *Front. Drug Discov.* 3:1201419. doi: 10.3389/fddsv.2023.1201419
- 2 Sun, D., Gao, W., Hu, H., & Zhou, S. (2022). Why 90% of clinical drug development fails and how to improve it? *Acta pharmaceutica Sinica. B*, 12(7), 3049–3062. <https://doi.org/10.1016/j.apsb.2022.02.002>
- 3 The Unbearable Cost of Drug Development: Deloitte Report Shows 15% Jump in R&D to \$2.3 Billion ([genengnews.com](http://genengnews.com))
- 4 Sun et al. (2022). *Acta Pharm Sin B*. 12(7): 3049–3062. doi: 10.1016/j.apsb.2022.02.002
- 5 The Business Times. (2023). Local life sciences sector faces gap in late stage capabilities
- 6 Liu A., (2023). Fierce Biotech. Special Report. 20 years in, Singapore still searches for its biotech success story.
- 7 Jumper, J., Evans, R., Pritzel, A. et al. Highly accurate protein structure prediction with AlphaFold. *Nature* 596, 583–589 (2021). <https://doi.org/10.1038/s41586-021-03819-2>
- 8 Abramson, J., Adler, J., Dunger, J. et al. Accurate structure prediction of biomolecular interactions with AlphaFold 3. *Nature* 630, 493–500 (2024). <https://doi.org/10.1038/s41586-024-07487-w>
- 9 <https://www.biopharmatrend.com/ai-drug-discovery-pipeline/>
- 10 Rodriguez, A. et al. (2023). Unlocking the Potential of AI in Drug Discovery. BCG report commissioned by Wellcome Trust
- 11 Carreras-Puigvert, J., & Spjuth, O. (2024). Artificial intelligence for high content imaging in drug discovery. *Current opinion in structural biology*, 87, 102842. Advance online publication. <https://doi.org/10.1016/j.sbi.2024.102842>
- 12 <https://endpts.com/first-ai-designed-drugs-fall-short-in-the-clinic-following-years-of-hype/>
- 13 <https://www.biopharmatrend.com/ai-drug-discovery-pipeline/>
- 14 Visan, A. I., & Negut, I. (2024). Integrating Artificial Intelligence for Drug Discovery in the Context of Revolutionizing Drug Delivery. *Life (Basel, Switzerland)*, 14(2), 233. <https://doi.org/10.3390/life14020233>
- 15 <https://www.genengnews.com/topics/drug-discovery/ai-in-drug-discovery-trust-but-verify/>
- 16 Pun, F. W., Ozerov, I. V., & Zhavoronkov, A. (2023). AI-powered therapeutic target discovery. *Trends in pharmacological sciences*, 44(9), 561–572. <https://doi.org/10.1016/j.tips.2023.06.010>

## APPENDIX V. HEALTHCARE AND IMAGING

### AUTHORS:

Dr. Rosa Qi Yue So  
(I2R, A\*STAR)

Assoc. Prof Daniel Ting Shu Wei  
(Singapore National Eye Centre,  
Singhealth Duke-NUS)

Prof. Ngiam Kee Yuan  
(NUS, National University Hospital,  
National University Cancer Institute)

Assoc. Prof Tan Cher Heng  
(NTU, Tan Toeh Seng Hospital)

Dr. Rick Goh Siow Mong  
(IHPC, A\*STAR)

## Executive Summary

AI has emerged as a transformative force in healthcare, aiming to optimize care delivery amidst increasing demands from aging populations and workforce shortages. By optimizing clinical workflows, AI automates routine tasks, allowing healthcare professionals to focus on complex decision-making and

patient care. New multimodal AI tools also hold the promise of a shift towards personalized, community, and home-based care, with the expectation of extended healthspan through early detection, personalized treatments, and chronic disease management.

## Introduction

In recent years, the integration of Artificial Intelligence (AI) into healthcare has emerged as a transformative force, addressing critical needs and augmenting the capabilities of traditional healthcare systems. As the demands on healthcare providers intensify due to an ageing population and decreased workforce, the urgency to optimize care delivery becomes paramount. AI offers multifaceted solutions that not only alleviate the strain on the existing healthcare workforce but also facilitate a shift from hospital-centric models to personalized, community, and home-

based care. Moreover, AI-driven innovations promise to extend healthspan by enabling early detection, personalized treatment plans, and proactive management of chronic conditions. The following grand challenges aims to target ways that healthcare can be further transformed, enhancing efficiency, accessibility, and patient outcomes.

The availability of AI solutions in healthcare has seen remarkable growth (figure 1), offering promising advancements in diagnostics, personalized medicine, and operational

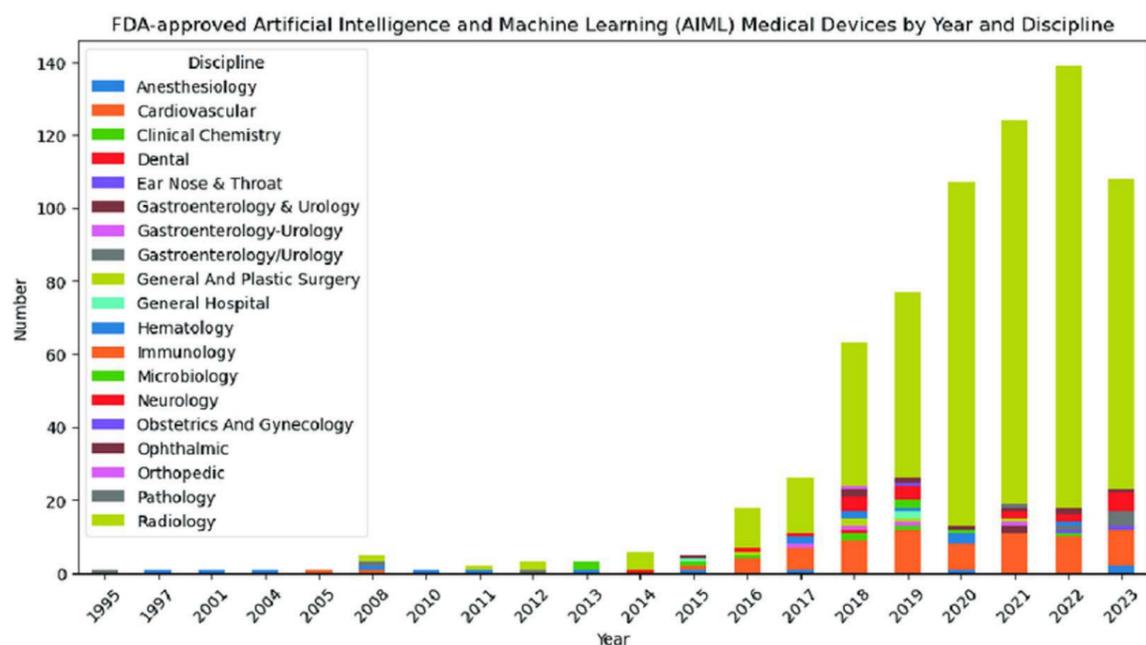


Figure 1: Increase in FDA approved AI tools for medicine, with radiology as the leading use case. Data taken from fda.gov.

efficiency. However, widespread adoption faces significant hurdles. Challenges such as data fragmentation, interoperability issues among different systems, and the need for rigorous validation of AI algorithms in clinical settings create barriers to seamless integration. Addressing these gaps is crucial to realizing the full potential of AI in healthcare, ensuring that innovations translate effectively into reduction in healthcare costs and improved outcomes for patients globally.

### AVAILABILITY OF QUALITY DATA

In the realm of AI for healthcare, one of the most pressing challenges lies in the availability and integration of healthcare data. Currently, valuable datasets are often fragmented across various healthcare providers, research institutions, government agencies and industry, existing in isolated silos that hinder comprehensive analysis. Bridging these silos is essential for enabling robust research and development in healthcare AI, as it requires

seamless access to diverse data sources ranging from hospital records to cohort studies and lifestyle information.

Incorporating patient-generated data such as wearable device metrics, survey inputs and lifestyle photos (e.g. diet) adds another layer of complexity and richness to analyses. However, a significant hurdle is the inconsistency in data quality across these sources, which complicates accurate and reliable insights.

The annotation of healthcare data demands specialized expertise and is a labour-intensive process, often lacking standardized protocols, further impeding progress in AI-driven healthcare solutions.

Another challenge is the lengthy process of accessing, cleaning, and de-identifying healthcare data to make it available to researchers. Most of the time, researchers only have access to a subset of the data, making it difficult to derive meaningful outcomes.

### DEPLOYMENT CHALLENGES

The deployment of AI healthcare solutions faces multifaceted challenges that span technological, ethical, and operational dimensions. Trustworthiness is a critical concern, as AI algorithms must not only deliver accurate diagnoses and treatment recommendations but also be explainable to healthcare providers who must incorporate new guidelines and diagnostic approaches. Model drift, where AI performance degrades over time due to changing data or environments, complicates long-term reliability and interoperability across different healthcare systems. The potential for errors in AI predictions introduces risks, highlighting the need for robust error management

### Background

The workshop commenced with an overview of the transformative potential of AI in healthcare, emphasizing its role in enhancing patient care through data-driven insights and advanced technologies. AI's ability to integrate into clinical workflows was highlighted, showcasing its potential to revolutionize diagnostics, patient management, and treatment personalization.

A recurring theme across the sessions was the integration of AI into healthcare systems to streamline processes and improve efficiency. Key areas included the use of AI for real-time data analysis, predictive modeling for patient outcomes, and support for clinical decision-making. The discussions underscored the need for robust data infrastructure and the importance of addressing biases in AI algorithms to ensure equitable healthcare delivery.

The workshop identified several challenges in the implementation of AI, such as data silos, the need for large-scale data sets, and the complexities of integrating AI into existing healthcare systems. However, these challenges were also seen as opportunities for innovation. The importance of developing AI tools that are transparent, trustworthy, and capable of handling real-world clinical data was emphasized.

strategies and clear delineation of responsibilities between AI systems and human healthcare providers.

Integrating AI tools into clinical workflows effectively may impose additional workload burdens in the short term on healthcare professionals. For example, deploying a model that forecasts patient deterioration requires additional manpower to respond whenever the risk exceeds a certain threshold. Such additional responsibility on the clinical staff can be burdensome, especially when a certain level of false positive is to be expected. Ensuring competency in AI among healthcare workers is crucial to foster awareness and promote buy-in, facilitating the responsible integration of AI technologies into healthcare delivery for improved patient outcomes.

Specific applications of AI were discussed, including its use in mental health, chronic disease management, and radiology. AI's role in monitoring patient conditions, providing personalized coaching, and supporting mental wellness through digital platforms was highlighted. In radiology, the focus was on improving diagnostic accuracy and addressing the challenges of AI integration into radiological workflows.

Ethical considerations were a significant focus, particularly the need for responsible AI governance and addressing biases in AI models. The discussion emphasized the importance of human oversight in AI applications and the ethical implications of AI decision-making in healthcare. Ensuring patient data security and developing frameworks for the ethical use of AI were also critical points.

Through the workshop on 'AI in Healthcare and imaging', several existing gaps were identified, such as **data fragmentation**, and deployment challenges such as **lack in reliability** and **interoperability**. Such difficulties slow development of next generation AI models and also hinder widespread adoption of AI tools.

## Grand Challenges

A few grand challenge statements from collective discussions at the workshop are as follows.

### GRAND CHALLENGE 1:

#### CAN MULTIMODAL FOUNDATIONAL AI FOR MEDICINE BE USED TO GENERATE EXPLAINABLE AND PERSONALIZED ADVICE FOR INCREASED COMMUNITY-BASED OR SELF-MANAGEMENT OF HEALTH?

##### OBJECTIVE

Development of multimodal medical foundational models is a fast-moving research area, with new state-of-the-art models emerging every few months<sup>1,2,3,4,5,6</sup>. These large foundation models are generic and versatile, with ability to interpret various data types (text, images, videos, etc) and provide tailored response to diverse clinical needs and queries. The most recent medical foundation models (Med-GEMINI<sup>4,5</sup>, BiomedGPT<sup>6</sup>) are built upon large volumes of clinical data from the electronic medical record (EMR) system, as well as biomedical literature databases (e.g. pubmed).

One current gap is that data outside of clinical systems, especially those related to population health and wellness, are not incorporated in these foundation models. This grand challenge aims to harness diverse data sources and advanced algorithms to generate personalized, explainable advice that enhances community-based or self-management of health. By developing foundation models based on comprehensive population health data, including clinical data, omics data and lifestyle data from in between clinic visits, the models can better predict individual health outcomes and tailor interventions accordingly (figure 2).

This concept of a digital twin, which simulates an individual's health profile based on real-time data inputs, facilitates proactive forecasting of health conditions and enables early intervention strategies. A personalized approach ensures that healthcare advice and interventions are not only relevant but also actionable and comprehensible to individuals, promoting active engagement in health management.

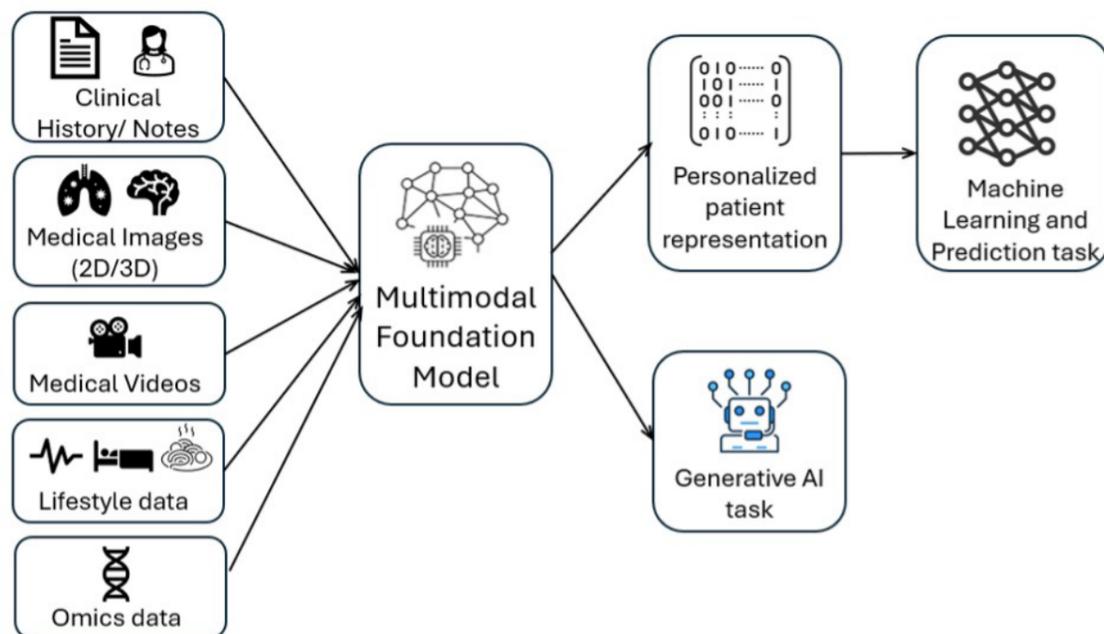


Figure 2: Building of a multimodal foundation model with clinical and population data, for downstream prediction or generative AI task.

### DATA REQUIREMENTS

Building or fine-tuning a general multimodal medical foundation model requires high-quality datasets that contain diverse data types, such as text, images, and other modalities like demographics, lifestyle and multiomics data. These datasets could come from public datasets and also proprietary datasets unique to Singapore. The following contains some examples of data that is present within the ecosystem.

International Public datasets:

- **MIMIC-III and MIMIC-IV:** Comprehensive datasets containing de-identified health records, including clinical notes, from ICU patients.
- **PubMed Central (PMC):** A free archive of biomedical and life sciences journal literature, providing a large collection of research articles.
- **Medical Literature Analysis and Retrieval System Online (MEDLINE):** Offers access to research papers, abstracts, and medical articles.
- **BioBank UK:** Contains various medical data types, including imaging, genetic data, and clinical reports. It's a large-scale, comprehensive dataset for studying the interaction between biological factors and disease.
- **The Cancer Genome Atlas (TCGA):** A comprehensive database containing genomic data linked to various types of cancer, which can be used in conjunction with medical imaging and text data.

Public and proprietary datasets within Singapore

- **National registries,** such as registry of kidney disease, stroke, cancer and cardiac disease
- **National diabetes database**
- **HealthierSG data**

- **Lifestyle data** including wearables data such as from HPB's step challenge, and Screen4Life program.

- Genomics data from **PRECISE SG100k**
- Data from **cohorts studies** such as SEED, HELIOS, ATTRACT, CADANCE, GUSTO

### AI METHOD REQUIREMENTS

Training a general medical foundation model involves using a variety of AI methods that enable the model to learn patterns from diverse data types such as text, medical images, sensor data, and genomic information. Below are some common AI methods used:

- Transfer Learning
- Self-Supervised Learning (SSL)
- Contrastive Learning
- Multi-Task Learning (MTL)
- Attention Mechanisms and Transformers
- Vision Transformers (ViTs)
- Generative Models (GANs and VAEs)
- Few-Shot and Zero-Shot Learning
- Reinforcement Learning (RL)

### COMPUTE REQUIREMENTS

The amount of compute requirement would depend on the size of the multimodal medical foundation model. A small to medium-scale model with 500M to 1B parameters would approximately need 100,000 to 500,000 GPU hours; 32-64 GB of GPU memory per GPU; 32 to 128 GPUs, depending on batch size and data parallelism. For example, training a BERT-based multimodal model with 1 billion parameters takes approximately 100-200K GPU hours when fine-tuned on large datasets.

A large-scale model with 10B+ parameters would require 1 to 10 million GPU hours; 80-128 GB of GPU memory per GPU; 512 to 1,024 GPUs for several weeks of training. For example, GPT-3 (175 billion parameters) required an estimated 3640 petaFLOPS-days of compute.

**GRAND CHALLENGE 2:  
HOW CAN AI BE USED TO OPTIMIZE ALL ASPECT OF CLINICAL WORKFLOW (E.G RADIOLOGY) TO REDUCE WORKLOAD BY 50%?**

This grand challenge is to reduce workload significantly by using AI to optimize all aspects of a clinical workflow like radiology. By leveraging AI's ability to automate routine tasks such as data analysis, scheduling, and administrative duties, clinicians can focus more on complex tasks such as decision-making and patient interaction.

Using radiology as an example, currently AI models for medical imaging are mainly built around automatic image segmentation and diagnostic inference, interpretation and reporting. However, having only an AI diagnostic model is not sufficient to realize the promise of increased productivity. Other "less challenging" AI functions such as machine scheduling optimization, post-inference

evaluation and results communication are often ignored but plays a large role in improving overall productivity (figure 3). Other important considerations include efficient storage of AI outputs (e.g. image heatmaps, image annotations).

Another key to integrating AI into clinical workflows requires productive human-AI interactions, which can only result by building trust in the AI system. This involves demonstrating the reliability and accuracy of AI algorithms through continuous AI monitoring processes and ensuring transparency in AI-driven insights and interpreted. Automatic updating of an AI model after model drift is still under research but an essential feature to ensure continual reliability of the AI output.

**OBJECTIVE**

Having an end-to-end AI solution that is sufficiently trusted would be key to achieving a significant (e.g. 50%) reduction in workload for clinical workflows such as radiology.

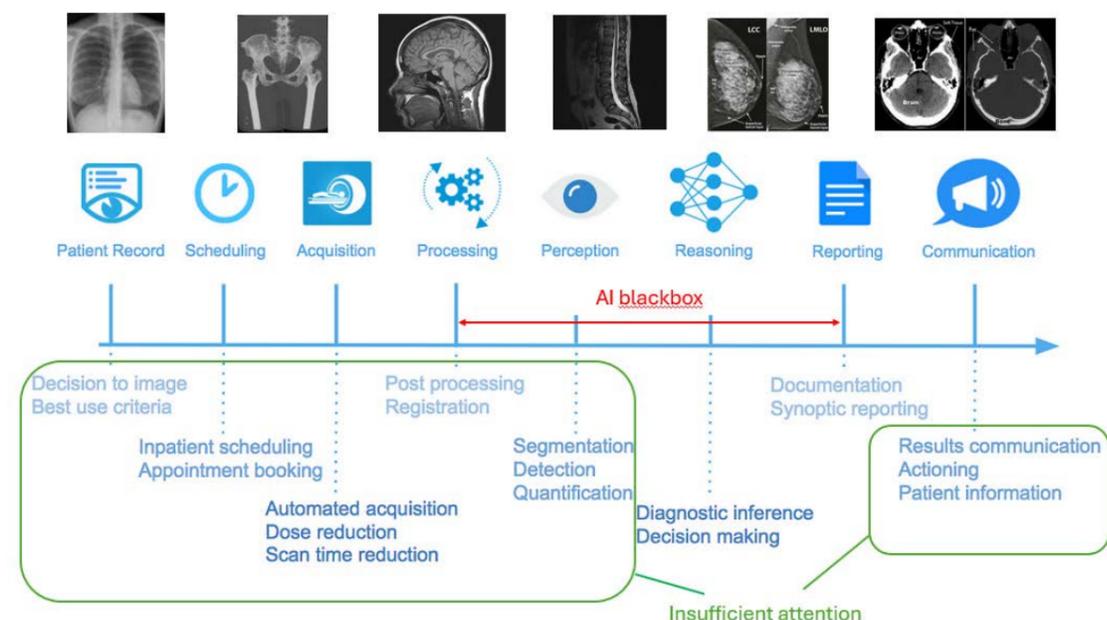


Figure 3: adapted from presentation by Dr. Steven Wong, highlighting areas of radiology workflow that is often neglected in AI research (green), as well as need for more transparent AI model that can be adjusted for model drift during implementation.

**DATA REQUIREMENTS**

There are several high-quality radiology datasets available for research, including modalities such as X-rays, CT scans, MRIs, and other medical imaging data, along with radiology reports. Some examples include:

- MIMIC-CXR
- CheXpert
- OpenI Chest X-ray Dataset
- IU X-ray Dataset (Indiana University)
- PEIR Digital Library
- COVID-19 Radiography Database (with Reports)
- RadQA Dataset

In order to develop models for other segments of radiology workflow, proprietary data, such as appointments and scheduling data from various imaging machines and EMR would need be extracted.

**AI METHOD REQUIREMENTS**

Frontier methods for report generation include advanced **transformer-based models** and **multi-modal learning** approaches. Transformers, such as GPT and BERT enhances natural language understanding and generation, allowing for more sophisticated and contextually accurate report writing. These models are integrated with vision transformers (ViTs) or CNN-based feature extractors to process and interpret radiological images. Multi-modal frameworks combine image analysis with text generation by aligning visual features with textual descriptions through attention mechanisms and cross-modal learning.

To build trustworthy medical radiology AI tools, several key methods can be employed: **Explainable AI (XAI)** helps make AI decisions transparent through tools like saliency maps, while **uncertainty quantification** and **model calibration** ensure predictions are reliable. Robustness to data shifts is enhanced through data augmentation and domain adaptation, and bias mitigation ensures fairness across diverse populations. Integrating human-in-the-loop systems allows radiologists to oversee

and validate AI predictions, enhancing trust. Additionally, continuous learning allows AI models to stay updated with real-world data, and multi-modal AI integrates contextual patient information for more accurate predictions.

**COMPUTE REQUIREMENTS**

Training of state-of-the-art models like transformers (e.g., GPT, BERT) or vision transformers (ViTs) requires significant computational power. For example, training a large transformer model for text generation may require multiple GPUs (e.g., 8-16 NVIDIA V100/A100) and significant memory (e.g., 128-256 GB) for a few days to several weeks. Once trained, the inference for a model may be done with a single GPU (e.g., NVIDIA A100) or TPU, with low latency and minimal memory needs (e.g., 16-32 GB RAM).

**GRAND CHALLENGE 3:  
HOW CAN AI BE USED TO SHORTEN R&D CYCLES FOR HEALTHCARE INNOVATION BY 50%?**

**OBJECTIVE**

The grand challenge of using AI to shorten the research and development (R&D) cycle for healthcare innovation involves leveraging its capabilities across multiple critical stages. For the development of AI-based tools, AI such as LLM can automate labour-intensive tasks such as data cleaning, quality checks, annotation, and feature engineering, accelerating the preparation and refinement of datasets essential for research.

Numerous studies have also shown that AI can be used to accelerate clinical trials, from optimising study design to determining patient recruitment and reducing drop-out rates<sup>7,8,9</sup>. It is possible to use AI in future to reduce validation and impact assessment cycles, leveraging retrospective data to perform first stage assessment of real-world effectiveness of healthcare interventions. These approaches enhance the agility of R&D cycles, allowing for iterative improvements based on continuous feedback and data-driven decision-making.

## DATA REQUIREMENTS

For accelerating clinical trials, some example data used could be:

- Patient Records and Electronic Health Records (EHRs), e.g. MIMIC-III
- Clinical Trial Registries (e.g. ClinicalTrials.gov, EU Clinical Trials Register)
- Pharmacological Databases that track adverse drug reactions and other safety-related data, which can inform trial designs and optimize safety monitoring.
- Health Economics and Outcomes Research (HEOR) Data for assessing the broader impact of clinical interventions.
- Patient-Reported Outcomes (PROs) as part of clinical trials.

## AI METHOD REQUIREMENTS

AI methods are used to optimize clinical trial design, patient recruitment, and reducing

drop-out rates. For study design, AI-driven simulation models, predictive analytics, and Bayesian optimization help refine trial parameters, such as sample size and endpoints, by predicting outcomes based on various design configurations. Patient recruitment is enhanced by matching patients to trials using electronic health records and predictive models. Optimization algorithms further improve recruitment strategies by balancing patient demographics and trial timelines. To address drop-out rates, risk stratification models can be used to identify high-risk patients and personalized engagement strategies to improve adherence.

## COMPUTE REQUIREMENTS

For the running simulations and predictive analytics, high-performance GPUs or a cluster of CPUs (multiple NVIDIA V100 or A100 GPUs) can be used. Training may take days to weeks depending on the compute power and size of data used.

## AI Methods and Data - Challenges and Opportunities

With increasing collection of healthcare data, increasingly large foundation models for AI Healthcare have been developed that requires significant computational resources for customization and fine-tuning. Efficient methods for fine-tuning large medical foundation models present significant opportunities to enhance healthcare outcomes by optimizing performance while managing computational resources. Techniques such as **transfer learning**, which leverages pre-trained models and adapts them to specific medical tasks, can accelerate model development and reduce the need for extensive data and computation. **Few-shot learning** and **active learning** further improve efficiency by allowing models to learn from smaller, more relevant datasets or dynamically selecting the most informative examples. Additionally, innovations in **model pruning** and **quantization** reduce the computational burden of deploying large models without sacrificing accuracy. These methods enable more accessible, cost-effective, and scalable AI solutions.

Second, data silos in healthcare arise when patient data is fragmented across various systems and institutions, hindering comprehensive analysis and integrated care. Privacy-preserving methods, such as **encryption and federated learning**, address these issues by allowing AI models to be trained across decentralized data sources. Secure multi-party computation (SMPC) allows multiple parties to jointly compute functions over their data without revealing the underlying data to each other. These techniques enable the effective utilization of fragmented data while maintaining stringent privacy and security standards.

Finally, AI healthcare tools requires a close collaboration between artificial intelligence and human expertise to achieve optimal clinical outcomes. **Human-in-the-loop (HITL)** approaches integrate AI systems with real-time input and oversight from healthcare professionals, ensuring that AI predictions are interpreted and validated by experts,

which helps in refining model accuracy and contextual relevance. **Interactive AI tools**, such as decision support systems and intelligent diagnostic assistants, augment human decision-making by providing actionable insights while allowing clinicians to guide the AI's learning process based on practical experience. In addition, **adaptive**

**learning systems** adjust AI models dynamically based on user feedback and evolving clinical scenarios, ensuring continuous improvement and alignment with real-world needs. This synergy between AI and human expertise enhances diagnostic precision, treatment personalization, and overall efficiency in healthcare delivery.

## Singapore's Role

In general, Singapore is well positioned to take the lead in addressing these grand challenges in AI for Healthcare due to **availability of Asian-centric clinical and population datasets**, strong **ecosystem collaboration and alignment towards innovation** and **reputation for good governance**.

Singapore is well positioned to lead in AI for healthcare research due to a few factors. First, Singapore's commitment to building a robust digital infrastructure and integrated healthcare system ensures availability of high-quality clinical data. With initiatives like the National Electronic Health Record system and development of the TRUST platform, quality healthcare datasets are made available for training and validating AI algorithms. Singapore also holds a number unique datasets collected from healthy citizens, such as various programs under the Health Promotion Board that has collected lifestyle data, and the PRECISE program that is one of Asia's leading genome database. The combination of such population datasets together with high quality real-world clinical data puts Singapore in a good position to lead the development of deep healthcare AI capabilities that is targeted at developing health strategies catered to the Asian phenotype.

Second, Singapore has a strong ecosystem of collaboration between healthcare providers (NUHS, Singhealth, NHG), academia (NUS/NTU/SUTD/SIT/A\*STAR) and industry. These institutions are at the forefront of AI research, with multidisciplinary teams working together to push the frontiers of AI healthcare and imaging research. With government led initiatives like HealthierSG and AgeWell, together with supportive grants and incentives, these research teams are well positioned to lead the development of AI solutions for the transformation of healthcare systems in Singapore, and beyond.

Lastly, Singapore has earned a strong reputation for its good governance and forward-thinking regulatory framework, which is conducive to the growth of AI in healthcare. Key initiatives include the development of clear guidelines for AI use in clinical settings and a commitment to data privacy. Such regulatory support has positioned Singapore to take a global lead in AI healthcare implementation, while ensuring that advancements are both effective and ethically sound.

## Conclusions

In conclusion, AI holds immense promise in revolutionizing healthcare by addressing current challenges and enhancing patient-centred outcomes. Despite hurdles such as data privacy and deployment complexities, AI-driven innovations are poised to reshape care delivery by enabling personalized, community-focused approaches and proactive health management. AI not only alleviates workforce burdens but also enhances the precision and

efficiency of medical decision-making. Moving forward, continued investment in AI research, coupled with robust regulatory frameworks, interdisciplinary collaboration and user education, will be essential to fully harness AI's potential in transforming healthcare into a more accessible and responsive system that meets the evolving needs of patients and providers.

## REFERENCES

- 1 Wornow, M., Xu, Y., Thapa, R. *et al.* The shaky foundations of large language models and foundation models for electronic health records. *npj Digit. Med.* **6**, 135 (2023).
- 2 Moor, M., Banerjee, O., Abad, Z.S.H. *et al.* Foundation models for generalist medical artificial intelligence. *Nature* **616**, 259–265 (2023).
- 3 Yuting He, Fuxiang Huang, Xinrui Jiang, *et al.* Foundation Model for Advancing Healthcare: Challenges, Opportunities and Future Directions. arXiv:2404.03264v1 [cs.CY] 04 Apr 2024
- 4 Khaled Saab, Tao Tu, Wei-Hung Weng, *et al.* Capabilities of Gemini Models in Medicine. arXiv:2404.18416 Apr 2024
- 5 DeepMind. Advancing Multimodal Medical Capabilities of Gemini. arXiv: 2405.03162 May 2024
- 6 Zhang, K., Zhou, R., Adhikarla, E. *et al.* A generalist vision–language foundation model for diverse biomedical tasks. *Nat Med* (2024).
- 7 Zhang, B., Zhang, L., Chen, Q. *et al.* Harnessing artificial intelligence to improve clinical trial design. *Commun Med* **3**, 191 (2023).
- 8 Matthew Hutson. How AI is being used to accelerate clinical trials. *Nature* **627**, S2-S5 (2024)
- 9 Chopra H, Annu, Shin DK, Munjal K, Priyanka, Dhama K, Emran TB. Revolutionizing clinical trials: the role of AI in accelerating medical breakthroughs. *Int J Surg* 109(12):4211-4220. Dec 2023.

## APPENDIX VI. GENOMICS



### AUTHORS:

Assoc. Prof Joanne Ngeow  
(NTU, National Cancer Centre Singapore)

Prof. Mile Sikic  
(Genome Institute of Singapore, A\*STAR)

## Executive Summary

This whitepaper, a result of the AI for genomics workshop, outlines potential grand challenges in the field. The challenges span from data generation for assembling complex genomes and developing DNA language models tailored for genome annotation, to data utilization that prioritizes privacy and leverages AI for the early detection of chronic diseases such as cardiovascular and cancer. While many challenges primarily concern human genomics, this paper also highlights the need

for assembling and annotating plant genomes. Additionally, apart from the foundation models focused on human genomes, we explore smaller, more specific models for the vast diversity of microbial genomes. In drafting this whitepaper, we have meticulously utilized existing data, such as that from the National Precision Medicine Programme, incorporated insights from leading experts in AI and genomics based in Singapore and existing computational resources.

## Introduction

Integrating Artificial Intelligence (AI) into genomics heralds a transformative era in biotechnology and medicine, redefining what is possible in healthcare and biological sciences. In Singapore, a nation renowned for its robust healthcare system and cutting-edge technology initiatives, the relevance of AI in genomics is particularly significant. As the volume of genomic data expands exponentially, traditional analytical methods have become increasingly insufficient to harness this information fully. With its powerful capabilities in machine learning, data processing, and pattern recognition, AI emerges as a pivotal technology to meet these challenges.

Singapore's unwavering commitment to becoming a smart nation positions it at the forefront of adopting advanced technologies in healthcare. AI applications in genomics, from enhancing genetic sequencing processes to revolutionising disease prediction, prevention, and treatment, are a testament to Singapore's leadership in this field. The integration of AI can dramatically improve the accuracy and speed of genomic analysis, enabling personalised medicine and advancing our understanding of genetic influences on disease phenotypes. Furthermore, AI-driven tools are crucial in interpreting the functional significance of genetic variations, aiding drug

discovery, and tailoring healthcare treatments to individual genetic profiles, showcasing Singapore's prowess in the healthcare and technology sectors. Besides human genomics, it is important to highlight the potential of genomics in understanding, utilising and protecting the unique biodiversity of South-eastern Asia.

## Background

The rapid advancements in genomic technologies and data science have brought to the forefront the immense potential of Artificial Intelligence (AI) in transforming genomics research and its application in medicine. Recognising this potential, we organised a multidisciplinary workshop, bringing together leading experts in biology, data science, and computer science IT experts from the government and policymakers. Each provided valuable perspectives that enriched the dialogue and contributed to a holistic understanding of the current landscape and future needs.

The workshop was structured to foster deep discussions and facilitate comprehensive analysis across several key areas: scalability of computational resources, data privacy and security, ethical considerations in genetic research, the development of AI models capable of handling complex genomic datasets and focus on fields beyond human genomics including plant genomics and metagenomics.

## Grand Challenges

As a result of the organised workshop, we identified four potential grand challenges:

- Developing graph neural network approach for de novo assembly of complex genomes
- Developing and fine-tuning DNA language models for automated genome annotation.

This white paper proposes several potential grand challenges that emerge from integrating AI into genomics, specifically focusing on their implications in Singapore. These include data privacy and security, the need for robust AI models that can accurately interpret vast genomic datasets, and the requirement for significant computational resources.

Through a series of presentations and breakout sessions, the workshop addressed the critical need for robust, scalable AI tools to manage and interpret the increasing influx of genomic data. Discussions also covered the technological and ethical frameworks necessary to integrate tools safely into clinical and research settings.

This white paper consolidates the insights and recommendations from the workshop, aiming to guide future research and funding priorities in AI for genomics. It is a foundational document to stimulate further exploration and innovation in this dynamic field, ensuring that developments align with scientific possibilities and societal needs. The identified grand challenges are set to catalyse significant breakthroughs in how we create, understand, protect, manipulate, and leverage genomic information through AI, promising to reshape the future of personalised medicine and biodiversity exploration and conservation.

- Developing models for preserving genomic data privacy during model training
- Developing predictive models of early detection of chronic diseases (e.g. cancer, diabetes, chronic kidney disease, cardiovascular diseases).

## Grand Challenges

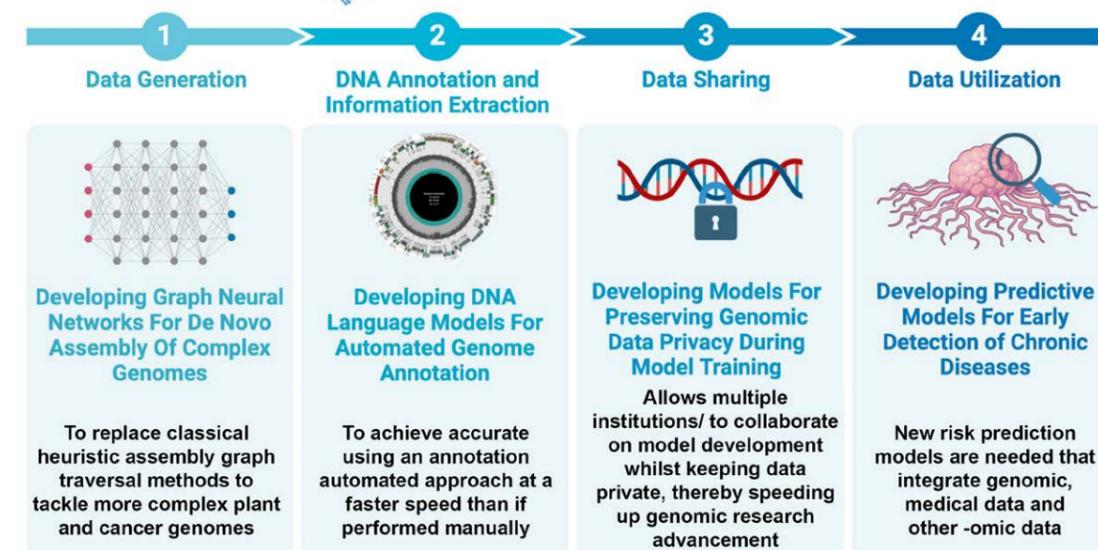


Figure 1: Summary of grand challenges in AI for genomics

The proposed challenges follow a stream beginning with the generation of genome assemblies, followed by deciphering the genomic language to extract critical information. This data, while preserving privacy, is then used to develop predictive models for chronic disease detection. Notably, this stream applies not only to human genomes but also to plant and microbial genomes. Moreover, all its components/challenges can be developed in parallel, starting with the existing data.

### GRAND CHALLENGE 1: DEVELOPING GRAPH NEURAL NETWORKS APPROACH FOR DE NOVO ASSEMBLY OF COMPLEX GENOMES

De novo genome assembly is a crucial tool in genomics, particularly useful for creating genome maps from scratch without relying on reference genomes. This technique<sup>1</sup> has primarily focused on human, vertebrate, and microbial genomes, but recent efforts aim to

tackle the more complex plant and cancer genomes where there are still a lot of space for the improvement.

Plant genomes are diverse, varying significantly in size and complexity across species, which impacts key agricultural traits like yield, disease resistance, and drought tolerance. De novo assembly allows researchers to explore this genetic diversity in high detail, aiding in breeding programs and helping to understand evolutionary relationships and adaptations among different plant species.

In cancer research, de novo assembly offers unique insights into the genomic alterations that distinguish various cancers. It is particularly valuable for identifying somatic mutations unique to cancer cells, especially in regions poorly represented by reference genomes. This method also enables the study of tumour subclones to track tumour evolution and progression, which is essential for crafting personalised treatment strategies that target specific genetic changes in tumours.

## OBJECTIVE

The current de novo assembly methods are focused on solving one genome in time by handcrafting parameters. The goal is to utilise newly telomere-to-telomere reconstructed genomes and new long, highly accurate reads to develop a graph neural networks framework, which will replace classical heuristic assembly graph traversal methods and show improvement in implementing it for plant and cancer genomes. The new genome assembler based on machine learning can resolve various genomes, from diploid to plant and cancer.

## DATA REQUIREMENTS

Many high-quality telomere to telomere human and vertebrate genomes and a few smaller plant genomes such as *A. thaliana*.

## AI METHOD REQUIREMENTS

Graph neural networks

## RESOURCE REQUIREMENTS

4 H100 cards for 36 months

---

### GRAND CHALLENGE 2:

## DEVELOPING AND FINE-TUNING DNA LANGUAGE MODELS FOR AUTOMATED GENOME ANNOTATION

Genome data annotation is a crucial process in genomics that marks specific genome regions to clarify their functions, essential for understanding the biological implications of genetic sequences. This process significantly contributes to genetics, molecular biology, and bioinformatics, enhancing our understanding of biological systems and fostering innovations in medicine and agriculture.

One possible way to improve annotation is by using DNA language models. Although there are already a few of them<sup>2,3,4,5,6</sup>, they still struggle to understand the complex language of DNA due to their size. One of the biggest challenges is what alphabet to use. Protein and RNA language models use single nucleotide resolution. However, this is not feasible for DNA due to input size constraints. Instead of single nucleotides, the models use non-overlapping kmers, which might lead to losing important information in case of insertion or

deletion of a single nucleotide, which might change all kmers in a sequence.

In addition to genomic data, there is an opportunity for multimodal models which incorporate other helpful information, including epigenomic information (DNA modifications) and transcriptomes.

The process includes two main types of annotations: (i) structural annotation, which identifies genomic elements such as genes, exons, introns, regulatory motifs, and non-coding RNA structures, and (ii) functional annotation, which predicts the functions of these genomic regions using methods like gene ontology and pathway mapping.

Despite its advances, genome annotation faces several challenges. The complexity of many organisms' genomes, filled with repetitive sequences, hinders accurate annotation. Rapidly evolving genomic technologies can make existing annotations obsolete, necessitating continuous updates. Moreover, integrating diverse data types to improve annotation accuracy remains a logistical and technical challenge.

The solution should be modular enough to support adaptation to various types of genomes, from short microbial genomes to long plant genomes.

## OBJECTIVE

Today, an experienced curator spends one day annotation a single genome at the National Precision Medicine program. The goal is to achieve an automated approach that might do this in an hour. A similar need exists to produce models for annotation of plant genomes and extracting insights from microbial genomes.

Expected outcomes:

- Multimodal DNA language model that can handle genomes of different scales from tens of thousands to tens of billions of nucleotides.
- A finetuned model for fast annotation of human genomes might also be used for other genome types. It should achieve human genome annotation in one hour.

## DATA REQUIREMENTS

Numerous reconstructed and annotated genomes exist. The haplotype-resolved human genome has 6 billion nucleotides. PRECISE has recently started a project to sequence and assemble 100 human genomes with high accuracy. The sequencing will include the detection of epigenomic information and RNA transcripts.

## AI METHOD REQUIREMENTS

DNA language models based on transforms, or other architectures such as Hyena<sup>6,7</sup> or Mamba<sup>8</sup>.

## RESOURCE REQUIREMENTS

Training and finetuning DNA language model – 18 months, 16 H100 cards.

---

### GRAND CHALLENGE 3:

## DEVELOPING MODELS FOR PRESERVING GENOMIC DATA PRIVACY DURING MODEL TRAINING

The National Precision Medicine programme has gathered a lot of useful information about the genomes of the local population. This information can further improve human biology's understanding and open new avenues in diagnosis and finding essential biomarkers. However, this data should be used while keeping privacy. Data can be used by researchers, but it might also provide a stream of revenue while providing access to pharma companies. One of promising approaches to access this data to train AI models while keeping data private is federated learning.

Federated learning is a transformative approach in machine learning that enables the development of predictive models using decentralised data while ensuring privacy. This is particularly crucial in genomics, where sensitive genetic data can be utilised without compromising individual privacy. In this method, raw data stays on local servers, and only model updates are shared for aggregation, significantly reducing privacy risks and regulatory challenges.

This technique allows multiple institutions, such as hospitals and biotech companies, to collaborate on model development while keeping their data private. Such collaboration not only improves model performance but also speeds up genomic research advancements. Federated learning is applicable in various genomic studies, including disease risk prediction and drug treatment personalisation, by training models on genetic markers without accessing raw data.

## OBJECTIVE

Develop federated learning approach together with data adaptation which will take into account genomics faces challenges, including handling data heterogeneity and coordinating updates across different IT systems. Ensure model robustness against biases and training errors. The expected outcome is for a federated learning model incorporated in the TRUST platform.

## DATA REQUIREMENTS

1000 genomes; gnomad; PRECISE

## AI METHOD REQUIREMENTS

Federated learning

## RESOURCE REQUIREMENTS

Storage space, 4 H100 GPU cards for 18 months.

---

### GRAND CHALLENGE 4:

## DEVELOPING PREDICTIVE MODELS OF EARLY DETECTION OF CHRONIC DISEASES (E.G. CANCER, DIABETES, CHRONIC KIDNEY DISEASE, CARDIOVASCULAR DISEASES)

The advancement of AI models holds great potential for identifying personalised genomic and non-genomic markers that inform specific aetiologies within cohorts, enabling personalised medicine. However, the integration of biodata with AI requires addressing challenges such as high dimensionality, explainability, and appropriate knowledge-representation modelling.

Traditional AI models, like deep convolutional neural networks (CNNs), have limitations as they assume a static representation of information, hindering their ability to capture temporal effects and establish causality. New models will be needed to integrate genomic and medical data as well as other -omic data for risk prediction models. One example is clonal hematopoiesis of indeterminate potential (CHIP), a common ageing-related phenomenon that serves as a biomarker for leukaemia risk as well as increased risk of cardiovascular disease. The presence of CHIP is further influenced by underlying genetic, lifestyle and environmental factors. Given the complexity of the problem, AI can potentially help us understand what causes CHIP and

also how the presence and absence of CHIP can influence/modify other genomic, lifestyle and environmental factors.

#### **OBJECTIVE**

Model for early detection that can be finetuned to at least two chronic diseases.

#### **DATA REQUIREMENTS**

PRECISE-SG100K; UK Biobank; ALL of Us Program

Disease databases such as TCGA; ICGC

#### **AI METHOD REQUIREMENTS**

Neural networks, Diffusion models

#### **RESOURCE REQUIREMENTS**

8 H100 cards for 18 months.

## AI Methods and Data - Challenges and Opportunities

### DATA

#### **DATA DIVERSITY AND QUALITY**

There is a critical need for well-labeled and high-quality multi-omic datasets covering everything from genomics to phenotypic outcomes. There are several types of possible omics data such as short and long DNA reads and RNA reads. It is important to note that DNA reads might hold information about (i) DNA epigenetic changes, ie. DNA modifications such as 5mC and 5hmC (ii) chromatin interactions, ie. Hi-C, Omni-C, micro-C and Pore-C, and (iii) chromatin accessibility, ie ATAC-seq. Now single cell DNA sequencing is still not mature enough, but this technology might be critical for more precise genomics information.

#### **DATA INTEGRATION**

It is important to take care of data integration of complex dataset which includes omics data, but also critical health information such as disease outcomes. During the data integration following precise bioinformatics pipeline should of high importance.

#### **DATA UTILISATION**

There are two types of data available for these challenges. Publicly available genomics data which is abundant, and data related to local smaller cohorts. From the other perspective there is a need to develop models on smaller size precisely labelled dataset and then scale to larger datasets and cohorts.

### AI

Challenges should focus on utilizing self-supervised learning for foundational model development, particularly where labelled data is insufficient. In cases with ample labelled data, a supervised approach can be used to create robust models. Self-supervised learning should be employed for feature extraction and dimensionality reduction from large datasets, helping to establish causal links from genetic data to patient outcomes. Additionally, while many challenges necessitate building models

from scratch, leveraging pre-trained models for fine-tuning with specific datasets is recommended wherever feasible. This approach can not only enhance the current models but also reveal their limitations, providing critical insights for developing new models.

Although we expect the most of proposed solutions for prepared challenges will be based on graph neural networks and large foundation models, we argue that solutions should be inspired with genomics using genomics specific inductive biases.

## Singapore's Role

In creating these challenges, we took in consideration that there are experts in Singapore with content and technical expertise in genomics and AI, including de novo genome assembly<sup>9,10,11,12</sup> for microbial<sup>13</sup>, plant<sup>14,15,16</sup>, human, and cancer genomes<sup>17,18,19,20,21,22</sup>, development of language models (ie. project SEALD- Southeast Asian Languages in One Network Data), language models in biology<sup>23</sup> and graph neural networks<sup>24,25,26,27</sup>, federated learning (Singapore's Federated Learning platform Synergos), work with various omics data<sup>28</sup>, and chronic diseases. Finally these challenges plan to intensively use human

genome data produced by the National Precision Medicine programme<sup>29,30,31</sup>. Soon there will be more than 100.000 genomes sequenced by short read technologies, but also high-quality genomes sequenced by long reads for which epigenomic DNA changes and RNA transcripts will be provided too.

Finally, there are huge investments in GPU cards at National Supercomputing Centre (NSCC) Singapore and local universities and ASTAR. The deep infrastructure will allow us to translate research data for direct application in the clinics.

## Conclusion

In conclusion, this whitepaper on AI for genomics has identified several grand challenges that need to be addressed to harness the full potential of AI in genomics. These include the development of graph neural network approaches for de novo assembly of complex genomes, the creation and fine-tuning of DNA language models for genome annotation, the formulation of strategies to preserve genomic data privacy during model training, and the construction

of predictive models for early detection of chronic diseases. This work is foundational, setting the stage for significant advancements in personalized medicine and biodiversity conservation by leveraging AI. It also highlights Singapore's strategic positioning and expertise in genomics and AI, aiming to catalyze breakthroughs in these critical areas. This whitepaper serves as a call to action for researchers, policymakers, and industry leaders to collaborate and invest in these challenges to drive innovation and improve health outcomes.

## REFERENCES

- Li, H. & Durbin, R. Genome assembly in the telomere-to-telomere era. *Nat. Rev. Genet.* (2024) doi:10.1038/s41576-024-00718-w.
- Dalla-Torre, H. *et al.* The Nucleotide Transformer: Building and evaluating robust foundation models for human genomics. *bioRxiv* (2023) doi:10.1101/2023.01.11.523679.
- Zhou, Z. *et al.* DNABERT-2: Efficient foundation model and benchmark for multi-species genome. *arXiv [q-bio.GN]* (2023).
- Sanabria, M., Hirsch, J., Joubert, P. M. & Poetsch, A. R. DNA language model GROVER learns sequence context in the human genome. *Nat. Mach. Intell.* (2024) doi:10.1038/s42256-024-00872-0.
- Ji, Y., Zhou, Z., Liu, H. & Davuluri, R. V. DNABERT: pre-trained Bidirectional Encoder Representations from Transformers model for DNA-language in genome. *Bioinformatics* **37**, 2112–2120 (2021).
- Nguyen, E. *et al.* HyenaDNA: Long-Range Genomic Sequence Modeling at Single Nucleotide Resolution. *arXiv [cs.LG]* (2023).
- Poli, M. *et al.* Hyena Hierarchy: Towards Larger Convolutional Language Models. *arXiv [cs.LG]* (2023).
- Gu, A. & Dao, T. Mamba: Linear-time sequence modeling with selective state spaces. *arXiv [cs.LG]* (2023).
- Vaser, R. & Šikić, M. Time- and memory-efficient genome assembly with Raven. *Nature Computational Science* **1**, 332–336 (2021).
- Stanojevic, D., Lin, D., Florez De Sessions, P. & Sikic, M. Telomere-to-telomere phased genome assembly using error-corrected Simplex nanopore reads. *bioRxiv* (2024) doi:10.1101/2024.05.18.594796.
- Nagarajan, N. & Pop, M. Sequence assembly demystified. *Nat. Rev. Genet.* **14**, 157–167 (2013).
- Vrček, L. *et al.* Geometric deep learning framework for de novo genome assembly. *bioRxiv* (2024) doi:10.1101/2024.03.11.584353.
- Bertrand, D. *et al.* Hybrid metagenomic assembly enables high-resolution analysis of resistance determinants and mobile elements in human microbiomes. *Nat. Biotechnol.* **37**, 937–944 (2019).
- Lim, A. H. *et al.* Genome assembly and chemogenomic profiling of National Flower of Singapore Papilionanthe Miss Joaquim 'Agnes' reveals metabolic pathways regulating floral traits. *Commun. Biol.* **5**, (2022).
- Salojärvi, J. *et al.* The genome and population genomics of allopolyploid *Coffea arabica* reveal the diversification history of modern coffee cultivars. *Nat. Genet.* **56**, 721–731 (2024).
- Teh, B. T. *et al.* The draft genome of tropical fruit durian (*Durio zibethinus*). *Nat. Genet.* **49**, 1633–1641 (2017).
- ICGC/TCGA Pan-Cancer Analysis of Whole Genomes Consortium. Pan-cancer analysis of whole genomes. *Nature* **578**, 82–93 (2020).
- Toh, M. R. *et al.* Global epidemiology and genetics of hepatocellular carcinoma. *Gastroenterology* **164**, 766–782 (2023).
- Sharma, A. *et al.* Onco-fetal reprogramming of endothelial cells drives immunosuppressive macrophages in hepatocellular carcinoma. *Cell* **183**, 377–394.e21 (2020).
- Ng, A. W. T. *et al.* Disentangling oncogenic amplicons in esophageal adenocarcinoma. *Nat. Commun.* **15**, 4074 (2024).
- Zang, Z. J. *et al.* Exome sequencing of gastric adenocarcinoma identifies recurrent somatic mutations in cell adhesion and chromatin remodeling genes. *Nat. Genet.* **44**, 570–574 (2012).
- Tam, W. L. & Weinberg, R. A. The epigenetics of epithelial-mesenchymal plasticity in cancer. *Nat. Med.* **19**, 1438–1449 (2013).
- Penić, R. J., Vlašić, T., Huber, R. G., Wan, Y. & Šikić, M. RiNALMo: General-Purpose RNA Language Models Can Generalize Well on Structure Prediction Tasks. *arXiv [q-bio.BM]* (2024).
- Defferrard, M., Bresson, X. & Vandergheynst, P. Convolutional neural networks on graphs with fast localized spectral filtering. *arXiv [cs.LG]* (2016).
- Bresson, X. & Laurent, T. Residual Gated Graph ConvNets. *arXiv [cs.LG]* (2017).
- Liu, J., Kawaguchi, K., Hooi, B., Wang, Y. & Xiao, X. EIGNN: Efficient Infinite-Depth Graph Neural Networks. *arXiv [cs.LG]* (2022).
- He, Y. & Hooi, B. UniGraph: Learning a cross-domain graph foundation model from natural language. *arXiv [cs.LG]* (2024).
- Goh, W. W. B., Wang, W. & Wong, L. Why batch effects matter in omics data, and how to avoid them. *Trends Biotechnol.* **35**, 498–507 (2017).
- Wu, D. *et al.* Large-scale whole-genome sequencing of three diverse Asian populations in Singapore. *Cell* **179**, 736–749.e15 (2019).
- Chan, S. H. *et al.* Analysis of clinically relevant variants from ancestrally diverse Asian genomes. *Nat. Commun.* **13**, 6694 (2022).
- Wong, E. *et al.* The Singapore National Precision Medicine strategy. *Nat. Genet.* **55**, 178–186 (2023).

## APPENDIX VII. DIGITAL PHENOTYPING

### AUTHORS:

Dr Creighton Heaukulani  
(MOH office for Healthcare Transformation - MOHT)

Prof Robert J.T. Morris  
(MOHT, NUS)

Assoc. Prof Jimmy Lee  
(Institute of Mental Health Singapore, NTU)

### Executive Summary

Machine learning is transforming the scientific discovery process in fields such as biomolecular modeling and materials science by accelerating exploratory processes with algorithmic solutions. The health sciences are similarly amenable to being transformed by this meta-strategy through the automated discovery and optimization (via machine learning methods) of digital bio- and behavioral-markers, the optimization of downstream interventions, and the model-based transfer of parts of these strategies to adjacent conditions and analogous care settings. This approach is enabled by large, longitudinal, and exploratory cohort datasets that measure data across a broad range of determinants underlying all aspects of human health. This approach promises to dramatically accelerate the health research process, which traditionally iterates between exploratory studies that aim to discover and refine health markers and their associated

interventions. Singapore is poised to lead this strategy, by creating some of the world's most comprehensive exploratory cohort datasets (such as the GUSTO study and the National Steps Challenge), a comprehensive EHR system covering much of the population, and world-leading digital sensing and intervention capabilities that have already been proven and are now ready to be expanded to a broad set of conditions. Proposed work under this white paper includes the expansion and integration of these cohort datasets with behavioral, genomic, social, and environmental determinants of health. Existing interventions will be transferred across all mental health and cardiovascular diseases and will branch out to important new challenges including obesity and healthy ageing. Our strategy collectively builds toward a foundation model for all healthcare in Singapore, which drives holistic interventional strategies aiming to maximize human health and potential.

## Introduction

### MACHINE LEARNING IN HEALTHCARE

Data and technology have transformed healthcare (Stoumpou et al., 2023). Longitudinal records provide insights into disease progression, while omics data and precision medicine enable a comprehensive understanding of individual patients for customised treatments. Mobile devices further revolutionise healthcare data collection and delivery of treatment. Machine learning methods are used to analyse these data sources, providing clinicians with actionable insights for diagnosis and treatment, which contribute to improved healthcare outcomes and reduced costs.

### ACCELERATING THE HEALTHCARE RESEARCH PROCESS

Healthcare science research has a long tradition and maintains a meticulously high standard. As such, it is costly and time-consuming, requiring clinicians and researchers to formulate hypotheses for

new biomarkers and interventions based on existing theory, followed by validation and refinement through exploratory studies. This process is iterative and its expansion to new conditions or populations occurs through a sequential and speculative process.

This scientific process is being transformed by the application of machine learning to large and exploratory research cohort datasets, which replaces the need to conduct exploratory research trials in the pursuit of discovering biomarkers and interventions. Instead, machine learning algorithms automatically discover and engineer biomarkers and their associated interventions from both cross-sectional and longitudinal datasets that measure many determinants of health. See Figure 1. Research may therefore advance more quickly to a validation trial, reducing timelines from decades to years. Moreover, addressing adjacent conditions or populations can be parallelized. The time from hypothesis to validation of a biomarker or intervention can be measured and reduced.

### Digital sensing and intervention eases the research bottleneck in the healthcare impact life-cycle

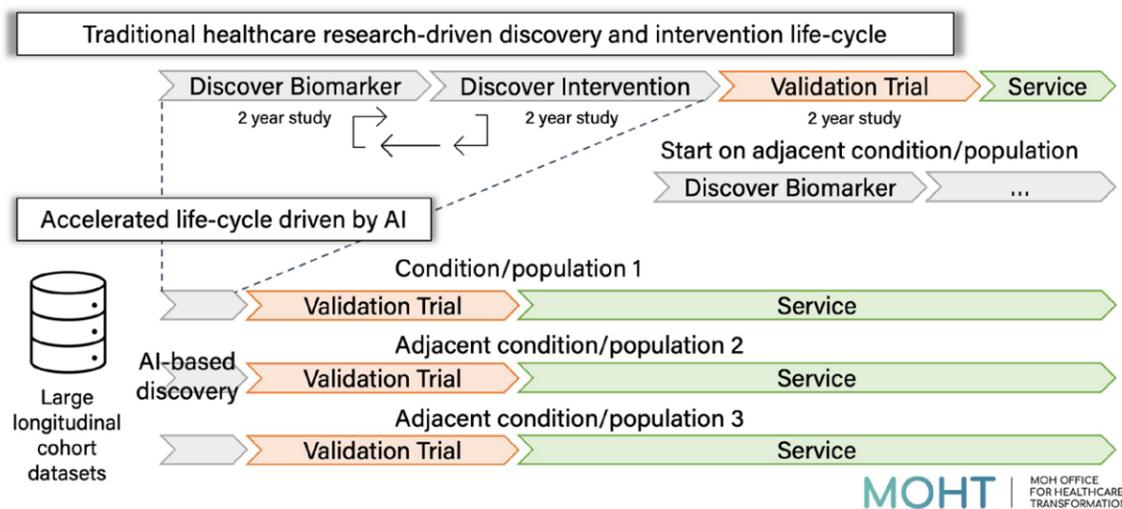


Figure 1: AI-based discovery from large longitudinal cohort datasets accelerates the research process by eliminating iterative and sequential exploratory research studies and parallelizes workstreams across adjacent conditions.

### MEASURING ALL DETERMINANTS OF HEALTH

This strategy requires large exploratory cohort datasets measuring a broad range of determinants of health that transcend diagnoses. While electronic health records (EHR) and genomic profiles have formed the bulk of the data historically studied, it is only recently that behavioral determinants, which form the largest share of healthcare determinants at 40% (Schroeder, 2007), have been comprehensively measured. Digital

phenotyping, which is the characterization of an individual's behaviors through the continuous analysis of data from digital devices, has emerged as the highest fidelity strategy to measure these behavioral determinants of health (Torous et al., 2016). MOHT and others in Singapore have expanded and integrated digital phenotyping to include other forms of telehealth monitoring and longitudinal datasets (including health records), which has created an expansive digital sensing and intervention strategy across the lifespan and population.

## Background

In May 2024, we conducted a workshop focusing on cohort datasets, digital sensing and interventions (with particular focuses on cardiometabolic disease management, behavioral factors, and digital mental health), and AI methods that accelerate research in healthcare. After an opening by Prof Robert Morris (Chief Technology Strategist, MOHT), who set out our challenge statements, leaders across the Singapore healthcare ecosystem were invited to give keynote talks, including A/Prof Daniel Fung (CEO, IMH), who showcased how mental healthcare in Singapore is a leading use-case for digital health interventions and highlighted the challenge of translating research into improved healthcare outcomes and reduced costs. Digital mental health innovations are best exemplified by the HOPES digital phenotyping project, which was presented by A/Prof Jimmy Lee (Psychiatry Senior Consultant, IMH). A/Prof Andy Khong (Head of Psychology, NTU) presented work on behavior change in the WellFeet study, which is a predecessor to the CADENCE digital health platform for behavior change in the

management of cardiovascular disease, and A/Prof Lian Leng Low (Duke-NUS) presented on the EMPOWER platform and study for behavior change in diabetes management. Prof Michael Meany (NUS and A\*STAR) presented findings from the GUSTO and S-PRESTO trials, now likely the largest birth and maternal cohorts in the world. Dr Praveen Deorani (MOHT) presented results on the PTEC and AMI-HOPE programs on the telehealth management of hypertension and post-discharge AMI patients. Finally, Dr Creighton Heaukulani (MOHT) presented some ideas on the application of transfer learning and other machine learning methodologies that are poised to accelerate healthcare research.

The workshop also conducted roundtable discussions on a variety of topics including the discovery of digital biomarkers, feature engineering for digital biomarkers and corresponding interventions, expansion of strategy to new patients and conditions, translational research and real-world evidence, and ethics and safety.

## Grand Challenges

### GRAND CHALLENGE: TOWARD A FOUNDATION MODEL FOR HEALTHCARE IN SINGAPORE

**Motivation 1:** Healthcare outcomes are determined by a range of determinants that are not currently holistically integrated in models for biomarkers and interventions in relation to disease and wellness.

**Motivation 2:** The healthcare research lifecycle is an iterative and speculative process searching for an optimizing biomarkers and interventions that spans many years.

#### OBJECTIVE

As a result of the organised workshop, we identified a potential grand challenge as building a foundation model for healthcare in Singapore. This white paper proposes a unifying strategy to improve healthcare outcomes and reduce costs by accelerating healthcare research through machine learning applied to foundation models for Singapore's health. Our strategy is to expand existing cohort datasets to include a broader range of health determinants, including behavioral measures from digital phenotyping. We will additionally launch new initiatives to generate important cohort datasets including the planned Brightline study at NTU, which will initially target 500 students and subsequently expand and will measure a broad range of behavioural determinants of importance to population health. We plan to join multiple longitudinal cohorts, both serially and

concurrently, to cover the entire lifespan, as well as to expand the conditions being investigated. We will integrate further data on genomics (from Singapore's PRECISE SG100K genotyping program), environmental determinants (from the National Environmental Agency), and social determinants (from the Silver Generation Office), among others. Combined with LLMs (Large Language Models) mined on the medical literature, and eventually large international longitudinal cohorts such as the Epic COSMOS dataset, we will arrive at a Foundation Model for Singapore's lifestyle, phenotype and optimal healthcare interventions that maximize health and human potential across the population. See Figure 2.

Our key performance indicators are the number of transformational care initiatives developed and deployed in clinical and community health settings; the number of patients and care providers on service; the improvement of health and wellness outcomes; and the reduction in the costs of care. Resulting acceleration of the healthcare research process will be measured in the time from project launch to deployment of the intervention, which may be compared to traditional healthcare research and development timelines. Under this whitepaper and the workplan therein, we expect three to five new disease programs to be launched and five to six expansions of existing programs, which all will lead to measurable improvements in health outcomes and reduced costs.

### Digital Phenotyping Collects all Human Experiences and Drives Models which Predict Health Outcomes

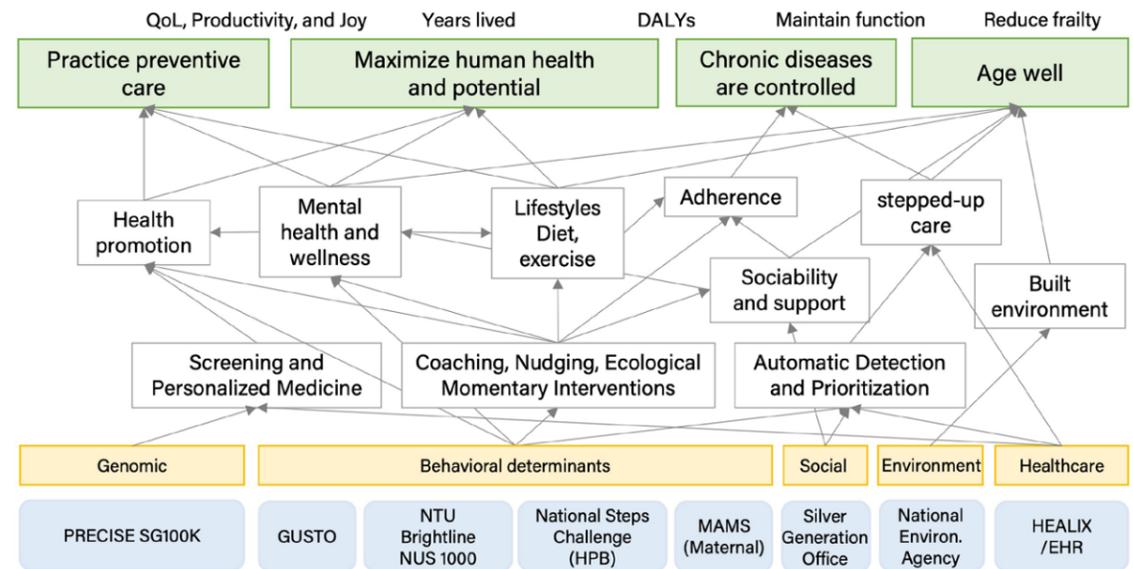


Figure 2: By assembling large cohort datasets that integrate measures of all determinants of health in Singapore, machine learning models may drive the discovery of interventional strategies that impact all of human health.

#### KEY DIFFERENTIATOR – INCLUDING DATA ON PATIENT BEHAVIORS AND HEALTH

Behavioral interventions are a key part of our strategy, as it is well known that patient behaviours control about 40% of outcomes, measured in premature deaths, yet have only recently been continuously and systematically measured in the healthcare context (via digital phenotyping (Torous et al., 2016)), as compared to genomic determinants (30% of outcomes) and healthcare determinants (10% of outcomes) (Schroeder, 2007). They directly impact health promotion and adherence in the management of chronic conditions. The EMPOWER Study and App at SingHealth and NUS uses AI to create personalised nudges for health coaching in chronic disease management (Kwan et al., 2022). The platform also leverages a localized GPT model trained on medical data. EMPOWER has seen statistically significant improvements in HbA1c levels after 3 months from over 1,000 diabetes patients.

Also in diabetes management, the WellFeet Program developed by NTU aims to reduce diabetic foot ulcers by offering a coaching app for patients, caregivers, and healthcare. It includes an AI-powered chatbot trained on extensive data and specialist articles. A feasibility study showed significant improvements in health literacy and self-care behaviours among users. WellFeet has paved the way for the development of CADENCE, a coaching platform targeting cardiovascular disease prevention and control through a holistic review of health status and a creative approach to long-term user engagement.

Both the EMPOWER team at SingHealth and NUS and the WellFeet and CADENCE teams at NTU and NUS will be key partners.

### KEY DIFFERENTIATOR - THE FOCUS ON INTERVENTIONS AND TRANSLATION

A key feature of our approach is that it is applied simultaneously to the discovery and optimization of both biomarkers and interventions. For example, in the HOPES project, the discovery of important biomarkers was a by-product of the machine learning algorithm developed to predict relapses, which drives our intervention of AI-based stepped-up care (Heaukulani et al., 2024). Similar kinds of interventions have been shown in PTEC and AMI-HOPE to be safe, as they provide an extra layer of care and do not inadvertently de-escalate care. High specificities in detection helps to ensure that there are minimal and acceptable increases in interventional costs, and in fact may result in fewer clinic visits. Being driven by interventions ensures that our strategy has translational impact beyond the confines of research and has high probability of resulting in cost-effective care improvements. MOHT operates in a manner where new techniques are always concurrently developed and deployed with clinical staff, with an emphasis on cost-effective care transformation.

### DATA REQUIREMENTS

Access to the following existing datasets will be sought:

- electronic medical records from HEALIX/TRUST;
- SG100K genomics dataset from PRECISE;
- National Steps Challenge from HPB;
- GUSTO and S-PRESTO cohorts from NUS and A\*STAR;
- social indicator data from the Silver Generation Office;
- environmental data from the National Environment Agency.

Additionally, new datasets will be created, including further research studies in the HOPES mental health programme and interventional phases of the Brightline digital phenotyping program for student wellness at NTU. Additionally, multiple new datasets from studies will be generated from cohorts amongst the following settings: ageing, familial hypercholesterolemia, obesity, mental disorders, and/or cardiometabolic diseases.

### AI METHOD REQUIREMENTS

Our strategy is to “stitch together” different cohort datasets in order to cover the entire human life span, but different cohort datasets will not commonly contain the same patients. A semi-supervised learning approach must therefore be taken, in which classes of patients (described by distinct collections of model parameters and which are often referred to as patient archetypes) are automatically inferred from the data. State-space models are one such approach (Alaa & van der Schaar, 2019; Krishnan et al., 2015). While individual patients are rarely shared between cohorts, these archetypes (if inferred jointly across all cohort datasets) will always be shared as factors determining all human health and disease. These approaches are semi-supervised because the automatic inference of archetypes occurs in an unsupervised learning fashion, and the training of model parameters to optimally predict health or disease measurements is regressed onto the inferred archetypes in the traditional way. While many machine learning models are inherently learning latent representations of individual patients that are regularized amongst one another, these semi-supervised learning methods often make the model structure for different archetypes explicit and sometimes even aim to interpret, or identify, those archetypes – a task referred to as inference in the medical literature (which is *not* how the machine learning community, and we in this article, use the term).

Generative models, as recently popularized by large language models (LLMs) like OpenAI’s GPT, and *federated learning*, where model training is performed on multiple datasets that are kept decentralized during training, have both been considered with hope and trepidation for healthcare applications (Rieke et al., 2020) and will both play key roles in our strategy. This will be exemplified by an upcoming project to develop a clinical decision support (CDS) system to assist GPs within the Healthier SG initiative, in which every Singapore resident will enrol with a GP that will holistically manage their long-term health. GPs must currently integrate a patient’s historical data (if it is even available to the GP) with their current presentation, national care guidelines (which may be difficult to navigate in the event of comorbidities), and the latest findings from international sources in order to make an optimal recommendation to a patient, all within a very limited time frame and under the pressure of an increasing patient load. We aim to develop machine learning-based tools derived from our foundation model to allow GPs to tackle this challenge and focus on more critical elements of care. This GP CDS Project will at first provide the care team tools that ingest patient health records and make care recommendations (for example, in preventive care such as vaccinations and in the management of chronic diseases) based on Ministry of Health guidelines. We will then research and develop methods to integrate LLM based recommendations that incorporate the medical literature and international cohort datasets. Because our methods will learn from data from different private GP practices and other national patient data sources, we will need to research and develop federated learning methods to train our models without ever agglomerating datasets.

In the traditional health research process, branching into new conditions usually begins its development following success on an existing strategy for an adjacent population. (See Figure 1.) Instead of launching a new sequence of exploratory trials, our accelerated strategy utilizes concurrent observation and machine learning to discover features within the cohort datasets that are transdiagnostic and the new condition can be approached in parallel, provided the dataset measures a broad enough range of the determinants of health (for example we now derive over 200 patient features and clinical scales measuring many symptoms and functioning beyond psychosis in HOPES). *Transfer learning*, which is a well-established technique in machine learning in which knowledge inferred from a particular task is re-used to boost performance on a related task, will enhance this process. It should be noted that transfer learning will also be applied to interventions. For example, an LLM-based therapist’s assistant in mental health is currently being used by mental health professionals in MOHT’s Let’s Talk forum for youth mental wellness (Heaukulani et al., 2024), and techniques from this tool will be transferred to the development of the LLM-based features in the GP CDS Project described earlier.

### COMPUTE REQUIREMENTS

Two million GPU hours on an A100 machine will be required, along with 1pb of storage for 2 years. This will be required for the training of foundation models on the integrated datasets.

## Singapore's Role

### SINGAPORE'S WORLD-LEADING RESEARCH COHORT DATASETS

Singapore is well-positioned to lead this AI for Science strategy, because it already possesses some of the world's most comprehensive exploratory cohort datasets. For example, the Growing Up in Singapore Towards healthy Outcomes (GUSTO) (Soh et al., 2014) and Singapore Preconception Study of Long-Term Maternal and Child Outcomes (S-PRESTO) (the S-PRESTO Study Group et al., 2021) studies are now among the world's largest birth cohorts. Coordinated by the wide-coverage public healthcare system, Singapore's EHR is amongst the world's most comprehensive, especially for chronic disease care. The National Steps Challenge run by the Health Promotion Board (HPB) measures physical activity at the population level and is one of the largest population health projects worldwide (Yao et al., 2020). Our strategy is to expand these cohorts to include the broadest possible range of healthcare determinants and to integrate these data sources.

### SINGAPORE'S WORLD-LEADING DIGITAL SENSING AND INTERVENTION STRATEGIES

Singapore programs on digital sensing and intervention are world leading. In the HOPES Project from the Institute of Mental Health (IMH) and MOHT, the most comprehensive dataset of schizophrenia phenotypes in the world has been collected, now comprising over 220 million events from wrist devices and the smartphones (Abdul Rashid et al., 2021). Machine learning algorithms discovered significant behavioral markers from among over 200 features and the resulting predictive algorithms can anticipate psychosis relapses 21 days in advance with a sensitivity of 91% and a specificity of 95%. These algorithms drive the automated delivery of digital therapeutics and stepped-up care continuously and beyond the clinic based on determined need. HOPES is currently deployed as a clinical service at IMH for

patients with psychosis and mood disorders and is now expanding to further mental health conditions based on AI extrapolations from our large cohort dataset. (See section on transfer learning below.)

In MOHT's PTEC (Teo et al., 2023) and AMI-HOPE initiatives, telehealth devices are used to continually collect vital signs in the management of hypertension and in the post-discharge care of heart attack patients. Key biomarkers are derived by machine learning models from electronic health records and will soon be derived from this developing longitudinal cohort dataset. Interventions include behavioral and adherence coaching using care coordinators, nurses and pharmacists and automatically stepped-up care when needed. PTEC is now in all polyclinics across Singapore and achieves a 300% increase the number of controlled hypertension cases compared to standard care. The program is now expanding to diabetes management with the inclusion of at-home blood glucose monitors.

### KEY PARTNERSHIPS

The partnership between IMH and MOHT, which has developed world-leading competence in digital phenotyping and interventions over the past 5 years, will lead this initiative. Other key partners include the custodians of existing cohort datasets and those responsible for the health and wellbeing within those cohorts. This includes NUS and A\*STAR (the GUSTO and S-PRESTO studies), NTU (Brightline), HPB (National Steps Challenge), PRECISE (SG100K), HEALIX and TRUST (EHR, Genomics and other national data), the National Environment Agency, and the Silver Generation Office. MOH and the three health clusters will also be key partners, having responsibility for the healthcare services for residents in Singapore. Tech partners include Synapxe and commercial companies such as PheBe Health (the team behind the latter are the originators of one leading digital phenotyping platform), as well as the technology developed by collaborators such as the EMPOWER and CADENCE platforms, to mention just a few.

## Conclusion

Machine learning is already accelerating the healthcare research process, but this is only being realized separately, disease-by-disease, and rarely integrating comprehensive data across multiple health determinants. Our strategy is to take advantage of Singapore's world leading cohort studies and datasets, integrating them along with international data and new cohorts that we will collect, resulting in the development of a foundation model for all of human health in Singapore. Using this strategy of broad-spectrum data collection, data integration, and interventional discovery, we will address a large portion of Singapore's burden of disease, initially focusing on

cardiometabolic diseases, mental disorders, obesity, familial hypercholesterolemia, healthy ageing, and the maintenance of wellness. The research and development requirements of this initiative are immense, requiring the development of new methodology and tasteful application of generative models, semi-supervised learning, transfer learning, federated learning, and more. Despite the intense research agenda, the nature of MOHT and clinical partners ensures that research is translated and that outcomes are always measured in terms of improved healthcare and reduced costs amongst the national population as a result of deployed and scaled interventions.

## REFERENCES

- Abdul Rashid, N. A., Martanto, W., Yang, Z., Wang, X., Heaukulani, C., Vouk, N., Buddhika, T., Wei, Y., Verma, S., Tang, C., Morris, R. J. T., & Lee, J. (2021). Evaluating the utility of digital phenotyping to predict health outcomes in schizophrenia: Protocol for the HOPE-S observational study. *BMJ Open*, 11(10), e046552. <https://doi.org/10.1136/bmjopen-2020-046552>
- Alaa, A. M., & van der Schaar, M. (2019). Attentive State-Space Modeling of Disease Progression. *Advances in Neural Information Processing Systems* 32. NeurIPS 2019.
- Heaukulani, C., Phang, Y. S., Weng, J., Lee, J., & Morris, R. J. (2024). Deploying AI Methods for Mental Health in Singapore: From Mental Wellness to Serious Mental Health Conditions. Machine Learning for Cognitive and Mental Health Workshop (ML4CMH), AAAI 2024, Vancouver, BC, Canada.
- Krishnan, R. G., Shalit, U., & Sontag, D. (2015). *Deep Kalman Filters* (arXiv:1511.05121). arXiv. <http://arxiv.org/abs/1511.05121>
- Kwan, Y. H., Yoon, S., Tan, C. S., Tai, B. C., Tan, W. B., Phang, J. K., Tan, N. C., Tan, C. Y. L., Quah, Y. L., Koot, D., Teo, H. H., & Low, L. L. (2022). EMPOWERing Patients With Diabetes Using Profiling and Targeted Feedbacks Delivered Through Smartphone App and Wearable (EMPOWER): Protocol for a Randomized Controlled Trial on Effectiveness and Implementation. *Frontiers in Public Health*, 10, 805856. <https://doi.org/10.3389/fpubh.2022.805856>
- Rieke, N., Hancox, J., Li, W., Milletari, F., Roth, H. R., Albarqouni, S., Bakas, S., Galtier, M. N., Landman, B. A., Maier-Hein, K., Ourselin, S., Sheller, M., Summers, R. M., Trask, A., Xu, D., Baust, M., & Cardoso, M. J. (2020). The future of digital health with federated learning. *Npj Digital Medicine*, 3(1), 119. <https://doi.org/10.1038/s41746-020-00323-1>
- Schroeder, S. A. (2007). We Can Do Better—Improving the Health of the American People. *New England Journal of Medicine*, 357(12), 1221–1228. <https://doi.org/10.1056/NEJMs073350>
- Soh, S.-E., Chong, Y.-S., Kwek, K., Saw, S.-M., Meaney, M. J., Gluckman, P. D., Holbrook, J. D., Godfrey, K. M., & on behalf of the GUSTO Study Group. (2014). Insights from the Growing Up in Singapore Towards Healthy Outcomes (GUSTO) Cohort Study. *Annals of Nutrition and Metabolism*, 64(3–4), 218–225. <https://doi.org/10.1159/000365023>
- Stoumpos, A. I., Kitsios, F., & Talias, M. A. (2023). Digital Transformation in Healthcare: Technology Acceptance and Its Applications. *International Journal of Environmental Research and Public Health*, 20(4), 3407. <https://doi.org/10.3390/ijerph20043407>
- Teo, S. H., Chew, E. A. L., Ng, D. W. L., Tang, W. E., Koh, G. C. H., & Teo, V. H. Y. (2023). Implementation and use of technology-enabled blood pressure monitoring and teleconsultation in Singapore's primary care: A qualitative evaluation using the socio-technical systems approach. *BMC Primary Care*, 24(1), 71. <https://doi.org/10.1186/s12875-023-02014-8>
- the S-PRESTO Study Group, Loo, E. X. L., Soh, S.-E., Loy, S. L., Ng, S., Tint, M. T., Chan, S.-Y., Huang, J. Y., Yap, F., Tan, K. H., Chern, B. S. M., Tan, H. H., Meaney, M. J., Karnani, N., Godfrey, K. M., Lee, Y. S., Chan, J. K. Y., Gluckman, P. D., Chong, Y.-S., ... Cheng, Z. R. (2021). Cohort profile: Singapore Preconception Study of Long-Term Maternal and Child Outcomes (S-PRESTO). *European Journal of Epidemiology*, 36(1), 129–142. <https://doi.org/10.1007/s10654-020-00697-2>
- Torous, J., Kiang, M. V., Lorme, J., & Onnela, J.-P. (2016). New Tools for New Research in Psychiatry: A Scalable and Customizable Platform to Empower Data Driven Smartphone Research. *JMIR Mental Health*, 3(2), e16. <https://doi.org/10.2196/mental.5165>
- Yao, J., Tan, C. S., Chen, C., Tan, J., Lim, N., & Müller-Riemenschneider, F. (2020). Bright spots, physical activity investments that work: National Steps Challenge, Singapore: a nationwide mHealth physical activity programme. *British Journal of Sports Medicine*, 54(17), 1047–1048. <https://doi.org/10.1136/bjsports-2019-101662>

# APPENDIX VIII. MATERIALS/CHEMISTRY



## AUTHORS:

Assoc Prof Kedar Hippalgaonkar  
(NTU/IMRE, A\*STAR)

Prof Kostya Novoselov  
(Institute for Functional Intelligent Materials -  
IFIM, NUS)

Prof Andrey Ustyuzhanin  
(NUS IFIM, Constructor University)

Dr Beatrice Soh  
(IMRE, A\*STAR)

## Executive Summary

The whitepaper underscores a transformative intersection between artificial intelligence (AI) and the fields of materials science and chemistry, highlighting a pioneering frontier ripe for breakthroughs in the design, development, and deployment of novel materials. It delineates the critical role of AI in bridging the theoretical and experimental dichotomy, thereby facilitating the development of materials with unprecedented functionalities. The whitepaper discusses the significant challenges that stand in the way, including data generation and accessibility, multiscale phenomena modelling, materials inverse design, accelerated computational

modelling, and robust experimental validation. It emphasizes Singapore's strategic position as a global R&D hub blessed with robust governmental support, advanced computational resources, and a culture of international collaboration, which collectively provides a conducive ecosystem for overcoming these challenges. Through a multidisciplinary approach that integrates AI with theoretical insights and experimental validation, the document envisions a future where materials science and chemistry are significantly advanced, leading to smarter, sustainable, and precise materials.

## Introduction

The overarching goal in materials science and chemistry is to rationally design materials with specified functions. One of the significant challenges in achieving this is to understand the relationship between the atomic and molecular composition of matter and its macroscopic properties and to leverage this knowledge to achieve desired functionality through chemical synthesis.

Much of materials science cannot be conducted in-situ, hence there is a critical

need for an integrated approach that combines both theory and experiments. Currently, a significant gap exists between theoretical predictions and experimental results in materials science and chemistry. Although digital twins have shown great promise in other fields, such as modelling the urban heat effect through the Digital Urban Climate Twin by Cooling Singapore 2.0, similar models are lacking for material systems. This gap presents a unique opportunity for AI to make a transformative impact.

In the heart of transforming global industries and technological capabilities, the intersection of artificial intelligence (AI) with materials science and chemistry represents a pioneering frontier poised for breakthrough advancements. Bridging the profound complexities of materials at atomic levels to their macroscopic counterparts, this domain encapsulates the ambition to develop materials with unprecedented functionalities and properties. The challenges encountered in this endeavour – encompassing data generation and accessibility, multiscale phenomena modelling, materials inverse design, accelerated computational modelling, and robust experimental validation – define a roadmap for revolutionary discoveries and innovations.

As we venture into this exploration, the role of AI emerges not merely as a tool but as an integral collaborator in bridging the longstanding gap between theoretical models and experimental outcomes. This collaboration promises a transformation in how materials are designed, developed, and deployed across various sectors. The current state, while ripe with potential, highlights specific areas for growth and development that can usher in a new era of material science and chemistry – one that is faster, more efficient, and predictively accurate than ever before.

## Background

With the rapid progress in AI and computational methods, the potential to revolutionize materials science and chemistry has become increasingly evident. To explore this potential, we hosted a multidisciplinary workshop that brought together leading experts in materials science, physical sciences, data science, and computer science from various research institutions and IHLs in Singapore. The diverse expertise contributed to a rich and comprehensive discussion, offering a broad perspective on the current challenges and future directions in the field.

Singapore, with its strategic positioning as a global R&D hub, robust support from the government for AI and advanced materials, and profound collaborative networks both locally and internationally, stands at the forefront of this transformative journey. With dedicated efforts towards enhancing high-performance computing resources, fostering talent in AI and material science, and spearheading innovative public-private partnerships, Singapore is poised to address these grand challenges head-on. This mission, deeply embedded in the nation's strategic initiatives, reflects not only a commitment to scientific advancement but also a vision for a future where the materials that shape our world are smarter, more sustainable, and intricately tailored to meet the demands of tomorrow.

This white paper delves into these grand challenges, outlining the current landscape, potential paths forward, and the pivotal role that AI for Materials Science and Chemistry can play in this dynamic field. Through this exploration, we aim to catalyse a global dialogue on overcoming these challenges, fostering collaboration, and highlighting Singapore's emerging leadership in driving the next wave of innovations in materials science and chemistry.

The workshop was designed to encourage in-depth discussions and thorough analysis of key topics. It began with brief presentations by experts in the areas of functional chemistry, intelligent functional materials, soft materials and devices, two-dimensional materials, solid-state materials, computational materials and foundational models for materials. These presentations were followed by panel discussions that allowed for further exploration and exchange of ideas among the participants.

After the presentations, attendees engaged in breakout sessions focused on structural materials, functional materials, computational materials, soft materials, and machine learning methods. These sessions provided a platform for more specialized discussions, enabling participants to address specific challenges and opportunities within their respective domains and paving the way for future collaborative efforts and innovation.

## Grand Challenges

As a result of the workshop, we identified five potential grand challenges:

- Data Generation and Accessibility
- Multiscale Phenomena Modelling
- Materials Inverse Design
- Accelerate Computational Modelling
- Robust and Efficient Experimental validation

---

### GRAND CHALLENGE 1: DATA GENERATION AND ACCESSIBILITY

---

A critical and emerging grand challenge in the integration of AI in materials science and chemistry is enhancing data generation and accessibility. This challenge encompasses several key aspects, including synchronization across research groups, synthetic data generation, experimental data collection and curation, and the establishment of standardized knowledge representation and semantics for generated data and processes. Addressing these challenges is pivotal for advancing the field and maximizing the potential of AI-driven discoveries.

#### OBJECTIVE

To enhance data generation and accessibility across computational and experimental studies.

This white paper captures the key discussions and recommendations from our workshop, serving as a foundational document to guide future research priorities in AI for materials science and chemistry in Singapore. The document is structured by the identified grand challenges, which will fundamentally change the way research in materials science and chemistry is carried out as progress in overcoming them is made.

### EXPECTED OUTCOME 1: SYNCHRONIZATION ACROSS RESEARCH GROUPS

Achieving seamless synchronization across research groups presents a significant technical and organizational challenge. Research in materials science and chemistry generates vast amounts of data, often stored in disparate databases with varying formats. This diversity hampers the ability to share, compare, and integrate data across research initiatives. There is a compelling need for centralized platforms and interoperable databases that promote data sharing and communication between groups. These platforms must support a variety of data types and origins, ranging from experimental results to computational simulations, facilitating a holistic approach to materials research.

### EXPECTED OUTCOME 2: SYNTHETIC DATA GENERATION

Synthetic data generation stands as a promising solution to augment scarce experimental datasets, particularly in scenarios where obtaining real-world data is impractical or too expensive. Through techniques like Diffusion Networks, synthetic data that mimics real experimental outcomes can be created, providing an invaluable resource for training AI models. However, the challenge lies in ensuring the validity and reliability of synthetic data, such that it accurately represents the complex phenomena of materials science and chemistry.

### EXPECTED OUTCOME 3: EXPERIMENTAL DATA COLLECTION AND CURATION

The collection and curation of experimental data in a standardized, high-quality format is another critical obstacle. Consistency in data collection protocols and metadata documentation is essential for the robust validation of AI models. Moreover, the development of high-throughput experimental tools, aided by AI, can revolutionize data collection, enabling rapid generation of vast datasets. Nevertheless, these datasets must be carefully curated and annotated to ensure their usefulness, necessitating advanced data management strategies and technologies.

### EXPECTED OUTCOME 4: KNOWLEDGE REPRESENTATION AND SEMANTICS

A foundational challenge in enhancing data accessibility and utility is the establishment of standardized knowledge representation and semantics for generated data and processes. Effective representation involves developing a universal language or ontology for materials science and chemistry data, encompassing both experimental and computational findings. This standardized framework would facilitate the integration of data from various sources, supporting more complex and accurate AI-driven analyses and predictions. Additionally, it would enable more intuitive querying and interpretation of data, enhancing collaboration among researchers and accelerating the pace of innovation.

### EXPECTED OUTCOME 5: SOLICITATION AND CURATION OF UNSTRUCTURED KNOWLEDGE

The solicitation and curation of this unstructured knowledge using advanced language models and ontologies represents a complementary frontier in the quest to unlock novel insights and data. By employing sophisticated AI-driven tools like transformer-based language models tailored to decipher complex scientific texts, researchers can automate the extraction of valuable data and information previously inaccessible. Coupled with the development of powerful ontologies

that structure this extracted knowledge, they orchestrate a more comprehensive and accessible data landscape. This concerted effort not only amplifies the depth of data available for AI models but also bridges the gap between theoretical research and practical application, positioning the field to leapfrog traditional research methods and foster a richer, more interconnected scientific ecosystem.

### AI METHOD REQUIREMENTS

Interoperable AI platforms to integrate disparate datasets, diffusion models to generate synthetic data, AI-enabled high-throughput experimental tools, knowledge representation models for universal data framework, transformer-based language models to extract knowledge.

---

### GRAND CHALLENGE 2: MULTISCALE PHENOMENA MODELLING

---

A fundamental goal in materials science and chemistry is to rigorously understand the connections between processing, structure, properties and performance. Achieving this understanding would enable the rational design of materials with desired properties and functionalities. The major challenge lies in the fact that chemical composition alone does not determine final properties; phases, defects, morphology, structure and processing also play crucial roles. Currently, no multiscale modelling approach can capture these complexities. AI presents a promising solution to integrate these parameters and link them to material properties and functionalities. This challenge is critical because the behaviour of materials cannot be fully understood or predicted without considering the complex interactions that occur across different scales, from atomic to macroscopic levels.

#### OBJECTIVE

To understand the fundamental connections between processing, structure, properties and performance.

**EXPECTED OUTCOME 1:  
MULTISCALE FOUNDATION MODELS TO  
ADDRESS NATURE'S COMPLEXITY**

The main technical bottleneck to achieving this is a need to bridge multimodal data across different time and length scales for a given material system. Material systems span multiple length scales. For example, polymers consist of covalently bonded atoms at the atomic scale, to long chains of repeating units at the molecular scale, to crystalline or amorphous regions at the nanometre scale, to different morphologies at the micrometre scale, to a material with unique properties at the macroscopic scale. Different simulation techniques are used to explore phenomena at different length scales, from quantum simulations for atomic-level detail to continuum models for larger scales. There is a trade-off between the achievable length and time scales for the various simulation techniques, therefore datasets from each simulation span a wide range of scales. So, one of the primary hurdles in modelling multiscale phenomena is capturing the inherent complexity of nature's processes. To address this, there's a push towards developing Multiscale Foundation Models that can seamlessly integrate data across varying scales. These models aim to encapsulate the myriad interactions between processing, structure, properties, and performance in materials science. A significant challenge here lies in the technical bottleneck of bridging data that spans from atomic to macroscopic scales, which is further complicated by the different types of data involved, some of which cannot currently be accurately represented in a machine-readable form. One potential AI solution is to develop multiscale foundational models that are universally applicable across various applications or even domains. These models would integrate data from different length and time scales to provide a holistic understanding of material behaviour. Moreover, these multiscale models could be adaptable and applicable to a wide range of materials.

**EXPECTED OUTCOME 2:  
THEORETICAL FOUNDATION FOR MULTISCALE  
PHENOMENA ANALYSIS**

A strong theoretical foundation is crucial for the analysis of multiscale phenomena. This involves not only understanding the individual phenomena at each scale but also how these scales interact and influence each other. Current computational modelling methods are too slow and often too narrow to adequately predict material performance across scales. Accelerating computational modelling through AI can provide insights that are impossible to obtain through experimentation alone, due to practical constraints such as time, cost, and the physical limitations of observing phenomena at the atomic or molecular level.

**EXPECTED OUTCOME 3:  
INTERPRETABILITY AND KNOWLEDGE  
EXTRACTION FROM TRAINED MODELS**

Enhancing the interpretability of AI models and the extraction of knowledge from them are vital for advancing multiscale phenomena modelling. The ability to interpret what an AI model has learned can lead to deeper insights into material behaviour across scales. However, current AI models often act as "black boxes," providing little insight into how they arrive at their predictions. Developing AI solutions that can offer explanations for their outputs is crucial for the model's acceptance and trust by scientists and engineers. Furthermore, knowledge extraction from these models could inform future research directions and practical applications, making AI a valuable tool for innovation in materials science.

**EXPECTED OUTCOME 4:  
HIGH-THROUGHPUT COMPUTATIONAL TOOLS  
AND AI-DRIVEN EXPERIMENTATION**

An additional item to consider is the development of high-throughput computational tools and the integration of AI with experimental approaches. This involves leveraging AI to enhance experimental designs and outcomes, thereby increasing the speed and reducing the costs of discovering and testing new materials. By combining high-throughput experimental tools with AI's predictive capabilities, researchers can collect and analyse data more efficiently, leading to faster iterations in the material design process.

**EXPECTED OUTCOME 5:  
SUSTAINABILITY AND EFFICIENCY IN  
MATERIAL DESIGN**

Lastly, the emphasis on sustainability and efficiency in material design emerges as a crucial aspect of addressing multiscale phenomena modelling. AI models that can predict the environmental impact and energy efficiency of materials across their lifecycle will be invaluable. Incorporating sustainability criteria into multiscale models will enable the design of more environmentally friendly materials without compromising performance.

Integrated multiscale models that draw rigorous connections from the microscale to macroscale of a material will revolutionize material design, leading to more efficient and innovative applications. As an example, consider the development of high-performance battery materials. Understanding how the atomic and molecular structures of electrode materials influence their electrochemical properties would enable the design of batteries that are more efficient, have higher energy densities, and longer lifespans. This would allow for the precise tuning of materials through controlled synthesis and processing techniques, ultimately leading to batteries that charge faster, last longer, and are safer for consumer use.

**AI METHOD REQUIREMENTS**

Multiscale foundation models to integrate data from multiple time and length scales

---

**GRAND CHALLENGE 3:  
MATERIALS INVERSE DESIGN**

---

This challenge fundamentally inverts the traditional materials discovery paradigm by beginning with desired properties or functionalities and working backward to discover the underlying structures and processes that can manifest those attributes.

**OBJECTIVE**

To rationally design a material with specified properties and functionality.

**EXPECTED OUTCOME 1:  
INVERTIBLE GENERATIVE MODELS MAPPING  
COMPLEX FUNCTIONAL MATERIALS INTO  
STRUCTURES AND PROCESSES**

The cornerstone of the Materials Inverse Design challenge is the development of invertible generative models. These advanced AI models are tasked with mapping desired functional outcomes back to their structural and process origins. Unlike traditional forward generative models that predict outcomes based on given inputs, invertible models aim to decipher the complex lattice of possibilities leading to a specified material property or function. This involves a deep understanding of the intricate relationship between material structures, their processing conditions, and the resulting properties.

**EXPECTED OUTCOME 2:  
ADVANCED SAMPLING TECHNIQUES FOR  
GUIDED MATERIALS GENERATION**

Integral to effective inverse design is the implementation of advanced sampling techniques. These techniques enable the exploration of the vast, often underexplored, design space of materials to identify novel material compositions and configurations. Leveraging AI-driven approaches such as active learning and Bayesian optimization, these sampling methods focus on efficiently navigating the design space. By guiding the generation of materials towards regions with the highest potential for success, they significantly reduce the time and resources required for discovery.

**EXPECTED OUTCOME 3:  
FOUNDATION MODELS TUNABLE TO  
LOW-DATA CASES**

The challenge of materials inverse design is exacerbated in low-data scenarios, common in the exploration of new materials spaces where experimental data is scarce. To this end, foundation models that are tunable to low-data cases are crucial. These models utilize transfer learning and few-shot learning techniques to adapt to new tasks using minimal data. By leveraging the underlying patterns learned from extensive pre-training on related tasks, these models can effectively make predictions or generate insights for new material systems with sparse data inputs.

**EXPECTED OUTCOME 4:  
ADVANCED MULTICRITERIA OPTIMIZATION  
MODELS**

A vital component of the Materials Inverse Design challenge is the development of advanced multicriteria optimization models. These models tackle the multifaceted nature of materials design, where trade-offs between different properties and functionalities often exist. By employing sophisticated optimization algorithms capable of navigating these trade-offs, such models aim to identify materials that optimally balance multiple criteria. This multi-objective approach is essential for tailoring materials to complex applications requiring a delicate balance of properties.

**AI METHOD REQUIREMENTS**

Invertible AI models to map material function back to structure and process, advanced sampling techniques to explore vast design space, foundation models for sparse datasets, multi-objective optimization models for materials design.

---

**GRAND CHALLENGE 4:  
ACCELERATING COMPUTATIONAL  
MODELLING**

---

In situ experiments in materials science and chemistry are often limited by practical constraints such as time, cost and the ability to observe phenomena at the atomic or molecular level. Computational simulations play a crucial role in bridging this gap, providing insights that are difficult or impossible to obtain through direct experimentation alone. However, current computational modelling is too slow to predict material performance reliably. For example, simulating the folding process of a medium-sized protein (about 100 amino acids) might require a molecular dynamics simulation running at a speed of approximately 10 nanoseconds per day of computation on a high-performance computing cluster. This process, which can naturally take milliseconds to seconds, would require years of computation time to simulate accurately. Faster and more accurate computational modelling is crucial for accelerating multiscale simulations and maximizing the output from the computational resources available. There is an opportunity to integrate AI with computational simulations

to minimize computational resources without compromising accuracy. To tackle this challenge, several key areas must be developed and enhanced.

**OBJECTIVE**

To accelerate computational modelling of material systems.

**EXPECTED OUTCOME 1:  
INTEGRATED SIMULATION SOFTWARE WITH AI  
MODELS CAPABLE OF ROBUST UNCERTAINTY  
ESTIMATION**

Integrating simulation software with AI models offers a robust framework for addressing the inherent uncertainties in computational modelling. These AI-enhanced tools can predict outcomes and simultaneously estimate the uncertainties of these predictions, providing a more comprehensive view of the possible ranges of behaviour for new materials under various conditions. This approach not only accelerates the pace of computational modelling by offering more accurate predictions but also instils greater confidence in the results, facilitating informed decision-making in the experimental validation phase.

**EXPECTED OUTCOME 2:  
DIFFERENTIABLE COMPUTATIONAL  
ALGORITHMS FOR FASTER MULTICRITERIA  
OPTIMIZATION**

The creation and implementation of differentiable computational algorithms are essential for enabling faster multicriteria optimization. These algorithms allow for the direct optimization of complex models with respect to multiple criteria, facilitating a quicker convergence to optimal solutions. By computationally deriving gradients for all model parameters, it becomes possible to efficiently navigate the vast space of material properties and synthesizing processes, speeding up the search for materials that meet specific multi-faceted requirements.

**EXPECTED OUTCOME 3: NOVEL (NON-VON  
NEUMANN) COMPUTING ARCHITECTURES  
INCLUDING PROBABILISTIC COMPUTING TO  
REDUCE COMPUTATIONAL COMPLEXITY**

Exploring novel computing architectures, such as those based on non-Von Neumann paradigms, presents a revolutionary approach to reducing computational complexity. Probabilistic computing, quantum computing,

and neuromorphic computing are examples of such architectures that offer the potential to process complex computations more efficiently than traditional computing paradigms. Interestingly enough, those approaches are constantly searching for better underlying platforms and materials. By leveraging these advanced computing frameworks, the field can overcome the current limitations imposed by traditional computational resources, significantly speeding up the modelling process and reducing CO2 footprint related to extensive computations.

**AI METHOD REQUIREMENTS**

AI models for uncertainty estimation in computational predictions, differential algorithms for multi-objective optimization, non-Von Neumann architectures for more efficient simulations

---

**GRAND CHALLENGE 5:  
ROBUST AND EFFICIENT  
EXPERIMENTAL VALIDATION**

---

Computational models provide detailed insights at various scales, from atomic to macroscopic levels, but their predictions need to be validated against real-world data to ensure their reliability. This integration necessitates robust validation against experimental data, as simulations alone cannot account for all the complexities and variables present in actual material systems. Accurate experimental data serves as a benchmark to verify and refine computational models, ensuring their predictions are applicable to real-world scenarios. This challenge resolution is foundational to bridging the gap between theoretical models and real-world applications.

However, the scarcity of high-quality experimental data, particularly for dynamic processes, poses a significant challenge. One major technical bottleneck is the absence of standardized data collection formats, which leads to variations in experimental results across different settings. This inconsistency complicates the integration and comparison of data from diverse sources. Additionally, varied experimental conditions can significantly affect the results, making it challenging to build a comprehensive and consistent database for model validation and training. To address this multifaceted challenge, several key areas require development and refinement.

**OBJECTIVE**

To generate high-quality experimental data for validation of computation models

**EXPECTED OUTCOME 1:  
AUTOMATED EXPERIMENTAL FACILITIES  
CONNECTED TO AI MODELS**

The advent of automated experimental facilities that seamlessly interface with AI models marks a revolutionary stride toward robust validation. By integrating AI into the experimental workflow, these facilities can dynamically adjust experimental parameters in real-time based on incoming data and model predictions. This level of automation and AI integration not only accelerates the experimental phase but also enhances the precision and relevance of the validation efforts, ensuring that theoretical models undergo testing under the most informational conditions.

**EXPECTED OUTCOME 2:  
SYNTHESIS AND ANALYSIS PROTOCOLS  
ENHANCED BY AI**

To further bolster experimental validation, tailored synthesis, and analysis protocols enhanced by AI are paramount. These protocols leverage AI to optimize material synthesis processes and analytical techniques, ensuring that the generated materials and the data derived from their analysis are of high quality and relevance. AI-driven optimization aids in identifying the most effective experimental conditions, reducing the time and resources spent on trial-and-error approaches.

**EXPECTED OUTCOME 3:  
AI MODELS INTEGRATED IN COMPUTATIONAL  
SOFTWARE FOR EXPERIMENTAL DESIGN**

Integrating AI models directly into computational software used for designing experiments represents another critical avenue. These AI-enhanced tools can predict the outcomes of various experimental setups, guiding scientists in selecting the most promising ones for actual testing. By forecasting potential complications or highlighting promising areas of exploration, these tools play a pivotal role in streamlining the validation process, making it both more efficient and directed.

#### EXPECTED OUTCOME 4: UNCERTAINTY ESTIMATION IN EXPERIMENTAL VALIDATION

A robust experimental validation framework must inherently account for the uncertainties associated with both experimental and computational predictions. Implementing sophisticated uncertainty estimation mechanisms is crucial for assessing the confidence in experimental validation outcomes. These mechanisms allow researchers to quantify the reliability of the data and the predictions derived from AI models, providing a clearer understanding of the materials being studied and the fidelity of the models used.

## Role of Singapore in AI for Materials Science and Chemistry

In the evolving landscape of AI for Materials Science and Chemistry, Singapore stands at the cusp of global leadership, buoyed by its strategic advantages and commitment

The ability to generate large, standardized experimental datasets for model validation will enable the establishment of a centralized database in Singapore to provide a consistent and comprehensive data repository. This, in turn, can facilitate the development of robust foundational models. Moreover, developing high-throughput experimentation tools will facilitate the scaling out of experiments from the research environment to industrial scale.

#### AI METHOD REQUIREMENTS

AI-enabled automated experimental systems, AI methods to optimize synthesis protocols and suggest experiments

to technological innovation. Singapore has several strategic advantages that put it in pole position to become a leader in AI for materials science and chemistry.

TABLE 1: PROJECTED HIGH-PERFORMANCE COMPUTING (HPC) REQUIREMENTS FOR RIE 2025

Year	Computational need	Estimated HPC hours (hours/year)
2025	Development of multiscale AI models	200,000
	AI-driven material discovery (phase 1)	250,000
2026-2027	Integration of experimental and computational data	500,000
	High-throughput virtual screening	700,000
2028-2029	AI for sustainable materials	1,000,000
	Real-time data-driven manufacturing optimization	1,500,000
2030	Autonomous AI laboratories	2,000,000
	Predictive AI for material lifespan and performance	2,500,000

First, situated as a global R&D hub in Asia, Singapore benefits from its central location, facilitating collaborations and knowledge exchange across international borders. Within Singapore, there are strong collaborations between academia and industry (e.g. between NUS, NTU and A\*STAR) and strong government-

industry initiatives (e.g. Startup SG to support startup development). Outside of Singapore, we also have strong partnerships with top overseas universities (e.g. MIT, Cambridge, UC Berkeley), industrial partners (e.g. Siemens and IBM Research) and public-private partnerships (e.g. Google AI and BASF). Such

strong collaborations can serve to enhance data availability and accessibility, which is a hindrance particularly with industrial data.

Second, the government's strong support for AI and advanced materials further enhances Singapore's competitive edge, with dedicated funding, policy frameworks, and incentives attracting top talent and companies to the region. Strong governance can be leveraged to implement policies and infrastructure that provides a competitive edge in addressing the Materials Inverse Design challenge. For example, highly skilled researchers trained in both materials science and AI are needed to carry out the research. We can envision educational policies put in place at the national level to facilitate such training. As another example, computational resources at the national level can be efficiently distributed to meet the differing computational needs of various research groups and areas. Projected computing resources needs for AI in materials science and chemistry for RIE 2025 are listed in Table 1.

Educational policies and initiatives aimed at cross-training researchers in both materials science and AI are fundamental to cultivate the expertise needed for Robust and Efficient

Experimental Validation. By establishing high-throughput experimental facilities and integrating AI into the experimental design and validation processes, Singapore can lead in efficiently bridging the gap between theoretical models and real-world applications.

The envisioned centralized database for standardized experimental datasets, as highlighted in the document, will be critical for supporting the development of foundational models and facilitating the Materials Inverse Design and Robust and Efficient Experimental Validation challenges. Singapore's strong governance and infrastructure can play a pivotal role in overcoming technical bottlenecks, such as data accessibility and integration across different research platforms.

Moreover, the development of high-throughput experimentation tools, boosted by Singapore's advanced manufacturing and technological capabilities, will be essential in scaling experiments from research environments to industrial applications. This directly contributes to the Accelerating Computational Modelling challenge by enabling faster iteration and integration of computational predictions and experimental data.

## Conclusion

Advancements in AI offer ground-breaking opportunities across the spectrum of materials science and chemistry, a realm rich with complexity and promise. As highlighted in the document, the path forward is marked by pressing grand challenges: foregrounding Data Generation and Accessibility, the pivotal role of Multiscale and Inverse Design models, fuelled by Accelerated Computational Modelling, and imperative Robust and Efficient Experimental Validation. These challenges underscore the necessity for a cohesive, integrated approach that marries the theoretical and experimental under the expansive umbrella of AI's capabilities.

Singapore is strategically positioned to spearhead this exploration, drawing

upon its robust technological ecosystem, governmental advocacy for AI and advanced materials research, and a deep-rooted culture of international collaboration and innovation. By focusing on harnessing its strengths — including advanced computational resources projected to meet the rising demands of AI-driven materials science and a concerted push towards high-throughput experimental methodologies — Singapore aims to transcend traditional barriers in materials research. This concerted approach not only promises to catalyse breakthroughs in material science but also to position Singapore at the vanguard of global research and innovation, redefining the intersection of AI and material science for future generations.

# APPENDIX IX. CHEMICAL AND BIOMANUFACTURING

## AUTHORS:

Prof Saif A. Khan  
(NUS)

## Executive Summary

Singapore is one of the world's largest energy, chemical and bio-pharmaceutical manufacturing hubs, and has positioned itself as a destination for next-generation sustainable, digitalized and net-zero manufacturing solutions in response to global trends<sup>1,2,3</sup>. There is also significant momentum globally in the domain of AI for (accelerated) chemical and materials discovery for a variety of intended applications in pharmaceuticals, energy and sustainability. The ultimate success of these discoveries depends on the *manufacturability* of new molecules and

materials (or products containing them), *sustainably* and at scale.

This whitepaper articulates a set of overarching scientific grand challenges for chemical and biological manufacturing that can be addressed with the growing power of AI methods. The paper also outlines specific directions within these grand challenges with tremendous potential for direct and tangible impact on the Singaporean advanced chemical, biological and materials manufacturing landscape.

## Background

The whitepaper consolidates the outputs of a multidisciplinary workshop on AI for chemical and biological manufacturing organized on April 24, 2024, with >100 attendees from academia, industry, government agencies and A\*STAR<sup>1</sup>. The workshop was designed around five key topics focusing on the pivotal role of AI in accelerating science and innovation in: (i) chemical and (bio)chemical manufacturing (eg. small-molecule drugs made via thermochemical, biocatalytic/biochemical and photo/electrochemical routes), (ii) biologics manufacturing (eg. antibodies made from the culture of animal cells), (iii) materials manufacturing (eg. polymers, nanoparticles and their assemblies for advanced consumer products), (iv) process systems engineering,

control and optimization and (v) modeling, simulation and design of multi-scale processes. Each topic was covered by a keynote talk and had a dedicated discussion group composed of industrial and academic experts to provoke deep discussions and collective brainstorming.

In the following, we outline grand challenge statements that arose from the workshop, which have been distilled and synthesized from the voluminous workshop discussions. These challenge statements cut across the individual topics covered and have the potential to revolutionize advanced chemical and biological manufacturing.

## Grand Challenges

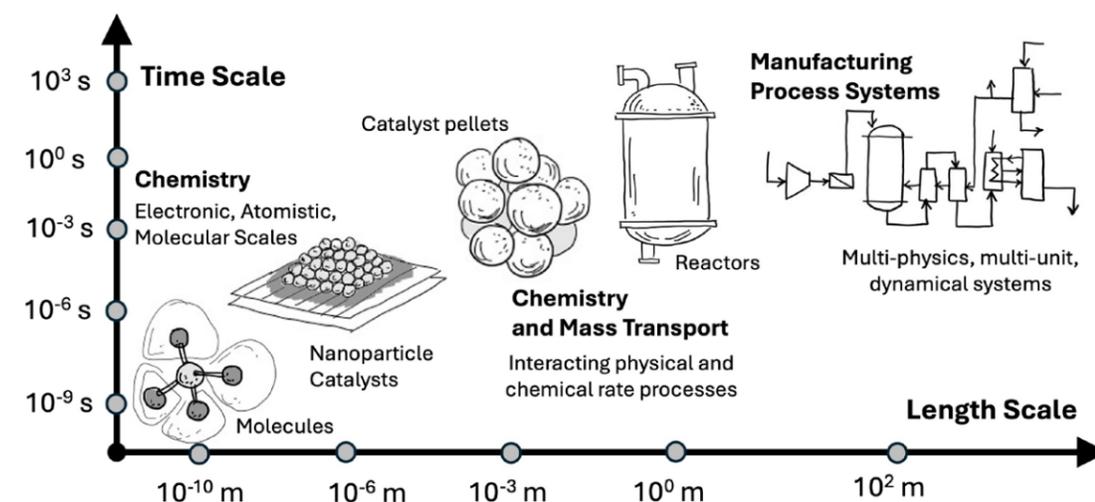


Figure 1: Chemical and biological manufacturing is an inherently *multi-scale* and *multi-physics* challenge, involving *multi-unit process systems*.

Chemical and biological manufacturing processes are inherently complex multi-unit, multi-physics *systems*, involving spatial and temporal dynamics across *multiple scales* (Fig. 1). As such, the synthesizability of a molecule or material at the lab scale typically does not guarantee the viability of a *scaled-up* manufacturing process. This complexity also deeply impacts *process operation* and *control*, and has traditionally imposed harsh limits on the flexibility, agility and sustainability of industrial processes.

### GRAND CHALLENGE 1: RADICAL ACCELERATED DESIGN

Scaled-up process design in chemical and biological manufacturing has traditionally been an iterative, time- and resource-intensive combination of empiricism, past experience and, increasingly, mathematical modeling and simulation, with no *a priori* guarantee of success. As such, **the ability to achieve one-shot de novo design of scaled-up, economically viable and environmentally sustainable manufacturing processes for entirely new molecules or materials is a scientific grand challenge** across the various molecular manufacturing sectors, including

energy, (bio)pharmaceuticals, specialty chemicals, agrochemicals, etc. An example of the paradigm-changing *disruptive impact* enabled by solutions to this grand challenge would be the AI-enabled discovery of a new drug with the accompanying *a priori complete specification* of a scalable and modular net-zero manufacturing process for the final drug product that can be prototyped and deployed in a time scale of *weeks* instead of years. This is *not possible* today.

### OBJECTIVES

The main objective is to leverage AI to reduce reliance on human experience, ensuring one shot optimal process design for new problems across various industrial sectors and for plants of different scales and complexity, thus significantly reduce overall R&D costs and lead times.

Deliverables include systematic machine learning-based process design and optimization protocol(s), and a generalized software package(s) that can leverage limited experimental data and interact with process simulators like Aspen HYSYS to generate optimal process designs for chemical and biological manufacturing across sectors.

#### DATA REQUIREMENTS

Data needed for scaled-up process design are typically acquired from computer simulations using process simulators for plants of multiple scales, along with experiments on lab-scale and pilot-scale plants. Specifically, while it is desirable to achieve a high level of generalizability of the resulting process design framework(s), the data for energy systems and data for (bio)pharmaceutical plants will need to be collected/prepared separately (see Section 4.2 below).

#### AI METHOD REQUIREMENTS

Dense neural networks and recurrent neural networks will form the foundation for modeling both static and dynamic systems. Feature selection and extraction techniques, including principal component analysis and autoencoders, will identify and extract the most relevant features for further analysis and modeling. To expedite the scale-up process, transfer learning can be employed to apply insights gained from lab-scale experiments to large-scale plant operations. Additionally, modular learning approaches can be utilized to integrate various modules developed for individual unit operations (e.g., foundation models, data-driven models, gray-box models, and first-principles models) into a cohesive global model for complex process networks to support process scale-up design with improved modeling accuracy and flexibility.

#### COMPUTE REQUIREMENTS

2 H100 cards for 12 months

#### GRAND CHALLENGE 2:

### RADICAL 'OUT-OF-DISTRIBUTION' DYNAMICS

The second grand challenge is the (AI-enabled) **dynamic operation and control** of (existing and new) manufacturing processes **in non-intuitive 'out-of-distribution' process windows** that enable ultra-precise manufacturing with embedded quality control while allowing agility in response to varying feedstocks, energy sources and supply-chain disruptions. **Real-time monitoring and control** is widely considered to be a **'holy grail' problem**, and is severely challenged by the same complexities (i.e. multi-unit, multi-

physics, multi-scale spatiotemporal dynamics) that make scaled-up process design such a daunting challenge. Addressing this problem would be particularly *critical* in next generation manufacturing processes that start from net-zero feedstocks that could come from a variety of sources of varying quality and driven by a mix of energy sources with temporally varying supply.

#### OBJECTIVES

Most of the existing dynamic modeling and control frameworks predominantly depend on mechanistic models grounded in physical laws, such as mass and energy balances. This dependency has significantly restricted the application of these solutions to complex chemical and biological processes. *Therefore, key objectives are (i) to seamlessly integrate machine learning and data analytics into control and optimization theories and (ii) develop innovative, practical data-enabled dynamic modeling and advanced control frameworks and algorithms.* These solutions will be generalizable to various chemical and biological processes and scalable to achieve autonomous and resilient plant-wide system operations and manufacturing that can provide a higher level of safety, production consistency, energy efficiency, and sustainability.

Deliverables include emerging machine learning-based dynamic modeling and control methods that leverage state-of-the-art AI concepts, validated in terms of efficient computation and reduced energy consumption.

#### DATA REQUIREMENTS

Training a machine learning-based process model and a learning-enabled process controller often requires vast amounts of data. However, chemical and biological industry players typically do not share their datasets publicly. Additionally, real-time control necessitates online data collection. To address these challenges will require the development and utilization of benchmark chemical and biological process simulators for data generation and control performance evaluation, while also leveraging existing benchmark lab-scale setups to generate experimental data (see Section 5.2 below).

#### AI METHOD REQUIREMENTS

Traditional AI methods, such as long short-term memory (LSTM) networks and convolutional neural networks (CNNs) can initially be developed for process modeling and used as baseline models for performance comparison with more advanced AI techniques. Meta-learning and foundation models can then be developed to create unified process models that allows for rapid adaptation to new processes with minimal training and data requirements. Additionally, advanced neural network architectures, such as liquid neural networks and neural ordinary differential equations, can be developed for dynamic systems. Lastly, online machine learning and continual learning can be employed to enhance model performance in real-time, addressing process disturbances and model uncertainties effectively. For machine learning-based control, model predictive control (MPC) and

reinforcement learning (RL) are two prominent advanced control methods. Specifically, machine learning models can be integrated into existing MPC frameworks to enhance their performance. Novel approaches, such as physics-informed machine learning models, Koopman operators, and input-convex neural networks, can be explored to improve control performance and computational efficiency. Reinforcement learning techniques can also be employed to derive optimal control policies from process operational data, where safe reinforcement learning methods will be developed to ensure safe data collection and learning.

#### COMPUTE REQUIREMENTS

2 H100 cards for 24-36 months (GPUs are needed for machine learning controller design and improvements)

## Impact of AI Methods and Data - Challenges and Opportunities

The complexity of the above grand challenges is particularly well suited to a **multi-modal data-driven and AI-enabled approach**, with potential for game-changing economic impact.

#### AI METHODS

Conventional *de novo* process design is a sequential, stage-wise workflow, with each stage traversing a hierarchy of length and time scales, up from molecules to process systems. The workshop discussions highlighted the role of AI methods in **enabling fundamentally new, accelerated and disruptive integrative workflows** that allow single-shot full-scale process design. New processes are designed for specified starting materials, via a specific synthetic *route* (i.e. sequence of chemical or biochemical reactions) and under *constraints* such as yield and purity of the target molecule(s). AI methods have already made great strides in the prediction of viable synthetic routes for target molecules from specified starting materials<sup>4</sup>.

#### AI FOR THERMODYNAMICS AND (BIO) CHEMICAL KINETICS

The first stage of manufacturing process design involves the measurement or estimation of critical manufacturing-relevant thermodynamic properties, particularly those of molecular mixtures (such as phase equilibrium and solubility) and chemical/kinetic properties (such as reaction rate coefficients and the distribution of possible byproducts) in a given synthetic route. These are *intrinsic, scale-independent properties and rate processes, agnostic* to the specific processing equipment used. Work in this stage conventionally involves low-throughput, labor-intensive *lab-scale experimentation* in limited parameter spaces. **AI methods can enable orders-of-magnitude accelerations in the generation of such information, via learnt and interpretable molecular structure-thermodynamic/kinetic property mappings.** Information from this stage is used to identify possible superstructures for the final manufacturing system (i.e. a sequence of reactor(s) and separation units, for example), also known as process 'flow sheets' (see below).

## AI FOR PROCESS MODELING AND SCALE-UP

The science of scale-up for process design involves understanding and accounting for *physical rate processes*, such as fluid mechanics, mass and energy transport within the units comprising the process flow sheet, which are typically strongly coupled to *chemical rate processes*, such as reaction kinetics, in highly non-linear *scale-dependent* fashion. Therefore, the next stage in process design is to quantify this *scale-dependence*, which dictates how far away from ideality (i.e. the thermodynamic limits identified in 4.1.1) any full-scale manufacturing process operates. This ultimately decides the energy consumption and environmental sustainability of any given process at scale. **AI methods can enable orders of magnitude accelerations in quantifying and modeling scale-dependence**, which is typically described by systems of partial differential equations for each process unit, parameterized by the information gathered in 4.1.1 (these are computationally intensive to solve even for simple, single fluid phase systems). **AI methods can enable the facile development of high fidelity, interpretable dynamic models for entire process flow sheets**, which allow inverse flowsheet design and enable process operation in 'out-of-distribution' process windows, ultimately facilitating fully automated 'lights-off' manufacturing environments [4].

## INTEGRATION

AI-enabled acceleration of *de novo* process design will enable the flexible, modular design of *bespoke* process hardware (reactors, separation units, etc.) for *intensified*, energy efficient and sustainable molecular transformations. **This would completely invert the conventional (and conservative) paradigm**, which has been to manufacture molecules and materials in *standardized* large-scale equipment, often leading to vast material and energy inefficiencies. Furthermore, these methods have significant potential to be **integrated at the scale of foundation models for manufacturing**. Such foundation models will be able to map an evolving space of novel functional molecules/materials to spaces of viable, scalable manufacturing processes - a *fundamental and unaddressed* problem that falls squarely within the domain of AI4SCI.

## DATA

The workshop discussed the issue of data at length; it was recognized that **these grand challenges are faced with the paradox of simultaneously having too much and too little data**. Given the multi-scale nature of manufacturing, which spans molecules to supply chains, data in the domain tend to be heterogeneous (discrete/continuous, scalar/vectorial, static and/or spatio-temporally varying), high dimensional, fragmented, with varying levels of noise and bias, and are constrained by known physical laws. There thus exist open challenges of *data ontology* and especially *representation*, particularly at the level of interconnected process units that make up process *systems*.

## DATA ABUNDANCE

There is already an abundance of historical measured process-scale data across industry partners, spanning a very wide variety of molecules manufactured and processes deployed. The volume of available *multi-modal* data across chemical and biological manufacturing is also poised to expand exponentially as companies deploy digitalized and distributed 'smart' real-time sensing of key process parameters. There is thus considerable potential to tap into these data via a *federated learning* model, for which Singapore is a unique hub for cross-sector efforts. However, these data tend to be clustered around nominal, baseline operating conditions (for safety and operational stability reasons) in full-scale process systems, and as such may not be as information rich as the data volume would imply. Next, in parallel with observational data, there exist voluminous multi-scale *simulation* data across academia, research institutes and industry, spanning broad variety of methods, from *ab initio* molecular design to computational fluid dynamics (CFD) and process flow sheet simulations.

## DATA SCARCITY

A primary challenge in building AI-enabled multi-scale molecule-to-process models that address the grand challenges is the comparative *scarcity* of data that *links* the various scales. Starting at the molecular scale, there exist numerous international consortia that are mapping molecular structure to

*chemical* properties, for example, with large datasets comprising literature, experimental and simulation data. However, *process-relevant* physical/thermodynamic properties, particularly those of molecular mixtures, such as phase equilibria, solubility and solid-state form, and chemical properties, such as reaction rate coefficients, are comparatively scarce and unsystematized; these are critical in making choices for and sizing process equipment like reactors and separators. There is an even greater scarcity of systematized information/data on scale-dependent, far-from-equilibrium *physical* rate processes, such as hydrodynamics, mass and heat transport.

## Singapore's Role

The biopharmaceutical sector in Singapore has seen substantial growth, with the country being home to leading pharmaceutical companies such as Pfizer, GlaxoSmithKline, and Roche. These companies have chosen Singapore for its advanced manufacturing capabilities, regulatory excellence, and strong intellectual property protection. The biopharmaceutical sector is supported by state-of-the-art facilities, including small molecule and biologics manufacturing

## DATA FACTORIES

Given the above considerations regarding the availability of data, the workshop recognized the potential for the development of **'Data Factories' - self-driving lab platforms for high quality multi-scale process-relevant dataset generation**, to augment multi-partner data. This would leverage considerable expertise already existing in Singapore, built via programmes such as the Accelerated Materials Development for Manufacturing (AMDM) project.

plants and research institutions. Additionally, the chemical manufacturing industry is bolstered by Singapore's strategic location, which facilitates efficient access to raw materials and global markets. Thus, this concentration and connectivity of basic science, AI expertise, industrial/academic R&D, policy and manufacturing in ~750 km<sup>2</sup> confers a significant competitive advantage to Singapore in playing a leading position in this domain.

## Conclusion

AI methods promise significant and fundamental paradigm shifts in chemical and biological manufacturing, enabling hitherto unforeseen, radical acceleration and

opportunities in the design and operation of state-of-art chemical and biological manufacturing processes for both known and novel molecules and materials.

## REFERENCES

- 1 <https://www.edb.gov.sg/en/our-industries/energy-and-chemicals.html>
- 2 <https://www.edb.gov.sg/en/our-industries/pharmaceuticals-and-biotechnology.html>
- 3 <https://projects.gbreports.com/singapore-chemicals-ingredients-and-materials-2023/welcome-letter-sections>
- 4 B. Decardi-Nelson *et al.*, Computers and Chemical Engineering 187 (2024) 108723.

# APPENDIX X. FINANCIAL SERVICES

## Financial Services

Registration, Coffee/Tea & Light Breakfast  
Opening Remarks by the AI for Science team – Profs. Yang Zhang and Kedar  
Hoppagankar  
Motivation and Drivers of the Thematic Workshop – Profs. Bo An and Ke-wei  
Huang  
Tea / Coffee Break  
Sharing by Breakout Leads  
Summary, Way Forward (Bo & Ke-Wei)  
Networking Lunch



### AUTHORS:

Prof. Bo An  
(NTU, Artificial Intelligence Research Institute)

Prof. Ke-Wei Huang  
(NUS, Asian Institute of Digital Finance)

## Executive Summary

The financial industry, an ecosystem with billions of users and over 80 trillion market capitalization, has become a paramount pillar for many countries (e.g., Singapore). Recent advances in AI are rapidly transforming the finance sector by offering powerful tools and models for data analysis, risk management, customer service and trading. In this white paper, we point out three promising directions

to explore: i) generative AI and theory-guided AI for financial service; ii) AI governance in financial services; iii) LLM and multimodal AI empowered financial services. Discussions on challenges, approaches and potential impact are provided. We believe Singapore is well positioned to take the lead on AI-powered financial applications in the next decade.

## Introduction

Over the last decade, Artificial Intelligence (AI) has dramatically transformed the financial services industry, driving innovation and reshaping operational dynamics. Globally, AI has become integral to the development of fintech solutions, risk management, fraud detection, and personalized customer experiences. In the United States, for instance, AI-driven robo-advisors like Wealthfront has

gained significant traction, managing over \$50 billion in assets as of 2023.<sup>17</sup> Meanwhile, Singapore's AI funding soared by 77 percent to US\$481.21 million in 2023 across 24 deals, signifying increasing faith in AI's potential.<sup>18</sup> As a global financial hub, Singapore's commitment to integrating AI into financial services is not merely an opportunity but a strategic imperative.

Recognizing the critical role of AI in maintaining its competitive edge, the Monetary Authority of Singapore (MAS) recently committed up to S\$100 million under the Financial Sector Technology and Innovation Grant Scheme (FSTI 3.0) to bolster quantum and AI capabilities in the financial sector. This significant investment underscores Singapore's dedication to becoming a global leader in AI and quantum technologies in financial services. The funding aims to support financial institutions in building and deploying AI models, developing AI platforms for industry-wide use cases, and establishing AI innovation centers that can anchor AI excellence in Singapore.

Despite the rapid advancements and significant investments in AI within the financial services sector, there are still notable limitations that must be addressed to fully harness the potential of these technologies. One of the primary challenges lies in the lack of theory-guided AI models that integrate established financial principles. While AI has excelled in data-driven decision-making, many existing models lack the foundational understanding of financial theories, leading to outcomes that may be optimal from a computational perspective but are misaligned with core financial principles. This gap underscores the need for theory-guided AI for financial services, where AI models are not only data-centric but also deeply rooted in the financial domain's theoretical frameworks.

Another significant limitation is the ongoing challenge of ensuring the interpretability, transparency, and fairness of AI models, as well as the responsible use and regulation of these technologies. As AI becomes more deeply integrated into financial services, questions about bias, accuracy, and ethical implications are increasingly critical. Current research often falls short in addressing these complexities, resulting in AI systems that may produce accurate predictions but lack transparency and accountability. The financial sector, characterized by stringent regulatory requirements, demands robust governance

frameworks to ensure that AI technologies are deployed ethically and in compliance with regulations. However, there is still a significant gap in developing comprehensive governance models that can keep pace with the rapid evolution of AI technologies. This limitation highlights the critical need for research and development in AI governance to build trust and accountability in AI-driven solutions.

Lastly, while Large Language Models (LLMs) are among the most prominent advancements in AI, their application within the financial sector is still emerging. LLMs have the potential to revolutionize customer interactions, automate compliance processes, and extract insights from vast amounts of textual data. However, research and implementation in finance are in the early stages, with limited exploration of their full capabilities and associated risks. Additionally, the financial industry is beginning to recognize the value of integrating diverse data types, leading to the potential combination of LLMs with multimodal AI, which incorporates text, images, audio, and numerical data. This integration presents an even greater opportunity for innovation, yet it remains underdeveloped. The current gap in fully leveraging LLMs and multimodal AI to enhance decision-making, customer service, and operational efficiency underscores the need to advance LLM and Multimodal AI empowered financial services, where the synergy of multiple data modalities can drive more comprehensive, accurate, and impactful AI-driven solutions in the financial sector.

This white paper proposes three key themes that encapsulate the grand challenges and opportunities presented by AI in financial services in Singapore. These themes collectively represent the forefront of AI innovation in financial services, offering a roadmap for Singapore to lead in the responsible and effective integration of AI into its financial ecosystem. This white paper aims to outline the grand challenges within these themes, guiding future research, policy development, and industry practices in the realm of AI and finance.

17 <https://www.reuters.com/business/finance/roboadvisor-wealthfront-reaches-50-billion-client-assets-2023-11-16/>

18 <https://kpmg.com/sg/en/home/media/press-releases/2024/02/singapore-ai-funding-skyrockets-fintech-remains-resilient.html>

## Background

NRF organized a research workshop to foster collaboration between experts in the field of financial services and AI, from both academia and industry. The discussions were structured around three key themes, which are crucial for advancing AI's role in financial services in Singapore.

### GENERATIVE AI AND THEORY-GUIDED AI FOR FINANCIAL SERVICES

The first theme explored the transformative potential of generative AI and theory-guided AI in financial services. These advanced AI tools are making significant strides in simulating complex financial environments with high accuracy, surpassing traditional models. By employing generative AI, financial institutions can enhance market analysis, decision-making, and scenario testing, improving overall efficiency and resilience. Key applications include simulating market responses for investors and regulators, managing investment risks, and serving as educational tools for students and researchers.

The workshop highlighted several innovative uses of generative AI and theory-guided AI:

- **Synthetic Data Generation:** Valuable for training AI models when real data is scarce, or privacy concerns limit data availability.
- **Scenario Planning and Stress Testing:** Generating economic scenarios to test model resilience and assess potential risks.
- **Fraud Simulation:** Enhancing fraud detection systems through the simulation of synthetic fraudulent activities.
- **Personalized Financial Products:** Customization of financial products based on individual risk profiles.

However, the deployment of generative AI in financial services faces several challenges:

- **Processing Large Multimodal Data:** Integrating diverse data types into coherent models remains difficult.

- **Building High-Fidelity Simulators:** Accurate market simulators are challenging to create, especially with limited data.
- **Trading Decision Complexity:** Developing reliable AI systems for trading decisions under varying conditions is complex.
- **Output Quality Evaluation:** Ensuring the accuracy and reliability of AI-generated insights requires robust validation frameworks.
- **System Integration:** Incorporating generative AI into existing financial systems presents compatibility and implementation challenges, necessitating careful consideration of both technological and organizational factors.

Generative AI and theory-guided AI have the potential to revolutionize financial services by enabling sophisticated financial products, enhancing risk management, and improving regulatory compliance. However, addressing challenges related to data processing, model validation, and system integration is crucial for realizing these benefits. The workshop discussions emphasized the importance of responsible AI development to ensure positive and equitable outcomes in the financial sector.

### REGULATING AI USAGE IN FINANCIAL SERVICES

The second theme focused on the critical regulatory challenges associated with the growing integration of AI in financial services. As AI becomes increasingly embedded in the industry, robust regulatory frameworks are essential to ensure safe, ethical, and transparent AI usage.

The workshop highlighted several key challenges and necessary actions:

- **Bias in Datasets:** AI models can inherit biases from their training data, potentially leading to unfair treatment of certain customer segments. Addressing these biases is crucial for maintaining trust in AI-driven financial services.

- **Accountability and Trust:** Establishing clear lines of accountability is essential, especially when multiple parties are involved in the AI model supply chain.
- **Fraud and Deepfake Detection:** As AI-driven fraud and deepfake techniques become more sophisticated, there is an urgent need for advanced detection mechanisms to combat these growing threats.
- **Vendor-Sourced Models:** Ensuring third-party AI models comply with local regulations and are resilient against fraud is critical.
- **Cross-Border Data Transfers:** Varying data privacy laws across countries, particularly within ASEAN, complicates the transfer of financial data, creating a fragmented regulatory landscape.

The workshop also emphasized the need for:

- **Explainable AI and Transparency:** AI models must be transparent in their decision-making processes to ensure compliance and prevent bias.
- **Traceability of Datasets:** Ensuring data traceability is critical for accountability and auditing.
- **Global Collaboration on Standards:** Harmonizing regulations through international collaboration, especially in data privacy, is essential.
- **Liability and Risk Management:** Clear guidelines are needed to determine who is responsible when AI systems fail.
- **Ethical AI Deployment:** AI systems should be deployed fairly and equitably, avoiding the exacerbation of existing inequalities.

To fully realize the benefits of AI in financial services, key regulatory issues must be resolved. Research should focus on developing standards for explainable AI, improving data traceability, and fostering international collaboration. By addressing these challenges, Singapore can lead in the responsible deployment of AI in financial services, ensuring that these technologies enhance the efficiency, security, and fairness of global markets.

### LLM-EMPOWERED FINANCIAL SERVICES

The third theme explored the transformative impact of Large Language Models (LLMs) in financial services, highlighting their potential to revolutionize decision-making and extract value from unstructured data. LLMs can enhance multimodal forecasting, analyze ESG reports, and support applications like automated customer service and personalized financial advice. However, challenges such as reducing prompt dependency, addressing language biases, and overcoming hardware limitations were identified as critical areas for future research.

The workshop identified several key challenges and potential applications of LLMs:

- **Customer Servicing:** LLMs could revolutionize customer service, but challenges in handling complex queries and ensuring transparency need to be addressed. Solutions include integrating LLMs internally with human experts in the loop.
- **Asian LLMs:** Developing LLMs that can handle the linguistic diversity of Asia is essential. A dual-model approach, where a smaller model handles personalization and a more powerful LLM addresses complex queries, was proposed.
- **Risk Management:** Understanding and managing LLMs' risk profiles is key for their use in high-stakes financial decisions.
- **Knowledge Graphs:** LLMs could automate the construction of domain-specific knowledge graphs, though challenges in interpreting tabular data remain.

The workshop also emphasized the challenges and needs:

- **Accuracy and Reliability:** Ensuring the robustness of LLM outputs is crucial for critical financial decisions.
- **Adaptability:** LLMs must be able to adapt to rapidly changing market conditions and regional variations, such as those in ASEAN.
- **Strategic Planning:** Developing LLMs capable of long-term strategic planning, beyond immediate responses, remains challenging.
- **Transparency and Compliance:** LLM decision-making processes must be transparent to ensure regulatory compliance and maintain trust.

## Grand Challenges

### GRAND CHALLENGE 1: GENERATIVE AI AND THEORY- GUIDED AI FOR FINANCIAL SERVICE

Generative AI like SORA<sup>1</sup> has achieved remarkable performance and holds the potential to act as general-purpose simulators of the physical world. Under the context of finance, adopting advanced generative AI techniques with financial data (e.g., stock price, financial news and company earnings reports) can provide tools to model the complex financial systems more accurately and offer a deeper understanding of investors behavior. In addition, building financial world simulator would enable market analysis, informed decision-making, robust scenarios testing, thereby enhancing the efficiency and resilience of the finance sector. A wide range of users can benefit from financial generative AI applications. For finance regulators, they can provide insights into how markets react to changes in policy. For professional investors, they can improve risk management by analyzing potential risk and black swan events. For business school students, they can provide an easy-to-use education tools to understand financial markets. For economic researchers, they can offer a realistic testbed of new economic hypothesis and theories.

#### OBJECTIVE

To build a high-fidelity financial market simulator, three major challenges are as follows: i) how to efficiently process large amount of multimodal financial data; ii) how to achieve generalization ability and robustness with limited financial data; iii) how to systematically evaluate the quality of the simulator and use it to enhance downstream decision-making. Existing work on market simulator<sup>2</sup>, time series generation<sup>3</sup> and world model<sup>4</sup> has shown the potential for financial world simulators. Technically, the approach can be based on two mainstream generative AI methods (generative adversarial network and diffusion models). Novel design and components with financial insights are required to further improve the performance and address the above challenges for building realistic financial simulators.

The expected outcome is to develop a realistic financial market simulator.

#### DATA REQUIREMENTS

Long-term price and volume data of different financial markets.

#### AI METHODS REQUIREMENTS

Diffusion model and GAN.

#### RESOURCE REQUIREMENTS

Storage space, financial data purchase, 8 H100 GPUs for 12 months.

### GRAND CHALLENGE 2: AI GOVERNANCE IN FINANCIAL SERVICES: INTERPRETABILITY, TRANSPARENCY, AND ETHICAL COMPLIANCE

The rapid adoption of AI systems in the finance sector brings significant opportunities but also introduces complex challenges related to data privacy, ethical considerations, and the potential for biased decisions. As AI becomes increasingly capable of analyzing real-time financial information and making high-stakes decisions with minimal human oversight, the risks associated with its use intensify. Professor Michael Wellman highlighted in his US Senate testimony that financial regulations are struggling to keep pace with the advancements in AI technology. This gap underscores the urgent need for robust AI governance frameworks in financial services, particularly in Singapore, where the financial sector is a key pillar of the economy and a global financial hub.

In Singapore, the importance of regulating AI systems in financial services cannot be overstated. The financial sector here is deeply interconnected with global markets, making it a prime target for AI-driven innovations. However, this also exposes it to risks that could have far-reaching implications. Ensuring data privacy, preventing fraud, and promoting the ethical use of AI are critical to maintaining trust and stability in the financial system. Singapore's position as a leader in fintech and its commitment to innovation place it at the



Figure 1: US senate testimony: AI risks to the financial sector

forefront of addressing these challenges. To mitigate the risks associated with AI in finance, there is a pressing need to build AI models that are not only powerful but also transparent and explainable. This line of research is vital for fostering a secure and ethical financial environment, both in Singapore and globally.

#### OBJECTIVE

The integration of AI into financial services presents several grand challenges that must be addressed through rigorous scientific research and practical industry applications to ensure responsible and ethical use.

AI models, particularly those deployed in high-stakes financial decision-making, often function as "black boxes," making it difficult for stakeholders to understand the reasoning behind specific decisions. This opacity not only undermines trust but also complicates efforts to align AI-driven decisions with regulatory standards and ethical norms. A critical research challenge is the development of advanced methodologies that enhance the interpretability of AI models while preserving their predictive accuracy. This involves exploring novel approaches such as explainable AI (XAI) techniques, model-agnostic interpretability frameworks, and hybrid models that combine human expertise with machine learning. Additionally, a key area of focus is developing models that can be explained and validated through established finance and economics theories, rather than relying solely on variable importance scores. By grounding AI decisions in well-understood theoretical frameworks, these models can provide more meaningful and contextually relevant explanations, fostering deeper trust

and understanding among financial institutions and their clients. The goal is to create AI systems that are not only powerful but also transparent, accountable, and theoretically sound, enabling their seamless integration into regulatory frameworks and boosting confidence across the financial sector.

PROF MICHAEL WELLMAN, 2023

AI models are frequently trained on historical data that may contain inherent biases, leading to decisions that perpetuate or even exacerbate existing inequalities. This is particularly concerning in financial services, where biased algorithms can result in unfair outcomes, such as discriminatory credit scoring or loan approvals. A major research challenge is the development of robust techniques for bias detection, analysis, and mitigation in AI models. This includes exploring algorithmic fairness, adversarial debiasing, and the creation of synthetic datasets that better represent diverse populations. By advancing these methods, researchers can help ensure that AI systems in finance deliver fair and equitable outcomes, contributing to the broader goal of social justice in financial services.

As AI systems increasingly assume decision-making roles in finance, the ethical implications of their use become paramount. This challenge encompasses the need for comprehensive frameworks that guide the ethical deployment of AI in areas such as algorithmic trading, risk management, and customer interactions. The rapid evolution of AI technologies outpaces current regulatory frameworks, necessitating dynamic governance models that can adapt to new developments while maintaining ethical standards. Research in this area should focus on developing principles for AI ethics, creating

adaptive regulatory mechanisms, and ensuring that AI systems operate within the bounds of societal values and legal requirements. The establishment of these frameworks is essential for promoting the responsible and sustainable integration of AI into the financial sector, balancing innovation with public trust and accountability.

- **Enhanced Interpretability and Transparency:** Development of AI models that are not only accurate but also interpretable and grounded in finance and economics theories. This will enable transparent, accountable decision-making aligned with regulatory standards, fostering trust among stakeholders.
- **Mitigated Bias and Increased Fairness:** Creation of equitable AI models that detect and correct biases, ensuring fair outcomes across diverse populations. This will lead to more inclusive financial services, reducing disparities and promoting social justice.
- **Ethical Use and Dynamic Regulation:** Establishment of adaptive governance frameworks that ensure the ethical deployment of AI in finance. These frameworks will evolve with technological advancements, balancing innovation with public trust and regulatory compliance.

#### DATA REQUIREMENTS

- **Historical Financial Data:** Market prices, trading volumes, credit scores, transaction histories.
- **Textual Data:** Financial reports, news articles, analyst reports, customer communications.
- **Demographic and Socioeconomic Data:** Population group insights, bias detection.
- **Behavioral Data:** User interactions, transaction patterns, digital footprints.
- **Regulatory Data:** Laws, regulations, compliance guidelines.

#### AI METHODS REQUIREMENTS

- Explainable AI (XAI)
- Bias Detection and Mitigation
- Federated Learning & Differential Privacy
- Ethical AI Frameworks

#### RESOURCE REQUIREMENTS

2 H100 GPUs for 12 months.

### GRAND CHALLENGE 3: LLM AND MULTIMODAL AI EMPOWERED FINANCIAL SERVICES

AI agents<sup>5</sup> utilizing large language models (LLMs) are transforming the landscape of decision-making. The fusion of Large Language Models (LLMs) and multimodal AI represents a transformative opportunity for the financial services industry. LLMs, with their advanced Natural Language Processing (NLP) capabilities, can process and generate human-like text, making them invaluable for tasks such as customer interaction, compliance automation, and textual data analysis. When combined with multimodal AI—which involves integrating diverse data types such as text, images, audio, video, and structured numerical data—the potential for innovation is significantly amplified. Multimodal data mining allows for a more comprehensive understanding of complex financial scenarios by extracting insights that might not be apparent when analyzing each data type in isolation. This integrated approach can drive advancements in areas such as financial forecasting, fraud detection, portfolio optimization, and personalized financial advice, ultimately leading to more informed and effective decision-making within the financial sector.

#### OBJECTIVE

The integration of LLMs and multimodal AI into financial services presents several significant challenges that must be addressed to fully realize their potential. One of the primary challenges is the complexity involved in combining high-dimensional data from different modalities while ensuring accurate data representation and alignment. This complexity often leads to models that overfit the training data, resulting in poor generalization to unseen data—a critical issue in the dynamic and volatile financial sector. Additionally, LLMs, though powerful in processing textual data, are currently limited by their reliance on prompt engineering and their inconsistent outputs for similar prompts, raising concerns about reliability and consistency in financial applications. Another challenge is the need for real-time accuracy and adaptability, as financial markets are continuously evolving. Traditional LLMs struggle with maintaining up-to-date information and adapting to rapid market changes. To address these issues, FinAgent<sup>6</sup> has been proposed as a multimodal

foundation agent specifically designed for the finance domain. FinAgent incorporates a market intelligence module that processes a diverse range of data to provide accurate financial market analysis. Its unique dual-level reflection module enables rapid adaptation to market dynamics while integrating a diversified memory retrieval system, allowing the agent to learn from historical data and improve decision-making processes over time.

Building on the success of FinAgent, there is significant potential to extend LLM and multimodal AI agents into various financial scenarios by developing customized models that cater to specific financial sectors, such as banking, insurance, and investment companies, particularly within ASEAN countries. The financial environments in these regions are diverse, and a generic model may not suffice. Customized LLMs must incorporate localized financial knowledge, regulatory requirements, and cultural nuances, ensuring that the models are contextually relevant and effective. However, this customization introduces additional challenges in terms of data collection, model training, and evaluation. Creating an associated evaluation dataset that reflects the unique financial and economic conditions of these regions is crucial for ensuring the models' accuracy and reliability. This effort also demands careful consideration of the computational resources and infrastructure required to support such specialized models. While the potential benefits of these innovations are substantial, navigating the complexities of integrating advanced AI systems into existing financial infrastructure, ensuring regulatory compliance, and managing the computational demands will be essential to their success.

The successful implementation of LLM and multimodal AI technologies in financial services is expected to yield advanced AI agents capable of making accurate, efficient, and contextually aware decisions across a variety of financial scenarios, particularly those tailored to the unique needs of different financial sectors (e.g., banking, insurance, investment, and real estate.) and ASEAN countries. These agents would enhance operational efficiency by automating complex tasks, including compliance and risk management, while providing highly personalized financial services through the integration of diverse data sources. The customized LLMs, enriched with localized financial knowledge and regulatory insights, would improve fraud detection, financial forecasting, and portfolio management. Additionally, the development of a region-specific evaluation dataset would ensure that these AI agents are not only effective but also reliable and adaptable to the evolving financial landscape. Ultimately, the outcome would contribute to a more responsive, resilient, and regionally nuanced financial sector.

#### DATA REQUIREMENTS

To enable LLM and multimodal AI in financial services, access to a wide range of diverse datasets is essential, including:

- **Textual Data:** Financial news, official filings, and analysts' reports.
- **Numerical Data:** Financial asset prices, financial ratios, and economic indicators.
- **Alternative data in finance.**

#### AI METHODS REQUIREMENTS

Large language models.

#### RESOURCE REQUIREMENTS

Storage space, 16 H100 GPUs for 12 months.



Figure 2: Overview of FinAgent

## AI Methods and Data – Challenges and Opportunities

### GENERATIVE AI AND THEORY-GUIDED AI FOR FINANCIAL SERVICES

As AI technologies continue to evolve, Generative AI and Theory-Guided AI have emerged as powerful tools for modeling complex financial systems and enhancing decision-making processes. These approaches offer significant opportunities to improve financial simulations, risk assessment, and market predictions, but they also introduce challenges that need to be addressed to fully realize their potential in the financial sector.

- **Generative Adversarial Networks (GANs):** GANs, first introduced by Goodfellow et al. (2014)<sup>7</sup>, have become a cornerstone of generative AI, enabling the creation of realistic synthetic data that can mimic financial markets. In finance, GANs can generate synthetic financial time-series data, simulate rare market events, and provide robust testing environments for trading algorithms. A recent survey by Wiese et al. (2020)<sup>8</sup> in the financial domain explores the applications of GANs for generating financial data. However, challenges include ensuring the quality of the generated data, preventing mode collapse (where the model generates limited variations of the data), and balancing the trade-off between realism and computational complexity.
- **Diffusion Models:** Diffusion models, introduced by Sohl-Dickstein et al. (2015)<sup>9</sup>, have gained attention for their ability to generate high-quality data by

simulating the gradual diffusion process. In the financial sector, these models can be used for time-series generation and the creation of realistic market scenarios that capture the inherent uncertainties of financial systems. A recent study by Cao et al. (2023)<sup>10</sup> discusses the application of diffusion models in financial forecasting, emphasizing their potential to improve market predictions and scenario planning. The opportunity lies in the model's ability to produce fine-grained and controllable simulations. However, challenges include the high computational cost and the need for extensive tuning to capture the intricate patterns of financial data accurately.

- **Theory-Guided Neural Networks (TGNNs):** TGNNs integrate domain-specific financial theories into neural network architectures, as outlined by Karpatne et al. (2017)<sup>11</sup>. By embedding established financial principles into AI models, TGNNs offer the dual benefits of data-driven insights and theoretical grounding. This approach enhances model interpretability and ensures that predictions are consistent with known financial laws. The opportunity here is to improve the accuracy and reliability of financial models while providing clear explanations for model predictions. However, challenges include the complexity of integrating diverse financial theories into AI models and ensuring that these models remain flexible enough to adapt to new data and emerging financial paradigms.

### AI GOVERNANCE IN FINANCIAL SERVICES

As AI becomes deeply embedded in financial services, it brings both opportunities and significant challenges. The growing reliance on AI demands robust governance to ensure transparency, accountability, and ethical compliance. This section explores the essential aspects of AI governance in finance, focusing on the need for explainable AI, bias mitigation, and data privacy, to build a trustworthy and equitable financial ecosystem.

- **Explainable AI (XAI):** Ensuring that AI models in finance are interpretable is crucial for compliance and trust. Seminal works like LIME (Ribeiro et al., 2016)<sup>12</sup> and SHAP (Lundberg & Lee, 2017)<sup>13</sup> have laid the groundwork for XAI. Recent advances have focused on developing explainable deep models tailored for financial applications, addressing the unique complexities of financial data. However, challenges remain in balancing model accuracy with interpretability, particularly in high-stakes financial environments where decisions must be both explainable and reliable.
- **Bias in AI Models:** AI models, particularly those trained on historical data, can inherit biases that lead to unfair outcomes in financial services. Addressing these biases requires ongoing research into fairness-aware algorithms (Agarwal et al., 2018)<sup>14</sup> and the development of synthetic datasets designed to mitigate bias in financial applications. For example, Altman et al. (2024)<sup>15</sup> have explored the use of realistic synthetic financial transactions to enhance anti-money laundering models. The opportunity lies in creating more equitable financial systems, but the challenge is to ensure that bias mitigation strategies are robust and do not introduce new biases or worse performance.

- **Transparency and Accountability:** Regulatory frameworks must ensure that AI models are not only accurate but also transparent and accountable. This is particularly important in financial applications where decisions can have significant legal and financial implications. The work by Doshi-Velez & Kim (2017)<sup>16</sup> on the importance of interpretability in AI underscores the need for robust regulation. Additionally, emerging legal frameworks (Uzougbo et al., 2024)<sup>17</sup> are beginning to address the ethical considerations and accountability issues. The key challenge here is developing regulations that keep pace with rapid AI advancements while ensuring that these regulations are enforceable and do not stifle innovation.
- **Data Privacy and Security:** Ensuring data privacy while leveraging AI in finance is a significant challenge. Techniques such as federated learning (McMahan et al., 2017)<sup>18</sup> and differential privacy (Dwork et al., 2006)<sup>19</sup> offer potential solutions by allowing AI models to be trained on decentralized data without compromising privacy. However, implementing these techniques in financial services is complex and requires careful consideration of data governance policies. Zhang et al. (2021)<sup>20</sup> provide a comprehensive survey of these methods, highlighting both the opportunities for enhanced privacy and the challenges of ensuring data security.

---

## LLM AND MULTIMODAL AI EMPOWERED FINANCIAL SERVICES

---

The integration of Large Language Models (LLMs) and multimodal AI systems into financial services represents a significant leap forward. These technologies combine advanced NLP capabilities with the ability to process and integrate multiple data modalities—such as text, images, numerical data, and audio—into a unified model. However, while these advances hold great promise, they also present unique challenges.

- **Large Language Models (LLMs):** LLMs, including models like GPT-3 (Brown et al., 2020)<sup>21</sup> and BERT (Devlin et al., 2019)<sup>22</sup>, have revolutionized the generation of human-like text. In financial services, these models offer opportunities for automating report generation, sentiment analysis, and customer interaction. However, LLMs face significant challenges, including prompt dependency, overfitting, and the need for extensive computational resources. Additionally, LLMs often require fine-tuning, which introduces further complexities. Recent studies (e.g., Zhao et al., 2024)<sup>23</sup> provide a comprehensive review. A growing area of interest is the use of LLM agents in multi-agent collaboration, where multiple specialized LLMs interact to tackle complex tasks collaboratively, enhancing scalability and adaptability in financial applications. The challenge here lies in coordinating these agents to avoid conflicts and ensure coherent decision-making based on financial domain knowledge.

- **Multimodal AI:** The ability to integrate diverse data types into a single analytical framework allows for more comprehensive financial analyses. Techniques like cross-modal transformers (Tsai et al., 2019)<sup>24</sup> and contrastive learning (Radford et al., 2021)<sup>25</sup> have shown promise in effectively combining these modalities. In financial applications, multimodal AI can enhance risk assessment, improve fraud detection, and offer richer insights for portfolio management. The key opportunity here is the potential to uncover insights that would be missed if only a single data modality were considered. However, challenges include the increased complexity of models (and reduced transparency), the difficulty of aligning data from different modalities, and the need for robust pre-processing pipelines to handle noisy or incomplete data.
- **Time-Series Analysis:** Time-series models like LSTMs and newer attention-based models are crucial for forecasting in finance. When integrated with LLMs, these models can provide context-aware predictions that adapt to market changes. Innovations such as temporal fusion transformers (Lim et al., 2021)<sup>26</sup> can handle sequential data more effectively, offering enhanced performance in volatile market conditions. The opportunity lies in improving the accuracy and relevance of financial forecasts by leveraging the contextual understanding that LLMs provide. However, challenges include ensuring that these models remain responsive to sudden market shifts and avoiding the potential over-reliance on historical data, which may not always be indicative of future trends.

## Singapore's Role

Singapore is strategically positioned to lead the integration of artificial intelligence (AI) in financial services, not only within its own borders but across ASEAN, Asia, and the Middle East. This leadership is driven by a combination of advanced infrastructure, a rich talent pool, access to diverse financial data, and a forward-thinking regulatory environment.

---

### STRATEGIC INFRASTRUCTURE AND ECOSYSTEM

---

Singapore's robust digital infrastructure provides the foundation for its leadership in AI and finance. The country's state-of-the-art data centers, high-speed connectivity, and advanced cloud computing capabilities enable the processing and analysis of vast amounts of financial data in real-time. This infrastructure supports the development and deployment of sophisticated AI models that can be tailored to the dynamic and complex financial environments of ASEAN and broader Asian markets. Moreover, Singapore's strategic geographical location and its role as a global financial hub make it a natural gateway for AI-driven financial innovations to be disseminated across the Asia-Pacific and Middle Eastern regions. The country's well-established financial services sector, combined with its strong ties to international markets, ensures that AI innovations developed in Singapore are well-positioned for regional and global adoption.

---

### WORLD-CLASS TALENT AND ACADEMIC EXCELLENCE

---

Singapore's commitment to nurturing a deep pool of AI talent is another critical factor in its leadership. The nation's investment in education and research has cultivated a highly skilled workforce capable of pushing the boundaries of AI in finance. Singapore's National AI Strategy (NAIS 2.0) include programs such as the AI Visiting Professorship and the AI Accelerated Masters Programme for attracting top-tier researchers

and practitioners from around the world, creating a vibrant academic community that is deeply integrated with industry needs. This talent base is further supported by collaborations between Singapore's leading academic institutions and global research bodies. These partnerships foster the exchange of knowledge and expertise, driving innovation in AI applications for financial services. Singapore's ability to attract and retain top talent, coupled with its emphasis on continuous learning and development, ensures that it remains at the forefront of AI research and application.

---

### ACCESS TO DIVERSE DATA

---

Singapore's role as a financial hub grants it access to a vast array of financial data, crucial for building accurate and adaptable AI models. This diverse data pool, coupled with strong data governance policies, allows for the development of AI technologies that are both innovative and trustworthy, tailored to the specific financial challenges of ASEAN, Asia, and the Middle East.

---

### SUPPORTIVE REGULATORY FRAMEWORK

---

Singapore's regulatory environment is designed to balance innovation with stability. The proactive frameworks established by regulatory bodies like the Monetary Authority of Singapore (MAS) ensure that AI applications in finance are ethically sound and transparent. This regulatory leadership positions Singapore as a model for other countries in the region, promoting the adoption of best practices in AI governance.

Singapore's combination of infrastructure, talent, data access, and regulatory foresight positions it as a leader in developing AI models for financial services across ASEAN, Asia, and the Middle East. By focusing on regional needs, Singapore is not only advancing its own financial sector but also contributing to the global evolution of AI in finance.

## Conclusions

AI in financial services represents a transformative force with the potential to significantly enhance the efficiency, transparency, and resilience of the finance sector, particularly in Singapore. While challenges such as ensuring interpretability, mitigating bias, and maintaining ethical compliance remain, the pursuit of innovative

AI techniques, including generative AI, theory-guided models, and large language models, offers unprecedented opportunities. By addressing these challenges through rigorous research and practical applications, Singapore can solidify its position as a global leader in AI-driven financial innovation, ultimately benefiting both its local economy and the broader global financial community.

## REFERENCES

- 1 SORA. "General-purpose simulators of the physical world." OpenAI. Available at: <https://openai.com/index/sora/>.
- 2 Xie, Tao, et al. "Market-GAN: Adding Control to Financial Market Data Generation with Semantic Context." *AAAI Conference on Artificial Intelligence*, 2024.
- 3 Yoon, Jinsung, et al. "Time-series Generative Adversarial Networks." *Advances in Neural Information Processing Systems*, 2019.
- 4 Hao, Zheng, et al. "Reasoning with Language Model is Planning with World Model." *Conference on Empirical Methods in Natural Language Processing*, 2023.
- 5 Xi, Meng, et al. "The Rise and Potential of Large Language Model Based Agents: A Survey." *arXiv preprint arXiv:2301.01234*, 2023.
- 6 Zhang, Jing, et al. "A Multimodal Foundation Agent for Financial Trading: Tool-Augmented, Diversified, and Generalist." *Proceedings of the ACM SIGKDD International Conference on Knowledge Discovery & Data Mining*, 2024.
- 7 Goodfellow, Ian, et al. "Generative Adversarial Nets." *Advances in Neural Information Processing Systems* 27, 2014.
- 8 Wiese, Matthias, et al. "Quant GANs: Deep Generation of Financial Time Series." *Quantitative Finance* 20.9 (2020): 1419-1440.
- 9 Sohl-Dickstein, Jascha, et al. "Deep Unsupervised Learning Using Nonequilibrium Thermodynamics." *International Conference on Machine Learning*, PMLR, 2015.
- 10 Cao, Yiqing, et al. "Financial Forecasting with Diffusion Models." *arXiv preprint arXiv:2301.01234*, 2023.
- 11 Karpatne, Anuj, et al. "Theory-Guided Data Science: A New Paradigm for Scientific Discovery from Data." *IEEE Transactions on Knowledge and Data Engineering* 29.10 (2017): 2318-2331.
- 12 Ribeiro, Marco Tulio, Sameer Singh, and Carlos Guestrin. "Model-Agnostic Interpretability of Machine Learning." *arXiv preprint arXiv:1606.05386*, 2016.
- 13 Lundberg, Scott M., and Su-In Lee. "A Unified Approach to Interpreting Model Predictions." *Advances in Neural Information Processing Systems* 30, 2017.
- 14 Agarwal, Alekh, et al. "A Reductions Approach to Fair Classification." *International Conference on Machine Learning*, 2018.
- 15 Altman, Erik, et al. "Realistic Synthetic Financial Transactions for Anti-Money Laundering Models." *Advances in Neural Information Processing Systems* 36, 2024.
- 16 Uzougbo, Ngozi Samuel, Chinonso Gladys Ikegwu, and Adefolake Olachi Adewusi. "Legal Accountability and Ethical Considerations of AI in Financial Services." *GSC Advanced Research and Reviews* 19.2 (2024): 130-142.
- 17 Doshi-Velez, Finale, and Been Kim. "Towards a Rigorous Science of Interpretable Machine Learning." *arXiv preprint arXiv:1702.08608*, 2017.
- 18 McMahan, Brendan, et al. "Communication-Efficient Learning of Deep Networks from Decentralized Data." *International Conference on Artificial Intelligence and Statistics*, 2017.
- 19 Dwork, Cynthia. "Differential Privacy." *International Colloquium on Automata, Languages, and Programming*, Berlin, Heidelberg: Springer Berlin Heidelberg, 2006.
- 20 Zhang, Chen, et al. "A Survey on Federated Learning." *Knowledge-Based Systems* 216 (2021): 106775.
- 21 Brown, Tom B., et al. "Language Models are Few-Shot Learners." *Advances in Neural Information Processing Systems* 33, 2020.
- 22 Devlin, Jacob, et al. "BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding." *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, 2019, pp. 4171-4186.
- 23 Zhao, Huaqin, et al. "Revolutionizing Finance with LLMs: An Overview of Applications and Insights." *arXiv preprint arXiv:2401.11641*, 2024.
- 24 Tsai, Yao-Hung Hubert, et al. "Multimodal Transformer for Unaligned Multimodal Language Sequences." *Proceedings of the Conference of the Association for Computational Linguistics*, 2019.
- 25 Radford, Alec, et al. "Learning Transferable Visual Models from Natural Language Supervision." *International Conference on Machine Learning*, PMLR, 2021.
- 26 Lim, Bryan, et al. "Temporal Fusion Transformers for Interpretable Multi-Horizon Time Series Forecasting." *International Journal of Forecasting* 37.4 (2021): 1748-1764.

# APPENDIX XI. ELECTRONICS/ SEMICONDUCTORS



## AUTHORS:

Dr J. Senthilnath  
(I<sup>2</sup>R, A\*STAR)

Dr Li Xiaoli  
(I2R, A\*STAR)

Prof Yeo Yee Chia  
(NUS-IME, A\*STAR)

Prof Aaron Thean  
(NUS)

Prof Nagarajan Raghavan  
(SUTD)

Dr Sridhar Narayanaswamy  
(IHPC, A\*STAR)

## Executive Summary

The semiconductor industry, **contributing 7% to Singapore's GDP and accounting for 44% of its manufacturing output**, is a *cornerstone of the nation's economic vitality*. Integrating artificial intelligence (AI) into this sector significantly amplifies its capabilities, leading to smarter design processes, more efficient production, and enhanced defect detection. This fusion of AI and semiconductor technology could establish Singapore as a global leader in technological innovation and sustainable development. Given its critical importance, it is essential for government funding to support AI advancements in semiconductor design and manufacturing. Such investments not only push the boundaries of scientific knowledge and foster innovation but also address urgent global challenges. Strategic funding in this pivotal sector can elevate Singapore to the forefront of semiconductor manufacturing, driving economic growth and improving quality of life worldwide.

This whitepaper elaborates on the significance of AI in semiconductor design and manufacturing, underscoring the need

for targeted government research funding. It explores the intersection of AI with the semiconductor industry, highlighting key scientific challenges and potential solutions in both research and development (R&D) and manufacturing. Institutions such as A\*STAR, NRF, and local universities are well-positioned within Singapore's ecosystem, from foundational research to applied manufacturing, with AI-for-Science grants serving as an initial catalyst.

Addressing major scientific challenges in AI for semiconductors involves: i) **AI-guided Semiconductor Manufacturing Processes:** Integrating digital twins, physics-guided simulations, and design of experiments (DoE) to generate data and optimize manufacturing processes, ii) **AI-enabled 3D-IC packaging:** Using AI techniques to optimize 3D-IC package configuration and enhance manufacturing outcomes, iii) **AI-accelerated Semiconductor Failure Analysis:** Employing sensor analytics and visual inspection for rapid and precise defect prediction and localization.

## Introduction

The semiconductor industry has leveraged artificial intelligence (AI) to swiftly tackle challenges in both research and development as well as manufacturing. Figure 1 illustrates the schematic representation of the semiconductor manufacturing process, showcasing the transition from real data to virtual data and the application of algorithms to scientific advancements.

Current data reveals that the number of nanometer-scale transistors approximately doubles every two years, in line with Moore's Law. This rapid growth underscores the urgent need for innovative approaches in semiconductor manufacturing processes. Conventional methods face increasing challenges at various levels—from micro (atomistic) to macro (device-level)—including process optimization for semiconductor fabrication<sup>1</sup> and thin films<sup>2</sup>, packaging of integrated circuits (2D/2.5D/3D)<sup>3</sup>, and multi-level failure analysis (device-level<sup>4,5</sup> and circuit-level<sup>6</sup>). Additionally, the rising complexity of integrating diverse materials into nanoscale

structures exacerbates issues related to defect formation and isolation. Similarly, in the AI domain, the exponential growth in model parameter sizes each year results in increased modeling complexity, necessitating advanced computational resources.

Exploring AI for semiconductor manufacturing not only accelerates the process but also enhances production efficiency by circumventing traditional trial-and-error methods<sup>7</sup>. The advancement of semiconductors depends on discovering efficient materials and optimizing fabrication processes [1]. AI's effectiveness is significantly increased when it learns from diverse data distributions. In certain scenarios, AI models must make predictions with limited semiconductor wafer resources or conduct experiments without compromising accuracy and performance. Furthermore, collaboration between humans and machines is essential for advancing semiconductor process development, with human experts providing critical guidance and insight<sup>8</sup>.

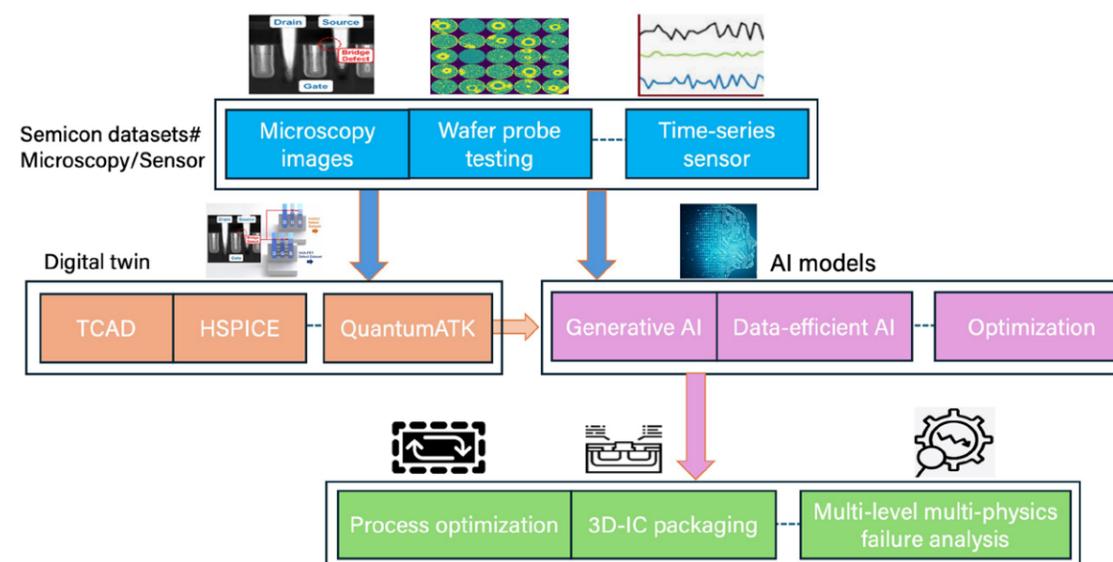


Figure 1: A schematic representation for the semiconductor manufacturing process.

## Background

On 6<sup>th</sup> May, 2024, A\*STAR organized a workshop on AI for semiconductor manufacturing. The event, led by Dr J. Senthilnath (I<sup>2</sup>R ASTAR) and Dr Li Xiaoli (I<sup>2</sup>R ASTAR), attracted over 150 participants, with a dynamic mix of in-person and online attendees. With nine expert speakers presenting, Prof Yeo Yee Chia (NUS-IME ASTAR) on AI for semiconductor R&D and manufacturing, Dr Li Xiaoli (I<sup>2</sup>R ASTAR) on AI for semiconductor DoE and failure analysis, Prof Aaron Thean (NUS) on physics-guided AI defect isolation in microelectronics, Ms Bernice Zee (AMD) on the application of machine learning in semiconductor failure analysis, Prof Nagarajan Raghavan (SUTD) on AI-enabled foundry yield estimation, Prof Chen Zhi Ning (NUS) on AI reshapes RF and EM device technology, Dr Jason Png Ching Eng (IHPC ASTAR) and Dr Sridhar Narayanaswamy (IHPC ASTAR) on physics guided AI for design and analysis of semiconductor: devices, packaging and manufacturing, and Dr J. Senthilnath (I<sup>2</sup>R, ASTAR) on AI-guided microscopy for semiconductor. The workshop delved into the intersection of AI and semiconductor design and manufacturing, highlighting key scientific challenges and potential solutions.

The workshop discussed the challenges faced by traditional methods in semiconductor manufacturing, and the potential for innovative AI tools to tackle these complex problems. The discussion emphasized the crucial role of data curation in training AI models. It was highlighted that curating the right datasets is essential for developing AI solutions that can contribute to breakthroughs in semiconductor

manufacturing. A structured approach to data collection and curation across academia and industry was recognized as necessary.

Robust technological infrastructure was noted as crucial to support AI initiatives in Singapore. This infrastructure must be capable of handling complex computations and large volumes of data to meet the demands of advanced AI applications in semiconductor manufacturing.

Integrating physics constraints into AI models was proposed as a way to improve semiconductor scanning techniques. By embedding these constraints, AI models can become more efficient and accurate in simulating and predicting semiconductor behaviors. Developing efficient surrogate models was also emphasized to speed up simulation processes without sacrificing accuracy.

The workshop also conducted roundtable discussions on various topics including i) AI-guided semiconductor manufacturing process, ii) AI-enabled yield estimation, and iii) AI-accelerated semiconductor failure analysis. The workshop concluded by emphasizing the importance of increased collaboration between academia and industry, as well as the need for funding in research and development projects addressing identified challenges and technological needs in Singapore. The goal is to drive innovation and collaborative efforts in the semiconductor industry, particularly within the context of Singapore.

## Grand Challenges

In the realm of semiconductor design and manufacturing, it is vital to address several sub-challenges related to semiconductor design and manufacturing.

---

### GRAND CHALLENGE 1: AI FOR SEMICONDUCTOR DESIGN AND MANUFACTURING

---

**Process parameters:** (a) Discovery and Development – The identification of superior semiconductors requires exploration of a vast parameter space. This demands innovative methods for material discovery and development, (b) Chemical Mechanical Polishing (CMP) – It is crucial to fine-tune CMP parameters to enhance uniformity across different thickness ranges. This includes adjustments for within-wafer (WiW), wafer-to-wafer (WTW), and wafer extreme edge (WEE) variations.

**IC packaging:** (a) 3D packaging – The 2D IC packaging approach has limitations in integration, power consumption, and reliability. On the contrary, 2.5D IC packaging offers improvements over 2D packaging. However, it still necessitates a further reduction in interconnect lengths and power consumption through 3D packaging, (b) 3D inspection – The existing method mainly relies on 2D Fast Fourier Transform (FFT) to analyze current flow patterns. To circumvent this there is a need to build a 3D to understand current flow patterns.

**Failure analysis:** (a) Time Constraint – Depending on the nature of the failure, the failure analysis (FA) process can necessitate expert evaluation for weeks to months, (b) Data Imbalance – There is often a disproportion in the number of non-defective samples compared to defective ones, requiring oversampling of the defect class to achieve class balance.

### OBJECTIVE

During the workshop, we organized three breakout sessions to understand the impact of the selected topics, and participants from both the academy and industry actively shared their views. This white paper proposes a unified strategy to improve the overall semiconductor design and manufacturing process in,

i) **AI-guided Semiconductor Manufacturing Processes:** Integrating digital twins, physics-guided simulations, and design of experiments (DoE) to generate data and optimize manufacturing processes, ii) **AI-enabled 3D-IC packaging:** Using AI techniques to optimize 3D-IC package configuration and enhance manufacturing outcomes, iii) **AI-accelerated Semiconductor Failure Analysis:** Employing sensor analytics and visual inspection for rapid and precise defect prediction and localization.

Expected outcomes in Semiconductor Manufacturing Processes:

i) **Enhanced Speed:** Achieving up to 10 times faster optimization in overall process efficiency, integrated circuit (IC) packaging, and defect localization.

ii) **Design Discovery:** Identifying candidate semiconductor designs with a success rate exceeding 40%, based on experimental or analytical validations.

iii) **AI-Assisted Process Tuning:** Utilizing AI to fine-tune process parameters, resulting in more stable yields with variance reduced to less than 0.5%, even amidst noise.

### DATA REQUIREMENTS

Data Management Systems – Robust databases and data management platforms for storing, organizing, and analyzing semiconductor manufacturing data, including tools for data integration and sharing.

## SEMICONDUCTOR DATABASE MANAGEMENT AND MODEL SHARING

*i) Open Datasets:* Several publicly available semiconductor datasets are valuable for research and development: *i) Time-Series Datasets:* These include manufacturing operation data and semiconductor quality data collected from various sensors, such as, UCI SECOM dataset<sup>9</sup>, Wafer dataset<sup>10</sup>, and Chemical Mechanical Planarization (CMP) of wafer dataset<sup>11</sup>, *ii) Single/Multi-Label Defect Datasets:* During wafer fabrication, different types of defects can co-occur on a single wafer. Examples include Mixed-type defects dataset<sup>12</sup>, *iii) Visual Inspection Datasets:* These datasets support defect localization and type classification, such as, PCB dataset<sup>13</sup>.

*ii) Industrial datasets:* Semiconductor companies often hesitate to share confidential data due to its sensitive nature. Effective data collection and curation are crucial for developing AI models that generalize well. To address these challenges, consider the following solutions: *i) National Semiconductor Translational and Innovation Center* could provide data for Singapore ecosystem to drive the Co-innovation, *ii) Integration of Digital Twin and Simulation Data:* Using digital twins and simulation data can help overcome limitations associated with experimental data, leading to more robust AI-driven solutions, *ii) Data Augmentation:* AI methods can augment existing data to mitigate limitations posed by the scarcity of real datasets, *iii) Federated Learning and Privacy-Preserving Analytics:* Exploring federated learning and privacy-preserving data analytics can facilitate the sharing of parameters and models rather than raw data, helping to establish a comprehensive semiconductor database or shared models within Singapore.

### AI METHOD REQUIREMENTS

In order to accomplish the previously mentioned objectives, addressing limited data issues and expected outcomes, it is necessary to develop advanced AI methods. These methods include the following:

**Data augmentation:** Use reinforcement learning to search for a better data augmentation process.

**Generative oversampling:** Address the imbalance where non-defect samples are more abundant than defect samples by oversampling the defect class.

**Transfer learning:** Transfer knowledge from a source to a target domain with varying distributions, such as from simulation to experimental datasets or from current generation components to next generation components.

**Online learning:** Enable AI methods to adapt and generalize when learning from less data and when the distribution of the defect class varies.

**Generative inverse design:** Create realistic surrogate physical simulators.

### COMPUTE REQUIREMENTS

Computational Infrastructure – GPUs and cloud-based resources to handle large-scale data processing, AI model training, and simulations. The proposed computational resource requirements for large scale computational tasks are as follows: EPYC 7702P CPU – 3 M core hours, A100/V100 GPU – 0.05 M card hours, Storage – 2x 4 TB, RAM – 256 GB and Implementation – python (pytorch and other necessary libraries).

### OTHER REQUIREMENTS

*i) Human Resources:* (a) Domain Experts – Specialists in semiconductor design, manufacturing processes, and AI/data science to guide project development and ensure accurate implementation, (b) Data Scientists and Machine Learning Engineers – Professionals skilled in developing and applying machine learning models, including generative models and AI-driven optimization techniques, (c) Research Scientists: Experts to contribute to fundamental research and development of new technologies, such as advanced microscopy techniques and AI integration with physics laws.

*ii) Financial Resources:* (a) Funding – Grants and investments to support research and development activities, purchase of equipment, and hiring of specialized personnel. (b) Budget for Equipment – Investment in advanced microscopy tools, digital twin technologies, and other critical hardware necessary for semiconductor design and manufacturing.

*iii) Collaborative Resources:* (a) Partnerships – Collaboration with academic institutions, research centers, and industry partners to leverage expertise, share resources, and drive co-innovation, (b) Industry Networks – Engagement with semiconductor industry networks to stay updated on emerging trends, technologies, and best practices.

## AI METHODS AND DATA – CHALLENGES AND OPPORTUNITIES

### AI METHODS

Emerging data-driven AI methods present promising alternatives to the traditional trial-and-error approach. Addressing challenges related to semiconductor database management and model sharing is crucial, particularly in data-scarce environments and for preserving privacy.

*i) Generative inverse design:* This approach involves using generative models such as Variational Autoencoder (VAE), Generative Adversarial Network (GAN) and Diffusion Models (DM) to discover and design semiconductors with desired properties or to predict properties by exploring the search space. Although generative inverse design is still under development and has not yet surpassed the accuracy of first-principle calculations [8], it holds significant potential for advancing semiconductor design.

*ii) Integration of Physics Laws into AI Models:* Embedding fundamental physics laws and constraints into AI models can enhance the accuracy and efficiency of simulating semiconductor behaviors. For instance, in magnetic field microscopy when current flows through a sample, changes in resistance induced by the local magnetic field are detected using physics law such as Ampere's circuital law. This approach enhances the reliability of predictions and optimizations, avoiding the pitfalls of relying solely on AI models that may generate unreliable data<sup>14</sup>.

*iii) Development of AI Surrogate Models:* Advanced AI models, such as deep learning, have the potential to transform traditional semiconductor design and manufacturing methods. They are particularly valuable in scenarios where traditional simulations or experiments are too costly or time-consuming. Efficient surrogate models<sup>15</sup> can expedite

simulation processes while maintaining high accuracy, thereby streamlining operations and reducing time-to-market.

### VISION AND TRANSFORMATIVE ASPECTS

*i) Technical role:* The impact of AI on semiconductor design and manufacturing is limited when projects lack expertise in either the semiconductor domain or AI/data science. This gap often leads to challenges in implementing solutions due to insufficient domain knowledge, physics understanding, or experience, which can result in excessive data requirements, slow processes, and inflated team sizes. However, projects that integrate expertise in both AI and semiconductor domains can effectively implement innovative solutions. Leveraging domain knowledge, physics, and experience allows for optimal data utilization, algorithm development, and efficient implementation.

*ii) Scientific data:* a) **National Semiconductor Translational and Innovation Center (NSTIC):** NSTIC aims to provide valuable data to the Singapore ecosystem, fostering co-innovation, b) **Big Data in Semiconductor Manufacturing:** Generating large volumes of metrology data, such as After-Develop Inspect (ADI) and After-Etch Inspect (AEI), enhances tool-to-tool and chamber-to-chamber matching, process stability, and tool productivity. This data helps reduce lithography overlay errors, enables fast analysis of statistical process control charts, ensures line stability, improves yield, and supports preventive maintenance.

*iii) Next Era of Microscopy:* An intelligent scanning scheme that employs multi-source microscopy can harness the strengths of various microscopy types. For example, high-end microscopy offers superior image quality but requires extensive scanning time. In such cases, a hybrid scanning sensor can quickly locate defects, which can then be precisely identified using high-resolution microscopy scanners.

*iv) Advanced AI models:* Generative AI, data-efficient learning, resource-efficient learning, and optimization techniques hold significant potential for transforming semiconductor design and manufacturing. Applying these advanced AI models to areas such as process optimization, 3D-IC packaging, and failure analysis could provide groundbreaking advancements.

## SINGAPORE'S ROLE

Singapore plays a crucial role in the global semiconductor supply chain, producing **10% of all chips** and **around 20% of semiconductor manufacturing equipment**. The country hosts a diverse semiconductor ecosystem with research and development (R&D) and manufacturing activities from IC design, wafer fabrication, packaging and testing<sup>16</sup>. The semiconductor industry accounts for **7% of Singapore's GDP and 44% of its manufacturing output**, making it a cornerstone of the nation's economic strength<sup>17,18</sup>.

Out of the top 15 semiconductor firms, around 9 have established operations in Singapore. These firms cover all aspects of the industry, including IC design, assembly, packaging and testing, wafer fabrication,

## Conclusions

This whitepaper addresses a pivotal challenge at the intersection of AI and semiconductor design and manufacturing. Its impact on the Singapore ecosystem is broad and influential, providing guidance to researchers, informing industry needs, and fostering the advancement of smarter, more efficient, and sustainable semiconductor manufacturing processes. AI is already accelerating the semiconductor R&D and manufacturing process, primarily through process optimization, and there is limited integration of multi-level physics for failure analysis. Our strategy aims to utilize Singapore's world-leading cohort studies and datasets, integrating them with open database and new cohorts that can be collected,

and equipment/raw material production for the semiconductor value chain. Expanding on the significance of AI for Semiconductor Design and Manufacturing, the workshop emphasized its relevance for government research funding. The semiconductor R&D and industry play a pivotal role in our nation's economy, contributing significantly to our GDP and dominating our electronics and manufacturing output. Recognizing its importance, government funding in AI for semiconductor design and manufacturing becomes crucial. Such investment not only advances scientific knowledge and fosters innovation, but also addresses pressing global challenges. By strategically investing in this critical sector, the government can propel our country to the forefront of the semiconductor industry, thereby driving economic growth and enhancing the quality of life for people worldwide.

resulting in the development of an AI model for all of manufacturing industries in Singapore. By employing a strategy of data collection, data integration, and solution finding to address a large portion of Singapore's economy through the development process as well as the maintenance. AI for semiconductor R&D requirements of this initiative are immense, necessitating the development of new technology. Despite the intense research agenda, the semiconductor industries ensure that research is translated and that outcomes are always measured in terms of improved yield and reduced costs among the production line as a result of deployed and scaled interventions.

## ACKNOWLEDGEMENTS

We would like to express our gratitude to the I<sup>2</sup>R team (Dr Rajdeep, Dr Ji Wei, Dr Ashish, Dr Zhuoyi, Dr Aye Phyu, and Dr Chuan Sheng), IME team (Dr Kart Leong, Dr Hyunsu, Dr Navab), Dr Riko (IMRE), Ms Bernice (AMD), Prof Zhining (NUS), Dr Jason (IHPC), and Dr Shengkai (NMC) for the valuable discussion during the workshop and drafting of this white paper. We also extend our thanks to the workshop participants from both the academic and industry sectors, who actively engaged in the roundtable discussions.

## REFERENCES

- 1 Xie, J., Zhou, Y., Faizan, M., Li, Z., Li, T., Fu, Y., Wang, X. and Zhang, L., 2024. Designing semiconductor materials and devices in the post-Moore era by tackling computational challenges with data-driven strategies. *Nature Computational Science*, pp.1-12.
- 2 Dutta, R., Tian, S.I.P., Liu, Z., Lakshminarayanan, M., Venkataraj, S., Cheng, Y., Bash, D., Chellappan, V., Buonassisi, T. and Senthilnath, J., 2022. Extracting film thickness and optical constants from spectrophotometric data by evolutionary optimization. *Plos one*, 17(11), p.e0276555.
- 3 Lau, J.H., 2022. Recent advances and trends in advanced packaging. *IEEE Transactions on Components, Packaging and Manufacturing Technology*, 12(2), pp.228-252.
- 4 Pan, J., Low, K.L., Ghosh, J., Senthilnath, J., Ferdous, M.M., Lim, S.Y., Zamburg, E., Li, Y., Tang, B., Wang, X. and Leong, J.F., 2021. Transfer learning-based artificial intelligence-integrated physical modeling to enable failure analysis for 3 nanometer and smaller silicon-based CMOS transistors. *ACS Applied Nano Materials*, 4(7), pp.6903-6915.
- 5 Zhou, B., Jieming, P., Sivan, M., Thean, A.V.Y. and Senthilnath, J., 2023, June. Quantile Online Learning for Semiconductor Failure Analysis. In *ICASSP 2023-2023 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)* (pp. 1-5). IEEE.
- 6 Yang, J., Wang, B., Wu, Y. and Ivanov, A., 2005. Fast detection of data retention faults and other SRAM cell open defects. *IEEE Transactions on computer-aided design of integrated circuits and systems*, 25(1), pp.167-180.
- 7 Chen, Y.L., Sacchi, S., Dey, B., Blanco, V., Halder, S., Leray, P. and De Gendt, S., 2024. Exploring Machine Learning for Semiconductor Process Optimization: A Systematic Review. *IEEE Transactions on Artificial Intelligence*.
- 8 Kanarik, K.J., Osowiecki, W.T., Lu, Y., Talukder, D., Roschewsky, N., Park, S.N., Kamon, M., Fried, D.M. and Gottscho, R.A., 2023. Human-machine collaboration for improving semiconductor process development. *Nature*, 616(7958), pp.707-711.
- 9 [Online]. Available: <https://www.kaggle.com/datasets/paresh2047/uci-semcom/code>
- 10 [Online]. Available: <https://www.timeseriesclassification.com/description.php?Dataset=Wafer>
- 11 [Online]. Available: <https://www.phmsociety.org/events/conference/phm/16/datachallenge>
- 12 [Online]. Available: <https://www.kaggle.com/datasets/co1d7era/mixedtype-wafer-defect-datasets>
- 13 [Online]. Available: <https://www.kaggle.com/datasets/akhatova/pcb-defects/data>
- 14 Gibney, E., 2024. AI models fed AI-generated data quickly spew nonsense. *Nature*, 632(8023), pp.18-19.
- 15 Rautela, M., Senthilnath, J., Huber, A. and Gopalakrishnan, S., 2022. Toward Deep Generation of Guided Wave Representations for Composite Materials. *IEEE Transactions on Artificial Intelligence*, 5(3), pp.1102-1109.
- 16 [Online]. Available: <https://www.edb.gov.sg/en/business-insights/insights/what-makes-singapore-a-prime-location-for-semiconductor-companies-driving-innovation.html>
- 17 [Online]. Available: Written reply to PQ on Singapore's semiconductor manufacturing industry (mti.gov.sg)
- 18 [Online]. Available: Semiconductors play major role in Singapore manufacturing decline and recovery | Plant Engineering

# APPENDIX XII. HYBRID QUANTUM COMPUTING



## AUTHORS:

Assoc. Prof. Dario Poletti  
(NUS Centre for Quantum Technology  
- CQT, SUTD)

Dr. Kishor Bharti  
(A\*STAR IHPC)

Dr. Ye Jun  
(A\*STAR Quantum Innovation Centre, IHPC)

Dr. Patrick Rebentrost  
(NUS CQT)

## Executive Summary

In the following, we are going to describe how AI can support the development of quantum computing, and also how quantum computing can further accelerate AI and could bring

it to deal with more complex systems in a more accurate, interpretable and energy effective manner.

## Introduction

Hybrid quantum computing is often used to describe a particular variational optimization procedure. In this case one aims to find the minimum of a complex function which would require significant time to evaluate on a classical computer. One would then evaluate it using a quantum processor, and then feed the outcome of this computation to a classical computer which will guide the choice of parameters to find the minimum of this function [Peruzzo2014]. As a number of problems can be turned into this computational framework, e.g. finance, logistics etc, this approach is considered by a portion of the community as the most promising for an application of quantum computing in the short term [Bharti2022]. To mention one highlight, recently 6400 nodes of the Fugaku supercomputer assisted a Heron quantum processing unit with up to 77 qubits and 10570 gates to simulate the  $N_2$  molecule [RobledoMoreno2024].

It is however a path which has uncertainties as two main problems are still affecting this direction: barren plateaus [McClean2018] and the amount of data required [Schuld2021]. Furthermore, there is yet no proven quantum advantage in using hybrid quantum computing. Nonetheless, the performance of this type of computations can still be significantly accelerated by improving the integration between classical and quantum components.

A sizeable portion of the community considers that it would be beneficial for Singapore's research agenda to consider a broader connotation for hybrid quantum computing. It is in fact hard to conceive that a useful quantum processor could function without the concurrent use of a classical computer, independent of the type of computation done on the quantum processor, for example to control the setup, to compile the codes in machine language, for the noise reduction

techniques and so on. In short, every computation we can currently envision will be, in this more general sense, a hybrid classical-quantum computation. Furthermore, the pace of advances in quantum computing hardware in the past few years is impressive, even with a recent realization of 48 logical qubits [Bluvstein2024], although still with limited practical use. We thus recommend that the community in Singapore refers to hybrid quantum computing in this more general sense.

Embracing the notion that quantum computers will rely on classical computers, it is useful to state the differences between classical and quantum computers.

Quantum algorithms rely on a completely different type of approach compared to classical ones, as interference and entanglement play a crucial role, and measurements are probabilistic. It thus requires a completely different type of thinking.

## Background

We organized a workshop to discuss the main research directions that should be studied and how AI can accelerate them, while also focusing on how quantum computing can boost the progress of AI. We had experts from academia, research institutes and industry. In the morning, the speakers addressed the current difficulties and potential advantages of integrating classical and quantum computing, how to proceed towards full-stack solutions (from simple codes agnostic of the hardware or whether the computation is classical or quantum, to the controls of the specific hardware), the advantages of different experimental implementations, and the potential of quantum-inspired algorithms whose exploitation is at its infancy.

Classical computers are considered to work better than quantum ones for problems with large amount of data and simpler computations, while quantum processors can have significant speed-up compare to classical ones for problems with few data and complex computations. A key aspect is that solving certain problems requires different computational complexity on classical and quantum processors. For some problems, the complexity of a quantum algorithm is up to exponentially smaller compared to the best-known classical algorithm, as for example, for Shor's algorithm [Shor1999] or the quantum linear systems algorithm for sparse matrices [Harrow2009]. Furthermore, quantum systems of  $L$  qubits are described by a vector with  $2^L$  complex numbers, while  $L$  bits are described by a binary vector of size  $L$  [Nielsen2000], and in that sense a quantum computer with  $L$  qubits processes vectors of much higher dimension compared to a classical computer with  $L$  bits. Another important difference is the processing speed of single operations, where classical processors outpace by orders of magnitudes the quantum ones, which results in latencies.

The afternoon of the workshop was organized in breakout sessions focused on near-term applications for AI in quantum computing and the bottlenecks to its integration. Topics included error correction, optimal control, and the potential for AI to design more efficient quantum algorithms. The discussion highlighted a few near-term directions such as application in many-body quantum physics and material science, followed by chemistry and biology.

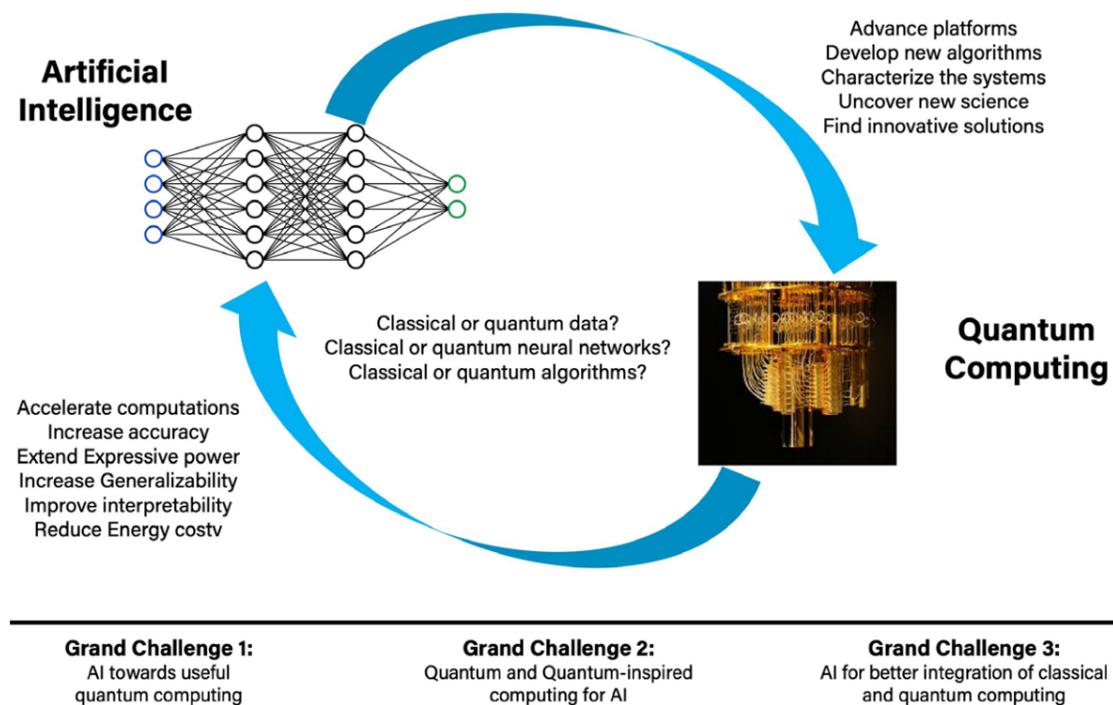


Fig.1 showing both directions of the expected outcomes and highlighting the three challenges.

## Grand Challenges

We highlight three Grand Challenges:

- how to have useful quantum computers.
- how to improve AI with quantum and quantum-inspired computing.
- how to integrate classical and quantum computing systems.

### USEFUL QUANTUM COMPUTERS

There are computations that can be exponentially or polynomially accelerated by quantum algorithms. These include linear algebra solvers [Harrow2009], Fourier transform [Coppersmith1994], and search in unsorted sets [Grover1996, Grover1997], which are pervasive to countless everyday applications. Quantum computers are also

the only viable way to digitally simulate quantum devices/materials/systems which is at the essence for the development of quantum technology. In the shorter term, applications can be foreseen first in the fields of physics, e.g. simulating strongly correlated materials and their dynamics, or quantum sensors leveraging on many-body quantum systems, and of chemistry, e.g. evaluating electro-chemical properties of molecules, or simulating their formation dynamics, and also biology in which larger molecules interact with each other giving rise to biological processes. This, in subsequent years, could also lead to applications in pharmacology and engineering. Optimization is also a very important area of research, with pervasive applications from finance, to logistics and more, in which quantum computing can bring acceleration.

The main obstacle in this grand challenge is noise. Noise reduces the coherence and entanglement in the system, brings in errors in the preparation, modification, and measurement of the quantum system. It has origins that depend on the materials used, on the fabrication of the components, and their utilization. In all these aspects, AI can potentially help in obtaining significant reductions of noise. For instance, recently in [Bausch2023] the authors used a transformer-based neural network to decode a surface code. An alternative to surface codes could be quantum low-density parity check codes [Breuckmann2021] whose error correction can be boosted by AI to achieve error rates more than 10 orders of magnitude smaller than the physical error rate.

There are also other bottlenecks in which AI could help the progress, for example the ability to simulate large quantum systems, which can be used, for instance, to benchmark, characterize and certify how computations are performed on quantum processors, or how to best implement sets of computations.

Another important bottleneck is the quantum algorithms available. As mentioned earlier, quantum algorithms rely on profoundly different principles compared to classical ones, and in general we have not developed a clear way of obtaining new algorithms and possibly new basic components such as linear algebra and Fourier transform. An AI could develop new approaches, for example using more local operations compared to many large entangling operations, which we would have not considered in a somewhat similar way to how AlphaGo has produced a completely new style of playing Go [Silver2016, Silver2017]. For example, an AI could propose codes which obtain similar performance while providing saving of quantum gates.

A major difficulty is also the certification of the quantumness of the computation and the characterization of the behavior of the system. For the latter, the use of neural networks can help simulate ever larger and deeper circuits [Carleo2017, Jonsson2018].

It should not be discarded that each implementation of quantum processors has its own specific bottlenecks and difficulties, on top of dealing with noise. For example, how to accelerate the gate operations, how to implement preparation and measurements in a more efficient manner, how to increase the scale while still being able to control them accurately and possibly reducing the size of the overall setup and increase the stability of the system. Last but not least, how to produce materials with less impurities so that noise can be reduced significantly. These are also places in which AI can accelerate progress in a way specific to the platform used.

### OBJECTIVE

Singapore could invest in the development and use of AI in works that cover the full spectrum from a better hardware, e.g. materials and fabrication, to middleware and algorithms. Many of these works are largely multidisciplinary, e.g. quantum computer design researchers working together with material and nano-fabrication scientists, or the use of machine learning for simulating the dynamics of quantum processors, while some can be very focused on a single, yet possibly impactful aspect, e.g. a novel quantum algorithm.

### DATA REQUIREMENTS

Unlike applications of AI in the classical domain, there is significantly less data availability in the quantum one. Starting from the materials used in each implementation, AI can help reduce the level of noise and improve the controllability of the system. Data about noise level in different conditions of preparation of the setup and use needs to be produced and standardized so that AI can be applied. Data on the effect of control protocols in the presence of noise also needs to be produced and made available to the community. For error correction, we can rely, at the small scale, on simulated data, yet for larger setups it would be necessary to collaborate with experimentalists and industries to have access to data which can be used by the community to uncover novel error correction techniques using AI. For compilers, one could

partially tap on simulations, for example with tensor networks, however new tools should be developed to find ways to simulate deeper circuits with long-range couplings. The advantage of AI for novel algorithms could be found with computation on small size systems first, for which data could be produced with different algorithms to solve similar tasks. However, similarly to compiling, the greatest advantage could occur in the case of large and deep circuits for which data is currently not available.

To represent the many-body wave-function of a quantum computation with neural networks, one can tap on unsupervised learning, for which no data is needed. However, for realistic, noisy, scenarios, one would still need to complement this process with data from experiments.

#### **AI METHOD REQUIREMENTS**

Reinforcement learning is the most natural approach for this grand challenge, together with the use of generative models such as transformers. For instance, reinforcement learning is ideal to improve the control of the quantum setup, and generative models could provide novel solutions for compilers or algorithms.

Neural networks can be used also to represent the wave-function of the system [Carleo2017] and for this, convolutional neural networks have shown to be effective, however the best models include information on the symmetries of the systems studied.

#### **COMPUTE REQUIREMENTS**

Between research on materials, noise mitigation, error correction, compilers, algorithms and neural network quantum states, 8 H100 cards for 36 months would be needed. However, some of this, especially the research on the materials, could be shared with other main directions of AI for Science.

---

### **IMPROVE AI WITH QUANTUM AND QUANTUM-INSPIRED COMPUTING**

---

In the last few years AI has progressed in a way that was unthinkable by the lay people and some of the experts in the field. However, it still suffers from some important problems: the lack of high accuracy performance, reduced interpretability, limited performance on quantum systems and large energy cost.

For instance, an accuracy of 99%, which is already difficult to achieve in several tasks, is insufficient for deployment of AI in systems like traffic control, management of supply chain or for scientific applications like characterization of systems' properties. In these cases, either the use of AI is not safe, or it is below the state of the art provided by traditional methods.

The output of neural networks is also notoriously difficult to interpret and thus still difficult to trust in a way that we are comfortable with [Chakraborty2017, Melis2018, Roscher2020]. More interpretable models are energy-based models [Ackley1985] as their behavior can be associated with physical models for which we have built an intuition. While energy-based models have performance limitations which other models have surpassed, the inclusion of interpretable (possibly physics-based) aspects, and the application in physics models can help build a key to interpret neural networks better [Karniadakis2021]. Quantum-inspired computing [Orus2019] also provides a more interpretable classical machine learning architecture, even for classical tasks [Aizpurua2024].

Furthermore, recent years have seen the application of AI to describe quantum systems [Carleo2017]. However, the complexity of quantum systems composed or accurately described by  $L$  qubits grows exponentially with  $L$ . It is thus expected that only a quantum-based neural network could accurately describe a large quantum system.

AI currently requires significant power, however current studies are investigating the energy cost of performing operations on quantum processors where there is a concrete possibility that significant energy savings can be obtained, for instance if the time to run an operation is significantly shorter, and/or the number of qubits required is orders of magnitude smaller than that of classical bits [FellousAsiani2023].

#### **OBJECTIVE**

For this challenge, one should do two things: use quantum, quantum-inspired or physics enhanced methods to solve classical problems. Then aim for better accuracy, generalization power, and shorter running time and map what limits their performance with the corresponding physical properties that these algorithms are known to be able to study. This can also increase their interpretability properties, either because of the use of physical principles or because these tools have been applied to prototypical physical models and it is thus clear in which conditions they best perform. In parallel, one should also consider quantum problems, for which quantum systems have a clear quantum advantage, e.g. an exponential speedup in the computation and exponential large space in which to store the data.

Then one can consider developing new neural network models which include physical aspects to increase the accuracy of its performance, its interpretability and reduce the fluctuations. For example, including energy, particle-number or momentum conservations.

In parallel, the community could invest in studying accurately the energy cost or running computations on a quantum processor, and the preparation and measurement of the quantum processor, as a function of the platform and algorithms used, including the classical computers and devices connected to the quantum one. This could evaluate under which circumstances a quantum computer can reduce, and to what extent, the energy requirements.

#### **DATA REQUIREMENTS**

For this challenge we can tap on many available standard data sets from different classical tasks, from image recognition to anomaly detection, or engineering problems. For quantum data, we will have to produce data in collaboration with experimentalists and/or industry.

#### **AI METHOD REQUIREMENTS**

Tensor-network based algorithms, and physics-informed models will be key in this challenge. Also, the use of quantum models, most importantly for quantum machine learning problems with quantum data.

#### **COMPUTE REQUIREMENTS**

3 H100 for 36 months.

---

### **INTEGRATION OF CLASSICAL AND QUANTUM COMPUTER SYSTEMS**

---

This challenge envisions a paradigm shift in computing by seamlessly integrating classical and quantum systems, creating a new class of hybrid quantum solutions. This would be done by addressing critical challenges at the quantum-classical interface, particularly in middleware, control electronics, and fundamental theoretical frameworks.

Currently, hybrid quantum systems face scalability, latency, and error correction challenges. Leading approaches, including IBM's Qiskit Runtime [Qiskit] and Quantinuum's System Model H1 [Quantinuum] provide some integration between classical and quantum components. However, there are significant improvements in middleware and control electronics to achieve seamless integration.

Recent innovations in control electronics are addressing these challenges. For instance, the Quantum Instrumentation Control Kit (QICK), developed by Fermi National Accelerator Laboratory engineers, has demonstrated improved quantum computer performance while reducing control equipment costs [Ding2024]. This compact system incorporates the capabilities of an entire rack of equipment into a single electronics board, addressing issues of space requirements and response times that are critical for qubit manipulation.

A notable contributor in this field is Qibo, an open-source framework for quantum computing that manages the full stack from low-level libraries to control quantum hardware, to high-level quantum algorithm implementation [Qibo]. Qibo's associated library, Qibolab, serves as an open-source hybrid quantum operating system, providing the software layer required to automatically execute circuit-based algorithms on custom self-hosted quantum hardware platforms [Efthiou2024].

Current research is also focusing on addressing fundamental theoretical challenges in classical-quantum integration. This includes developing new mathematical models to describe the interaction between quantum and classical systems and exploring the theoretical limits of hybrid quantum-classical computations.

#### OBJECTIVE

AI can thus help in:

- Developing a unified theoretical framework that bridges the gap between classical and quantum computation, addressing the inherent differences in information processing between the two paradigms.
- Creating formal models for hybrid quantum-classical algorithms that can fully exploit the strengths of both paradigms while accounting for their distinct operational principles.
- Optimising the interface between classical and quantum components to minimise information loss and maximise computational advantage.
- Establishing a rigorous mathematical foundation for describing and analysing hybrid quantum-classical systems, including their dynamics, computational complexity, and information flow.

#### DATA REQUIREMENTS

Quantum circuit execution logs provide invaluable insights into the performance and behaviour of quantum algorithms, allowing researchers to identify bottlenecks and optimise execution patterns. Error correction code performance metrics are crucial for improving the reliability of quantum computations, enabling AI models to learn and adapt error mitigation strategies. Latency measurements between classical and quantum components offer a quantitative basis for optimising the interface between these two paradigms, essential for achieving seamless integration. Resource utilisation statistics for various hybrid algorithms help understand the efficiency of different approaches, guiding the development of more optimised hybrid solutions. Quantum state tomography data provides a comprehensive view of the quantum system's state, enabling AI models to understand better and predict quantum behaviour.

#### AI METHOD REQUIREMENTS

Reinforcement learning and generative model can help provide novel, and better performing, solutions to integrate classical and quantum systems, despite their fundamental differences. Adversarial approaches could be explored too.

AI-assisted theorem proving bounds about quantum-classical integration can accelerate the development of formal frameworks for hybrid and quantum systems, providing a solid theoretical foundation for practical advancements.

#### COMPUTE REQUIREMENTS

min 2 x DGX H100 nodes with 16 H100 GPUs in total for 36 months

## Singapore's Role

Since 2007, Singapore is home to the Centre for Quantum Technologies which has catalyzed and still accelerates the progress of quantum technologies in the country. Furthermore, key competences are distributed in A\*STAR, NTU, NUS and SUTD, including classical and quantum algorithms, classical and quantum machine learning, quantum-inspired methods, many-body physics, sensing, error correction, quantum thermodynamics, high performance computing, foundations in physics, and all the major platforms such as neutral gases, trapped ions, superconducting qubits and photonics.

We have a National Quantum Office that steers and coordinates a national coherent effort in research directions to avoid ineffective competition and foster fruitful collaborations.

## Conclusions

Singapore is well equipped to tackle major challenges in quantum computing tapping on classical and quantum computer science, AI, material science and nanofabrication innovations.

What could help is building stronger links between HPC, Computer Science and Quantum Physics communities to tap on each other resources and expertise to produce innovative and higher impact solutions, and of course to have/develop locally multi-qubits quantum processors with different platforms.

The National Supercomputing Centre also supports significantly fundamental and applied research which needs major computing power.

Singapore is also home to reputed local quantum technologies start-ups which work on relevant topics including compilers, like Horizon Quantum Computing, error correction algorithms, like Entropica Labs, hybrid computing on noisy processors, like AngelQ. More companies are emerging too, e.g. RAQS Consulting.

Last, Singapore has a strong ecosystem in material science and engineering, nanofabrication, and computer science, and being a city state fosters cross-disciplinary research that can lead to significant impact on the progress of quantum computing.

And on top of this, also close collaborations with partners oversea to increase the expertise and produce more impactful results.

It could be challenging to compete directly with large tech corporations, and some academic institutions and partnerships, but investigating well in human capital and in both domain-specific and cross-disciplinary research it is possible for Singapore to be at the forefront of AI for quantum computing and quantum computing for AI research, development and enterprise.

## REFERENCES

- [Ackley1985] D. H. Ackley, G. E. Hinton, T. J. Sejnowski, A learning algorithm for Boltzmann Machines, *Cognitive Science*, 9 (1), 147 (1985)
- [Aizpurua2024] B. Aizpurua, R. Orus, Tensor Networks for Explainable Machine Learning in Cybersecurity, arXiv:2401.00867 (2024)
- [Bausch2023] J Bausch, et al., Learning to Decode the Surface Code with a Recurrent, Transformer-Based Neural Network, arXiv:2310.05900 (2023)
- [Bharti2022] K. Bharti, et al., Noisy intermediate-scale quantum algorithms, *Review of Modern Physics* 94, 015004 (2022)
- [Bluvstein2024] D. Bluvstein, et al., Logical quantum processor based on reconfigurable atom arrays, *Nature* 626 (7997), 58-65 (2024)
- [Breuckmann2021] N.P. Beckmann, J.N. Eberhardt, Quantum low-density parity-check codes, *PRX Quantum* 2 (4), 040101 (2021)
- [Carleo2017] G. Carleo, M. Troyer, Solving the quantum many-body problem with artificial neural networks, *Science* 355, 602 (2017)
- [Chakraborty2017] S. Chakraborty et al., Interpretability of deep learning models: A survey of results, 2017 IEEE SmartWorld, San Francisco, CA, USA, 1-6 (2017)
- [Coppersmith1994] D. Coppersmith, IBM Research Report No. RC19642 (1994)
- [Ding2024] C. Ding, et al., Experimental advances with the QICK (Quantum Instrumentation Control Kit) for superconducting quantum hardware, *Phys. Rev. Research* 6, 013305 (2024)
- [Efthymiou2024] Stavros Efthymiou, et al., Qibolab: an open-source hybrid quantum operating system, *Quantum* 8, 1247 (2024)
- [FellousAsiani2023] M. Fellous-Asiani et al., Optimizing Resource Efficiencies for Scalable Full-Stack Quantum Computers, *Physical Review X* 4, 040319 (2023)
- [Grover1996] L. K. Grover, A fast quantum mechanical algorithm for database search, *Proceedings of the twenty-eighth annual ACM symposium on Theory of computing*, ACM 212 (1996)
- [Grover1997] L. K. Grover, Quantum mechanics helps in searching for a needle in a haystack, *Physical Review Letters* 79 (2), 325 (1997)
- [Harrow2009] A. W. Harrow, A. Hassidim, S. Lloyd, Quantum algorithm for linear systems of equations, *Physical Review Letters* 103 (15) 150502 (2009)
- [Jonsson2018] B. Jónsson, B. Bauer, G. Carleo, Neural-network states for the classical simulation of quantum computing, arXiv:1808.05232 (2018).
- [Karniadakis2021] G. E. Karniadakis, et al., Physics-informed machine learning, *Nature Reviews Physics* 3, 422 (2021)
- [McClellan2018] J.R. McClellan, et al., Barren plateaus in quantum neural network training landscapes, *Nature Communications* 9, 4812 (2018)
- [Melis2018] D. A. Melis, T. Jaakkola, Towards Robust Interpretability with Self-Explaining Neural Networks, *Advances in Neural Information Processing Systems* 31 (2018)
- [Nielsen2000] M.A. Nielsen, I.L. Chuang, *Quantum computation and quantum information*, Cambridge University Press (2000)
- [Orus2019] R. Orus, Tensor networks for complex quantum systems, *Nature Reviews Physics* 1, 538 (2019)
- [Peruzzo2014] A Peruzzo, et al., A variational eigenvalue solver on a photonic quantum processor, *Nature communications* 5 (1), 4213 (2014)
- [Qibo] <https://github.com/qiboteam/qibo>
- [Qiskit] <https://www.ibm.com/quantum/blog/qiskit-runtime-capabilities-integration>
- [Quantinuum] <https://www.quantinuum.com/news/riken-selects-quantinuum-system-model-h1-for-large-scale-hybrid-quantum-supercomputing-platform-in-japan>
- [RobledoMoreno2024] J. Robledo-Moreno, et al., Chemistry beyond exact solutions on a quantum-centric supercomputer, arXiv:2405.05068 (2024)
- [Roscher2020] R. Roscher, et al., Explainable Machine Learning for Scientific Insights and Discoveries, *IEEE Access* 8, 42200 (2020)
- [Schuld2021] M. Schuld, Supervised quantum machine learning models are kernel methods, arXiv:2101.11020 (2021)
- [Shor1999] P. W. Shor. Polynomial-time algorithms for prime factorization and discrete logarithms on a quantum computer, *SIAM Review* 41 (2), 303 (1999)
- [Silver2016] D. Silver, et al., Mastering the game of Go with deep neural networks and tree search, *Nature* 529 (7587), 484 (2016)
- [Silver2017] D. Silver, et al., Mastering the game of go without human knowledge, *Nature* 550 (7676), 354 (2017)

## APPENDIX XIII. SUSTAINABILITY



### AUTHORS:

Prof Madhavi Srinivasan  
(NTU, Energy Research Institute @ NTU)  
Prof Yeoh Lean Weng  
(A\*STAR)

Assoc. Prof Arvind Easwaran  
(NTU)

### Executive Summary

In this whitepaper, AI as a necessary element in revolutionizing the nominal routine today by unearthing a science-informed sustainable future for planet earth from the fundamentals is discussed. It is imperative that Sustainability must be embraced at all levels which creates enough room for new discoveries, unimagined pathways, novel methodologies etc. The unprecedented momentum recently towards sustainability has fast-forwarded investigation into alternate options and management

measures with time as the essence. This paper also discusses on how AI could be leveraged as a Sustainability Pathfinder in increasing choices, widening routes, resolute complexities, aid phenomena, innovate solutions and ultimately improve sustainability transformation rates thereby creating new knowledge and exploring the unknown. The suitable methods, grand challenges, scientific enablement, other needs & resources are put-forth for capability building in this domain.

### Introduction

AI for Sustainability deals with AI's critical role in advancing sustainability goals from enterprises, organizations all the way to achieving national decarbonization agendas and beyond. Sustainability goals are best dealt with when the underpinning 'carbon' is well understood- which also paves the way for the appropriate strategies, routes, mechanisms etc., in its timely mitigation. This further necessitates to specifically tackle the Scope 1&2 emissions, while scope 3 from downstream activities is expected to be mitigated while each organization / user transforms towards sustainability by mitigating their individual scope 1&2 emissions.

For effective Scope 1&2 decarbonization efforts from upstream activities, the transition processes primarily involve or intersects with electricity consumption, alternate process methods, product analysis, effective monitoring & ecosystem footprints. Hence the catalyzing topics for Sustainability where AI shall play a key role include:

- Future Grid - Orchestration of the dynamic operating envelope in the envisioned Energy Transition.
- Electrification - Journey displacing fossil fuels by swapping traditional components & systems with electrified analogues.

- Carbon Management - Suite of innovations, mechanisms & processes aiding in facilitating the transition towards a Low-carbon future.
- Urban Nexus - Intersection of multi-vector initiatives in an urban sustainability transformation journey.

Each of the topics must be dealt with both individually and inclusively for a holistic sustainability transformation. The Ultimate Grand Challenge for Sustainability is in nominating plausible routes for a timely transformation of phenomena within a defined geography towards attaining Net Zero goals (or) effectively being Decarbonized across sectors.

The vision for AI for sustainability is effectively deciphering pathways dynamically to achieve a decarbonized future. More specifically, to reduce the carbon emissions (Scope 1 & 2) towards zero.

## Background

Decarbonization has taken main-stage as the priority initiative within organizations and nations with several pledges / commitments put-forth in the recent times. In parallel, governments are gradually implementing 'carbon tax' which impacts businesses and sustainability transformation becomes no more a choice, but a mandate.

Also, noteworthy that many targets remain ambitious, and organizations retort to relooking at the needs half-way for achieving targets on a timely manner. This is frequently due to misinformed goal setting or insufficient feasibility assessment towards decarbonizing pathways which is multi-faceted and complex. This is true to various sectors / industry segments.

While carefully looking at the carbon footprints, it is apparent that systemic changes had to take place across sectors to facilitate any transition. For e.g, the scope 2 'indirect electricity' emissions are not feasible to be

Decarbonization is indeed a transformative journey as it involves Systemic changes across multiple sectors. While doing so, there are serious economic implications, intersects policy & regulation, requisites adhering environmental justice & equity, maintaining climate resilience and necessitates climate adaptation, behavioral & cultural shifts, infrastructural needs & global co-ordination – all of it with progressing technological innovation & adoption and most importantly with a long-term commitment. Addressing the individual grand challenges in aforementioned topics aids in resolving the Ultimate Grand Challenge of Decarbonization. As the challenges are associated with a time-factor, the problem statements become dynamic in nature adding to its complexity.

mitigated at large by an organization, without a grid that facilitates clean electrons. While the clean electrons could be guaranteed only when an appropriate energy mix and its handling is plausible. Such intersects deem AI necessary to imagine / find / advice / invent / discover / solve and thereby catalyze decarbonization with both a macroscopic view and a microscopic lens.

In Singapore, the Singapore Green plan 2030 forms as a comprehensive national strategy aimed at advancing Singapore's sustainable development agenda over the next decade. The plan builds on past environmental initiatives and sets ambitious targets across various sectors to address climate change, enhance resource efficiency, and improve the living environment. Apart from the Green plan, individual government agencies and organisations also have specific targets in decarbonisation / attaining NetZero commonly ranging between 2030 to 2050.

## Grand Challenges

The grand challenges laid out across the sub-domains of Sustainability have a general demand for AI in pathfinder approaches which must be science-informed, from first principles if need-be and in parallel leverage AI in addressing the technical challenges enroute to path-finding trajectories as laid below.

However, the challenges are more quantifiable when a boundary in extent of surface area or timeframe is well defined. For the best use of research investments, the quantifiable challenges in this initiative shall be applied to Singapore and its allied NetZero or decarbonisation targets.

---

### GRAND CHALLENGE 1: FUTURE GRID

---

While AI imagines the path-ahead to Grid decarbonisation, it is also imperative for AI to address challenges in:

- Grid Resiliency by maintaining energy balance & network stability adhering to regulatory & policy needs.
- Forecasting patterns & Anomaly detection in Energy systems, Smart Grids & Grid applications such as EV charging.
- Assess / introduce Flexibility from various resources / infrastructure such as buildings, Storage methods and ancillaries.
- Asset Management in existing & new infrastructures with prediction.
- Overall Energy Costs Reduction promoting Grid interactive energy communities and further assess the economic impact.

#### OBJECTIVE

The quantifiable challenge is to develop a feasible time-oriented transformation plan for NetZero emissions fulfilment of Singapore's power sector by year 2050<sup>1</sup> with Guaranteed Grid Resiliency and Contingency mechanisms allowing for 15-25% Flexibility.

---

### GRAND CHALLENGE 2: ELECTRIFICATION

---

The Electrification grand challenge lies in effective switch over of the present fossil-fuel based Industrial practices / processes or sectoral consumption patterns to electrified operations maximizing renewable usage within a targeted timeline.

AI-guided electrified transformation across industries / sectors forms a key role of AI such as achieving clean Mobility (by 2040 at Singapore). It shall also address / assist industry enablers by:

- Effective planning by predicting demand and real-time intelligence.
- Optimizing infrastructural usage, timeliness of processes, energy footprint, overall costs etc.
- Improving process peripherals such as Safety, Reliability, Uptime etc.
- Energy Management in an electrified scenario.

#### OBJECTIVE

The quantifiable challenge is to develop an Electrification map for achieving 100% clean mobility in Singapore by year 2040<sup>2</sup> for successful displacement of all fossil fuel based transport off the-road.

---

### GRAND CHALLENGE 3: CARBON MANAGEMENT

---

Carbon management caters to the grand challenge of deciphering technical route(s) to NetZero emissions at a sectoral level in target time via improving efficiency, promoting circularity, lowering carbon intensity, enhancing sustainability and reducing waste; via both physical & digital methods.

Emissions reduction with AI considering choices such as Avoidance, Reduction, Substitution, Sequestration & Offsetting is a benefit at large. Incrementally, AI can aid in:

- Processing large variety of data-sets / Bigdata.
- Validation & Verification of both data and physical processes.
- Accounting with best-practices or evolutionary policies.
- Causal analysis in products / process / projects for Company's GHG.
- Automation of Scope 1 to 3 emissions monitoring / reporting.
- Reconciliation of different resources at different levels.
- Closing the loop phenomena / circulatory assessments.

#### OBJECTIVE

The quantifiable challenge would be to chart a transformational plan for national climate target reducing 2030 emissions to 60MtCO<sub>2</sub>e<sup>3</sup> while provisioning carbon-tax at 50\$/tCO<sub>2</sub>e and lead to net-zero emissions by 2050.

#### GRAND CHALLENGE 4: URBAN NEXUS

Urban Nexus models aim at resolving the grand challenge of ensuring sustainability across sectors with systemic shifts, multi-dimensional constraints, multi-vector challenges & targets in a phased manner.

Solving the multidimensional puzzle is beyond human interpretation where AI becomes an able candidate in smoothening the nexus accommodating the quotient of energy / carbon at the cross-roads of urban infrastructure. While doing so, AI can aid in:

- Infusing Sustainability across sectors by prescriptive means with environmental awareness & being judgemental of time.
- Productivity improvements across phenomena / environment importantly buildings.
- Yield enhancements such as in urban farming.
- Preserving heritage or in sustainability transformation.
- Uplifting human health & communities.
- Prioritizing resiliency across sectors.

#### OBJECTIVE

The quantifiable Challenge is to emulate a representative systemic environmental sustainability method aligning to '80-80-80 targets' of Green building masterplan<sup>4</sup> by achieving 80% of buildings go green, 80% improvement in energy efficiency over 2005 levels – all by 2030.

#### RESOURCE REQUIREMENTS

Computing needs: each sub-domain will have specific needs from distributed computing, High-performance computing, and cloud-based computing.

HPC needs per sub-domain is estimated to consume 1,500,000 hours per annum. Hence AI for Sustainability in totality will need on an average 6,000,000 hours per annum.

- CPU (in Millions of Core Hours, in M core hours): 6
- GPU (in Millions of Card Hours, in M card hours): 2
- V100/A100/H100 GPU cards or others in the projected GPU card hours: NVIDIA H100 + NVLink switch system
- Storage (Terabytes): 200
- Special Requirements (e.g. software, etc), if any:  
NVIDIA Omniverse or digital twinning frameworks, CodeAssist tools, Functional mock-up interfaces, Model representation environments

#### AI METHODS REQUIREMENTS

ML approach allows for the identification of models and observations based on short timescales, providing new ways to evaluate and interpret model differences.

Learning methods applicable on different datasets and inference needs include Supervised learning (Regression, Deep Learning, Support Vector Machines), Unsupervised learning (K-means, Principal Component Analysis, Anomaly Detection), Reinforcement learning (Q-learning, Actor-critic algorithm) & Transfer Learning.

ML techniques, such as Convolutional Neural Networks (or CNNs), are increasingly utilized in Climate Science to evaluate climate models; identify model characteristics; and assess model performance in comparison to observational data.

Variations of Neural networks, especially Physics informed Neural networks (PINN) or Hybrid AI form a key role in transformation trajectories and coupling scenarios. Often the critical dynamics are overlooked by other methods, but are imperative for sustainability science, especially while extending simulations to several-kms wide. PINNs are expected to reduce complexity and achieve coupling to a large extent. Nevertheless, existing frameworks such as NVIDIA Modulus and platforms such as NVIDIA Omniverse could be tapped where necessary to enrich the simulations / digital twinning. An example in case of predicting carbon injection methods and capacities for storage in a nested FNO architecture<sup>5</sup> can be handy in developing sub-routines in different strategies.

However, various Regression methods also play a significant role in data-based

approaches. Generative models help in instantiating different components for situational & evolutionary analysis.

'Carbon Models' across different scopes & sectors envisioned to be developed along this effort would form a baseline in evaluating NetZero pathfinders and a groundbreaking tool enabling the coupling of science with target maps for actionable insights. This could be helpful in examining specific events or broader phenomena related to Sustainability/ carbon mitigation. It could also help us understand the interactions between several subsystems in a science-informed manner and achieve process estimates with precision.

Research shall focus on tracking phenomena over time with high-resolution data and conducting regional studies in areas with unique characteristics or significant events. Furthermore, it will allow the performing of long-term simulations to explore different emissions scenarios. Collaboration with other projects and scientists will continuously refine the model, fostering interdisciplinary research within and beyond the AI4SCI initiative.

## AI Methods and Data – Challenges and Opportunities

#### DATA

While the domain specific grand challenges have to be overcome, it is equally important to look at the bottlenecks (or) challenges in implementing AI to advance science in Sustainability.

#### DATA FOR FUNDAMENTALS

Digitalization: or equivalently 'data availability' across different process or ecosystems is found crucial in developing AI approaches / training AI models. Given most of the transitional elements are futuristic in nature, the volume of data sets required is to be paid attention for appropriate synthesis or aggregation.

Quality: In case of established processes and consumption processes, the quality of data must be ensured by rigorous verification methods. This also spreads across the extent of data available for any sensible implementation with validation means and trust factors associated.

Sensitivity: As data shall involve behavioral patterns and high-risk process environments, the sensitivity of the data has to be respected / observed while implementing AI methods.

Platform needs: Data Interoperability and standardization or harmonization in a platform-based approach must be explored to suit AI needs solving challenges in integration and in-situ measurements.

## INTERDEPENDENCY

Constraints Decoupling / Trade-off: Individual optimization seems to be a straightforward implementation. However, complexity becomes several-folds moving upstream while challenges have to be addressed comprehensively. Trade-off in constraints shall form a continual process to improvise outcomes.

Local vs Global scaling / adaptability: The boundary of consideration in the case of Sustainability shall face drastic scaling challenges. This should be overcome by increasing the population of learning parameters from different environments and enriching awareness on compliances needs across regions.

## APTNESS / ADOPTION THRESHOLDS

Robust decision making: Several of the target applications are safety-critical in nature (e.g., the electrical grid), which necessitates that any novel decision making engine based on AI has been thoroughly verified and tested. However, verifying such black-box and data-driven models is non-trivial, particularly when training datasets are not representative (which is typical for real world use cases). Therefore, alternate design approaches that improve robustness of AI based decision making must be developed.

Interpretability: As agendas involve economic implications, environmental justice & resilience, the 'art-of-making-decisions' by AI must be understood for AI-based decisions to be convincing. This should be addressed from modelling practices to evolutionary computation processes to be explainable from roots.

Human-in-loop: Transition is a mission-critical process impacting negatively even due to slight changes. Hence, human intervention is essential in decisions along the execution. Compatibility of the AI-environment with such prefaces promotes adoption.

Vulnerability: the severity of attacks in certain processes would be detrimental and even irreversible in nature. There could also be situations in model drift or overfitting causing unintended operations. Adversarial defense mechanisms have to be in place in developing AI models for Sustainability.

Sustainable implementation: the required resources and environmental impact in doing so is a parameter of choice for individuals & organizations to set a degree of optimization.

## AI METHODS

Noteworthy that, the objectives of optimization and the contingent inputs / boundary conditions differ in time domain due to the dynamic nature of the problem and multi-vector perspectives involved at each step. This requires the constant looping of 'information' obtained by AI making sense of the training / user data back to the AI engine/ model to create new information based on evolving scenarios and observations and carry semantics of user responsiveness / preference.

Secondly, in AI based pathfinding tasks, the demand for 'years-ahead simulations' is rooted in the global optimization with a spectrum of local information, at greater fidelity, both to advance scientific understanding as well as to link to impacts and better integrate local knowledge, including observations.

Hence, to mitigate the dangers of full-throttled HPC implementation with comprehensive simulations with changing scenarios, the AI schema needs to develop mechanisms for critical inquiry. This is also true in improvising systems once physical realization of AI-suggested outcomes takes place to further improvise towards achieving goals (or) in scenarios when the task is classified a-priori into sub-process streams.

In certain cases, data forms the key whenever available. However, physical models being able to generalize matters also find a role outside the bounds of existing data and explore counterfactuals. Hence, it is foreseen that the AI/ML methods, the AI Engine/model and the computing necessary will have a hybrid combination of several techniques/models/frameworks to successfully addressing the scientific grand challenges for Sustainability.

## Singapore's Role

The ecosystem in Singapore is favorable in many aspects to catalyze AI based Science. In this case, firstly, to effectively operationalize AI for Sustainability and secondly in realizing scientific breakthroughs for fast-forwarding sustainability initiatives.

Codified for AI exploitation: when it comes to energy systems / industrial processes / integration methodologies under the transformation purview, most of the domain relevant processes are codified for easy migration into AI environments for model replication and scaling.

Some of the relevant capabilities developed include:

- For instance, for future grids initiative, a representative micro-grid environment<sup>6</sup> with different energy resources and demands with uncertainty prediction has been developed.
- For electrification tasks, schedule optimization<sup>7</sup> forms a key where-in the attributes are studied in prior art.

It is easier to interpret when the role of AI/ML is compartmentalized into 'ML-inside' and 'ML-on top' based on the intended outcomes of AI and its field of view with respective AI Engines utilized.

ML inside: shall relate to ML that projects trajectories / outcomes based on the scenarios and assimilation of outcomes.

ML on-top: shall relate to ML that operates on the data to create insights and produce new information or patterns.

- In the case of urban building decarbonization, Energy Management techniques in Building environments<sup>8</sup> with reinforcement learning could act as an enabler and scaled accordingly.

- For interconnected environments requiring coupling of different subsystems using hybrid AI, a campus-level representative method<sup>9</sup> has been developed for further scaling and Multiphysics integration. This has also demonstrated the effectiveness of the platform in convergence speeds and efficiency of resources.

Data Repositories / Semantics available for a starting point in determining present state of different operations such as power generation & demand pattern, carbon footprint of key activities, infrastructural types and sizes, weather and environmental parameters etc. For few instances / scenario, representative data however at different scale is available for appropriate synthesis to infer the semantics of the individual processes.

Key resources for both geography centric and representative models exist for most of the existing energy infrastructure. However, new infrastructure forms can be synthesized with precision due to the traits observed in pilot systems and processes and shall be adapted for the region.

- The Singapore Energy statistics<sup>10</sup> is available for over a 5-year period with much granularity segregated by its value chain.
- Also, individual power producers have well defined operational datasets<sup>11</sup> backed by asset types and asset spreads. This is also true for different geographies and semantics can be derived from a wide variety of operational data<sup>12</sup>.
- For electrification scenarios, IEA has played a key role in providing free datasets<sup>13</sup> on several outlooks including renewable penetration, EV adoption, energy infrastructures, electricity statistics, carbon tracking, etc.
- With several climate change modelling initiatives, such as Earth-Climate digital twin platform from NVIDIA<sup>14</sup>, other existing weather forecasting models developed, abundant of environmental data is prevalent. In Singapore, administered by NEA, several resources for weather-based data with near-term projections are available for research.

## Conclusions

For successful AI for Sustainability realization, given its depth & width, a focused research and development program, both to improve and implement models and to understand how and why they behave as they do is necessary. In addition to a strong scientific focus, the program must support basic advances in technical fields, such as informatics, numerical mathematics, and especially AI. Doing so will ensure that AI for Sustainability stimulates science, and knowledge gains to maintain transformation pace with demands of our changing world.

Further 'closing-the-gap' in the near-term:

- Mandated ESG Reporting: Singapore is one of the few countries which has mandated ESG reporting<sup>15</sup> progressively which enhances the volume of quality data for effective learning and enhancing precision in AI outcomes.
- Chartered NetZero Future: Recent multi-agency efforts referred as Singapore Green Plan<sup>16</sup> have outlaid specific targets or milestones to achieve across various sectors with timeframes which could form a basis for developing methods operating coherently to test science-informed routes phased sequentially simplifying interdependency constraints.
- Infrastructural investments pipelined: following recent national budget, S\$5B is allocated for Future Energy Fund<sup>17</sup> to propel Singapore's Transition. This not only brings alignment to the potential impact of AI for Sustainability but also helps in validating AI based outcomes in different sub-domains.
- Enriched Technical Innovation: Due to a structured research agenda in the past decade in preparing ourselves for the future, there exists a pool of innovations & solutions to be instantiated for transitional assessments / impacts.

While primarily intended as a telescope looking into the sustainable future of our own planet, it will also aid in incremental improvements towards set targets. It is also foreseen to interpret new observations across geographies enabling us to cope up with unimagined scenarios & beyond.

Having a shared vision for a Sustainable future, meeting the laid targets in the transformational agenda is imperative for attaining Nation's sustainability. However, frequently a physical realization for a trial-and-error approach strategy in different agendas becomes impractical. With the proposed AI techniques and quantifiable targets set forth in this paper along different streams

of decarbonization, when successfully developed and demonstrated, becomes a 'science-informed most-guaranteed pathway' discovered to achieve or realize a decarbonized future, with Singapore becoming a role-model in leveraging AI for Science meeting Sustainability demands. It would also be our responsibility for appropriate use of AI for embarking on a sustainable future.

## REFERENCES

- 1 <https://www.straitstimes.com/singapore/environment/importing-more-clean-energy-among-ways-to-help-singapores-power-sector-reach-net-zero-emissions-by-2050-report>
- 2 <https://www.straitstimes.com/singapore/environment/green-vehicles-add-power-to-the-fight-against-climate-change>
- 3 [https://www.nccs.gov.sg/singapores-climate-action/singapores-climate-targets/overview/#:~:text=On%2025%20October%2022%20C%20Singapore,Emissions%20Development%20Strategy%20\(LEDs\)](https://www.nccs.gov.sg/singapores-climate-action/singapores-climate-targets/overview/#:~:text=On%2025%20October%2022%20C%20Singapore,Emissions%20Development%20Strategy%20(LEDs))
- 4 [https://www1.bca.gov.sg/docs/default-source/docs-corp-buildsg/sustainability/20220726\\_singapore-green-building-masterplan-booklet.pdf?sfvrsn=151fba03\\_8](https://www1.bca.gov.sg/docs/default-source/docs-corp-buildsg/sustainability/20220726_singapore-green-building-masterplan-booklet.pdf?sfvrsn=151fba03_8)
- 5 Wen, Gege & Li, Zongyi & Long, Qirui & Azzadenesheli, Kamyar & Anandkumar, Anima & Benson, Sally. (2023). Real-time High-resolution CO2 Geological Storage Prediction using Nested Fourier Neural Operators. *Energy & Environmental Science*. 16. 10.1039/D2EE04204E.
- 6 Subrat Prasad Panda, Blaise Genest, Arvind Easwaran, Rémy Rigo-Mariani, Pengfeng Lin, Methods for mitigating uncertainty in real-time operations of a connected microgrid, *Sustainable Energy, Grids and Networks*, Volume 38, 2024, 101334, ISSN 2352-4677, <https://doi.org/10.1016/j.segan.2024.101334>.
- 7 J. P. E. Raja and A. Easwaran, "Event-Driven Real-Time Multi-Objective Charging Schedule Optimization For Electric Vehicle Fleets," 2022 IEEE 25th International Conference on Intelligent Transportation Systems (ITSC), Macau, China, 2022, pp. 3283-3289, doi: 10.1109/ITSC55140.2022.9922489.
- 8 Sharath Ram Kumar, Rémy Rigo-Mariani, Benoit Delinchant, Arvind Easwaran. Action Masking for Safer Model-Free Building Energy Management. ACM SIGEnergy Workshop on Reinforcement Learning for Energy Management in Buildings & Cities (RLEM), Nov 2023, Istanbul, Turkey. <hal-04299564>
- 9 K. Zhang et al., "Towards City-integrated Distributed Generation: Platform for Interconnected Micro-grid Operation (PRIMO)," IECON 2020 The 46th Annual Conference of the IEEE Industrial Electronics Society, Singapore, 2020, pp. 3791-3796, doi: 10.1109/IECON43393.2020.9255134
- 10 <https://www.ema.gov.sg/resources/singapore-energy-statistics>
- 11 [https://www.semcorp.com/media/z3ffbzyf/sci\\_fy2023\\_operational-data.pdf](https://www.semcorp.com/media/z3ffbzyf/sci_fy2023_operational-data.pdf)
- 12 <https://datasets.wri.org/dataset/globalpowerplantdatabase>
- 13 <https://www.iea.org/data-and-statistics/data-sets>
- 14 <https://nvidianews.nvidia.com/news/nvidia-announces-earth-climate-digital-twin>
- 15 <https://www.esgtoday.com/singapore-to-introduce-mandatory-climate-reporting-beginning-2025/>
- 16 <https://www.greenplan.gov.sg/>
- 17 <https://www.channelnewsasia.com/singapore/future-energy-fund-clean-fuel-lawrence-wong-budget-2024-hydrogen-nuclear-natural-gas-4128656>

# APPENDIX XIV. EDUCATION



## AUTHORS:

Assoc Prof Tan Seng Chee  
(NTU National Institute of Education)

Assoc Prof Ben Leong  
(NUS, AI Centre for Education Technologies)

## Executive Summary

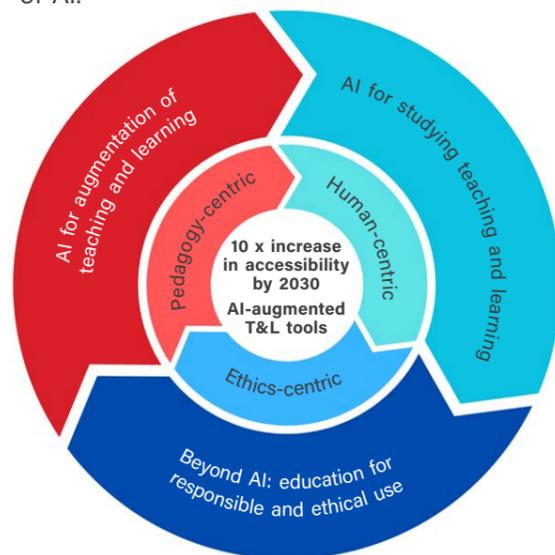
Recent developments in Artificial Intelligence (AI) have paved the way for possibilities previously unimaginable, with profound and disruptive impacts on many industries and domains. One such domain is education, which is of critical importance to Singapore. With its advanced level of development, Singapore has the opportunity to make a quantum leap by leveraging the advances in both pedagogy and technology. This White Paper proposes a new research programme to address the most urgent questions regarding the effective incorporation of AI into Singapore's education system, focusing on both the macro-level (systemic) and the micro-level (pedagogical).

We propose to:

- Comprehensively develop and investigate new pedagogies that leverage AI, including Generative AI (GenAI);
- Develop human-centric techniques and pedagogies to develop skills like critical thinking, creativity and collaboration that will become increasingly important in an AI-driven future;
- Assess the effectiveness of AI applications and pedagogies, so that teachers can be advised on how they should be responding to the large number of GenAI-based tools that are available; and

- Investigate the potential shortcomings of GenAI in teaching and learning and develop tools and techniques to mitigate the negative impact.

Our vision is to to achieve a tenfold increase in the number of Singapore students with access to effective AI-enhanced educational tools, which are designed and tested with a focus on human-centric, pedagogical, and ethical considerations by 2030. AI is used both to develop tools and platforms to augment teaching and learning, and for studying teaching and learning. Equally important are programmes to educate students and teachers on the responsible and ethical use of AI.



## Introduction

There have been many recent advances in the field of Artificial Intelligence, and the one development that singlehandedly captivated the world's imagination with its potential to transform how we live, learn, work, and play, is Generative AI (GenAI) with OpenAI's ChatGPT premiere in Nov 2022.

In the education domain, international organisations such as UNESCO, OECD and World Economic Forum have published several articles on GenAI and the future of learning. In an article for the World Economic Forum, Acar (2024) proclaimed that GenAI is an opportunity to address long-standing "flaws of current education systems" such as inequitable access and teacher burnouts." Pons (2023) writing for OECD, posited that GenAI could personalise learning, and that it has the potential to turn classrooms into spaces where learning is more dynamic, flexible, collaborative, and individualised. Similarly, Miao and Holmes (2023) suggested several ways that GenAI can be incorporated into teaching and the research of learning in their UNESCO report.

In a similar vein, Singapore's Ministry of Education issued a response in Parliament in February 2023, suggesting that GenAI will "become more pervasive over time". Furthermore, educators in schools and Institutes of Higher Learning will need to be provided with "guidance and resources to effectively harness it to enhance learning" (MOE Singapore, 2023). A similar view was expressed in a World Economic Forum article (Rayner, 2023), that cited current experts suggesting that many skilled and high paying positions in the future would be augmented by AI.

The excitement that GenAI stirred has compelled many organisations and individuals rushing to apply AI in the teaching and learning industry. Yet, many researchers are also urging caution in our experimentation and research. The challenge is to ensure that any transformation with AI, will, at the very least:

- adhere to the principle of pedagogy first;
- address key education needs and gaps;
- amplify teachers' capacity and agency rather than replace them; and
- develop critical knowledge and skills in students such that they are in control of the technology rather than be subservient to it.

A study by Felton et al. (2023) claimed that most teaching positions (with some exception such as special education) might be at risk because of advances in GenAI. We do not want a future where AI completely replaces human instructors. We do not want teachers to lose control over the teaching process. And we do not want teachers to be reduced to a mostly managerial role, merely overseeing learning systems and learners with limited to no opportunities to nurture and educate learners through the exercise of one's own autonomy and abilities. An undesirable scenario would be one where an AI sets the assessment for the teacher, a second AI completes the assessment for student, and a third grades the assignment on the teacher's behalf.

Instead of simply studying how AI might be used to improve education, we also need to ensure that our learners are equipped with the necessary skills needed to adapt to a future where AI will take over many of the tasks currently done by humans. Skills like critical thinking, creativity, communication, and collaboration will become vitally more important than ever before. As such, there is an urgency to study how AI itself could potentially be used to help learners develop these skills. In short, how can we use AI not as a substitute to human teaching, but to strategically augment teaching and learning?

While most would agree that it is important to investigate how best to address current education gaps with AI, this concern is often not the priority for teachers who are burdened with many immediate deadlines and an onerous workload. It is also not a possibility that many

teachers have considered because they are unfamiliar with such technologies. Therefore, what many teachers need is the support to understand the technology sufficiently before they can understand what they can do with it, and so use it effectively to enhance teaching and learning.

For teachers who devote time to investigate these issues, they often encounter institutional challenges such as data ethics (such as privacy, confidentiality and right of use), AI governance (such as potential bias, discrimination and misuse), cybersecurity risks (such as adversarial attacks and model poisoning) and data ownership (such as copyrights, intellectual property and misrepresentation). Navigating these rules and regulations is daunting for individual teachers. And the costs of failure for attempting explorations with AI in their teaching is enough to discourage many from even trying.

## Background

A workshop on AI for Education was held on 9 May 2024, from 9 am to 2 pm, and was attended by more than 60 in-person participants from NTU, NUS, SUSS, A\*STAR and other higher institutions. The workshop's agenda included an introduction to AI4SCI by Kedar Hippalgaonkar and introduction to the workshop by the two workshop organizers Tan Seng Chee (NIE/NTU) and Ben Leong (NUS); invited speakers presented current work and ideas on how AI has been used or can be used to support teaching and learning, followed by small groups discussion and group sharing. The broad themes of group discussion included new methods of teaching, learning and assessment in the era of generative AI, AI for neurobiological methods of understanding learning, new methods for studying teaching and learning, and AI supporting lifelong learning and knowledge creation. In each group, the participants explored various innovations, emphasizing the transformative role of AI in teaching, learning, and educational research. Speakers presented on the integration of AI tools, data-driven approaches, and the

Finally, the introduction of GenAI tools, like ChatGPT, has also created problems for teaching and learning. Plagiarism has become more common. On one hand, some students are unsure of differentiating the ambiguities of using GenAI as an assistant compared to using it as a shortcut; on the other hand, it has become significantly more convenient for students to plagiarise and avoid detection by their instructors. Students who rely on GenAI tools for quick answers are at risk of undermining their own learning as they often receive these answers without the appropriate context to make full sense of the information they receive. As such, undue reliance on GenAI tools can have an adverse effect on the impact on students in terms of the development of core skills and valuable learning objectives across their education journey.

potential of AI to enhance educational outcomes. Breakout sessions facilitated in-depth discussions on these topics.

A key focus was the application of AI to foster knowledge creation and critical thinking in educational settings. AI tools can support teachers by providing sustainable solutions for integrating new information and knowledge into classrooms, ultimately promoting long-term problem-solving abilities in students. Emphasis was placed on the need for institutional support to help teachers effectively adopt AI technologies.

The workshop highlighted the importance of student engagement in the learning process. AI technologies, such as deep learning models, were discussed as tools to automatically detect and categorize pedagogical features in lecture recordings, providing real-time feedback to enhance student learning. The challenge of engaging students was addressed, with AI proposed as a means to personalize learning experiences and maintain student interest.

Several talks focused on leveraging large-scale educational data to improve learning outcomes. The establishment of centers like ADIS was discussed, which provides data API services and supports research and application development in education. Projects using deep learning models to predict student success and skills development showcased the potential of AI to offer early interventions and tailored educational resources.

The workshop underscored the importance of privacy and ethical considerations in using AI for educational purposes. The generation of synthetic data was presented as a solution to protect sensitive student information while still enabling meaningful research. Ensuring data quality and adhering to privacy regulations were emphasized as critical for the successful implementation of AI in education.

Innovative AI-driven tools, such as EEG headsets, were introduced to provide real-time feedback on students' learning and engagement. These tools can help educators identify when students are struggling and implement timely interventions. Additionally,

## Grand Challenges

### GRAND CHALLENGE 1: ACCESS TO AI FOR EDUCATION

The grand challenge is to achieve a tenfold increase in the number of Singapore students who have access to effective AI-enhanced educational tools, which are designed and tested with a focus on human-centric, pedagogical, and ethical considerations, by 2030.

In light of the challenges and difficulties laid out in the preceding section, the grand challenge for AI in Education is to boldly yet wisely integrate AI into Singapore's education system – not as a substitute, but as an augmentation tool to every aspect of teaching and learning. This is to be achieved by researching and developing ways to simultaneously enhance the teachers' ability to deliver quality and effective teaching and

the use of multimodal AI tutors and generative AI for speech was demonstrated, showcasing how AI can personalize and enhance the learning experience.

Breakout sessions explored new methods of teaching, learning, and assessment in the era of generative AI. Discussions included the role of AI in motivating students, providing detailed feedback, and facilitating social learning. The potential of AI to support lifelong learning and knowledge creation was also examined, with a focus on ethical considerations and the need for continuous, sustainable learning opportunities.

The workshop concluded with an emphasis on the need for continued exploration and integration of AI in education. By addressing challenges such as data privacy, ethical considerations, and the need for institutional support, AI has the potential to significantly enhance educational outcomes and support lifelong learning. The discussions highlighted the importance of collaboration between educators, researchers, and institutions to effectively leverage AI for the benefit of all learners.

learning, while also enhancing students' ability to learn and create knowledge across their entire journey as lifelong learners.

Teaching and learning are inherently complex social activities. We cannot achieve the grand challenge by simply providing all Singaporean children with ChatGPT access. Many of the educational outcomes we care about are either intangible or are difficult to measure as they require a long period of time to develop within the learner. As such, this grand challenge is better defined as a set of questions that need to be addressed together to ensure that the outcomes from the grand challenge are (i) impactful; (ii) practical, (iii) effective, and (iv) relevant to the local context of education within Singapore. Only by answering these questions can we be sure that the AI-enhanced educational tools that our children are finally provided with are truly making a positive difference.

## OBJECTIVE

We propose a comprehensive multi-faceted programme comprising two categories of questions. The first category concerns the micro level or pedagogical perspective, focusing on the methods and practices used in teaching and learning. The second category concerns the macro level or systemic perspective, looking at the broader, overarching aspects of incorporating AI into the education system, focusing on the large-scale and long-term impacts of AI on education, and consequently how it will affect the structure and policies of AI in education.

The urgent questions of this grand challenge that will enable Singapore's education system to make the quantum leap forward are as follows:

- Micro-level (Pedagogical)
  - » What are the different ways where GenAI can be used to improve teaching in existing domains? How do we assess the effectiveness of these methods to provide practical advice to teachers on how they should be responding to the large number of GenAI-based tools that are available?
  - » GenAI is prone to hallucination. What tools and techniques can we use to eliminate (or reduce) hallucination, while simultaneously developing pedagogy that will enable our students to identify hallucination and work effectively with GenAI?
  - » What are the negative impacts of GenAI tools? What tools and techniques can be used to mitigate the negative impact of such tools?
  - » As GenAI becomes more human-like and more capable in its abilities, what traditional learning objectives are we willing to give up in our education or at least allow for a transformation in light of AI-aided assistance/augmentation?

For example, the invention of search engines led to changes in how we teach students to search for resources in the library using index cards to teach students how to discern and verify information found online. Given the rapid developments in GenAI, we may have to start pre-empting these changes if we wish to stay ahead of the game.

- Macro-level (Systemic)
  - » How do we help students develop the agency, competency and tools to become knowledge creators, innovators, and critical thinkers who can shape our future (OECD, 2018)? Can we use AI to track and measure the efficacy of such efforts over a student's entire course of education?
  - » How can AI be used to study the effectiveness of teaching with AI?
  - » What are the long-term negative impacts of using AI tools for both teaching and learning (to both teachers and students)?

Through this endeavour, we envision the ubiquitous presence of AI-augmented learning in our Singapore schools that are designed and tested with a focus on human-centric, pedagogical, and ethical considerations. To this end, we will be able to achieve the capabilities required to nurture and prepare our learners well for an AI-driven future, cultivating within our learners the agency, competency and tools they need to become knowledge creators, innovators, and critical thinkers who can shape our future. This will enable Singapore to make the quantum leap forward by avoiding the pitfalls of using technology merely as a substitute for human abilities, but instead to educate a population capable of using AI tools to enhance how we live, learn, teach, work, and most importantly of all, how we make meaning in our day-to-day interactions with humans and machines.

## DATA REQUIREMENTS

The following is a potential list of data requirements for AI-augmented learning:

- Student Profiles: Demographic data, learning preferences, prior knowledge, and educational history. High-fidelity textual data (survey forms, historical academic records)
- Interaction Data: Clickstream data, time spent on tasks, response patterns, and engagement metrics. Moderate to high fidelity, capturing millisecond-level interaction data to analyse student engagement patterns
- Performance Data: Test scores, assignment grades, quiz results, and feedback. High fidelity, precise scoring and grading data, including detailed breakdowns of performance across different skill areas
- Behavioural Data: Logs of classroom behaviour and peer interaction records. Moderate to high fidelity, capturing enough detail to analyse student participation and collaboration without overwhelming the system
- Content Usage Data: Moderate to high fidelity, capturing enough detail to analyse student participation and collaboration without overwhelming the system.
- Student Submissions: Text (essays, reports), digital artefacts (presentations, projects). High fidelity, with clear and detailed capture of submissions, including formatting and metadata (e.g., time of submission).
- Peer Review Data: Textual feedback, numerical ratings. Moderate fidelity, sufficient to capture meaningful peer insights while managing subjectivity in reviews.
- Teacher Feedback: Textual comments, annotated submissions. High fidelity, capturing nuanced teacher feedback, including tone and context.

- Plagiarism Detection Data: Textual comparison data, source matching logs. High fidelity, ensuring precise identification of potential plagiarism, including similarity scores and source references.
- Classroom behavioral data: Textual logs and video recordings. Moderate fidelity, sufficient to track significant trends without infringing on privacy

## AI METHOD REQUIREMENTS

Various machine learning methods will be employed to implement personalized user experiences and optimize interactions. Techniques such as linear regression, decision trees, k-Nearest Neighbours (k-NN), and neural networks will be utilized to predict user behaviour, tailor content, and enhance decision-making processes based on individual user data.

Natural Language Processing (NLP) will be leveraged to classify, translate, and generate textual input and output, enabling the system to understand and respond to user queries, perform language translation, and create coherent and contextually relevant text.

Computer vision techniques will be applied to handle media-related tasks, including image classification, object detection, image segmentation, and face recognition. These techniques will enable the system to analyse visual content, identify and locate objects within images, segment images into meaningful parts, and recognize faces, thereby enhancing user interactions with media.

In addition, expert systems techniques will be integrated into the intelligent tutoring system to provide rule-based decision-making capabilities, particularly in areas requiring specialized educational expertise.

## COMPUTE REQUIREMENTS

Listed here is an estimation of the hardware cost, running cost, GPU hours and GPU core, assuming the development of 8 systems over 5 years:

	Estimate for five years		
	Hardware cost	Development	Running Cost GPU
Adaptive Learning Platforms	300000	1296000	8 A100
AI-driven assessment tools	200000	1296000	4 A100
Intelligent Tutoring Systems	300000	72000	8 H100
AI powered educational research tools	3600000	200000	36000 2 A100
AI Enabled Classroom management tools	300000	36000	4 H100
AI-driven Longitudinal Learning analytics	300000	72000	4 H100
Multi-modal learning systems	300000	72000	8 H100
AI ethical use and digital citizenship systems	100000	12000	4 A100
<b>Total</b>	<b>\$ 3,600,000.00</b>	<b>\$10,000,000.00</b>	<b>\$14,460,000.00</b>
<b>Grand Total</b>			<b>\$18,060,000.00</b>

Assumption: 360 schools include 10 IHLS

Total of 8 systems

Running at 18 hours per day on a 10% utilisation per system

## AI Methods and Data - Challenges and Opportunities

### STATE-OF-THE-ART EXPLORATIONS INTO THE USE OF GENAI FOR TEACHING AND LEARNING

In the use of genAI technologies for teaching, Ali et al. (2023) have leveraged OpenAI's large language models to guide students' interactions that support various learning methods like inquiry learning. It differed from typical interactions with chatbot that operate in a conventional, teacher-centric instructions. The initial study indicated adherence to the educational strategies, cognitive direction, and socio-emotional assistance. These elements were all incorporated generatively by the AI, without relying on fixed rules or explicit programming. Lee, Tan and Teo (2024) explored how GenAI can be leveraged to support students in their knowledge building efforts, such as to "summarize the existing pool of ideas generated by student discourse" and "prompt students to consider diverse perspectives, help them appreciate the value of varying viewpoints that are essential for collaborative knowledge building." A pilot study demonstrated the ability of ChatGPT in supporting these processes, but it also

showed invalid references. Moving forward, ways to enhance the accuracy of GPT-generated information are needed. Putting human-in-the-loop into the process of learning is also critical, and ways to enhance the accuracy of GPT-generated information are needed. Tan has also developed and is evaluating knowledge building learning companion for teachers and students using large language models.

In an environment where the evolution of Large Language Models (LLMs) is rapidly advancing, new methods such as:

- Design of human-in-the-loop feedback where human feedback is not just used to train models to improve the accuracy of machine response, but human feedback is considered in the front-end design. For example, where curriculum or benchmarking text can be inputted according to users' needs.
- Multimodal learning: Integrating text with other modalities like images, audio and video to deepen understanding of learners and learning context.

- Building on the concept of explainable AI (XAI) to develop models that produce analytical models alongside tools that interpret the machine learning model. This increased transparency will allow users to have a better understanding of the AI models if need to.

At NTU, the Centre for the Applications of Teaching and Learning Analytics for Students (ATLAS) has been coordinating the GenAI efforts in the university for teaching and learning. ATLAS has built prototypes for a Socratic chatbot for statistics and conducted faculty validation (Qiu et al., 2024) and experimental study. The chatbot will be expanded to a Mathematics course in the coming semester. ATLAS has also been supporting other faculty who are interested in building other types of GenAI applications to personalise learning in their classes. Examples that are featured in recent media publications include Waai which is a writing companion Gen-AI application based on design thinking principles and Prof Leodar which is a Singlish study buddy. These applications aim to enhance student learning outcomes by providing them with immediate individualised feedback and *feedforward* which are difficult to accomplish in courses with large enrolments and traditionally challenging core courses. More importantly, these efforts are meant to give agency back to faculty who can then control the response from the LLMs through prompt engineering and retrieval augmented generation and have access to student conversation histories.

At NUS, the AI Community-of-Practice (COP) has been featuring use cases of GenAI in teaching and learning as a way to expand both educators' and students' perspectives on the possibilities of GenAI uses in teaching and learning. Many GenAI users cannot see GenAI's capabilities beyond that of an answer generator or a cheating tool. These use cases demonstrate how GenAI can be employed in a variety of contexts, and in ways that creatively facilitate and enhance learning or working: from using it as an idea generator for

assignments or discussions, as a model for role-play to practice soft skills, as a personal assistant for students to help students with their learning, as a teaching assistant to guide and assist students with their assignments, to using it as a facilitator for self-reflection and team building, and more.

GenAI is not the only form of AI that is applicable to teaching and learning. Discriminative models are relevant in solving problem statements that require predictive and recommendation solutions. Examples from ATLAS at NTU would include the use of deep learning approaches to predict students who might require academic support in the upcoming semester (Qiu et al., 2022; Qiu et al., 2023; Qiu et al., in press) and efforts to leverage machine learning to map course curriculum to SkillsFuture taxonomies to identify skills gaps (Lai et al., 2024). These AI solutions help schools carry out early intervention and revise their curriculum so that their students have a better chance at being academically successful and have opportunities to develop needed skills.

### IMPROVING PEDAGOGY WITH AI

Tan (2024) has developed a pedagogical framework for the use of AI for teaching and learning for both students and teachers: (1) Students learning from, learning with, and learning about AI. Learning from AI is exemplified in intelligent tutoring systems, where AI takes on the roles of a tutor. Learning with AI is about AI augmenting or orchestrating learning, while students retain their agency and responsibility in learning. Learning about AI focuses on developing AI literacy and readiness. (2) Teachers play critical and irreplaceable roles in working in partnership with AI to help students achieve the learning explained in the preceding statement. Teachers, also, can learn from AI about a topic or about teaching approaches, learn with AI in designing lessons and getting feedback about their lesson design, and learn about AI when getting ready to adopt AI in their classrooms.

At the same time, the education sector would need to proactively think about how teaching, learning and assessment could be revolutionised in the next five years. For example, with the ability to animate a person using a single image, the days of meta-humans as opposed to text-based chatbots are nigh (Xu et al., 2024). Robotics and virtual/augmented reality technologies are integrating with GenAI at an incredible pace such that in the foreseeable near future, we could create experiential learning programmes that are far more immersive and responsive to actions by the learners. Finally, AI could enable self-regulated problem-based learning and authentic assessment on a massive scale. Every student gets to propose their own unique problem statements on issues they have personal interests in, leverage AI in collaboration with others with similar problem statements and be assessed on the quality of their solutions as part of their undergraduate programme.

The design, development, and evaluation of these and other next generation pedagogical ideas would most certainly be key research areas to advance teaching and learning in the age of AI.

### AI METHODS FOR STUDYING LEARNING

Besides how we could improve teaching and learning with AI, we should also consider how we can *study* teaching and learning in the age of AI *with* AI. AI could enhance current research methods and offer new methods to study the impact of learning interventions. For example, AI could potentially:

- Massively accelerate qualitative research with its ability to (a) summarise and (b) annotate large quantity of text in a short time.
- Synthesise a huge amount of multimodal student learning behaviours, interactions, outcomes and experiences data.
- Support real-time statistical, graphical and network analysis of multimodal data in combination with machine learning approaches.

In other words, the toolkits that we can build with AI to study proposed solutions for key pedagogical questions and education challenges can be greatly expanded, and this might lead to new methods for evaluating pedagogical methods, interventions, and programmes.

One exciting new direction that AI affords us is the ability to now conduct longitudinal studies on a far greater scale, either across a student's entire formal education or to compare trends across generations. It is possible to study and even measure educational outcomes that previously could not have been easily measured with such fine-grain details. With this new direction now made possible, we can use AI to potentially:

- Investigate how education policies may relate to learning outcomes
- Investigate on a large scale how students' performance and cognitive abilities develop over time to make inferences and changes in policies
- Identify shifting trends across generations of students to design more effective policies

It is also possible to use such data to identify specific learner needs, or keep a record of the types of scaffolding techniques that have been previously found effective for them – such information can then be used to better inform/advise teachers in the future on how they can best tailor their educational approaches to cater to each student's unique learning needs.

There are two possible approaches to analyze the datasets. First, a foundation model for tabular data can be employed (van Breugel & van der Schaar, 2024). Second, neural network based graph (causal) modeling can be considered for datasets with a mix of observational and interventional data (Lagemann et al., 2023). Both approaches are highly suitable to the high-dimensional (many variables) nature of our large datasets and the need for rapid, automated mining for insights that traditional statistical approaches lack.

### RESPONSIBLE USE OF AI FOR EDUCATION

As a final note, the responsible use of AI is of utmost importance. Where the temptation to take shortcuts to learning is now stronger than ever, how can we ensure that students do not consider this as their first or even their last resort, but instead practice better habits that will enable them to do well later in the workforce? Therefore, studies around the ethical use of AI and how the use of AI will impact human capabilities are necessary. Some research areas could include:

- how we should test the usefulness and fairness of AI interventions for teaching and learning to ensure that they truly benefit both the teachers and the students;

### Singapore's Role

We believe that Singapore is well-positioned to tackle the questions in the grand challenges raised above. Our schools are well-resourced. We have the Student Learning Spaces (SLS) that collects learning analytics data and the MOE SLS Office has indicated that it is willing to share the data collected in the SLS for research. AI-based applications like adaptive learning of mathematics and AI-assisted assessment of short answers have been added to SLS. It is also willing to share the current research problems that MOE is grappling with. This translates to access of data from at least 420,000 students, just counting the K-12 students.

At present, every secondary school student has a Personal Learning Device (PLD) as part of MOE's National Digital Literacy Programme. A number of primary schools are piloting a similar initiative. It is likely that every student in Singapore will eventually be equipped with a PLD. This presents an opportunity for researchers to use these

- how we should evaluate the effectiveness of the AI tools used to study teaching and learning practices to gain cogent insights;
- how we should go about studying what is lost in terms of human knowledge, skills and abilities when AI is adopted into education; and
- how we should develop students' digital citizenship to be responsible users contributing to safe, ethical and sustainable digital environments.

Appendix 1 lists some potential systems that can benefit teaching and learning in Singapore schools and higher education.

PLDs to capture more fine-grain data that could not be captured through SLS. Beyond data collection, the infrastructure is ready (at present, at least from secondary school onwards) to experiment with new methods for teaching and learning with AI.

As for universities, both NUS and NTU each have their own respective data lakes that are already capturing huge amounts of data about students and their learning activities through their learning management system and more. This infrastructure can be potentially expanded to capture additional data if necessary.

With sufficient resources, we could potentially build the infrastructure to conduct even bigger longitudinal studies, tracking learners throughout their entire formal education, and the means to expand the data collection through personal learning devices or the respective learning management systems. This data can then be used to determine the effectiveness of learning interventions.

To ensure that the learner is tracked when they move from one academic institution to the next, we could tap on the existing SingPass infrastructure which is already a well-secured platform, or design something

similar to SingPass, albeit constrained only to the context of education. A mechanism will also need to be put in place to ensure that we preserve data privacy while analyzing the learning data.

## Conclusion

It is without doubt that education remains critical to the future competitiveness of Singapore's economy. Given the recent advent of GenAI and the release of many other

advances in AI, it is imperative that Singapore invest in exploiting AI for education to ensure that we do not fall behind.

## REFERENCES

Acar, O. A. (2024, February 19). Generative AI opens up vast opportunities for education. World Economic Forum. <https://www.weforum.org/agenda/2024/02/with-generative-ai-we-can-reimagine-education-and-the-sky-is-the-limit/>

Ali, F., Choy, D., Divaharan, S., Tay, H. Y., & Chen, W. (2023). Supporting self-directed learning and self-assessment using TeacherGAIA, a generative AI chatbot application: Learning approaches and prompt engineering. *Learning: Research and Practice*, 9(2), 135–147. <https://doi.org/10.1080/23735082.2023.2258886>

Felten, E. W., Raj, M., & Seamans, R. (April 10, 2023). Occupational Heterogeneity in Exposure to Generative AI. SSRN. <https://dx.doi.org/10.2139/ssrn.4414065>

Holmes, W., & Miao, F. (2023). Guidance for generative AI in education and research. UNESCO Publishing. [unesdoc.unesco.org/ark:/48223/pf0000386693](https://unesdoc.unesco.org/ark:/48223/pf0000386693)

Lagemann, K., Lagemann, C., Taschler, B., & Mukherjee, S. (2023). Deep learning of causal structures in high dimensions under data limitations. *Nature Machine Intelligence*, 5(11), 1306–1316.

Lai, J.W., Zhang, L., Chan, Y.S., Sze, C.C., & Lim, F.S. (2024) Compelling Educational Offerings: A Study on the Efficacy of Skills Identification Platforms with Course Descriptions, *Inted2024 Proceedings*, pp. 2553-2561. <https://doi.org/10.21125/inted.2024.0709>

Lee, A. V. Y., Tan, S. C. & Teo, C. L. (2023). Designs and practices using generative AI for sustainable student discourse and knowledge creation. *Smart Learning Environment*, 10, 59. <https://doi.org/10.1186/s40561-023-00279-1>

MOE Singapore (2023, February 6). Artificial Intelligence Technologies/ChatGPT. [Parliamentary Replies]. <https://www.moe.gov.sg/news/parliamentary-replies/20230206-artificial-intelligence-technologies-chatgpt>

OECD. (2018). The future of education and skills – Education 2030. Retrieved from <https://www.oecd.org/education/2030/E2030%20Position%20Paper%20%2805.04.2018%29.pdf>

Pons, A. (2023). Generative AI in the classroom: From hype to reality. *Generative AI in the classroom: From hype to reality?* OECD.

[https://one.oecd.org/document/EDU/EDPC\(2023\)11/en/pdf](https://one.oecd.org/document/EDU/EDPC(2023)11/en/pdf)

Qiu, W., Khong, A. W. H., Supraja, S., & Tang, W. (2023). A Dual-Mode Grade Prediction Architecture for Identifying At-Risk Students. *IEEE Transactions on Learning Technologies*. vol. 17, pp. 803-814. <https://doi.org/10.1109/TLT.2023.3333029>

Qiu, W., Khong, A.W. H., & Lim, F. S. (in press). Enhanced Student-graph Representation for At-risk Student Detection. *Proceedings of International Symposium on Circuits and Systems (ISCAS)*.

Qiu, W., Supraja, S., & Khong, A. W. H. (2022). Toward better grade prediction via A2GP - an academic achievement inspired predictive model. *15th International Conference on Educational Data Mining (EDM 2022)*, 195-205. <https://dx.doi.org/10.5281/ZENODO.6852984>

Qiu, W., Su, C.L., Jamil, N.B., Ng, S.S.H., Chen C.M., & Lim, F. S. (2024). "I Am Here To Guide You": A Detailed Examination of Late 2023 Gen-AI Tutors Capabilities in Stepwise Tutoring in An Undergraduate Statistics Course, *INTED2024 Proceedings*, pp. 3761-3770. <https://doi.org/10.21125/inted.2024.0984>

Tan, S. C. (2024). *Using Generative AI for Intelligent Augmentation of Teaching and Learning*. Invited presentation at Murdoch University. June 2024.

Rayner, M. (2023, August 14). AI: 3 ways artificial intelligence is changing the future of work. World Economic Forum. <https://www.weforum.org/agenda/2023/08/ai-artificial-intelligence-changing-the-future-of-work-jobs/>

van Breugel, B., & van der Schaar, M. (2024). Why tabular foundation models should be a research priority. *arXiv preprint arXiv:2405.01147*.

Xu, S., Chen, G., Guo, Y., Yang, J., Li, C., Zang, Z., Zhang, Y. & Guo, B (2024). VASA-1: Lifelike Audio-Driven Talking Faces Generated in Real Time (arxiv.org). <https://doi.org/10.48550/arXiv.2404.10667>

## Annex

A list of potential systems for the AI education white paper

### ● AI-Enhanced Adaptive Learning Platforms

» **Description:** Develop platforms that personalize learning experiences for students by adapting content, difficulty, and pacing according to individual learning styles and progress. These systems could incorporate generative AI to create dynamic content, provide real-time feedback, and suggest additional resources tailored to each student's needs.

» **Use Case:** A platform that adjusts its teaching strategy based on real-time analysis of a student's interactions, helping them with difficult concepts by offering alternative explanations or practice problems.

### ● 2. AI-Driven Assessment Tools

» **Description:** Implement AI systems that automatically assess student work, including written assignments, projects, and even creative work, while providing detailed feedback. These tools would need to address challenges like bias, hallucination, and the accurate grading of diverse types of student submissions.

» **Use Case:** A tool that can grade essays with a focus on both content and style, offering students insights into how they can improve their arguments or writing techniques.

### ● Intelligent Tutoring Systems

» **Description:** Create AI-based tutors that can provide personalized instruction and support for students in various subjects. These systems could offer explanations, answer questions, and guide students through problem-solving processes, mimicking the role of a human tutor.

» **Use Case:** A virtual tutor that helps students learn specific topics by identifying weak areas in their knowledge and providing targeted exercises to improve.

### ● AI-Enabled Classroom Management Tools

» **Description:** Develop AI systems to assist teachers in managing classrooms by tracking student engagement, understanding classroom dynamics, and providing insights into student behavior and performance. These systems can also help in personalizing learning experiences based on the collective progress of the class.

» **Use Case:** A dashboard for teachers that highlights students who might need additional support or those who are excelling and could benefit from more challenging materials.

- **AI-Powered Educational Research Tools**

- » **Description:** Build AI tools that can assist researchers in studying educational outcomes, pedagogical methods, and the impact of AI in education. These tools would analyze large datasets from educational activities to generate insights that can guide future teaching strategies.
- » **Use Case:** A research platform that uses AI to analyze the effectiveness of different teaching methods across schools, helping educators identify best practices.

- **AI-Based Ethical Use and Digital Citizenship Systems**

- » **Description:** Develop systems that help students and teachers understand and practice ethical AI use, focusing on digital citizenship and responsible AI practices. These systems would include modules for learning about data privacy, AI bias, and the impact of AI on society.
- » **Use Case:** An interactive learning tool that guides students through real-world scenarios involving AI, teaching them how to navigate ethical dilemmas and make responsible decisions.

- **AI-Driven Longitudinal Learning Analytics**

- » **Description:** Create systems that track and analyze student learning over long periods, potentially from primary through tertiary education. These systems would provide insights into how educational interventions impact long-term learning outcomes and help refine educational strategies.
- » **Use Case:** A system that tracks a student's academic journey, identifying key moments where interventions had the most impact, and suggesting adjustments for future teaching methods.

- **Multimodal Learning Systems**

- » **Description:** Develop AI platforms that integrate text, images, audio, and video to create rich, immersive learning experiences. These systems could cater to different learning preferences and make complex concepts more accessible.
- » **Use Case:** A learning platform that uses virtual and augmented reality combined with AI to create interactive history lessons where students can experience historical events in a more engaging way.

## APPENDIX XV. SCIENCE, SOFTWARE AND SECURITY

### AUTHORS:

Prof. David Lo  
(SMU Information System and Technology Cluster, Centre for Research for Intelligent Software Engineering)

Assoc. Prof. Reza Shokri  
(NUS, Microsoft)

### Executive Summary

This whitepaper outlines findings from the NRF workshop on AI for Science, Software, and Security (AI4S3). The event provided a platform for over 100 participants from various sectors—including government agencies like the Home Team Science and Technology Agency (HTX) and Cyber Security Agency (CSA), universities such as Singapore Management University (SMU), National University of Singapore (NUS), Nanyang Technological University (NTU), and Singapore University of Technology and Design (SUTD) and research institutes including the Institute of High-Performance Computing (IHPC), Institute of Infocomm Research (I2R), Cambridge Centre for Advanced Research and Education

(CARES). Attendees engaged through a series of presentations and roundtable discussions. The workshop featured seven domain experts from HTX, SMU, NUS, and NTU, who delivered key insights on integrating AI at the intersections of Science, Software, and Security, specifically focusing on the Software of Science, Security of Science, Science of Software, and Science of Security. The event facilitated in-depth discussions on visions and grand challenges, the development of roadmaps, and the identification of roadblocks, setting the stage for future research and practical implementations in these critical areas.

### Introduction

AI has the potential to serve as a transformative catalyst across multiple domains. Building on the momentum of previous AI4SCI workshops, which showcased AI's capacity to revolutionize various scientific fields such as genomics and physics, this workshop explores AI's role at the convergence of Science, Software, and Security. Our discussions are structured around four primary areas: Software of

Science, Security of Science, Science of Software, and Science of Security. These are depicted in **Figure 1**, while **Figure 2** contrasts this workshop's unique focus relative to earlier AI4SCI thematic workshops, highlighting new horizontal (Software and Security of Science) and vertical dimensions (Science of Software and Security).

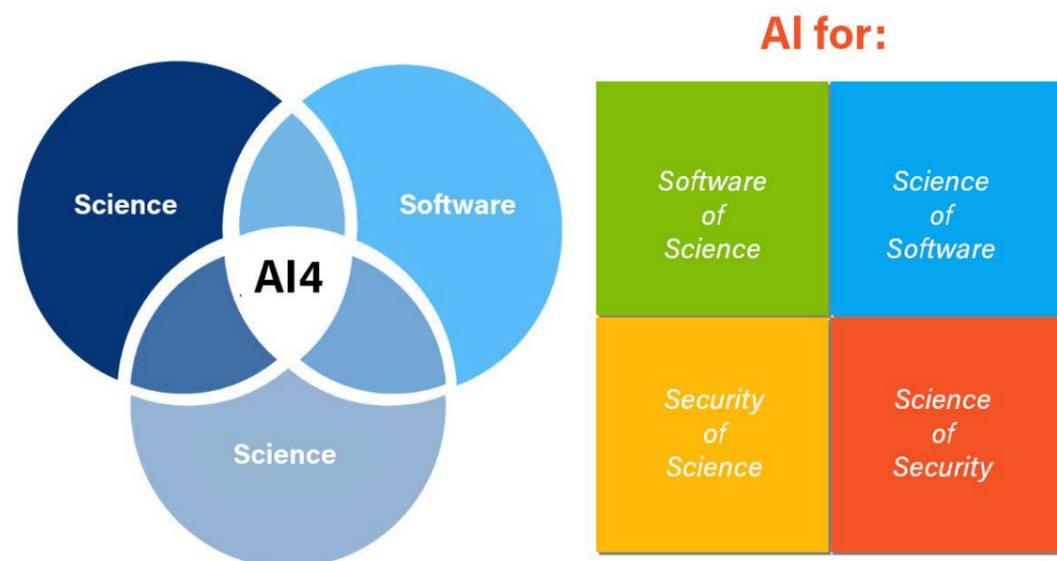


Figure 1: AI4S3 Workshop Scope

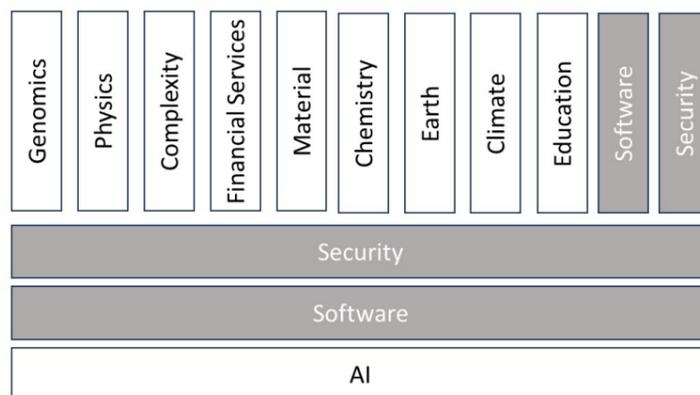


Figure 2: Unique Aspects (shaded in gray) of the AI4S3 Workshop among AI4SCI Thematic Workshops

## KEY AREAS OF FOCUS

**Software of Science:** This area involves using software to support scientific work across disciplines such as physics, biology, chemistry, and astronomy<sup>1</sup>. It encompasses activities from experimental design to data analysis and simulation. It has evolved into what is often referred to as Computational Science<sup>2</sup>. AI's integration into this field promises substantial advancements in how scientists conduct research and analyze data.

**Security of Science:** This area focuses on protecting scientific systems, data, and intellectual property from unauthorized access and applies security principles to safeguard the integrity and availability of scientific systems and data. Recent discussions<sup>3,4</sup>, notably those

led by the NATO Science for Peace and Security Programme, emphasize the growing importance of securing advancements in fields like synthetic biology and biotechnology against emerging threats.

**Science of Software:** This involves applying scientific methods to software development<sup>5</sup>. This includes both formal approaches using logical inference and mathematical modelling and empirical approaches that derive insights from real-world data and experimentation.

**Science of Security:** An interdisciplinary field that uses scientific methods to enhance the security of information systems and address challenges such as cyber threats and data breaches<sup>6</sup>. It integrates diverse disciplines, including cryptography, network security, and human factors in security.

## Background

The workshop facilitated a deep dive into how AI can catalyze advancements in these areas. Over 100 participants from various sectors attended, though space limitations

required us to select attendees for in-person participation. **Figure 3** captures a subset of these participants during the networking lunch.

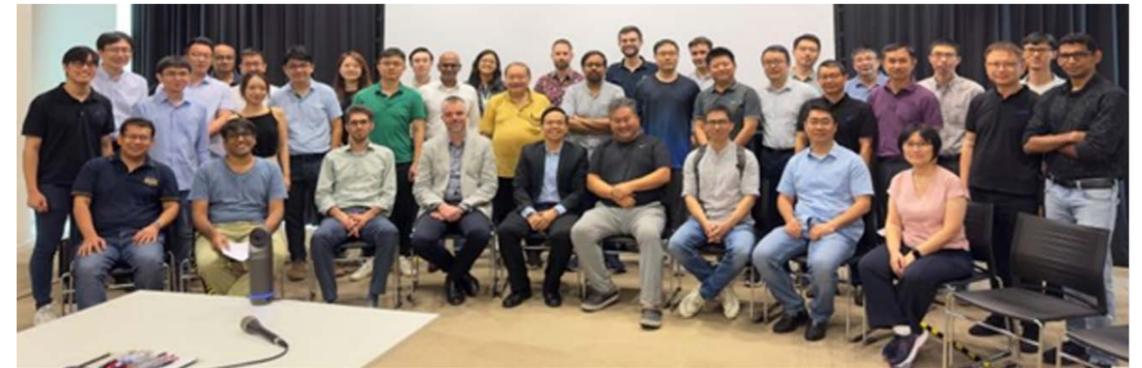


Figure 3: A Subset of In-Person Participants of the AI4S3 Workshop

The workshop explored the transformative potential of AI across the previously mentioned domains in two phases. Initially, a series of brief talks, each approximately 15 minutes long, were presented by domain experts. These sessions served to introduce the workshop's thematic areas and facilitate initial engagement among the participants. We were privileged to host a distinguished panel of speakers from various institutions in Singapore:

- AI for Catalyzing Software and Security of Science
- AI for Revolutionizing Science of Software (Formal)
- AI for Revolutionizing Science of Software (Empirical)
- AI for Transforming Science of Security

These discussions were facilitated by Swee Liang Wong (HTX) & Lwin Khin Shar (SMU), Jun Sun (SMU), Christoph Treude (SMU), and Prateek Saxena & Reza Shokri (NUS) respectively. Each group tackled three key questions, which are illustrated in **Figure 4**. The insights and conclusions drawn from these discussions are documented in Section 4, authored by the respective discussion leads with some editing from organizers.

- Computational Sciences in the Era of Generative AI, by Yi Di Yuan (HTX)
- Cybersecurity for Biotechnology Advanced by AI, by Lwin Khin Shar (SMU)
- AI for Software Science (Formal), by Jun Sun (SMU)
- AI for Software Science (Empirical), by Christoph Treude (SMU)
- AI Privacy and Science, by Reza Shokri (NUS)
- Automatically Creating Secure Code, by Prateek Saxena (NUS)
- Towards Building an AI Studio and Dataset for Security and Software Engineering Research, by Yang Liu (NTU)

Following the presentations, the workshop featured roundtable discussions focusing on four Grand Challenges:

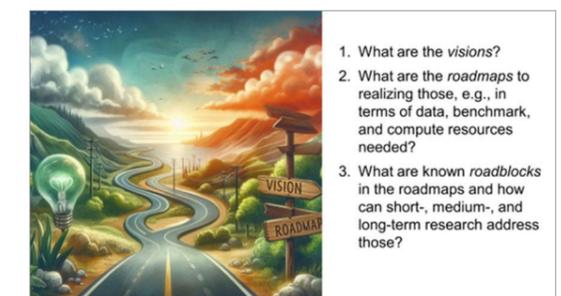


Figure 4: Three Key Questions Discussed in Each Breakout Group

## Grand Challenges

### GRAND CHALLENGE 1: AI FOR CATALYZING SOFTWARE AND SECURITY OF SCIENCE

#### OBJECTIVE

Our vision is to realize effective domain-specialized AI-powered framework, methodology, and tools for: (1) *constructing secure and safe software for physical sciences*, e.g., software for analyzing and detecting chemical signatures in drugs; (2) *detecting and mitigating security threats at the interfaces that intersect physical sciences and cyber*. The latter has been a subject of active recent research by US Army R&D and top universities around the globe (e.g., CMU, Yale, Columbia, etc.).

The research questions that we aim to address include:

- Software for computational sciences is primarily created by respective domain experts and rarely by software engineers. This often results in buggy/insecure code due to a lack of adherence to best practices. How can domain experts develop secure/safe code?
- AI models are primarily black box models. For end-user usage and trust, how to develop AI models that are explainable?
- A low-code/no-code environment would be highly desirable for experimentalists to tap into computational sciences, given that each domain is niche and would require extensive training to develop relevant software. Can generative AI help?

#### ROADMAPS

To achieve the vision discussed above, the following steps are needed:

- Acquisition of compute resources for embedding AI into the construction of software for physical sciences. The compute resources must fulfill security requirements, especially for homeland security purposes.
- Collection, curation, and secure sharing of data and benchmarks, especially in the computational sciences domain.

- Building an AI-powered collaborative platform for software engineers and computational scientists to interact to minimize the possibility of poor coding practices when developing software for science.
- Integration of AI to detect code smells, vulnerabilities, and threats into CI/CD pipelines of computational science software development.
- Design and development of AI assurance techniques to ensure that AI models used in computational sciences are trustworthy and explainable.
- Design and development of security mechanisms, as well as curation and sharing of best practices to secure both sensitive computational science data and models from leakage and backdoors, hereby preventing unauthorized or malicious access.
- Clear mapping of IP ownership and responsibility.

#### ROADBLOCKS

We see the following roadblocks that need to be overcome to realize the vision:

- Scientific domains are diverse, and their software packages are diverse, too. Construction of a shared platform and method to develop a low code environment requires careful research.
- Explainable AI models are rare for physical sciences, novel methods may be required.
- Assessment methodology for the security of open-source scientific software is nascent and thus much research is needed.
- When developing security solutions for emerging software such as computational science software, false alarms are often a huge problem for domain experts. We need to develop methods that achieve a good trade-off between precision and recall, increasing the “signal-to-noise” ratio.

Considerations on how to scale up the developed AI-powered solutions so that they are practical and useful for practitioners are essential.

### GRAND CHALLENGE 2: AI FOR REVOLUTIONIZING SCIENCE OF SOFTWARE (FORMAL)

#### OBJECTIVE

The vision is that AI should be able to solve all the problems that formal method research aims to solve, i.e., automatically synthesize software systems that are not only correct (i.e., the software system precisely implements what the user wants) but also secure (i.e., no vulnerabilities) and efficient (i.e., in terms of both time and energy consumption). Furthermore, it should do so in a way such that the correctness of the synthesized program is apparent and/or checkable, i.e., what is generated is not only the code but also different artifacts (such as documentation, contracts, tests, or formal proofs) that establish its correctness. In effect, it should realize trustworthy autonomous software engineering and automatic programming.

Such a vision is certainly demanding. A more achievable goal is (1) *to utilize AI techniques, particularly large models, to make formal method techniques and tools more accessible to ordinary users*; and (2) *to integrate AI and formal method techniques and tools such that the AI produced outputs (such as the code) are more friendly for formal methods (such as formal code verification techniques)*.

#### ROADMAPS

To achieve the vision discussed above, we aim to train dedicated large code models to assist the applications of formal methods and for generating verification-friendly code, either from scratch or based on existing open-source pre-trained models such as Llama3. To do that, we plan to take the following concrete steps:

- First, we would like to collect a high-quality dataset of verifiable code samples. Each of the data records in the dataset should contain all the artifacts of a complete software development process, including the requirement documentation, the high-level/low-level system design, the code, the test cases, the bug reports, the correctness proofs, and so on. Such data records may be more available in domains such as security protocol development.

- Second, we aim to design a model training (or fine-tuning) process that not only focuses on accuracy but also pays attention to the quality, security, and verification-friendliness of generated code.
- Third, we aim to train a large language model dedicated to the generation of correct software systems, either from scratch or by fine-tuning an existing open-source model such as Llama3, depending on the amount of computational resources.

#### ROADBLOCKS

The following challenges must be overcome to realize our vision:

- Training or fine-tuning large language models takes a lot of computational resources, and thus, we may be limited in terms of what kinds of models we can work with. Ideally, we would like to (1) train a large language model from scratch with a focus on generating correct and secure code, (2) or fine-tune an existing powerful open-source large model with enhanced capability of generating correct and secure code.
- Collecting a sufficiently large set of quality software projects with the above-mentioned artifacts may not be easy. While there is an overwhelming number of software projects on GitHub, many of them do not have a proper set of documentation, and most do not have any correctness proofs. An additional challenge is that it may not be easy to judge the quality of a GitHub project considering our requirements. One remedy is to work with software companies that have strict in-house code practices to collect such data.
- Reasoning and improving the trained/fine-tuned large language model may be challenging. That is, we anticipate that the model may not have good performance right away, and thus we must find a way to improve it iteratively. However, knowing exactly how to do so may be challenging. Possible remedies include data augmentation, systematic prompt engineering, trying different training methods, and applying model interpretation techniques such as causality analysis.

- Scalability may be challenging to achieve. That is, even with the model, it may not be feasible to generate a complete software project in one go. Our remedy would be to rely on the compositionality of software systems, i.e., to generate the software project step-by-step and component-by-component. For instance, we would first generate initial designs, system tests, and high-level code skeletons, before filling in the details of each class/function, and before generating the corresponding correctness proofs.

---

**GRAND CHALLENGE 3:**  
**AI FOR REVOLUTIONIZING SCIENCE OF SOFTWARE (EMPIRICAL)**

---

**OBJECTIVE**

Vision and challenges in the application of AI to the empirical science of software focus on deepening our understanding and improving the methodologies within software engineering research. The goal is to: (1) *enable AI to help in data analysis, pattern recognition, and even in formulating new research hypotheses*. This involves not just automating routine tasks, but *evolving AI capabilities to interactively collaborate with researchers*; (2) *develop AI systems that can interpret complex and varied data not only from the software itself but also from software development processes and provide actionable insights*. Interaction with the field of mining software repositories is crucial as it allows for the analysis of large amounts of data, which can reveal patterns that are instrumental for empirical studies. AI can play an important role by extracting and learning from these data patterns, contributing to the advancement of software engineering as a science.

**ROADMAPS**

To realize the vision, we plan to take the following concrete steps:

- Construction of extensive and diverse datasets that capture a broad spectrum of software development settings, ensuring that data collection is standardized and compatible across systems.
- Design and development of AI tools that are customized to navigate the complexities of software development data. This data is particularly challenging for AI to analyze effectively due to its diverse origins—humans, machines, and a variety of tools all contribute different types of data, such as code changes, bug reports, user feedback, and performance metrics.
- Design and development of adaptable AI systems that are quick to integrate new patterns as the field of software engineering is highly dynamic, with frequent changes in technologies and practices.
- Acquisition of computational resources needed to power AI capable of handling and interpreting the above-mentioned complex and changing datasets.

**ROADBLOCKS**

Several roadblocks could impede the aforementioned roadmaps:

- Privacy concerns and the proprietary nature of software development data can restrict access to the necessary datasets.
- The diverse nature of software practices complicates the standardization of the data and benchmarks that are needed for analysis.
- Limited access to software development practitioners for empirical data gathering poses another challenge.
- There is a risk that AI inadvertently produces misleading results or ‘fake’ scientific insights.

Addressing these challenges requires a long-term research strategy. In the short term, improving data anonymization techniques will be critical to increasing data sharing while respecting privacy. In the medium term, efforts should focus on developing AI tools capable of deriving actionable insights from non-standardized, complex, and even noisy data. Long-term initiatives should aim to establish industry-wide standards for data and benchmarks, fostering an ecosystem that supports AI-driven empirical science of software.

---

**GRAND CHALLENGE 4:**  
**AI FOR TRANSFORMING SCIENCE OF SECURITY**

---

**OBJECTIVE**

In an era of rapidly evolving technology, the security of information systems has become a paramount concern. The Science of Security seeks to address this challenge through an interdisciplinary approach that integrates scientific methods to enhance cybersecurity. The vision for AI for Science of Security is rooted in leveraging the power of artificial intelligence (AI) to transform the landscape of cybersecurity.

- *AI as a Proxy for Human Experimentation and Decision-Making in Security*: AI has the potential to act as a proxy for humans in conducting experiments and making decisions. Many of the security vulnerabilities need to be analyzed in situations where humans are involved. Conducting human studies can be costly and time consuming. We can use AI to simulate humans. Besides, by inferring rules, specifications, and policies, AI can significantly accelerate the process of identifying and addressing security vulnerabilities. This capability can lead to more efficient and effective security measures, reducing the reliance on human expertise and minimizing human error.

- *Usability of AI in Security*: For AI to truly revolutionize the field of security, it must be more usable. This involves developing intuitive interfaces and systems that allow security professionals to interact seamlessly with AI tools. Enhanced usability will ensure that AI can be integrated into existing workflows, making it a practical and indispensable tool for cybersecurity.
- *Rigorous Analysis of Security Algorithms*: AI can serve as a powerful tool for conducting rigorous analysis of algorithms, particularly in cryptographic protocols. By automating the process of algorithmic analysis, AI can provide deeper insights and identify potential weaknesses that might be overlooked by human analysts. This capability is crucial for developing robust cryptographic solutions that can withstand sophisticated cyber threats.
- *Safety Mechanisms for AI*: While AI offers immense potential, it is essential to incorporate safety mechanisms to handle scenarios of catastrophic failures. A ‘turn off’ switch for AI systems ensures that in the event of unexpected behavior or security breaches, AI can be quickly deactivated to prevent further damage. This safety measure is critical for maintaining control over AI systems and ensuring their reliability in security applications.

**ROADMAPS**

To achieve the vision for Science of Security, a strategic roadmap is necessary, addressing short-term, medium-term, and long-term goals:

- In the short term, feasible goals include high-level simulation of user behavior, automated patching, testing, and creation of benchmarks. These initiatives can provide immediate benefits by improving the detection and mitigation of security vulnerabilities.

- Medium-term efforts should focus on addressing the lack of data in certain domains and developing scientific measures of software security. Additionally, improving data-sharing practices, while respecting privacy and intellectual property concerns, is vital for advancing AI's capabilities in security.
- In the long term, addressing the "last mile problem" is crucial. While AI is proficient in general-purpose tasks, it struggles with specialization to specific tasks. Ensuring AI can reliably identify all bugs, not just some, is a significant challenge. Building trust in AI's ability to spot and mitigate all potential security issues is essential for its widespread adoption.
- The effectiveness of AI in security tasks remains a contentious issue, with divided opinions among experts. While some are optimistic about AI's potential, others remain skeptical, highlighting the need for further research and development to prove AI's capabilities in real-world security scenarios.

---

### DATA NEEDS

---

As described earlier; in addressing data needs for computational science, the availability of high-quality data remains a persistent challenge. Many essential datasets are not readily accessible, necessitating collaborative efforts with agencies and companies involved in developing such software tools as a partial solution. Furthermore, the task of amassing a comprehensive collection of quality software presents another hurdle. Despite the abundance of projects on platforms like GitHub, many lack essential documentation or validation proofs, complicating efforts to identify reliable resources. To mitigate this issue, partnerships with companies adhering to stringent code practices can offer some relief. Similarly, acquiring a robust dataset of quality security information, including vulnerabilities, proves challenging beyond what is typically available from sources like the National Vulnerability Database (NVD).

Singapore emerges as a promising hub uniquely equipped to tackle these challenges. Singapore has extensive data collection and refinement expertise for software and security purposes. Additionally, many researchers in Singapore are at the forefront of innovations in data privacy technologies, positioning the nation as a leader in effectively addressing the aforementioned data challenges.

### ROADBLOCKS

Despite the promising vision for Science of Security, several roadblocks hinder the full realization of its potential:

- One of the primary challenges is that AI can exhibit human-like mistakes, making similar errors in judgment and decision-making. This limitation undermines the reliability of AI systems, as they may not consistently provide accurate or effective security solutions.
- AI also possesses unique fundamental flaws, such as vulnerabilities to adversarial inputs and issues with memorization. These flaws differ from human errors and present new challenges in ensuring the robustness and reliability of AI systems in security contexts.
- Currently, AI is not adept at planning and system design. This limitation restricts its ability to autonomously develop comprehensive security strategies and solutions, necessitating significant human intervention and oversight.

---

### HPC NEEDS

---

In assessing the high-performance computing (HPC) requirements, our estimates indicate substantial annual computing needs across various resources. Specifically, our projections suggest a demand of 66.4 million core hours for CPUs, complemented by 2.8 million card hours for GPUs (including A100, H100, RTX A5000, and RTX A6000 models), alongside a requisite storage capacity of 2400 terabytes.

To contextualize these figures, we operate under certain assumptions. We anticipate the involvement of two teams dedicated to each of

the four areas illustrated in Figure 1 (rightmost diagram), resulting in a total of eight teams. Each team is comprised of approximately 10 members actively engaged in HPC activities, primarily focused on training and fine-tuning AI models, including small- and medium-sized (e.g., up to 20B parameters) AI models.

Our estimations are grounded in insights gleaned from prior research collaborations. For instance, a recent partnership with GovTech's Cybersecurity Group<sup>2</sup> that constructs a foundational model optimized for GovTech's specific operational needs, involving a team of 10 researchers.

## AI Methods and Data – Challenges and Opportunities

---

### AI METHODS

---

To advance the intersection of science, software, and security, various AI methods can be leveraged across the above-mentioned four areas: (1) Software of Science, (2) Security of Science, (3) Science of Software, and (4) Science of Security.

In the first two areas, the application and adaptation of AI are still emerging, presenting numerous opportunities for diverse AI methods. Predictive and generative AI techniques can enhance tasks related to building and securing software for physical sciences. Traditional AI methods, such as random forests and support vector machines, remain valuable due to their simplicity and interpretability, which can facilitate communication with scientists. Additionally, more advanced AI approaches, including specialized Large Language Models (LLMs) and Vision-Language Models (VLMs) fine-tuned on computational science code,

hold significant promise. Reinforcement Learning (RL), particularly Reinforcement Learning from Human Feedback (RLHF), offers potential for creating adaptive systems that learn from domain experts, such as scientists and software engineers, to collaboratively build and secure essential software for physical sciences.

In the last two areas, traditional AI methods have been employed to address various software and security challenges. Despite this, persistent issues in software quality continue to impact the economy and society, e.g., the recent global disruptions to airlines, banks, hospitals and government offices.<sup>1</sup> The advent of generative AI introduces new possibilities for tackling these problems, while also presenting challenges such as privacy concerns and inherent weaknesses in AI algorithms. This underscores the need for developing more privacy-preserving and robust AI solutions.

## Singapore's Role

There are many reasons why Singapore is an excellent place to conduct deep research in the AI for Science, Software, and Security (AI4S<sup>3</sup>) direction.

First, Singapore has a distinctive culture of strong government-academia partnership that is not readily replicable in many places. This is also true in the intersection of AI, Science, Software, and Security. For example, Singapore has a unique government agency, the Home Team Science and Technology Agency (HTX), that focuses on the intersection of Science, Software, and Security for real applications and is rapidly harnessing the power of AI. "HTX is the first of its kind Science and Technology Agency in the world that brings together science and engineering capabilities across the Home Team Departments to transform the homeland security landscape and keep Singapore safe."<sup>3</sup> Beyond HTX, there are many other government agencies (e.g., CSA, GovTech, DSTA, DSO, etc.) that are also actively working on various intersections of science, AI, software, and security. These government agencies have worked closely with academia for years, resulting in many examples of concrete impact on society. Singapore is thus an excellent place to invest in AI4S<sup>3</sup> direction as any deep research findings uncovered can be readily translated to impact practice and society.

Second, Singapore has strong science, AI, software, and security ecosystems. For example, Singapore is ranked highly among

the world research leaders in AI, Software, and Security. To illustrate, an authoritative Computer Science ranking put NTU as 7 worldwide, SMU as 4 worldwide, and NUS as 21 worldwide in AI<sup>4</sup>, Software<sup>5</sup>, and Security<sup>6</sup>, respectively.

An African proverb mentions, "If you want to go far, go together.". The same applies to advancement in this AI4S<sup>3</sup> direction. For us to make headway in this important direction, there is a need for collaborations between universities (e.g., SMU, NUS, NTU, SUTD, etc.) and research institutes (e.g., A\*STAR) in Singapore. Moreover, there is a need for active partnerships with relevant government agencies (e.g., HTX, CSA, etc.) that can help bridge basic research into societal impact. Of course, industry involvement is also needed. And many companies in Singapore have started to spend research efforts in the AI4S<sup>6</sup> direction; these include ST Engineering, Accenture, Salesforce, SAP, etc.

We need to also look beyond Singapore. There will be great value if collaborations are established with well-known research centers from academia and industry who have also started to invest much in the AI4S<sup>3</sup> direction. These include the Software Engineering Institute (SEI)<sup>7</sup>, Security Privacy Institute (CyLab)<sup>8</sup> at Carnegie Mellon University, Ivado<sup>9</sup> (Canadian academia-industry AI consortium led by Turing laureate Yoshua Bengio), Microsoft Research, Meta, etc.

## Conclusion

The workshop has effectively underscored the immense potential of AI as a transformative catalyst across four critical intersections of Science, Software, and Security (S3): Software of Science, Security of Science, Science of Software, and Science of Security. Each area can be mapped to a grand challenge and presents unique visions, the achievement of which could significantly impact both

the economy and society. Roadmaps were proposed to realize these visions. Additionally, roadblocks were identified, identifying key issues to be addressed in the short, medium, and long term. Importantly, the discussions underscored the necessity of incorporating AI4S<sup>3</sup> into Singapore's strategic plans for future AI research and infrastructure investments in the Science sector.

## REFERENCES

- 1 William H Press, William T Vetterling, Saul A Teukolsky, and Brian P Flannery. 1988. *Numerical recipes*. Cambridge University Press.
- 2 Angela B Shiflet and George W Shiflet. 2014. *Introduction to computational science: modeling and simulation for the sciences*. Princeton University Press.
- 3 Dov Greenbaum. 2023. *Cyberbiosecurity: A new field to deal with emerging threats*. Springer Nature.
- 4 Benjamin D Trump, Marie-Valentine Florin, Edward Perkins, and Igor Linkov. 2021. *Emerging threats of synthetic biology and biotechnology: addressing security and resilience issues*. Springer Nature.
- 5 Maurice H Halstead. 1977. *Elements of Software Science (Operating and programming systems series)*. Elsevier Science.
- 6 JASON Program Office. 2010. Science of Cybersecurity (JASON Report JSR-10-102). <http://fas.org/irp/agency/dod/jason/cyber.pdf>. (Accessed on May 20, 2024).

# APPENDIX XVI. ROBOTICS



## AUTHORS:

Dr Dongkyu Choi  
(Samsung Electronics, ex-A\*STAR)  
Prof David Hsu  
(NUS, Smart Systems Institute)

Prof Lee Wee Sun  
(NUS)

## Executive Summary

Today, with the explosive growth of AI, large language models (LLMs), e.g., GPT, can summarize business reports, answer questions about common illnesses, compose legal arguments, converse with humans on philosophy, etc. Yet, an embodied AI agent, such as a robot, equipped with all the power of GPT, cannot tidy up a messy kitchen or build a house. This disparity highlights the need for more *general embodied intelligence* within the context of robotics. Embodied intelligence integrates cognitive capabilities and physical embodiment to create robots that can perceive, understand, and interact with the *physical world and humans*. It goes

beyond traditional artificial intelligence (AI) by emphasizing the symbiotic relationship between a robot's physical presence and its cognitive functions. This integration is crucial for developing human-centered robots that perform complex tasks autonomously and adapt to the dynamic, physical world.

We propose to bring the research community in Singapore closer to the goal of general embodied intelligence by leveraging the latest developments in AI. With the participation from industry and government stakeholders, our efforts will lead to real-world use cases that make an impact on daily lives of our citizens.

## Introduction

The recent popularity of AI, especially those techniques related to large language models (LLMs), brought forth a lot of rosy promises in various areas of academic research and industrial applications. But the reality is that these models are inadequate for physically grounded tasks. The *general embodied intelligence* that humans possess enables them to understand the physical world and take actions in it for a variety of tasks. Having this type of intelligent capabilities in robots will transform industries such as healthcare, manufacturing, transportation, environmental sustainability, facilities management, and others. It will enable robots to perform tasks beyond the reach of conventional AI systems,

by providing solutions for assisting the daily life of the elderly or collaborating with human workers in automating complex manufacturing processes<sup>1</sup>.

To realize this enormous potential, however, we must tackle two "disconnects" that currently exist. First, while LLMs are believed to capture all past experiences symbolically in the abstract, robots must understand and act on the physical world *in situ*. Second, to operate in natural human environments, robots must understand human intentions and preferences, while communicating with humans effectively.

## Background

Recognizing enormous opportunities and challenges, we have organized a series of three workshops on robot foundation models. The first two, in March 2024 and July 2024, respectively, attracted over 300 participants from Singapore universities, research institutes, government agencies, and the industry. The third and the last workshop is planned for the beginning of 2025.

- **First Workshop on Survey of Foundation Models for Robotics:** This inaugural workshop in the three-event series focused on brainstorming ideas from Singapore's members of academia. The event was well attended by about *100 researchers from A\*STAR and universities* with some participants from various government entities. The core team presented their initial results from the NRP-funded survey and suggested seven categories to organize previous works. They also introduced their ongoing efforts on selective replications of reported systems. Other participants who submitted their short slide decks in advance described summaries of their ideas to jumpstart the group discussions.
- **Second Workshop on Foundation Models for Robotics:** The second event aimed at capturing industry use cases, around which several teams of researchers can start

designing their proposals for follow-up projects. There were many more registrants than the anticipated 100 attendees, resulting in the expansion of the venue. The final number of attendees were *155 from relevant industries/government agencies and 77 from the research community*. After keynote speeches by Prof Lee Wee Sun from the core survey team and a researcher from NVIDIA, the leads from the four proposal teams presented an overview of their respective research plan. The proposal topics included fundamental research (led by Prof David Hsu, NUS), healthcare (led by Prof Ang Wei Tech, NTU), manufacturing (led by Dr Zhang Jingbing, ARTC), and facilities management (led by Dr Yau Wei Yun, I2R).

- **Third Workshop on Foundation Models for Robotics:** The planned final workshop around the year end will summarize the core teams survey results and their recommendations for NRP and the Singaporean ecosystem for robotics research.

This whitepaper is informed by the outcomes from these events and serves as the core team's understanding of the challenges and potential solutions based on their survey and engagements to date.

## Grand Challenges

The grand vision for general embodied intelligence is to connect robots and humans in *mind* and *body* to augment and expand human capabilities in the physical world.

### GRAND CHALLENGE 1: ROBOT BUTLERS FOR HOMES

Home is where we spend long hours of our daily lives. As we know, maintaining one is a challenging mission that often requires dedicated manpower. Typical tasks include cleaning rooms, organizing stuff, doing laundry, preparing meals, receiving deliveries, walking dogs, and many others. There are automated solutions for some of these

tasks, like robot vacuums and smart washing machines, but they are specialized machines that perform limited functions and often require human preparation or intervention for successful operation.

### OBJECTIVES

In simple terms, we would like to have a robot butler that can manage our homes even when we are away on vacation. For this ambitious goal, of course, we need two important capabilities at the minimum:

- natural communication and common-sensical understanding of everyday tasks
- understanding and executing physical actions for everyday tasks.

---

**GRAND CHALLENGE 2:  
INTERACTIVELY TEACHABLE ROBOTS  
FOR MANUFACTURING**

---

Robots currently used for manufacturing are pre-programmed to work in isolation for the purpose of automation. But there are tasks that are not automatable, like those that require frequent changes in product design or case-specific decision making. For example, high-mix, low-volume manufacturing of consumer goods and remanufacturing of defective products require human expertise to adapt to new situations. In such cases, it will help to have robots that are capable of learning through interactive teaching as if they were human teammates.

**OBJECTIVES**

To be effective in fast changing manufacturing scenarios, the robot needs to:

- Be teachable through a combination of language instructions, visual observations, and possibly remote-controlled demonstrations.
- Learn extremely quickly, preferably from one or a small number of demonstrations together with a small amount of language instructions.

## AI Methods and Data – Challenges and Opportunities

---

**ROADBLOCKS**

---

**DATA-SCARCITY**

While natural language processing and computer vision have benefited enormously from the internet-scale data readily accessible, robotics faces a fundamental difficulty: the high cost of gathering large amounts of data from the physical world<sup>2</sup>.

---

**GRAND CHALLENGE 3:  
ROBOT GUIDE DOG FOR VISUALLY  
IMPAIRED**

---

A robot guide dog surpasses the ability of its well-trained animal counterpart and empowers the visually impaired to lead an independent life: travel to a community center by bus, explore a new shopping center, find an empty seat in a busy hawker center, and visit a hospital clinic. Current AI technology would already enable a robot guide dog to communicate better with the visually impaired than a real guide dog, but the embodied capabilities remain a challenge.

**OBJECTIVES**

An effective robot guide dog would at least need the capabilities to:

- Navigate and guide humans in an intuitive manner through diverse environments, some of which it has never seen before.
- Have a good understanding of the needs of the human it is guiding in situ, and be trusted to guide the human safely.

**DATA REQUIREMENTS**

The need for multimodal data is key (images, volatile organic compounds, sounds), including 3-D physical data which is paramount for robot. Both real-life data and simulators to approximate the physical world and generate synthetic data will be necessary.

**GENERALIZATION**

With sufficient data, robots today can succeed on many narrow, specific tasks seemingly difficult for humans. Generalization is, however, critical as it enables robots to apply learned knowledge to new, unseen tasks and environments, ensuring adaptability and robustness in diverse real-world scenarios. Generalization is key to intelligence and a ubiquitous requirement of new opportunities for robot deployment, from flexible manufacturing to homecare robots.

**ROBUSTNESS AND SAFETY**

Robustness ensures reliable performance and resilience in the face of unpredictable and dynamic real-world conditions, thereby enhancing safety and operational effectiveness.

**TASK MODEL SPECIFICATION**

Specifying (formal) models for robot tasks is difficult because of the inherent complexity and variability of real-world environments, which require robots to handle a wide range of unpredictable situations. Further, accurately modeling all the physical, cognitive, and sensory interactions involved makes it hard to create comprehensive and precise formal models.

**TRUST IN ROBOTS**

Robots are very complex engineering systems. Incorporating foundation models further increases this complexity. Users need to understand how robots make decisions and be confident in their consistent performance and adherence to safety standards. Building this trust requires clear communication, robust testing and validation, and addressing ethical and privacy concerns, while also providing users with positive, seamless interaction experiences with the robots.

**ROADMAP**

The success of these systems hinges on advances in a common set of core technology thrusts that can help us address the roadblocks listed above.

- *Data*. One major obstacle to robot learning is the lack of data. Compared with the internet-scale data for text, images, and video, the amount of 3-D physical data required for robot learning is miniscule. We need new mechanisms to scale up data gathering efforts vastly at a low cost. We need simulators to approximate the physical world and generate synthetic data. We also need novel ways to combine limited 3-D physical data with the more abundantly available textual and visual data for robot learning.

- *Simulation*. An orthogonal way of mitigating the requirements of real-world data is to build digital twins autonomously: leverage foundation models to create real-time dynamic replicas of diverse and complex physical environments. These digital twins are augmented with differentiable simulators (e.g.,<sup>3</sup>), allowing robots to receive instant feedback and adapt accordingly. The simulator may be differentiable to further enhance performance. This approach facilitates extensive training without real-world constraints and allows virtual rehearsal before actions are executed in the physical world.
- *Knowledge*. Foundation models, such as large language models, are trained on vast amounts of diverse data. They provide broad but potentially ungrounded knowledge of the physical world and the ability to communicate with humans fluently in natural languages<sup>1</sup>. By integrating these models into robot systems, we can exploit their extensive knowledge to enhance robotic system performance through a multitude of different roles: a commonsense knowledge base, a natural interface of communication, a “black-box” control policy, a general-purpose heuristic for reasoning,...
- *Learning*. Foundation models excel in few-shot in-context learning for specific task domains<sup>4</sup>, whether healthcare, manufacturing, etc.. By incorporating these models, robots can quickly learn new skills or adapt to new tasks with limited training data. This capability is particularly beneficial in dynamic environments where robots must be versatile and responsive to changing requirements.

- *Reasoning.* Systematic reasoning is a hallmark of intelligence. One popular approach today is to use LLMs directly as a black-box control policy<sup>5</sup> and query it for the next robot actions. It is fast and works surprisingly well for simple tasks, by exploiting LLMs' vast past experiences and commonsense knowledge. With increasing environmental variability and task complexity as well as limited data, the black-box policy quickly fails. We need a more sophisticated, powerful approach: use LLMs' knowledge to build a world model and apply a planning algorithm to reason about the model systematically<sup>6</sup>. While systematic reasoning is slow, it generalizes to arbitrary complex tasks by composing finite elemental knowledge. We need new robot architectures, possibly neuro-symbolic, to integrate both the fast

and the slow reasoning modes in order to achieve general-purpose robot intelligence.

- *Human interaction and collaboration.* Foundation models, particularly LLMs, have revolutionized human-robot interaction by enabling natural and intuitive communication between humans and robots, making it easier for users to give commands, ask questions, and receive responses in everyday language [1]. This improves user experience, reduces the learning curve for operating robots, and increases accessibility. Additionally, LLMs can interpret context and intent more accurately, allowing robots to perform complex tasks and provide more personalized assistance, ultimately fostering greater trust and collaboration between humans and robots.

## Singapore's Role

Recent dramatic advances have enabled AI agents to become useful assistants in white collar jobs. Embodied Intelligence will have an equally dramatic impact in providing assistance in the physical world. AI capabilities in embodied intelligence are currently far from its capabilities in the virtual world. Singapore should invest appropriately in embodied AI research so that it is well prepared when the technology starts to make substantial economic impact.

Singapore is highly suitable as a testbed for embodied intelligence technologies. The Singapore society is ageing, resulting in difficulties in obtaining sufficient physical manpower and incentivizing experimental

deployment of robotics technologies. As an indication of this, Singapore has the second highest robot density per employee in the world, behind Republic of Korea<sup>7</sup>. The strong industry interest is also reflected in the large number of attendees in our industry workshop in July 2024.

There is a strong research ecosystem in robotics in the universities (NUS, NTU, SUTD) and A\*STAR, supported by the National Robotics Programme. There is currently an opportunity for rapid research progress in embodied intelligence building upon the recent disruptive advances in AI foundation models and Singapore is well positioned to take advantage of the opportunity.

## Conclusions

Our proposed roadmap aims to bring the research community in Singapore closer to the goal of general embodied intelligence. With the participation from industry and

government stakeholders, our efforts will lead to real-world use cases that make an impact on daily lives of our citizens.

## REFERENCES

- 1 R. Bommasani et al. On the opportunities and risks of foundation models. arXiv:2108.07258, July 2022.
- 2 Y.i Hu et al. Toward general-purpose robots via foundation models: A survey and meta-analysis. arXiv:2312.08782, 2023.
- 3 Z. Huang, Y. Hu, T. Du, S. Zhou, H. Su, J.B. Tenenbaum, and C. Gan. PlasticineLab: A soft-body manipulation benchmark with differentiable physics. In Proc. Int. Conf. on Learning Representations, 2020.
- 4 Q. Dong et al. A survey on in-context learning. arXiv:2301.00234, 2022.
- 5 Brohan, A., Brown, N., Carbajal, J., Chebotar, Y., Chen, X., Chormanski, K., Ding, T., Driess, D., Dubey, A., Finn, C. and Florence, P., 2023. RT-2: Vision-language-action models transfer web knowledge to robotic control. arXiv preprint arXiv:2307.15818.
- 6 J. Sun et al. A survey of reasoning with foundation models: Concepts, methodologies, and outlook. arXiv:2312.11562, Jan 2024.
- 7 International Federation of Robotics. Global Robotics Race: Korea, Singapore and Germany in the Lead, <https://ifr.org/ifr-press-releases/news/global-robotics-race-korea-singapore-and-germany-in-the-lead>, Jan 10, 2024

# APPENDIX XVII. AI METHODS AND MATHEMATICS



## AUTHORS:

Assoc. Prof Bryan Low  
(NUS AI Institute)

Prof Ong Yew Soon  
(NTU, A\*STAR)

Prof Ng See Kiong  
(NUS)

Assoc. Prof Xavier Bresson  
(NUS, NTU)

Asst. Prof Li Qianxiao  
(NUS, IFIM)

## Executive summary

This report outlines important directions for methodology research that supports emerging applications of artificial intelligence (AI) in science and engineering. Five complementary directions are discussed, including how to combine AI with existing simulation and

experimental pipelines, and how to interface physical knowledge and AI to complete the cycle of learning and experimentation for scientific discovery. We also discuss how these methodology research topics can be used to support the local research scene.

## Introduction

In the past decade, the development of AI methods has been primarily driven by its application in computer vision, robotics and natural language processing. The associated techniques, including Image/video analysis using deep learning<sup>1</sup>, deep reinforcement learning<sup>2</sup>, and training or fine-tuning large language models<sup>3</sup>, have since matured. However, with the growing interest in AI for science and engineering, new challenges arise in the development of AI theoretical frameworks and algorithms.

Several new issues in the application of AI to scientific problems necessitates new methodology research. First, compared to vision and language applications, the scale of

existing datasets in science and engineering research is much smaller. To resolve this, two approaches may be taken simultaneously: 1) we need to develop new algorithms that can operate on small datasets, and 2) we may also establish AI-augmented workflows to speed up the generation of computational and experimental data. Second, compared with existing AI applications where prediction and other quantitative fidelity is of utmost importance, in AI for science one also desires a certain level of interpretability of the resulting data-driven models, either to enable human understanding of the physical phenomenon being studied, or for bridging models or workflows across related application domains.

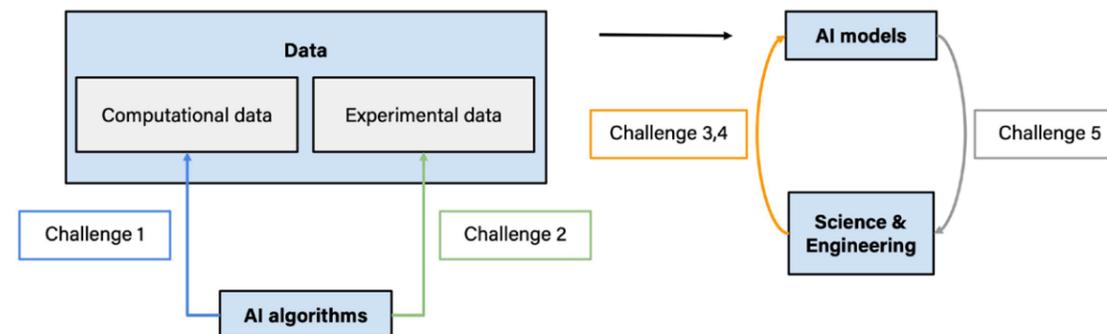


Figure 1: Key challenges in AI methods in science.

To address these issues, in this report we identify some key challenge statements for the methodology development of AI for science (Fig.1), outline the main problems and possible solutions. We also discuss the potential impact these methodological breakthroughs may enable in a variety of applications across science and engineering.

In the next section, we outline several methodological challenges that arise from AI applications in science and engineering, including a discussion of the current state of the art, possible approaches, their potential impact and expected outcomes.

## Grand Challenges

### GRAND CHALLENGE 1: AI-AUGMENTED SIMULATION OF COMPLEX PHYSICAL PROCESSES

#### BACKGROUND

Riding on the rapid advancements of machine learning model design in AI research, there is little shortage of large-scale, highly performant machine learning model architectures and training algorithms for the analysis of complex physical processes. Consequently, the key challenge shifts towards how to supply high quality, high throughput datasets to these learning algorithms, and how to ensure that the predictive results can lead to advancements in conceptual understanding – a problem of vital importance in the advancement of scientific knowledge. A promising line of attack to resolve this problem is to turn to high-fidelity computational data, e.g. from computational fluid dynamics (CFD) and density functional

theory (DFT) calculations, which provide accurate simulated environments to perform data-driven scientific discovery.

However, a key roadblock here is that these simulations are extremely expensive, and it is almost impossible to use them as a tool to probe the properties, let alone design, complex physical systems of scale. The grand challenge is thus: *can we use AI to accelerate simulations in scientific computing for knowledge extraction?* A concrete measure of this would be the ability to simulate and understand physically realistic systems at scale, e.g. to design a functional material from microscopic simulations and/or reduced order modeling. Resolving this challenge will have far reaching consequences for both computational science itself and the dataset supply for AI for science subsequently.

## STATE OF THE ART

Currently, there are many works on accelerating classical scientific computing using AI and machine learning. One can roughly classify them into two types: complete replacement or partial replacement. The first completely replaces, typically with supervised learning or generative modeling, a classical algorithm by a machine learning one. An example is learning energy predictions from DFT data<sup>4,5</sup>. The main issue with this approach is the general lack of data for training high-quality, generalisable models. The second approach is to interface data-driven and classical algorithms, often using AI to replace part of the expensive process of its solution<sup>6</sup>. Here the issue to resolve is how to do so while guaranteeing the stability and convergence of the algorithms. Moreover, these often require interfacing learning algorithms that require automatic differentiation (and some even require fast integration in moderately high dimensions), and traditional scientific computing ones which do not. This is another issue to be resolved that will greatly expand the applicability of AI in computational science. Finally, one still needs principled methods to summarize, interpret and exploit the vast computational data that will be generated. Currently AI-augmented reduced order modeling<sup>7</sup>, inverse design<sup>8</sup> and control<sup>9</sup> is an active area of research that needs to be further developed in order facilitate scientific understanding from data<sup>10</sup>.

## SINGAPORE'S ROLE

Singapore is well equipped to tackle this challenge. There are a number of research groups across university departments (Engineering, mathematics departments in NUS and NTU, NUS IFIM, etc.), public research entities (IHPC, IMRE in A\*STAR, etc.) and industry (SEA, Nvidia, etc.) who are interested in machine learning augmented scientific computing, as well as methodologies for subsequent analysis, e.g. reduced order modeling, inverse design and control. National computing resources provided by NSCC are also a natural platform that enable this

effort. Finally, this research closely supports many relevant AI for science domains, such as AI for synthetic biology and drug design, AI for materials science and chemistry, AI for chemical and biological manufacturing, etc., which all will benefit from increased simulation efficiency and automatic knowledge extraction. Besides developing the research manpower, to achieve this goal it is also required to expand national computing capabilities in NSCC, especially accelerated parallel computing (GPU) platforms.

## OBJECTIVE

We expect this research endeavor to serve as a cornerstone support for AI for science research efforts across many domains of interest. The overarching goal is to provide an efficient computational software infrastructure to enable fast simulation and analysis in a variety of domains ranging from materials science, chemistry and engineering, and to effectively interface the computational data generated with data-driven analysis, design and control.

---

## GRAND CHALLENGE 2: SCIENCE-INSPIRED EXPERIMENTAL DESIGN BACKGROUND

---

Despite improvements of AI towards modeling and making predictions about physical systems, it is still necessary to conduct physical experiments to high-fidelity generate data and verify the actual parameters of interest. Since these experiments can often be costly to obtain, either from required resources or the timescales involved in running the experiments, it is important to carefully design which experiments are important in order to obtain relevant information about the system that we are studying. The problem of experimental design (ED) is important in all disciplines of science, from agriculture, geology, or engineering, with existing ED methods covering a large range of problem settings in adaptive or non-adaptive learning, aiming to solve inverse problems, perform design optimization and more.

There are two main issues in performing ED for a general scientific setting. The first issue is that the ED problem can involve complicated search spaces, constraints or objective functions which may vary across each application. In agriculture, for example, experiments may involve varying soil temperature, humidity, mineral components, and other environmental conditions, resulting in a high-dimensional search space which may be discrete or continuous. In manufacturing applications, there may not be a clear single objective function which should be optimized for, and researchers have to select optimal designs which tradeoff between different desired properties. Performing optimization over such complex problem settings remains a challenge.

The second issue is that there are often practical constraints, often arising from field scientists, which restricts how the experiments may be conducted. Even though there are many feasible settings that the experiments can be conducted in, not all of these can be conducted with the same ease, or result in measurements that are of the same fidelity. Additionally, due to the difficulty in running experiments, it may be preferable that fewer design parameters are tested, or that multiple design parameters are all tested at once instead of in an adaptive setting, which would be more convenient for the scientists.

These two issues highlight a fundamental problem that while there are proposed ED algorithms for different conditions and objectives, there is often still a disconnect between methods that work under toy settings and methods that would be more applicable in practice. The big challenge is therefore to bridge the gap between theoretically sound ED methods and practical ED settings which match the constraints and requirements of practical experimental settings. There is currently a lack of realistic, real-life scientific dataset or benchmark that allows researchers to test and develop ED algorithms for these settings.

## SINGAPORE'S ROLE

Singapore is able to lead this effort well due to existing expertise both on the machine learning side, as well as in domain-specific scientific applications where ED problems are applicable. Research groups in public research entities (e.g. A\*STAR), industry, and academia conduct research in relevant scientific domain areas such as material science, synthetic biology and drug design, and precision engineering, as well as in relevant computer science domains such as active learning and experimental design.

## OBJECTIVE

We believe that to accelerate the development of practical ED methods, an effort should be placed in building up realistic, real-life datasets and/or self-driving labs which are good reflections of the challenges and constraints that are required by scientists in practical ED problems, such as complexity of search space or of objective functions, or constraints on feasible design parameters. These datasets could be hosted in a platform ('gym') where researchers could use, to test and develop their experimental design algorithms. This would allow machine learning researchers to more easily develop methods that can cater to the real use case and settings.

---

## GRAND CHALLENGE 3: PHYSICS INSPIRED NEURAL ARCHITECTURES

---

## BACKGROUND

Improving the rate and quality of computational and experimental data generation does not completely resolve the challenge of the lack of large datasets. In fact, there remains widespread challenges in the use of many of these AI methods in the diverse, complex scientific problems of today as data is expensive to acquire for many real-world complex scenarios. Methods that are currently the main focus of development in the AI/ML community are not fit-for-purpose when it comes to scientific data-sets that tend to be orders of magnitude smaller than what is typical in domains like natural language, and have underlying structures and representations that are also vastly different.

Given the anticipated data limitations in many domains, multiple strategies must be explored to mitigate this key challenge. Viable potential solutions include developing capabilities to generate synthetic data through scientific computing (through our current theoretical understanding of the system) and ensuring data generated is maximally useful (e.g. through Bayesian methods).

Physics-inspired neural networks (or more generally, knowledge-guided neural networks) are a potential alternative solution where the goal is to improve the performance of AI models through incorporation of priors as learning or inductive biases (e.g. imposing known governing equations or conservation laws on a model)<sup>11</sup>.

A key advantage of this is to mitigate data sparsity, as it is anticipated that the insertion of such knowledge can improve model performance given a limited dataset. These models can thus also be synergistic with other challenges (e.g. assist in generation of a larger dataset for a particular problem through AI). Another key advantage to such models is the fact that these models are explicitly trained to mitigate violation of key constraints, thereby also providing some degree of robustness in the model predictions (e.g. PINNs can learn models that will avoid violations of conservation laws in their predictions).

#### STATE OF THE ART

While there is a lot of prior knowledge in each domain that can be exploited, it is not always easy to incorporate them into the AI model. A key advantage of PINNs is that various kinds of knowledge can all be flexibly incorporated (e.g. ranging from the governing PDEs to simple empirical laws). Inserting physics information in PINNs is conceptually simple, but we need more advances to reduce the training (optimization) cost. There also remains the potential for more innovative ways to use our prior physics knowledge or theory to reduce the training cost of such PINN models, whether it's in the context of transfer learning, physics-derived representations, or better model initializations for optimization<sup>12,13,14,15</sup>. Current methods also focus more on mathematical priors (e.g. governing equations). However,

we can look beyond including equations, but also incorporate numerical models in some instances (synergies between scientific computing and machine learning)<sup>16,17</sup>.

Physical models (e.g. theories) impart guardrails, and PINNs are a good means of utilizing such guarantees in ensuring model predictions are consistent with our physical intuition. PINNs have been demonstrated as a means of inverse inference for integrating physical models with data<sup>18</sup>, but the right way to balance theory with high-quality observations that may not match can still be improved. PINNs can also be a useful way to extract some physics which can be further utilized to derive explainability.

Another problem brought up was that current models rely too much on metrics from the AI/ML community (e.g. L2 norm). Physics-based metrics are important for modeling physical systems (e.g. spatial smoothness, physics-derived metrics like lift and drag), but these are not well-defined in the context of current AI methods.

#### SINGAPORE'S ROLE

There are several teams in Singapore spanning A\*STAR, NUS and NTU working on the area of physics-informed neural networks, or physics-AI in general. For example, Physics-AI is a focus area in IHPC in A\*STAR, with teams working on physics-informed neural networks for various applications such as fluid dynamics. There are also teams in NUS and NTU who are working on physics-informed neural network methodologies. These capabilities include both more theoretical analysis and novel algorithms to improve the learning of PINNs, and to apply them to different domain problems.

#### OBJECTIVE

This research can support many domain areas in the general theme of AI for science. Physics-informed ML models can provide more robust models, even in more data-scarce scenarios in different domains like fluid dynamics and systems biology modeling. This can be both for forward estimation for different scientific problems (e.g. in flow system modeling), and state-of-health estimation for digital twinning when integrated with sensing in an inverse inference framework.

---

### GRAND CHALLENGE 4: REPRESENTATION LEARNING FOR SCIENCE

---

#### BACKGROUND

A more general approach to resolve data sparsity and related issues is through the development of representation learning for AI for science. The goal is to develop interpretable, generalizable, privacy-preserving representations that can also be "factorized" or "separable" for the specific domains such as molecular structures and healthcare data. This transformative representation learning approach will enable significant advancements in scientific discovery and application by providing robust and domain-specific representations that enhance interpretability, efficiency, and privacy.

Traditional representation learning approaches are not tailored to scientific data. On the one hand, a desired representation learning approach should be generalizable and applicable to diverse scientific fields and applications. On the other hand, domain-specific knowledge and characteristics need to be accounted for, such as certain invariance or equivariant properties.

The development of advanced representation learning techniques faces several quantifiable challenges that must be addressed to achieve the desired outcomes:

- *Developing Interpretable Latent Spaces for Domains Like Molecular Structures:* One of the primary challenges is to create latent spaces that are not only accurate but also interpretable, especially in complex domains such as molecular structures. Interpretable latent spaces facilitate better understanding and insights, allowing researchers to draw meaningful conclusions from the data.
- *Ensuring Privacy-Preserving Representation Learning in Sensitive Fields Like Healthcare:* In fields like healthcare, where data privacy is paramount, ensuring that representation learning methods can protect sensitive information while still providing useful insights is a significant challenge. Techniques must be developed that can balance the need for data utility with stringent privacy requirements.

- *Incorporating Domain Knowledge to Achieve Invariance or Equivariance in Representations:* Effective representation learning requires the incorporation of domain-specific knowledge to ensure that the representations exhibit certain desired properties, such as invariance or equivariance. This incorporation is crucial for making the models more robust and applicable to the specific characteristics of the data in various scientific fields.
- *Addressing the High Computational Costs Associated with Simulations in Material Design:* The computational expense of simulations, particularly in material design, poses a major challenge. Developing efficient representation learning methods that can reduce these costs while maintaining high accuracy and reliability is essential for practical applications and large-scale deployment.

Developing proper representation methods can also enable the interpretability of the learned AI models, allowing one to accumulate scientific knowledge and complete a cycle of learning.

#### STATE OF THE ART

There are a number of current representation learning approaches, but each with their shortcomings that need to be addressed.

- *Diffusion Models:* While diffusion models are effective for certain applications, they are computationally expensive, particularly for large-scale material design. This high computational cost limits their practical utility and scalability.
- *Hyperbolic Embeddings for Tree-Structured Data:* Hyperbolic embeddings are beneficial for hierarchical data structures but struggle to handle other types of data. Their limited applicability reduces their effectiveness in diverse scientific domains.
- *Graph Neural Networks (GNNs):* GNNs are powerful tools for relational data, providing robust insights and analysis. However, they often lack interpretability and are computationally intensive, which can hinder their practical implementation and broad adoption.

## SINGAPORE'S ROLE

Singapore is uniquely positioned to lead the advancement of representation learning in scientific domains due to its strong emphasis and tight integration of researchers both in AI methods and AI applications in science. Research institutes in A\*STAR and laboratories in research universities have strong capabilities in both algorithm and applied research, and there is already much existing research efforts on advancing representation learning for healthcare (e.g. ophthalmology and computer vision, reduced order modeling and materials science).

## OBJECTIVE

Two domains of application that these methodological advancements can immediately impact are materials design and healthcare. The development of efficient and interpretable representations is crucial for material design. By reducing the computational costs associated with simulations, researchers can accelerate the discovery and optimization of new materials, leading to significant advancements in this field. On the other hand, ensuring privacy while creating robust representations of patient data is essential for healthcare applications. Improved representation learning techniques can facilitate better diagnosis and treatment, enhancing patient outcomes and advancing medical research.

---

## GRAND CHALLENGE 5: INTERPRETABLE/EXPLAINABLE AI MODELS FOR SCIENTIFIC DISCOVERY

### BACKGROUND

AI models are increasingly used to analyze and make predictions on scientific data. Recent results have shown impressive performance on complex tasks such as protein folding and weather prediction, suggesting that the AI models have captured salient patterns in the data to make accurate predictions. Given such performant models, being able to understand how and why the AI model made its predictions could potentially reveal new scientific principles that the model has learnt. In addition, such understanding of the AI model provides a sanity check to ensure

that the predictions are not based on spurious correlations (e.g. due to technical artifacts in the data), thus providing confidence that its predictions will generalize.

Explainable AI in scientific discovery faces several challenges. First, ensuring the faithfulness of explanations is crucial, as explanations must accurately reflect the model's actual reasoning to be trustworthy. This is particularly important for methods that provide explanations for models post-hoc rather than inherently interpretable models (like decision trees). Secondly, complex models like deep neural networks that provide high accuracy may be difficult to interpret, creating a tradeoff between performance and explainability. Thirdly, explaining a model's predictions may involve information beyond the model itself, for example, relating to the interpretation or provenance of the input data. To further derive novel hypotheses or knowledge from the models will also require the method for interpreting/explaining predictions to understand scientific concepts beyond the model. Finally, effectively communicating the explanation to non-experts who may not have a deep understanding of the model's internals is a challenge that needs to be overcome for explainable methods to be broadly useful for science.

### STATE OF THE ART

To address these challenges, various methods for explainability have been developed<sup>19</sup>. Model-agnostic techniques like LIME and SHAP can provide understandable insights for any model through an approximation, which may not always be faithful to the original model. Intrinsically interpretable models, such as decision trees and linear models, offer transparency by design, though this may be at the expense of higher predictive accuracy. Post-hoc explainability methods, such as activation map visualization in neural networks, can elucidate complex models' decision processes, but still require manual interpretation and expert knowledge of neural networks to derive useful insights. Incorporating causal inference methods can help explore cause-and-effect relationships, making AI findings more robust and

scientifically valuable; however inferring causality is still a challenging problem especially in limited data situations that are common in practice.

In terms of generating explanations, using formal logic for explanations and explicit knowledge representation can make AI reasoning more structured and comprehensible. Recent large language models also offer a potential solution to improve the communication of explanations via natural language, but suffer from issues like hallucinations. They also need to be adapted to associate text with other forms of scientific data and knowledge to provide useful explanations in the scientific context.

## SINGAPORE'S ROLE

Singapore has multiple, diverse groups working on explainable AI at the universities and research institutes. These groups approach the problem from different perspectives, including human-computer interaction and formal methods, beyond the AI perspective. As the need for explainable AI goes beyond the scientific domain, existing work has spanned the range from being domain-agnostic to focusing on domains such as healthcare and social media analytics. Of note, explainable AI has been one of the key research focus areas of AI Singapore (AISG) and several of the abovementioned efforts have been supported by grants from AISG.

## OBJECTIVE

In the context of AI for science, interpretable/explainable AI is a horizontal capability that has broad application across all domain areas<sup>20,21,22,23,24,25</sup>, though targeted methods development will likely be required to incorporate domain knowledge and support AI models used in specific domain areas. For example, the development of interpretable/explainable anomaly detection methods<sup>26</sup> encoding domain knowledge may help scientists in the initial stages of hypothesis generation, by drawing their attention to salient yet unusual patterns in the data. Interpretable/explainable AI is also relevant in the context of the other methodological challenges in this report (e.g. in developing interpretable representations or more interpretable models

based on physics). More broadly, improving the interpretability and explainability of AI models will improve trust and hence adoption of AI tools in the scientific discovery process.

## RESOURCE REQUIREMENTS

As the development of methodology must go hand-in-hand with applications, the standard resources required for the latter are required. This includes high performance computing for simulation and machine learning, large scale data processing and storage, are required. On top of these, it is worthwhile to discuss some other unique requirements for the research outlined in this report. The first, and perhaps the most important requirement is the need to form and support interdisciplinary teams working on developing common AI methodology for several application domains. Here, AI expertise (machine learning, mathematics, high-performance computing) and domain expertise in the respective fields are equally important. Research topics should be formulated to focus on the development of AI methods, instead of for a particular application. This can then attract attention from experts from AI algorithm research, who can bring the latest developments in AI to these emerging applications. Second, there should be pipelines for efficient validation of these methodologies, so that the practical relevance of the research can be maintained and evaluated.

One may take inspiration from a successful implementation of this in the area of reinforcement learning we mentioned briefly in the discussion on experiment design. The OpenAI gym<sup>27</sup> provides an environment for methodology development directly facing algorithm researchers, and the clean interface, together with fast and quantitative performance validation enables non-domain experts to contribute meaningfully to the development of effective methods to accelerate learning. An equivalent implementation of a gym for AI for science and engineering, where problem statements are abstracted as a unified interface and performance metrics are clearly and readily fed-back to the algorithm researcher can be an invaluable tool to drive methodology research.

## Conclusion

In this report, we outlined several worthy directions of research for the development of methodology to support the application of AI in science and engineering applications. The key challenges include systematic improvements of simulation and experimentation, as well as a principled interface between AI systems and scientific knowledge.

Finally, while this report focuses on the development of methodology, we must also emphasize the importance of developing the theoretical underpinnings of these new methodologies, including their working principles, performance guarantees and fundamental limits. These require research that lies on the interface of mathematics, machine learning and physical/biological sciences.

## REFERENCES

- 1 Voulodimos, A., Doulamis, N., Doulamis, A. & Protopapadakis, E. Deep learning for computer vision: A brief review. *Computational intelligence and neuroscience* **2018**, (2018).
- 2 Arulkumaran, K., Deisenroth, M. P., Brundage, M. & Bharath, A. A. Deep reinforcement learning: A brief survey. *IEEE Signal Processing Magazine* **34**, 26–38 (2017).
- 3 Han, Z., Gao, C., Liu, J., Zhang, J. & Zhang, S. Q. Parameter-Efficient Fine-Tuning for Large Models: A Comprehensive Survey. Preprint at <http://arxiv.org/abs/2403.14608> (2024).
- 4 DeePMD-kit: A deep learning package for many-body potential energy representation and molecular dynamics.
- 5 Friederich, P., Häse, F., Proppe, J. & Aspuru-Guzik, A. Machine-learned potentials for next-generation matter simulations. *Nat. Mater.* **20**, 750–761 (2021).
- 6 Arisaka, S. & Li, Q. Principled Acceleration of Iterative Numerical Methods Using Machine Learning. in *Proceedings of the 40th International Conference on Machine Learning* 1041–1059 (PMLR, 2023).  
Benner, P., Gugercin, S. & Willcox, K. A Survey of Projection-Based Model Reduction Methods for Parametric Dynamical Systems. *SIAM Rev.* **57**, 483–531 (2015).
- 7 Sanchez-Lengeling, B. & Aspuru-Guzik, A. Inverse molecular design using machine learning: Generative models for matter engineering. *Science* **361**, 360–365 (2018).
- 8 Brunton, S. L. & Kutz, J. N. *Data-Driven Science and Engineering: Machine Learning, Dynamical Systems, and Control*. (Cambridge University Press, 2019).
- 9 Krenn, M. *et al.* On scientific understanding with artificial intelligence. *Nat Rev Phys* **4**, 761–769 (2022).
- 10 Karniadakis, G. E. *et al.* Physics-informed machine learning. *Nat Rev Phys* **3**, 422–440 (2021).
- 11 Wong, J. C., Gupta, A. & Ong, Y.-S. Can Transfer Neuroevolution Tractably Solve Your Differential Equations? *IEEE Computational Intelligence Magazine* **16**, 14–30 (2021).
- 12 Chiu, P.-H., Wong, J. C., Ooi, C., Dao, M. H. & Ong, Y.-S. CAN-PINN: A fast physics-informed neural network based on coupled-automatic-numerical differentiation method. *Computer Methods in Applied Mechanics and Engineering* **395**, 114909 (2022).
- 13 Wong, J. C., Ooi, C. C., Gupta, A. & Ong, Y.-S. Learning in Sinusoidal Spaces With Physics-Informed Neural Networks. *IEEE Trans. Artif. Intell.* **5**, 985–1000 (2024).
- 14 Hu, Z., Kawaguchi, K., Zhang, Z. & Karniadakis, G. E. Tackling the Curse of Dimensionality in Fractional and Tempered Fractional PDEs with Physics-Informed Neural Networks. Preprint at <http://arxiv.org/abs/2406.11708> (2024).
- 15 Jessica, L. S. E., Arafat, N. A., Lim, W. X., Chan, W. L. & Kong, A. W. K. Finite Volume Features, Global Geometry Representations, and Residual Training for Deep Learning-based CFD Simulation. Preprint at <http://arxiv.org/abs/2311.14464> (2023).
- 16 Um, K. & Brand, R. Solver-in-the-Loop: Learning from Differentiable Physics to Interact with Iterative PDE-Solvers.
- 17 Raissi, M., Yazdani, A. & Karniadakis, G. E. Hidden fluid mechanics: Learning velocity and pressure fields from flow visualizations. *Science* (2020) doi:10.1126/science.aaw4741.
- 18 Dwivedi, R. *et al.* Explainable AI (XAI): Core Ideas, Techniques, and Solutions. *ACM Comput. Surv.* **55**, 194:1–194:33 (2023).
- 19 Roscher, R., Bohn, B., Duarte, M. F. & Garcke, J. Explainable Machine Learning for Scientific Insights and Discoveries. *IEEE Access* **8**, 42200–42216 (2020).
- 20 Zhong, X. *et al.* Explainable machine learning in materials science. *npj Comput Mater* **8**, 1–19 (2022).
- 21 Jiang, S. *et al.* How Interpretable Machine Learning Can Benefit Process Understanding in the Geosciences. *Earth's Future* **12**, e2024EF004540 (2024).
- 22 Dybowski, R. Interpretable machine learning as a tool for scientific discovery in chemistry. *New Journal of Chemistry* **44**, 20914–20920 (2020).
- 23 Esterhuizen, J. A., Goldsmith, B. R. & Linic, S. Interpretable machine learning for knowledge generation in heterogeneous catalysis. *Nat Catal* **5**, 175–184 (2022).
- 24 Novakovskiy, G., Dexter, N., Libbrecht, M. W., Wasserman, W. W. & Mostafavi, S. Obtaining genetics insights from deep learning via explainable artificial intelligence. *Nat Rev Genet* **24**, 125–137 (2023).
- 25 Li, Z., Zhu, Y. & Van Leeuwen, M. A Survey on Explainable Anomaly Detection. *ACM Trans. Knowl. Discov. Data* **18**, 23:1–23:54 (2023).
- 26 Brockman, G. *et al.* OpenAI Gym. Preprint at <http://arxiv.org/abs/1606.01540> (2016).

## APPENDIX XVIII. GLOBAL R&D EFFORTS IN AI FOR SCIENCE

### North America

1. The US is a leading player in efforts in AI for Science. Multiple government-led workshops, initiatives, resources and funding opportunities for AI for Science have been active over the past half decade. The US Department of Energy is custodian of high performance computing resources (including quantum), vast experimental data, and also supports research into AI/ML methods. In 2019–2020, it organised a series of town halls examining opportunities in AI for Science, that was followed by establishment of multiple programme offices and research funding supporting AI research across various scientific domains. A follow up series of workshops and report were organised in 2022–2023 capturing the rapid progress in AI [Annex A: 1]. The National Science Foundation has established 25 AI institutes in the US since 2020 [Annex A: 2]. The National Institutes of Health has focused efforts on enabling AI through maximising data utility [Annex A: 3,4]. The Government of Canada and Ontario launched the Vector Institute in 2017, which aims to facilitate AI development and adoption across businesses, governments, and science [Annex A: 23].
2. The above government-led efforts are complemented by community efforts. Since 2021, the NeurIPS conference convenes global experts at annual AI for Science workshops [Annex A: 6]. More recently in 2023, the US National Academy of Sciences convened an AI for Science workshop [Annex A: 7]. The Nobel Turing Challenge, hosted by the Systems Biology

Institute, organises a series of workshops engaging the global community [Annex A: 8]. These community efforts are sustained by a healthy academic research ecosystem. Some leading institutes in AI and/or AI for Science include MIT, UC Berkeley, and Stanford [Annex A: 8–12]. University of Toronto, University of Chicago, Caltech, and Cornell have initiatives directed specifically at AI for Science [Annex A: 13–16]. Domain areas with significant AI efforts include materials (Berkeley Lab ML for Science, UC Berkeley Bakar Institute of Digital Materials, MIT Materials Science and Engineering, UToronto Acceleration Consortium), and biomedical and healthcare (MIT Jameel Clinic, Stanford AI for Structure-based Drug Discovery, Stanford Centre for AI in Medicine and Imaging). Efforts are also directed towards developing foundational AI models across different scientific fields, with the potential to democratise AI in science [Annex A: 17].

3. Significant effort and advances in AI for Science reside in industry. Industry giants Microsoft, Google, and Google DeepMind are leveraging their AI expertise to address scientific challenges in materials and chemistry, biomedical sciences, climate, nuclear fusion, sustainability, and quantum computing [Annex A: 18–20]. On a smaller scale, multiple start-ups and biotechs in the biomedical industry aim to leverage AI through generation, collection, and analysis of data, building foundational models for biology, and building semi-autonomous AIs for scientific research [Annex A: 21–23].

## United Kingdom

4. The UK is strongly committed to research and implementation of AI, as articulated in their National AI Strategy and Quantum Strategy, and correspondingly demonstrated by large investments in compute infrastructure and research. <sup>[Annex B: 1-3]</sup> UKRI supports research on AI and development of AI for applications including the sciences through a national institute for data science and AI, The Alan Turing Institute, as well as nine AI hubs hosted by UK universities <sup>[Annex B: 3-4]</sup>. Funding has also been specifically earmarked for the application of AI for life sciences <sup>[Annex B: 5]</sup>. Institutions with research initiatives focusing on AI for Science include Oxford, Cambridge, Imperial College London, and University of Southampton <sup>[Annex B: 6-9]</sup>. Strategic areas of priorities and strengths include healthcare, climate and sustainability, governmental policy, as well as research on the fundamentals of AI.

## Asia-Pacific

5. Within Asia-Pacific, Japan has made significant investments in supercomputing infrastructure, the RIKEN research institutes, as well as the AI Research Centre <sup>[Annex C:1-3]</sup>. Within AI for Science, Japan's government announced in mid-2023 plans to develop a generative AI to produce medical and scientific hypotheses by learning from research papers and experiments, leveraging on rich accumulated research data and supercomputing capabilities available at RIKEN. China's Ministry of Science and Technology and National Natural Science Foundation also initiated in 2023 a programme to promote innovations in AI models and algorithms for major scientific problems, and to construct a national open innovation platform for AI <sup>[Annex C: 4]</sup>. Australia's national science agency, CSIRO, hosts an AI for Missions Programme which is focused on translation and accelerating impact in a number of Missions using AI, including a mission on Scientific AI <sup>[Annex C: 5-6]</sup>.

## European Union

6. The EU is investing in projects to establish a Network of Excellence Centres (NoEs) in the area of AI, Data and Robotics to advance European AI research and development, and translation to real-world impact <sup>[Annex D:1-2]</sup>. It will provide programmatic funding to support research at the convergence of AI and Science (eg. life sciences, earth systems observation, nuclear science, social sciences and humanities, among others) <sup>[Annex D:3]</sup>. The EU is also establishing an Open Science Cloud to provide European researchers, innovators, companies and citizens with a federated environment to publish, find and reuse data, tools and services for research, innovation and education <sup>[Annex D:4]</sup>, and is investing in a European High Performance Compute network, which will enable application of AI to the domain of scientific research <sup>[Annex D:5]</sup>.

In addition, some European countries have been investing, independently from the EU initiatives, into AI and AI for Science research efforts. For example, France's Agence Nationale de la Recherche (ANR) has been the AISSAI project at CNRS, launched in 2024 and aimed at advancing AI in scientific discovery. The 3IA network, established in 2022, includes multiple nodes across the country, such as 3IA Côte d'Azur at Université Côte d'Azur, ANITI at Université de Toulouse, and MIAI Grenoble Alpes at Université Grenoble Alpes, all driving AI innovation in areas like healthcare, environmental science, and data processing. The 3IA initiative also includes industry-academic collaborations like the 3IA Prairie network. These initiatives, alongside industry labs like EDF's Data Innovation Lab, are positioning France as a growing hub for AI-driven scientific research. Spain and the Netherlands have also invested and hosted a number of AI for science labs and initiatives. <sup>[Annex E:23-28]</sup>

## South America

7. South America has seen a rise in efforts aiming to strengthen different countries' position in the global AI landscape, with increasing investments in AI research and development for scientific progress. In Chile, the National Center for Artificial Intelligence Research (CENIA), established in 2022, is supported by the Agencia Nacional de Investigación y Desarrollo (ANID), fostering AI research across various scientific disciplines <sup>[Annex E:18]</sup>. In Brazil, key initiatives include the C4AI at the University of São Paulo, launched in 2022 with funding from IBM and the São Paulo Research Foundation (FAPESP), which aims to advance AI applications in science and technology. Similarly, the Brazilian Center for Physics Research (CBPF) supports the Litcomp – IA project, focusing on AI-driven innovations in physics, funded by Petrobras and the Support Foundation for the Development of Scientific Computing (FACC). With additional initiatives such as UFRJ-Coppe launched in 2024, Brazil is also pushing forward a national strategy to foster AI innovation. These efforts, alongside broader regional initiatives, highlight South America's ambitions to build a robust AI ecosystem that can support scientific advancement and technological leadership. <sup>[Annex E:16-17, Figure x]</sup>

## Africa

8. Africa is aiming to become a key player in the global AI landscape, with growing efforts to harness AI for scientific research and innovation. In Rwanda, the World Economic Forum's Center for the Fourth Industrial Revolution (C4IR), a collaborative initiative among multiple institutions, is focused on leveraging AI to drive progress in sectors such as healthcare and agriculture. Launched to support Rwanda's digital transformation, the center promotes innovation with a particular emphasis on improving data access and enhancing the potential of AI to address local challenges. In South Africa, Wits University has introduced CirrusAI, an AI research initiative backed by the Cirrus Foundry Fund and Cortex Group, to advance scientific research using AI. Additionally, the African Institute for Mathematical Sciences (AIMS) launched an AI for Science Master's program in 2023, supported by Google DeepMind, to nurture AI talent across the continent. While these efforts are significant, they reflect Africa's ambition to become a leading force in AI for science, with ongoing initiatives aimed at accelerating progress in this field. <sup>[Annex E:57-59]</sup>



# ANNEX A.

## GLOBAL R&D EFFORTS IN AI FOR SCIENCE – NORTH AMERICA

### USA

NO	CATEGORY	INITIATIVE / EFFORT	DOMAIN	ASPECT OF PIPELINE: DATA / AI / COMPUTE	ASPECT OF AI: AI FOR SCIENCE/ RESEARCH ON AI / AI FOR APPLICATIONS	DETAILS
<b>USA</b>						
1	National Funding Agency	AI for Science Initiative, US DOE Office of Science (SC) <a href="https://science.osti.gov/Initiatives/AI">https://science.osti.gov/Initiatives/AI</a>	Cross-domain	AI; Data; Compute	AI for Science; Research on AI	<ul style="list-style-type: none"> <li>• DOE hosts high performance computers, the Exascale Computing Program (Aurora supercomputer), and Quantum Information Science Research Centers. The DOE Office of Science also generates vast experimental data sets. To maximize the impact of data, DOE supports development of new methods and algorithms that increase the reliability, robustness, and rigor of machine learning algorithm and methods to support their use in scientific research. Individual research programs focus on enhancing the analysis of the data for their disciplines to maximize the scientific impact of data.</li> <li>• Programme offices supporting AI research: Advanced Scientific Computing Research; Basic Energy Sciences; Biological and Environmental Research; Fusion Energy Sciences; High Energy Physics; Isotope R&amp;D and Production; Nuclear Physics</li> <li>• Funded projects under AI for Science Initiative: Fusion Energy, USD\$16M; Advanced Scientific Computing, USD\$30M; Fusion Energy Sciences, USD\$29M; Complex Systems, USD\$16M; Autonomous Optimisation and Control of Accelerators and Detectors, USD\$16M; High Energy Physics, USD\$6.4M</li> <li>• DOE reports               <ul style="list-style-type: none"> <li>» Original 2020 report on DOE AI for Science Townhalls, and Recommendation for creation of DOE AI for Science Initiative. <a href="https://www.osti.gov/biblio/1604756">https://www.osti.gov/biblio/1604756</a></li> <li>» Updated 2023 report on a follow up series of workshops in 2022. <a href="https://www.anl.gov/ai-for-science-report">https://www.anl.gov/ai-for-science-report</a></li> </ul> </li> </ul>
2	National Funding Agency	National AI Research Institutes, NSF <a href="https://new.nsf.gov/funding/opportunities/national-artificial-intelligence-research">https://new.nsf.gov/funding/opportunities/national-artificial-intelligence-research</a> <a href="https://aiinstitutes.org/">https://aiinstitutes.org/</a>	Cross-domain	AI	All Details provided focus on AI for Science	<ul style="list-style-type: none"> <li>• Launched in 2020, the NSF-led National Artificial Intelligence Research Institutes program consists of 25 AI institutes that connect over 500 funded and collaborative institutions across the U.S. and around the world.</li> <li>• Institutes relevant to AI for Science:               <ul style="list-style-type: none"> <li>» AI Institute for AI and Fundamental Interactions (MIT) aims to advance knowledge of fundamental interactions across scales using innovations in AI built upon ab initio physics principles, while simultaneously advancing the foundations of AI. They are targeting opportunities for ab initio AI to improve theory calculations, to improve experiments, and to advance the field of AI.</li> <li>» AI Institute for Molecular Discovery, Synthetic Strategy, and Manufacturing (University of Illinois at Urbana-Champaign), focuses on development of new AI-enabled tools to accelerate automated chemical synthesis and advance the discovery and manufacture of novel materials and bioactive compounds.</li> <li>» AI Institute for Artificial and Natural Intelligence (Columbia University) connects major progress Made in AI systems to our understanding of the brain.</li> </ul> </li> </ul>
3	National Funding Agency	NIH Strategic Plan for Data Science, Office of Data Science Strategy at NIH <a href="https://datascience.nih.gov/nih-strategic-plan-data-science">https://datascience.nih.gov/nih-strategic-plan-data-science</a>	Biomedical	Data	AI for Science; AI for Applications	<ul style="list-style-type: none"> <li>• Initiative aims to maximise the value of data generated through NIH-funded efforts, NIH aims to address the: findability, interconnectivity, and interoperability of NIH-funded biomedical data sets and resources; integration of existing data management tools and development of new ones; universalization of innovative algorithms and tools created by academic scientists into enterprise-ready resources that meet industry standards of ease of use and efficiency of operation; growing costs of data management.</li> <li>• Active grants are eligible for Administrative Supplements to Support Collaborations to Improve the AI/ML-Readiness of NIH-Supported Data</li> </ul>

4	National Funding Agency	NIH Common Fund's Bridge to Artificial Intelligence (Bridge2AI) program <a href="https://commonfund.nih.gov/bridge2ai">https://commonfund.nih.gov/bridge2ai</a> <a href="https://bridge2ai.org">https://bridge2ai.org</a>	Biomedical	Data	AI for Science; AI for Applications	<ul style="list-style-type: none"> <li>• Bridge2AI aims to facilitate widespread adoption of AI for biomedical challenges. Bridge2AI will use biomedical and behavioral research grand challenges to drive the development of ethics, standards, tools, data sets, and skills and workforce development strategies for linking scientific workflows, protocols, and other information about the data collection process into computable knowledge.</li> <li>• Efforts include: Generate new flagship biomedical and behavioral data sets that are ethically sourced, trustworthy, well-defined, and accessible; developing software and standards to unify data attributes across multiple data sources and across data types; Create automated tools to accelerate the creation of FAIR and ethically sourced data sets; Provide Providing resources to disseminate data, ethical principles, tools, and best practices; Create training materials and activities for workforce development</li> </ul>
5	Community	AI for Science workshops, NeurIPS <a href="https://ai4sciencecommunity.github.io/">https://ai4sciencecommunity.github.io/</a>	Cross-domain	AI	AI for Science	<ul style="list-style-type: none"> <li>• Workshop, part of NeurIPS conference</li> <li>• Up-to-date (2023) primer on AI for Science developments: <a href="https://medium.com/@AI_for_Science/ai-for-science-in-2023-a-community-primer-d2c2db37e9a7">https://medium.com/@AI_for_Science/ai-for-science-in-2023-a-community-primer-d2c2db37e9a7</a></li> <li>• Key researchers in AI for Science: <a href="https://ai4sciencecommunity.github.io/neurips23.html">https://ai4sciencecommunity.github.io/neurips23.html</a></li> </ul>
6	Community	National Academy of Sciences	Cross-domain	AI	AI for Science	<ul style="list-style-type: none"> <li>• Workshop in Oct 2023 explored the future of AI as an autonomous researcher performing scientific discovery: where AI stands, where it needs to go, and which disciplines should have increased investment for the utilization of AI scientists. <b>Report pending.</b> <a href="https://www.nationalacademies.org/event/40455_10-2023_ai-for-scientific-discovery-a-workshop">https://www.nationalacademies.org/event/40455_10-2023_ai-for-scientific-discovery-a-workshop</a></li> </ul>
7	Community	Nobel Turing Challenge	Cross-domain	AI	AI for Science	The Nobel Turing Challenge is a grand challenge aiming at developing a highly autonomous AI and robotics system that can make major scientific discoveries. The Initiative has convened multiple workshops globally since 2022. <a href="https://www.nobelturingchallenge.org">https://www.nobelturingchallenge.org</a>
8	Academic	MIT	Cross-domain	AI	All	<p>MIT hosts one of the first NSF AI Institutes, for AI and Fundamental Interactions. IAIFI aims to advance knowledge of fundamental interactions across scales using innovations in AI built upon ab initio physics principles, while simultaneously advancing the foundations of AI. They are targeting opportunities for ab initio AI to improve theory calculations, to improve experiments, and to advance the field of AI. <a href="https://iaifi.org">https://iaifi.org</a></p> <ul style="list-style-type: none"> <li>• The Electrical Engineering and Computer Science department (EECS), and Computer Science and Artificial Intelligence Lab (CSAIL) <ul style="list-style-type: none"> <li>» The department develops algorithms for modeling biological and clinical data across a range of modalities including imaging, text and genomics with collaborators. <a href="https://www.eecs.mit.edu/research/explore-all-research-areas/ml-and-healthcare">https://www.eecs.mit.edu/research/explore-all-research-areas/ml-and-healthcare</a></li> <li>» Connor Coley: Combines expertise in chemical engineering, computer science, and chemistry to accelerate molecular discovery. Develops platform technologies and workflows with relevance to small molecule drug discovery, chemical synthesis, and materials science. <a href="https://coley.mit.edu/research">https://coley.mit.edu/research</a></li> <li>» Sara Beery: computer vision methods that enable global-scale environmental and biodiversity monitoring across data modalities. <a href="https://beerys.github.io">https://beerys.github.io</a></li> </ul> </li> <li>• Schwarzman College of Computing: Cross-cutting structure serving as a home for CS and AI education and research. Both brings together existing MIT programs in computing and develops new cross-cutting educational and research programs. <a href="https://computing.mit.edu">https://computing.mit.edu</a></li> <li>• The Jameel Clinic incubates research at AI and life sciences to improve early detection of disease and aid in the personalization of treatment for these diseases. Their research focuses on creating and commercialising high-precision, affordable, and scalable ML solutions across 3 key areas. <a href="https://jclinic.mit.edu">https://jclinic.mit.edu</a> <ul style="list-style-type: none"> <li>» Clinical AI: e.g. Single-cell RNA sequencing foundation models; interpretation AI model for recurrence prediction after surgery in GI tumour; ML for dynamic information retrieval of EHR notes</li> <li>» Drug discovery: e.g. Discovery of structural class of antibiotics with explainable deep learning; improved influenza vaccine strain selection through deep evolutionary models; de novo design of protein structure and function with Rfdiffusion</li> <li>» Epidemiology: e.g. Predictive performance of multi-model ensemble forecasts of COVID-19; Deep learning model to predict future lung cancer risk from single low dose chest tomography</li> </ul> </li> </ul>

						<ul style="list-style-type: none"> <li>• Department of Materials Science and Engineering <ul style="list-style-type: none"> <li>» Ju Li lab developed Copilot for Real-World Experimental Scientist (CRESt), powered by various types of AI. CRESt suggests experiments and guides researchers through a process workflow – it can also retrieve and analyze data, switch equipment on and off, and drive robotic arms to mix liquids, or prepare materials for experimentation. <a href="https://dmse.mit.edu/news/accelerating-research-with-ai-assisted-experiments">https://dmse.mit.edu/news/accelerating-research-with-ai-assisted-experiments</a></li> </ul> </li> <li>• Department of Physics <ul style="list-style-type: none"> <li>» Max Tegmark: Work on Physics – using physics-based techniques to better understand biological and artificial intelligence. Work on AI safety – focus on mechanistic interpretability (MI): given a trained neural network that exhibits intelligent behavior, how can we figure out how it works. Auto-discovery of knowledge representations, hidden symmetries, modularity and conserved quantities. <a href="http://tegmark.org">http://tegmark.org</a> <a href="https://physics.mit.edu/faculty/max-tegmark">https://physics.mit.edu/faculty/max-tegmark</a></li> </ul> </li> <li>• Medical Engineering and Sciences <ul style="list-style-type: none"> <li>» James Collins: AI to discover novel classes of antibiotics and understand how they work, deep learning approaches for the de novo design of new antibiotics and the development of combination treatments. <a href="https://www.collinslab.mit.edu">https://www.collinslab.mit.edu</a></li> </ul> </li> <li>• Workshop planned for June 2024: AI Science: Strengthening the Bond Between the Sciences and Artificial Intelligence. Workshop will foster discussion of how AI intersects with biology, chemistry, physics, and related fields. This intersection is a two-way street, comprising i) leveraging scientific insight to develop new machine learning methods, and ii) developing new machine learning methods to advance the sciences. Specific technical AI areas of interest of focus, common to many of the Sciences, include: encoding of symmetries into models such as rotational equivariance; enabling neural networks to operate in bases suitable for the sciences; incorporating differential equations into neural networks for regularization and beyond; representation learning; generalization, extrapolation, uncertainty quantification, domain shift; batch active learning; transfer learning; generative and conditional-generative modeling such as diffusion and flow-based models; the role of data augmentation, and evaluation. Participation is encouraged from researchers focused on specific scientific domain applications, those on focused methodological development, as well as those who span the two. <a href="https://simons.berkeley.edu/workshops/aiscience-strengtheningbond-between-sciences-artificial-intelligence">https://simons.berkeley.edu/workshops/aiscience-strengtheningbond-between-sciences-artificial-intelligence</a></li> </ul>
9	Academic	UC Berkeley	Cross-domain	AI	All Details provided focus on AI for Science	<p>ML for Science at Berkeley Lab (DOE National Laboratory): Develop and share algorithms, software, tools, and libraries for scientific machine learning. Gather, organize and store scientific datasets in materials, energy, environment, biology, genomics, and astronomy. Host of HPC optimized for machine learning and advanced networking capabilities. <a href="https://ml4sci.lbl.gov/home">https://ml4sci.lbl.gov/home</a></p> <ul style="list-style-type: none"> <li>• AI+Science: Group of faculty in Electrical Engineering and Computer Sciences focused on intersection of AI+Science. <a href="https://ai-science.eecs.berkeley.edu/">https://ai-science.eecs.berkeley.edu/</a> <b>[Unclear if active]</b></li> <li>• Berkeley AI Research Climate Initiative: Community that aims to unite AI and climate-related researchers. <a href="https://ai-climate.berkeley.edu/index.html">https://ai-climate.berkeley.edu/index.html</a></li> <li>• Bakar Institute of Digital Materials for the Planet: Cross-departmental initiative aiming to use AI for Science to accelerate breakthroughs in climate change. <a href="https://bidmap.berkeley.edu">https://bidmap.berkeley.edu</a></li> </ul> <p>Workshop planned for June 2024: AI Science: Strengthening the Bond Between the Sciences and Artificial Intelligence. Workshop will foster discussion of how AI intersects with biology, chemistry, physics, and related fields. This intersection is a two-way street, comprising i) leveraging scientific insight to develop new machine learning methods, and ii) developing new machine learning methods to advance the sciences. Specific technical AI areas of interest of focus, common to many of the sciences, include: encoding of symmetries into models such as rotational equivariance; enabling neural networks to operate in bases suitable for the sciences; incorporating differential equations into neural networks for regularization and beyond; representation learning; generalization, extrapolation, uncertainty quantification, domain shift; batch active learning; transfer learning; generative and conditional-generative modeling such as diffusion and flow-based models; the role of data augmentation, and evaluation. Participation is encouraged from researchers focused on specific scientific domain applications, those on focused methodological development, as well as those who span the two. <a href="https://simons.berkeley.edu/workshops/aiscience-strengthening-bond-between-sciences-artificial-intelligence">https://simons.berkeley.edu/workshops/aiscience-strengthening-bond-between-sciences-artificial-intelligence</a></p>
10	Community	Community for Autonomous Scientific Experimentation (CASE), hosted at LBL <a href="https://autonomous-discovery.lbl.gov/">https://autonomous-discovery.lbl.gov/</a>	Cross-domain	AI	AI for Science	<ul style="list-style-type: none"> <li>• Community of researchers on the application of AI/ML methods to experimental and computational sciences, in order to enable autonomy and accelerate discovery. A resource to facilitate connection with researchers and to find information about methods, applications, literature and software tools.</li> </ul>

11	Academic	Stanford	Cross-domain	AI	All Details provided focus on AI for Science	<ul style="list-style-type: none"> <li>Stanford Artificial Intelligence Lab (SAIL): AI centre, not specific to AI for Science <a href="https://ai.stanford.edu">https://ai.stanford.edu</a></li> <li>AI for Structure-Based Drug Discovery: Small group of faculty. Aims to facilitate exchange of ideas between Stanford researchers developing groundbreaking machine learning methods that leverage molecular structure and industry scientists. <a href="https://aisbdd.stanford.edu/about">https://aisbdd.stanford.edu/about</a></li> <li>Center for Artificial Intelligence in Medicine and Imaging: Aim to develop, evaluate, and disseminate artificial intelligence systems to benefit patients. Center conducts research that solves clinically important imaging problems using machine learning and other AI techniques. <a href="https://aimi.stanford.edu">https://aimi.stanford.edu</a></li> </ul>
12	Academic	University of Chicago	Cross-domain	AI	AI for Science	<ul style="list-style-type: none"> <li>AI+Science: Research initiative on transformational AI-enabled scientific discovery across the physical and biological sciences, advancing core AI principles and training a new generation of diverse interdisciplinary scientists. <a href="https://datascience.uchicago.edu/research/ai-science/">https://datascience.uchicago.edu/research/ai-science/</a> <ul style="list-style-type: none"> <li>AI+Science conference hosted by UChicago and Caltech in Apr 2023</li> <li>AI+Science Summer School <a href="https://datascience.uchicago.edu/events/aiscience-summer-school-2024/">https://datascience.uchicago.edu/events/aiscience-summer-school-2024/</a></li> </ul> </li> </ul>
13	Academic	Caltech	Cross-domain	AI	AI for Science	<ul style="list-style-type: none"> <li>Caltech AI4science: Initiative aiming to bring together computer scientists with experts in other disciplines. <a href="https://www.ist.caltech.edu/ai4science/">https://www.ist.caltech.edu/ai4science/</a></li> </ul>
14	Academic	<a href="https://Science.ai.cornell.edu/">https://Science.ai.cornell.edu/</a>	Cross-domain	AI	AI for Science	<ul style="list-style-type: none"> <li>Cornell AI for Science (CUAISci): Focused on AI for scientific discovery, organized around the following broad research themes: AI for materials discovery, AI for physics, AI for the biological sciences, AI for the sustainability sciences. <a href="https://science.ai.cornell.edu/">https://science.ai.cornell.edu/</a></li> </ul>
15	Academic	Polymathic	Cross-domain	AI	AI for Science	<ul style="list-style-type: none"> <li>Polymathic's goal is to accelerate the development of versatile foundation models tailored for numerical datasets and scientific machine learning tasks. They aim to build AI models which leverage information from heterogeneous datasets and across different scientific fields, hence potentially democratizing AI in science by providing off-the-shelf models. Their team brings together pure machine learning researchers with domain scientists. Participating institutions: Simons Foundation, NYU, Cambridge, Schmidt Futures, Princeton, Berkeley lab. <a href="https://polymathic-ai.org">https://polymathic-ai.org</a></li> </ul>
16	Industry	Microsoft Research AI4Science <a href="https://www.microsoft.com/en-us/research/lab/microsoft-research-ai4science/faqs/">https://www.microsoft.com/en-us/research/lab/microsoft-research-ai4science/faqs/</a>	Cross-domain: Drug discovery, Molecular modeling, Life cycle analysis, PDEs	AI	AI for Science	<ul style="list-style-type: none"> <li>Microsoft seeks to drive major advances with lasting influence in the tools for scientific discovery through the use of machine learning, with a focus on 'fifth paradigm' scenarios. It aims to empower real-world impact on the most pressing societal problems including climate change, green energy, sustainable materials, and the discovery of new drugs.</li> <li>Projects <ul style="list-style-type: none"> <li>Chemistry, Biomedical: Bio Embedding – learn meaningful representations for biomolecules, design bio-inspired pretraining techniques; Generative Chemistry – collaboration with Novartis to speed up discovery of medicines; Graphormer, a deep learning package that allows researchers and developers to train custom models for molecule modelling tasks, such as materials science, or drug discovery;</li> <li>Decarbonisation – Project Carbonix: Ab Initio Aperiodic Molecular Dynamics - speed up Density Functional Theory with neural network, for physical and chemical process simulation, Materials and systems engineering – simulation of crystallisation process of hydrate in CCS, Discovery of new cathode materials through crystal structure design; Life cycle assessment - ML model of technical and behavioural pathways to EV adoption; Ophysics, Modeling: Fast Neural PDE Solver – to understand and forecast the world, leverage Physics behind PDEs</li> </ul> </li> </ul>
17	Industry	Google AI: Science AI, Quantum AI <a href="https://ai.google/discover/scienceai">https://ai.google/discover/scienceai</a> <a href="https://research.google/research-areas/general-science">https://research.google/research-areas/general-science</a>	Cross-domain	AI, Compute	AI for Science, AI for Applications, Research on AI	<ul style="list-style-type: none"> <li>Google AI has two focus areas: Science AI and Quantum AI</li> <li>Science AI is addressing science through breakthroughs in machine learning, cloud infrastructure, and data processing and analytics. It collaborates with research institutions around the globe, and share progress in scientific publications and open source releases Projects include: <ul style="list-style-type: none"> <li>Neuroscience: Connectomics to map the mouse brain</li> <li>Genomics: Pangenomes to improve equity of genomics</li> <li>Climate: Weather and climate prediction with general circulation models</li> <li>Biodiversity: Monitoring and measuring health of forest ecosystem and biodiversity through bird song</li> </ul> </li> <li>Quantum AI research efforts aim to build quantum processors and develop quantum algorithms. Nuclear fusion research: Deep reinforcement learning to autonomously discover how to control and contain plasma in a tokamak</li> </ul>

18	Industry	InSistro	Biomedical	AI	AI for Science	<ul style="list-style-type: none"> <li>InSistro collects and aggregating massive data sets, and harnesses advanced machine learning to accelerate drug discovery and development. <a href="https://www.insistro.com">https://www.insistro.com</a></li> </ul>
19	Non-profit	Future House	Biomedical	AI	AI for Science	<ul style="list-style-type: none"> <li>Philanthropically-funded moonshot with mission to build semi-autonomous AIs for scientific research, to accelerate the pace of discovery and to provide world-wide access to cutting-edge scientific, medical, and engineering expertise. Focus on biology. <a href="https://www.futurehouse.org/articles/announcing-future-house">https://www.futurehouse.org/articles/announcing-future-house</a></li> </ul>
20	Government	FASST	Cross-domain / Energy	AI / Data	AI for Science	<p>The <b>FASST Initiative</b> (Fast-Track to the Advancement of Science and Technology) is a U.S. Department of Energy (DOE) initiative designed to accelerate the deployment of cutting-edge scientific advancements and technological innovations. Launched to streamline research and development processes, FASST aims to foster collaboration across federal agencies, national laboratories, and industry to address critical energy challenges. By providing rapid access to funding and resources, the initiative supports breakthrough technologies that can transform the energy sector and drive progress in areas like clean energy, sustainability, and energy efficiency. FASST is part of the DOE's broader efforts to fast-track scientific discoveries and enhance America's leadership in energy innovation. <a href="https://www.energy.gov/fasst">https://www.energy.gov/fasst</a></p>
<b>CANADA</b>						
21	Non-profit / Government	Vector Institute for AI <a href="https://vectorinstitute.ai">https://vectorinstitute.ai</a> Tours for government representatives <a href="mailto:info@vectorinstitute.ai">info@vectorinstitute.ai</a>	Cross-domain	AI	All	<ul style="list-style-type: none"> <li>The Vector Institute launched in 2017 with support from the Government of Canada, the Government of Ontario, private industry, and in partnership with universities in Ontario. It aims to empower researchers, businesses and governments, to develop and adopt AI responsibly.</li> <li>Strategic research priorities: ML, Deep learning, AI for Science, Trustworthy AI, AI for Health, Foundation models</li> </ul>
22	Academic	Acceleration Consortium, University of Toronto Led by Alan Aspuru-Guzik	Materials	AI	AI for Science	<ul style="list-style-type: none"> <li>Consortium builds self-driving labs, combining material science with AI, robotics, and advanced computing, to rapidly design and test new materials.</li> <li>Global network of government, academia, and industry: to exchange resources and knowledge, to co-develop standards and best practices, and to connect players in the value chain to expedite research translation and commercialization.</li> </ul>
	Academia	<b>MILA (Montréal Institute for Learning Algorithms) i</b>	Cross-domain	AI	AI for Applications, Research on AI	<p><b>MILA</b> (Montréal Institute for Learning Algorithms) is a leading AI research institute based in Montreal, Canada, renowned for its contributions to deep learning and artificial intelligence. Founded in 2018 and supported by <b>over CAD 125 million</b> in funding from both government sources and private sector partners, including <b>Element AI</b> and <b>Hydro-Québec</b>, MILA focuses on advancing AI research with applications across various industries, including healthcare, finance, and transportation. The institute is a key player in the <b>Pan-Canadian Artificial Intelligence Strategy</b>, which aims to position Canada as a global leader in AI innovation. MILA's work is centered on developing cutting-edge AI algorithms, fostering collaboration between academic researchers and industry, and training the next generation of AI experts. It also works closely with institutions like <b>Université de Montréal</b> and various international partners to drive both theoretical advancements and practical applications of AI. <a href="https://www.essex.ac.uk/research-projects/ai-policy-observatory-for-the-world-of-work/national-and-regional-cases/canada">https://www.essex.ac.uk/research-projects/ai-policy-observatory-for-the-world-of-work/national-and-regional-cases/canada</a></p>
<b>BRAZIL</b>						
24	Government, Academia	Brazil's <b>National Strategy for Artificial Intelligence (EBIA)</b> ,	Cross-domain	AI	AI for Science	<p>Brazil's <b>National Strategy for Artificial Intelligence (EBIA)</b>, launched by the <b>Ministry of Science, Technology, and Innovation (MCTI)</b> in 2021, is a comprehensive initiative aimed at positioning the country as a global leader in AI research and development. The strategy outlines key priorities, including the development of AI applications in healthcare, agriculture, and public administration, while fostering a robust AI ecosystem that supports innovation and scientific progress. The strategy is backed by significant governmental investments, focusing on collaboration between academia, industry, and the public sector to promote the responsible and ethical use of AI technologies. Additionally, Brazil's growing commitment to AI is further reflected in the <b>CETIC.br</b> report on AI development, which highlights key initiatives and funding mechanisms aimed at boosting AI research across the country. Together, these efforts represent a strategic push to integrate AI into Brazil's scientific infrastructure, with a focus on sustainable and inclusive technological growth. <a href="https://cetic.br/media/docs/publicacoes/6/20240514085413/iso-year-xvi-n-1-ia-development-in-brazil.pdf">https://cetic.br/media/docs/publicacoes/6/20240514085413/iso-year-xvi-n-1-ia-development-in-brazil.pdf</a> <a href="https://www.gov.br/mcti/pt-br/acompanhe-o-mcti/transformacaodigital/arquivosinteligenciaartificial/ebia-summary_brazilian_4-979_2021.pdf">https://www.gov.br/mcti/pt-br/acompanhe-o-mcti/transformacaodigital/arquivosinteligenciaartificial/ebia-summary_brazilian_4-979_2021.pdf</a></p>

## ANNEX B.

# GLOBAL R&D EFFORTS IN AI FOR SCIENCE – UK

NO	CATEGORY	INITIATIVE / EFFORT	DOMAIN	ASPECT OF PIPELINE: DATA / AI / COMPUTE	ASPECT OF AI: AI FOR SCIENCE / RESEARCH ON AI / AI FOR APPLICATIONS	DETAILS
1	Government	AI Strategy	Cross-domain	Compute	All	<ul style="list-style-type: none"> <li>National AI Strategy: <a href="https://www.gov.uk/government/publications/national-ai-strategy/national-ai-strategy-html-version">https://www.gov.uk/government/publications/national-ai-strategy/national-ai-strategy-html-version</a></li> </ul>
2	Government	Quantum Strategy	Cross-domain	Compute	-	<ul style="list-style-type: none"> <li>UK Quantum Strategy: £2.5 billion <a href="https://www.gov.uk/government/publications/national-quantum-strategy/national-quantum-strategy-accessible-webpage">https://www.gov.uk/government/publications/national-quantum-strategy/national-quantum-strategy-accessible-webpage</a></li> <li>Quantum missions, building on Quantum strategy, aim to galvanise academics, industry and private investors to commit time and resource towards hitting significant milestones. Missions to achieve by 2030-2035: (1) Accessible, UK-based quantum computers capable of running 1 trillion operations and supporting key sectors of the economy; (2) Deployed the world's most advanced quantum network at scale, pioneering the future quantum internet; (3) Every NHS Trust will benefit from quantum sensing-enabled solutions, helping those with chronic illness live healthier, longer lives through early diagnosis and treatment; (4) Quantum navigation systems will be deployed on aircraft, providing next-generation accuracy for resilience that is independent of satellite signals; (5) Mobile, networked quantum sensors will have unlocked new situational awareness capabilities, exploited across critical infrastructure in the transport, telecoms, energy, and defence sectors. <a href="https://www.gov.uk/government/news/new-quantum-missions-launched-as-science-minister-visits-new-advanced-quantum-lab">https://www.gov.uk/government/news/new-quantum-missions-launched-as-science-minister-visits-new-advanced-quantum-lab</a></li> <li>National Quantum Computing Centre: Core funding from EPSRC <a href="https://www.nqcc.ac.uk">https://www.nqcc.ac.uk</a></li> <li><a href="https://www.ukri.org/what-we-do/strategic-priorities-fund">https://www.ukri.org/what-we-do/strategic-priorities-fund</a> ; <a href="https://www.ukri.org/what-we-do/browse-our-areas-of-investment-and-support/ai-and-data-science-for-engineering-health-and-Government-asg">https://www.ukri.org/what-we-do/browse-our-areas-of-investment-and-support/ai-and-data-science-for-engineering-health-and-Government-asg</a></li> <li>Nine AI hubs supported by £80 million. <a href="https://www.ukri.org/news/100m-boost-in-ai-research-will-propel-transformative-innovations">https://www.ukri.org/news/100m-boost-in-ai-research-will-propel-transformative-innovations</a>. AI hubs relevant to AI for Science: <ul style="list-style-type: none"> <li>» AI for collective intelligence (AI4CI), University of Bristol</li> <li>» AI hub for causality in healthcare AI with real data, University of Edinburgh</li> <li>» ProbAI: a hub for the mathematical and computational foundations of probabilistic AI, Lancaster University</li> <li>» AI for Chemistry: alchemy, University of Liverpool and Imperial</li> <li>» AI hub in generative models, University College London</li> <li>Omathematical foundations of intelligence: an 'Erlangen Programme' for AI, University of Oxford</li> </ul> </li> </ul>
4	Institute / National Funding Agency	The Alan Turing Institute	Cross-domain	AI	AI for Science; AI for Applications	<ul style="list-style-type: none"> <li>The Alan Turing Institute was created as the national institute for data science in 2015, and in 2017, AI was added to their remit. The Turing aims to provide an end-to-end, interdisciplinary pathway in data science and AI that enables impact at scale and major progress against societal challenges. Member universities are Cambridge, Edinburgh, Oxford, UCL, Warwick, Leeds, Manchester, Newcastle, Queen Mary University of London, Birmingham, Exeter, Bristol, and Southampton, and in 2023 the Institute launched an open university network to engage all UK universities. Grand challenge areas are: Defence and national security; Environment and sustainability; Transformation of health. <a href="https://www.turing.ac.uk">https://www.turing.ac.uk</a></li> </ul>
5	Government	AI Life Sciences Accelerator Mission	Biomedical	AI	AI for Applications	<ul style="list-style-type: none"> <li>£100 million in new government investment announced in 2023, targeted towards areas where rapid deployment of AI has the greatest potential to create transformational breakthroughs in treatments for previously incurable diseases. Parties working together will include government, industry, the NHS, academia and medical research charities. <a href="https://www.gov.uk/government/news/new-100-million-fund-to-capitalise-on-ais-game-changing-potential-in-life-sciences-and-healthcare">https://www.gov.uk/government/news/new-100-million-fund-to-capitalise-on-ais-game-changing-potential-in-life-sciences-and-healthcare</a></li> </ul>
6	Academic	Oxford AI4Science	Cross-domain	AI	All	<ul style="list-style-type: none"> <li>UKRI AI hub: Mathematical foundations of intelligence: an 'Erlangen Programme' for AI, led by Professor Michael Bronstein. Focusing on using mathematical principles, this hub will use geometry topology and probability to enhance AI methods. <a href="https://www.ox.ac.uk/news/2024-02-06-new-oxford-research-hub-propel-transformative-ai-innovations">https://www.ox.ac.uk/news/2024-02-06-new-oxford-research-hub-propel-transformative-ai-innovations</a></li> <li>Oxford AI4Science Lab is a part of the Department of Computer Science at the University of Oxford, led by Atılım Güneş Baydin. The lab specializes in probabilistic machine learning and scientific discovery, and collaborates with experts in disciplines including particle physics, heliophysics, astrobiology, Earth science, and computational social science, to solve important problems in these domains through application and development of AI methods. Highlights of work include: supercomputing-scale Bayesian inference in simulators of the Standard Model of particle physics, new calibration techniques for the NASA Solar Dynamics Observatory, and deploying machine learning algorithms onboard spacecraft launched to Earth orbit. Research has been funded by NASA, European Space Agency, US Department of Energy, and UK Space Agency. <a href="https://oxai4science.github.io">https://oxai4science.github.io</a> Community of scientists and AI specialists passionate about the use of AI to benefit science and society. Supported by Schmidt Futures.</li> </ul>

8	Academic	Imperial College London	Cross-domain	AI	AI for Science	<ul style="list-style-type: none"> <li>I-X Centre for AI in Science is dedicated to using AI to disrupt and advance Science, Engineering and Mathematics and is underpinned by core support from Schmidt Futures. Research in foundational areas: Explainable, Safe, and Robust AI; Humans and AI; Quantum and AI; Sensing and AI; Systems and Infrastructure for AI. Research in application areas: Economy, Health, Life, Planet, Resilience, Space. <a href="https://www.imperial.ac.uk/ix-ai-in-science">https://www.imperial.ac.uk/ix-ai-in-science</a></li> </ul>
9	Academic	AI4SD Network+, hosted by University of Southampton	Cross-domain	AI	AI for Science	<ul style="list-style-type: none"> <li>Funded by EPSRC</li> <li>Network+ aims to bring together researchers to leverage AI for scientific discovery, with a focus on design and synthesis of chemicals and materials.</li> </ul>

## ANNEX C. GLOBAL R&D EFFORTS IN AI FOR SCIENCE – ASIA-PACIFIC

NO	CATEGORY	INITIATIVE / EFFORT	DOMAIN	ASPECT OF PIPELINE: DATA / AI / COMPUTE	ASPECT OF AI: AI FOR SCIENCE / RESEARCH ON AI / AI FOR APPLICATIONS	DETAILS
<b>JAPAN</b>						
1	Government	Government of Japan	Biomedical, Cross-domain	AI	AI for Science; AI for Applications	<ul style="list-style-type: none"> <li>Japan's Ministry of Education, Culture, Sports, Science and Technology announced plans in Jul 2023 to develop a generative AI that produces medical and scientific hypotheses by learning from research papers and images of experiments. The generative AI will be first used for medical and material research with other areas added in the future. <a href="https://asia.nikkei.com/Business/Technology/Japan-to-develop-generative-AI-to-speed-scientific-discovery">https://asia.nikkei.com/Business/Technology/Japan-to-develop-generative-AI-to-speed-scientific-discovery</a> <ul style="list-style-type: none"> <li>Developing generative AI for one research area is estimated to cost roughly 30 billion yen (\$212 million). The ministry will seek funds for initial development in the fiscal 2024 budget.</li> <li>The project is expected to last eight years, with the technology made available for researchers nationwide from 2031.</li> <li>The Riken research institute, has a trove of accumulated research data, will lead the effort.</li> <li>Japan's Ministry of Economy, Trade and Industry will introduce a new supercomputer for research in 2024. The education ministry will increase the computing power of Riken's Fugaku supercomputer to make it easier to use in generative AI research.</li> <li>NTT, SoftBank are working on generative AI models compatible with the Japanese language.</li> </ul> </li> <li>Japanese and U.S. governments will collaborate on the development of AI for scientific research through data sharing and joint use of supercomputers for AI development. The collaboration will be led by RIKEN and the DOE Argonne National Laboratory. <a href="https://japannews.yomiuri.co.jp/politics/politics-government/20240212-168320/">https://japannews.yomiuri.co.jp/politics/politics-government/20240212-168320/</a></li> <li>Data, applications of AI to applications in the sciences, climate, and on Quantum-HPD hybrid Computing <a href="https://www.riken.jp/en/research/labs/r-ccs">https://www.riken.jp/en/research/labs/r-ccs</a></li> <li>RIKEN Center for Quantum Computing (RQC) explores the frontier of quantum technologies. It takes a full-stack approach, performing R&amp;D from hardware to software and from basic science to applications. <a href="https://www.riken.jp/en/research/labs/rqc">https://www.riken.jp/en/research/labs/rqc</a></li> <li>The RIKEN Center for Advanced Intelligence Project (AIP) aims to develop technologies for the welfare of society and humanity. AIP also conducts research on ethical, legal and social issues caused by the spread of AI technology. Research groups include those focused on AI itself, and others on applications of AI to applications in healthcare and social sciences. <a href="https://www.riken.jp/en/research/labs/aip/">https://www.riken.jp/en/research/labs/aip/</a></li> </ul>

2	Government / Academic	AI Research Center (AIRC),Japan	Cross-domain	AI; Compute	AI for applications	<ul style="list-style-type: none"> <li>The Artificial Intelligence Research Centre (AIRC) was inaugurated in 2015 under the National Institute of Advanced Industrial Science and Technology (AIST), to bring together global research talent and technologies in AI. The AIRC performs goal-oriented basic research to advance AI towards implementation of AI in the real-world. Areas of applications include mobility, productivity in manufacturing and service industry, healthcare, and security. <a href="https://www.airc.aist.go.jp">https://www.airc.aist.go.jp</a></li> </ul>
4	Government	Japan Science and Technology Agency (JST)	Cross-domain	AI; Compute;Data	AI for	<p>The initiative announced by the <b>Japan Science and Technology Agency (JST)</b> in October 2023, focuses on advancing <b>AI for scientific research</b> through the creation of a new research program called <b>AI-Driven Scientific Discovery</b>. This program, backed by a substantial ¥10 billion (approximately \$68 million) in government funding, aims to harness AI technologies to accelerate scientific breakthroughs, particularly in fields such as drug discovery, climate science, and materials engineering. The program will support collaborative research between universities, national laboratories, and private-sector companies, emphasizing the integration of AI tools in traditional scientific workflows. <a href="https://sj.jst.go.jp/news/202310/n1013-01k.html">https://sj.jst.go.jp/news/202310/n1013-01k.html</a></p>
5	Government	Ministry of Science and Technology, National Natural Science Foundation of China, China	Cross-domain	AI	AI for Science	<ul style="list-style-type: none"> <li>MOST and NSFC jointly initiated in early 2023 a programme for AI-driven scientific research addressing key problems in basic disciplines, as well as research needs in key sci-tech fields, such as drug development, gene research and biology. The programme aims to promote innovations in AI models and algorithms for major scientific problems, develop platforms for research fields, and construct a national open innovation platform for AI public computing power. It will bring together interdisciplinary research and development teams, promote the establishment of an innovation consortium, and build international academic exchange platforms. <a href="https://english.www.gov.cn/statecouncil/ministries/202303/28/content_WS642224f1c6d0f528699dc4bd.html">https://english.www.gov.cn/statecouncil/ministries/202303/28/content_WS642224f1c6d0f528699dc4bd.html</a></li> <li>China's "New Generation AI Intelligence Development Plan" : <a href="https://digichina.stanford.edu/work/full-translation-chinas-new-generation-artificial-intelligence-development-plan-2017">https://digichina.stanford.edu/work/full-translation-chinas-new-generation-artificial-intelligence-development-plan-2017</a></li> <li>McKinsey report on AI for China: <a href="https://www.mckinsey.com/capabilities/quantumblack/our-insights/the-next-frontier-for-ai-in-china-could-add-600-billion-to-its-economy">https://www.mckinsey.com/capabilities/quantumblack/our-insights/the-next-frontier-for-ai-in-china-could-add-600-billion-to-its-economy</a></li> </ul>
6	Government	Shanghai Science and Technology Commission	Cross-domain	AI	AI for Science	<p>The <b>Shanghai Science and Technology Commission's (STCSM)</b> initiative, launched in August 2024, focuses on advancing <b>AI for scientific research</b> in Shanghai, with an emphasis on accelerating AI applications in critical fields such as healthcare, drug discovery, and materials science. The initiative is part of a broader effort to position Shanghai as a global AI hub by integrating cutting-edge AI technologies into the scientific process, streamlining innovation, and facilitating rapid application of research breakthroughs. It aims to foster stronger collaboration between universities, research institutes, and industry, particularly in leveraging AI to tackle complex scientific challenges. This initiative is supported by significant funding from the Chinese government and aligns with the country's strategic priorities to drive AI-driven advancements in key sectors. <a href="https://stcsm.sh.gov.cn/news/20240805/9478ffa965264a4d9482e7b333f03167.html">https://stcsm.sh.gov.cn/news/20240805/9478ffa965264a4d9482e7b333f03167.html</a></p>
7	Government	CSIRO (Australia's National Science Agency)	Cross-domain	AI	AI for Science	<ul style="list-style-type: none"> <li>CSIRO AI for Missions Programme is focused on translation and accelerating impact in Missions using a range of AI technologies. (1) Data centred AI: involves sensing technologies, data generation, information processing and modelling, and analytics; (2) Human centred AI: focusses on robotics, reasoning, natural language processing and augmented reality; (3) Scientific AI: develops novel AI methods to accelerate scientific discovery and engineering. <a href="https://research.csiro.au/ai4m">https://research.csiro.au/ai4m</a></li> <li>CSIRO 2022 report highlighted the importance and transformative potential of development and adoption of AI for scientific research, and identified development pathways for research organisations seeking to upgrade AI capability. Report also includes a bibliometric analysis of trends in AI adoption by research domains. <a href="https://www.csiro.au/en/research/technology-space/ai/artificial-intelligence-for-science-report">https://www.csiro.au/en/research/technology-space/ai/artificial-intelligence-for-science-report</a></li> </ul>
8	Academic	Australian Artificial Intelligence Institute	Cross-domain	AI	All	<ul style="list-style-type: none"> <li>The AAIL vision is to develop theoretical foundations and advanced algorithms for AI, and to drive significant progress in related areas like computational intelligence, business intelligence, computer vision, data science, machine learning, brain computer interface, bioinspired neural networks and information systems. <a href="https://www.uts.edu.au/research/australian-artificial-intelligence-institute">https://www.uts.edu.au/research/australian-artificial-intelligence-institute</a></li> </ul>

## ANNEX D.

# GLOBAL R&D EFFORTS IN AI FOR SCIENCE – EUROPEAN UNION

NO	CATEGORY	INITIATIVE / EFFORT	DOMAIN	ASPECT OF PIPELINE: DATA / AI / COMPUTE	ASPECT OF AI: AI FOR SCIENCE / RESEARCH ON AI / AI FOR APPLICATIONS	DETAILS
1	Pan-European Organisation	EU Network of Excellence Centres (NoEs) in AI and Robotics <a href="https://www.vision4ai.eu/community">https://www.vision4ai.eu/community</a>	Cross-domain	AI	AI for science; Research on AI; AI for Applications	<p>AI solutions in the future, which may be applied in areas such as the healthcare sector or Autonomous driving, robotics, cybersecurity, media and document security.</p> <ul style="list-style-type: none"> <li>• euROBIN – This project brings together leading experts from the European robotics and AI research community, with a goal to establish a European ecosystem of robots that share their data and knowledge and jointly learn to perform an endless variety of tasks in human environments. It aims to demonstrate advances in three application domains, i. robotic manufacturing for a circular economy, ii. personal robots for enhanced quality of life iii. outdoor robots for sustainable communities.</li> <li>• ELIAS (European Lighthouse of AI for Sustainability) – This project aims to create a network of excellence connecting researchers in academia with industry to drive AI research for sustainable innovation and economic development. Projects include: i. AI for building optimisation, ii. AI for monitoring virtual infrastructure, iii. Responsible, user-centric advertising, iv. Mitigating misinformed migrant perception in EU, v. AI for forecasting of vegetation state, vi. Open materials discovery.</li> <li>• dAIEDGE – This project aims to establish a network of excellence for distributed, trustworthy, efficient and scalable AI at the Edge, and to advance Europe's innovation and technology base by developing a comprehensive policy and governance approach to AI in order for the EU to become a world leader in innovation in the data economy and its applications.</li> <li>• ENFIELD – This project aims to establish a center of excellence focused on advancing fundamental research in Adaptive, Green, Human-Centric, and Trustworthy AI, and elevate research within key sectors like healthcare, energy, manufacturing, and space by attracting top talents, technologies, and resources from leading research and industry entities in Europe.</li> </ul>
2	Pan-European Organisation	ELLIS (European Laboratory for Learning and Intelligent Systems)  <a href="https://ellis.eu/">https://ellis.eu/</a>	Cross-domain	AI	AI for Science; Research on AI;	<ul style="list-style-type: none"> <li>• A pan-European AI network of excellence which focuses on fundamental science, technical innovation and societal impact. ELLIS programmes cover the areas of: <ul style="list-style-type: none"> <li>i. Theory algorithms and computations of modern learning systems,</li> <li>ii. Machine learning for health, Quantum and physics-based machine learning,</li> <li>iii. Geometric deep learning,</li> <li>iv. Machine learning for earth and climate sciences,</li> <li>v. Natural intelligence,</li> <li>vi. Human-centric machine learning,</li> <li>vii. Robust machine learning,</li> <li>viii. Machine learning and computer vision,</li> <li>ix. Natural language processing,</li> <li>x. Multimedia/multimodal information,</li> <li>xi. Robotic learning, Interactive learning and interventional representations,</li> <li>xii. Symbolic machine learning.</li> </ul> </li> <li>• EU Partnership in AI, Data and Robotics - Focuses on deploying AI research results to industry settings.</li> <li>• Horizon Europe Pillar II – Includes calls for funding at the convergence of AI and life sciences, earth systems observation, nuclear science, or social sciences and humanities, among others.</li> <li>• European Innovation Council - One-stop shop for breakthrough innovators, providing support from idea to market, from early advanced research to commercialisation and scale-up. Supports breakthrough technologies and disruptive start-ups in any field, as well as targeted support for strategic challenges.</li> </ul>

3	Pan-European Organisation	EU programmes offering funding for AI in science	Cross-domain	AI	AI for Science; AI for Application	<p>EU programmes offering funding for AI in Science include:</p> <ul style="list-style-type: none"> <li>• ERC Synergy Grants – Allows proposals to be submitted by 2 to 4 scientists to tackle ambitious scientific challenges.</li> </ul>
4	Pan-European Organisation	European Open Science Cloud (EOSC) and AI4EOSC <a href="https://research-and-innovation.ec.europa.eu/strategy/strategy-2020-2024/our-digital-future/open-science/european-open-science-cloud-eosc_en">https://research-and-innovation.ec.europa.eu/strategy/strategy-2020-2024/our-digital-future/open-science/european-open-science-cloud-eosc_en</a> <a href="https://ai4eosc.eu/">https://ai4eosc.eu/</a>	Cross-domain	AI	AI for Science; AI for Applications	<ul style="list-style-type: none"> <li>• The European Open Science Cloud (EOSC) aims to provide European researchers, innovators, companies and citizens with a federated and open multi-disciplinary environment where they can publish, find and reuse data, tools and services for research, innovation and educational purposes. It will operate under well-defined conditions to ensure trust and safeguard the public interest, providing seamless access, FAIR (Findability, Accessibility, Interoperability and Reusability) management, and reliable re-use of research data and other digital objects (eg. methods, software and publications).</li> <li>• AI4EOSC will deliver an enhanced set of services for the development of Artificial Intelligence (AI), Machine Learning (ML) and Deep Learning (DL) models and applications for the European Open Science Cloud (EOSC). It will base its activities on the technological framework delivered by the DEEP- Hybrid-DataCloud H2020 project, which delivered the DEEP platform to exploit computing resources from pan-European e-Infrastructures. Identified use-cases include i) Automated Thermography, ii) Agrometeorology, iii) Integrated Plant Protection.</li> </ul>
5	Pan-European Organisation					<a href="https://council.science/news/ai-in-science/">https://council.science/news/ai-in-science/</a>
6	Industry	Bioptimus	Biomedical	AI	AI for Science	<ul style="list-style-type: none"> <li>• New startup launched with seed funding of \$35M, founded by former Google DeepMind scientists. Aims to build the first universal AI foundation model for biology. <a href="https://www.bioptimus.com/">https://www.bioptimus.com/</a></li> </ul>

## ANNEX E.

# GLOBAL R&D EFFORTS IN AI FOR SCIENCE – WORLD INFOGRAPHICS AND TABLE OF REFERENCES

### USA

NO	ACADEMIC LAB	INSTITUTION	MAIN FUNDING SOURCE	YEAR STARTED	REFERENCE (LINK)
1 2	Artificial Intelligence and Fundamental Interactions (IAFAI)	Massachusetts Institute of Technology (MIT) – Lead	National Science Foundation (NSF)	2020	<a href="https://iaifi.org/">https://iaifi.org/</a>
3	AI Institute for Molecular Discovery, Synthetic Strategy, and Manufacturing	Illinois-Urbana Champaign	National Science Foundation (NSF)	2020	<a href="https://moleculemaker.org/">https://moleculemaker.org/</a>
4	AI Institute for Artificial and Natural Intelligence	Columbia	National Science Foundation (NSF)	2023	<a href="https://arni-institute.org/">https://arni-institute.org/</a>
5	Cornell AI for Science Institute	Cornell	National Science Foundation (NSF)	2022	<a href="https://science.ai.cornell.edu/">https://science.ai.cornell.edu/</a>
6	Machine Learning for Science	UC Berkeley	Department of Energy AI 4 Science Initiative	2022	<a href="https://ml4sci.lbl.gov/home">https://ml4sci.lbl.gov/home</a>
7	Data Science Institute (DSI) AI + Science	University of Chicago	National Science Foundation (NSF)	2022	<a href="https://datascience.uchicago.edu/research/ai-science/">https://datascience.uchicago.edu/research/ai-science/</a>
8	AI4Science – Information Science and Technology	California Institute of Technology (Caltech)	Various, Amazon Web Service (AWS)	2018	<a href="https://ai4science.caltech.edu/about.html">https://ai4science.caltech.edu/about.html</a>
NO	INDUSTRY LAB		COMPANY		REFERENCE (LINK)
9	AI4Science – Microsoft Research		Microsoft		<a href="https://www.microsoft.com/en-us/research/lab/microsoft-research-ai-for-science/">https://www.microsoft.com/en-us/research/lab/microsoft-research-ai-for-science/</a>
10	Science AI, Quantum AI		Google		<a href="https://ai.google/discover/scienceai/">https://ai.google/discover/scienceai/</a> <a href="https://quantumai.google/">https://quantumai.google/</a>
11	InSipro		InSipro		<a href="https://www.insipro.com/">https://www.insipro.com/</a>
12	Accelerated Materials Design and Discovery (AMDD) – Toyota Research Institute		Toyota		<a href="https://www.tri.global/">https://www.tri.global/</a>

## Canada

NO	ACADEMIC LAB	INSTITUTION	MAIN FUNDING SOURCE	YEAR STARTED	REFERENCE (LINK)
13	Vector Institute for AI	Cross-institution	Ontario's Critical Technology Initiative Program (CTI), Government of Canada, private	2017	<a href="https://www.canada.ca/en/innovation-science-economic-development/news/2022/06/government-of-canada-launches-second-phase-of-the-pan-canadian-artificial-intelligence-strategy.html">https://www.canada.ca/en/innovation-science-economic-development/news/2022/06/government-of-canada-launches-second-phase-of-the-pan-canadian-artificial-intelligence-strategy.html</a>
14	Acceleration Consortium	University of Toronto	Canada First Research Excellence Fund (CFREF)	2023	<a href="https://acceleration.utoronto.ca/news/u-of-t-receives-200-million-grant-to-support-acceleration-consortiums-self-driving-labs-research">https://acceleration.utoronto.ca/news/u-of-t-receives-200-million-grant-to-support-acceleration-consortiums-self-driving-labs-research</a>
15	Mila Quebec	Cross-institution	Quebec Research and Innovation Investment Strategy (SQRI2)	N/A	<a href="https://mila.quebec/en/about/about-mila">https://mila.quebec/en/about/about-mila</a>

## Brazil

NO	ACADEMIC LAB	INSTITUTION	MAIN FUNDING SOURCE	YEAR STARTED	REFERENCE (LINK)
16	C4AI	University of Sao Paulo	IBM, Sao Paulo Research Foundation (FAPESP)	2022	<a href="https://c4ai.inova.usp.br/about/">https://c4ai.inova.usp.br/about/</a>
17	Litcomp – IA	Brazilian Center for Physics Research (CBPF)	Petrobras, Support Foundation for the Development of Scientific Computing (FACC)	2022	<a href="https://litcomp.cbpf.br/ai/">https://litcomp.cbpf.br/ai/</a>

## Chile

NO	ACADEMIC LAB	INSTITUTION	MAIN FUNDING SOURCE	YEAR STARTED	REFERENCE (LINK)
18	National Center for Artificial Intelligence Research (CENIA)	Inter-University	Agencia Nacional de Investigacion y Desarrollo (ANID)	2022	<a href="http://www.cenia.cl">www.cenia.cl</a>

## UK

NO	ACADEMIC LAB	INSTITUTION	MAIN FUNDING SOURCE	YEAR STARTED	REFERENCE (LINK)
19	Oxford AI4Science Lab	Oxford University	UKRI, NASA, ESA, US Department of Energy, UK Space Agency	2022	<a href="https://www.ox.ac.uk/news/2024-02-06-new-oxford-research-hub-propel-transformative-ai-innovations">https://www.ox.ac.uk/news/2024-02-06-new-oxford-research-hub-propel-transformative-ai-innovations</a>
20	Cambridge Accelerate Science	Cambridge University	Schmidt Future	2020	<a href="https://science.ai.cam.ac.uk/news">https://science.ai.cam.ac.uk/news</a>
21	I-X Centre for AI in Science	Imperial College London	Schmidt Future	2021	<a href="https://www.imperial.ac.uk/ix-ai-in-science/">https://www.imperial.ac.uk/ix-ai-in-science/</a>

NO	INDUSTRY LAB	COMPANY	REFERENCE (LINK)
22	AI4Science Microsoft Research Lab	Microsoft	<a href="https://www.microsoft.com/en-us/research/lab/microsoft-research-ai-for-science/">https://www.microsoft.com/en-us/research/lab/microsoft-research-ai-for-science/</a>

## France

NO	ACADEMIC LAB	INSTITUTION	MAIN FUNDING SOURCE	YEAR STARTED	REFERENCE (LINK)
23	AISSAI	CNRS	Agence Nationale de la Recherche (ANR) – National AI Research Program (PNRIA)	2024	<a href="https://aissai.cnrs.fr/en/">https://aissai.cnrs.fr/en/</a>
24	3IA Network – 3IA Côte d'Azur	Université Côte d'Azur	Agence Nationale de la Recherche (ANR) – National AI Research Program (PNRIA)	2022	<a href="https://3ia.univ-cotedazur.eu/">https://3ia.univ-cotedazur.eu/</a>
25	3IA Network – ANITI	Université de Toulouse	Agence Nationale de la Recherche (ANR) – National AI Research Program (PNRIA)	2022	<a href="https://aniti.univ-toulouse.fr/en/">https://aniti.univ-toulouse.fr/en/</a>
26	3IA Network – MIAI Grenoble Alpes	Université Grenoble Alpes	Agence Nationale de la Recherche (ANR) – National AI Research Program (PNRIA)	2022	<a href="https://miai.univ-grenoble-alpes.fr/">https://miai.univ-grenoble-alpes.fr/</a>
27	3IA Network – Pririe	Inter-University / Industry	Agence Nationale de la Recherche (ANR) – National AI Research Program (PNRIA)	2022	<a href="https://prairie-institute.fr/founding-members/">https://prairie-institute.fr/founding-members/</a>
28	LaborIA	INRIA	Agence Nationale de la Recherche (ANR) – National AI Research Program (PNRIA)	2021	<a href="https://www.laboria.ai/">https://www.laboria.ai/</a>

NO	INDUSTRY LAB	COMPANY	REFERENCE (LINK)
29	Data Innovation Lab	EDF R&D	<a href="https://www.edf.fr/en/the-edf-group/inventing-the-future-of-energy/rd-global-expertise/rd-experience/data-science-ai-world/data-innovation-lab-at-edf-rd">https://www.edf.fr/en/the-edf-group/inventing-the-future-of-energy/rd-global-expertise/rd-experience/data-science-ai-world/data-innovation-lab-at-edf-rd</a>
30	Bioptimus	Bioptimus	<a href="https://www.bioptimus.com/">https://www.bioptimus.com/</a>

## Spain

NO	ACADEMIC LAB	INSTITUTION	MAIN FUNDING SOURCE	YEAR STARTED	REFERENCE (LINK)
31	Institut d'Investigació en Intel·ligència Artificial (IIIA)	Consejo Superior de Investigaciones Científica	EU, Ministerio de Ciencia e Innovación	2020	<a href="https://www.iiia.csic.es/en-us/">https://www.iiia.csic.es/en-us/</a>
32	AI4EOSC	Consejo Superior de Investigaciones Científica (Lead)	EU, Horizon 2020	2022	<a href="https://ai4eosc.eu/about/">https://ai4eosc.eu/about/</a>

## The Netherlands

NO	ACADEMIC LAB	INSTITUTION	MAIN FUNDING SOURCE	YEAR STARTED	REFERENCE (LINK)
33	AI4Science Lab	University of Amsterdam	Faculty of Science (FNWI)	2020	<a href="https://ai4science-amsterdam.github.io/index.html">https://ai4science-amsterdam.github.io/index.html</a>
34	Vision4AI	University of Leyden (Lead)	EU Horizon 2020+	2020	<a href="https://www.vision4ai.eu/about/">https://www.vision4ai.eu/about/</a>

## China

NO	ACADEMIC LAB	INSTITUTION	MAIN FUNDING SOURCE	YEAR STARTED	REFERENCE (LINK)
35	AI for Science Institute	Cross-institutions	Ministry of Science and Technology (MOST), National Natural Science Foundation of China (NSFC)	2021	<a href="https://www.aisi.ac.cn/?page_id=17228955">https://www.aisi.ac.cn/?page_id=17228955</a>
36	Zeijang lab	Zhejiang University	Ministry of Science and Technology (MOST), National Natural Science Foundation of China (NSFC), Alibaba	2017	<a href="https://en.zhejianglab.com/">https://en.zhejianglab.com/</a>
37	Chen Frontier Lab	Tianqiao and Chrissy Chen Institute (TCCI)	Tianqiao and Chrissy Chen Institute (TCCI)	2020	<a href="https://www.cheninstitute.org/news/now-or-never-tcci-founder-tianqiao-chen-announces-one-billion-investment-in-ai-brain-science">https://www.cheninstitute.org/news/now-or-never-tcci-founder-tianqiao-chen-announces-one-billion-investment-in-ai-brain-science</a>

NO	INDUSTRY LAB	COMPANY	REFERENCE (LINK)
38	Microsoft Research – AI4Science	Microsoft	<a href="https://www.microsoft.com/en-us/research/lab/microsoft-research-ai-for-science/">https://www.microsoft.com/en-us/research/lab/microsoft-research-ai-for-science/</a>
39	Xtalpi	Xtalpi	<a href="https://www.forbes.com/sites/zinnialee/2023/12/01/tencent-backed-ai-drug-discovery-startup-xtalpi-files-for-hong-kong-ipo/">https://www.forbes.com/sites/zinnialee/2023/12/01/tencent-backed-ai-drug-discovery-startup-xtalpi-files-for-hong-kong-ipo/</a> <a href="https://www.xtalpi.com/en/about">https://www.xtalpi.com/en/about</a>
40	Insilico Medicine	Insilico Medicine	<a href="https://insilico.com/">https://insilico.com/</a> <a href="https://www.prnewswire.com/news-releases/insilico-medicine-raises-60-million-in-series-d-financing-to-advance-pipeline-and-launch-ai-powered-drug-discovery-robotics-laboratory-301561736.html">https://www.prnewswire.com/news-releases/insilico-medicine-raises-60-million-in-series-d-financing-to-advance-pipeline-and-launch-ai-powered-drug-discovery-robotics-laboratory-301561736.html</a>
41	N/1	Wanhua Chemical	<a href="https://en.whchem.com/cmscontent/779.html">https://en.whchem.com/cmscontent/779.html</a> <a href="https://english.news.cn/20240731/6f779fdb16a74c78b66e020e28a5dae0/c.html">https://english.news.cn/20240731/6f779fdb16a74c78b66e020e28a5dae0/c.html</a>
42	CATL – Hong Kong R&D center	CATL	<a href="https://www.catl.com/en/news/6155.html">https://www.catl.com/en/news/6155.html</a> <a href="https://www.all-about-industries.com/artificial-intelligence-is-making-its-way-into-asias-battery-industry-a-b039b0554dd51bc0d00eb03076285fb9/">https://www.all-about-industries.com/artificial-intelligence-is-making-its-way-into-asias-battery-industry-a-b039b0554dd51bc0d00eb03076285fb9/</a>

## Japan

NATIONAL INITIATIVE:					
Ministry of Education, Culture, Sports, Science and Technology (MEXT) KAKENHI					
8.5B Japanese Yen, 2024					

NO	ACADEMIC LAB	INSTITUTION	MAIN FUNDING SOURCE	YEAR STARTED	REFERENCE (LINK)
43	Correspondance and Fusion of AI and Brain Science	Okinawa Institute of Science and Technology Grad. Uni. (OIST)	Ministry of Education, Culture, Sports, Science and Technology (MEXT) KAKENHI	2016-2023	<a href="https://kaken.nii.ac.jp/en/grant/KAKENHI-INTERNATIONAL-16K21738/">https://kaken.nii.ac.jp/en/grant/KAKENHI-INTERNATIONAL-16K21738/</a>
44	Next Generation AI Research Center	University of Tokyo	University of Tokyo	2023	<a href="https://www.ai.u-tokyo.ac.jp/en/about#s2">https://www.ai.u-tokyo.ac.jp/en/about#s2</a>
45	AI for Science Platform	RIKEN TRIP	Ministry of Education, Culture, Sports, Science and Technology (MEXT)	2023	<a href="https://www.riken.jp/en/research/labs/r-ccs/ai_sci_plat/index.html">https://www.riken.jp/en/research/labs/r-ccs/ai_sci_plat/index.html</a>
46	Center for Advanced Intelligence Project (AIP)	RIKEN	Ministry of Education, Culture, Sports, Science and Technology (MEXT)	2016	<a href="https://www.lu.emb-japan.go.jp/itpr_fr/11_000001_00278.html">https://www.lu.emb-japan.go.jp/itpr_fr/11_000001_00278.html</a> <a href="https://www.riken.jp/en/research/labs/aip/">https://www.riken.jp/en/research/labs/aip/</a>
47	Generative AI Accelerator Challenge (GENIAC)	N/A	Ministry of Economy, Trade and Industry (METI)	2024	<a href="https://www.meti.go.jp/english/policy/mono_info_service/geniac/index.html">https://www.meti.go.jp/english/policy/mono_info_service/geniac/index.html</a>

## South Korea

NO	ACADEMIC LAB	INSTITUTION	MAIN FUNDING SOURCE	YEAR STARTED	REFERENCE (LINK)
48	AI Research Hub Project	Multi	Ministry of Science and ICT (MIST), private	2024	<a href="https://sciencebusiness.net/network-updates/south-korea-announces-public-call-ai-research-hub">https://sciencebusiness.net/network-updates/south-korea-announces-public-call-ai-research-hub</a> <a href="https://www.msit.go.kr/eng/bbs/view.do?sCode=eng&amp;mId=4&amp;mPid=2&amp;pageIndex=&amp;bbsSeqNo=42&amp;nttSeqNo=964&amp;searchOpt=ALL&amp;searchTxt=">https://www.msit.go.kr/eng/bbs/view.do?sCode=eng&amp;mId=4&amp;mPid=2&amp;pageIndex=&amp;bbsSeqNo=42&amp;nttSeqNo=964&amp;searchOpt=ALL&amp;searchTxt=</a> <a href="https://elblog.pl/2024/05/13/south-korea-to-invest-36-billion-won-in-ai-research-hub-by-2028/">https://elblog.pl/2024/05/13/south-korea-to-invest-36-billion-won-in-ai-research-hub-by-2028/</a>
49	Smart Lab	Computational Science Research Center (CSRC), Korea Institute of Science and Technology (KIST)	Korea Advanced Institute of Science and Technology (KAIST), Ulsan National Institute of Science and Technology (UNIST)	2024	<a href="https://www.kist.re.kr/eng/research/materials-latest-research-news.do?mode=view&amp;articleNo=13112&amp;title=Smart+labs+for+bespoke+synthesis+of+nanomaterials+are+emerging">https://www.kist.re.kr/eng/research/materials-latest-research-news.do?mode=view&amp;articleNo=13112&amp;title=Smart+labs+for+bespoke+synthesis+of+nanomaterials+are+emerging</a> <a href="https://csrc.kist.re.kr/bbs/board.php?bo_table=m01_01&amp;wr_id=1">https://csrc.kist.re.kr/bbs/board.php?bo_table=m01_01&amp;wr_id=1</a>
50	Global AI Frontier lab	New York University – Korea Advanced Institute of Science and Technology (NYU-KAIST)	Ministry of Science and ICT (MIST), National Science Foundation	2024	<a href="https://www.nyu.edu/about/news-publications/news/2024/may/korea-and-nyu-establish-global-ai-frontier-lab.html">https://www.nyu.edu/about/news-publications/news/2024/may/korea-and-nyu-establish-global-ai-frontier-lab.html</a> <a href="https://sciencebusiness.net/network-updates/south-korea-announces-public-call-ai-research-hub">https://sciencebusiness.net/network-updates/south-korea-announces-public-call-ai-research-hub</a>

## Australia

NO	ACADEMIC LAB	INSTITUTION	MAIN FUNDING SOURCE	YEAR STARTED	REFERENCE (LINK)
51	Mixed Reality lab	CSIRO Data61	Department of Industry, Science and Resources	2019	<a href="https://www.csiro.au/en/research/technology-space/data/mixed-reality-lab">https://www.csiro.au/en/research/technology-space/data/mixed-reality-lab</a>
52	AI for Missions Program	CSIRO Data61	Department of Industry, Science and Resources	2019	<a href="https://research.csiro.au/ai4m/">https://research.csiro.au/ai4m/</a>
53	Google Partnership	CSIRO Data61 -	Google, Department of Industry, Science and Resources	2021	<a href="https://www.csiro.au/en/news/All/News/2024/August/CSIRO-and-Google-accelerating-AI-for-science">https://www.csiro.au/en/news/All/News/2024/August/CSIRO-and-Google-accelerating-AI-for-science</a>
54	Australian Alliance for Secure Genomics and AI in Rare Disease (AASGARD)	Centre for Population Genomics (Lead)	Medical Research Future Fund	2024	<a href="https://www.garvan.org.au/news-resources/news/significant-mrff-funding-to-spearhead-ai-in-rare-disease-diagnosis">https://www.garvan.org.au/news-resources/news/significant-mrff-funding-to-spearhead-ai-in-rare-disease-diagnosis</a>

## New Zealand

NO	ACADEMIC LAB	INSTITUTION	MAIN FUNDING SOURCE	YEAR STARTED	REFERENCE (LINK)
55	Center for Artificial Intelligence Research (CAIR)	University of Auckland	ministry of Innovation Business and Employment	2000	<a href="https://www.aut.ac.nz/study/study-options/engineering-computer-and-mathematical-sciences/research/centre-for-artificial-intelligence-research-cair">https://www.aut.ac.nz/study/study-options/engineering-computer-and-mathematical-sciences/research/centre-for-artificial-intelligence-research-cair</a>

NO	INDUSTRY LAB	COMPANY	REFERENCE (LINK)
56	N/A	Litmaps	<a href="https://www.litmaps.com/">https://www.litmaps.com/</a>

## South Africa

NO	ACADEMIC LAB	INSTITUTION	MAIN FUNDING SOURCE	YEAR STARTED	REFERENCE (LINK)
57	CirrusAI	Wits University (Lead)	Cirrus Foundry fund, Cortex Goup	2019	<a href="https://cortexlogic.com/2019/08/06/school-of-chemistry-launches-new-ai-research-initiative-for-africa-cirrus-ai/">https://cortexlogic.com/2019/08/06/school-of-chemistry-launches-new-ai-research-initiative-for-africa-cirrus-ai/</a>
58	AI for Science Master's program	African Institute for Mathematical Science South Africa	Google Deepming	2023	<a href="https://ai.aims.ac.za/">https://ai.aims.ac.za/</a>

PAN AFRICA					
AI Africa Consortium – Cirrus foundry fund					

## Rwanda

NO	ACADEMIC LAB	INSTITUTION	MAIN FUNDING SOURCE	YEAR STARTED	REFERENCE (LINK)
59	World Economic Foundation Center for the Fourth Industrial Revolution	Inter-Institutions	N/A	N/A	<a href="https://c4ir.rw/press-release/">https://c4ir.rw/press-release/</a>

