

EU fellowship program, 2022-2023 academic year

University of California San Diego

**"Better law-making on AI
technologies in a digital age in a
complex global geopolitical context"**

V1.6 - 8 June 2023

Author:

**Peter Ide-Kostic, visiting scholar at the School of
Global Policy and Strategy**

Reviewer:

**Renee Bowen, Pastor Faculty Fellow, Professor;
Director, Centre for Commerce and Diplomacy of the
School of Global Policy and Strategy**



European Parliament

UC San Diego

SCHOOL OF GLOBAL POLICY AND STRATEGY

Contents

Thanks, and Acknowledgements	7
Part I: Introduction	8
1.0 Research purpose	9
1.1 Objectives	9
1.2 Structure of the report.	9
2.0 Executive summary	11
2.1 Artificial General Intelligence (AGI) reached by 2030	11
2.2 Prioritise national security over economic interests	13
2.3 A Democratic Technology Alliance for better coordination	14
2.4 Better lawmaking on AGI	15
Part II: What is AI	18
3.0 Simple examples of what is AI and what is not AI	19
3.1 Adding two numbers.	19
3.2 Recognising handwritten numbers	20
3.3 Analysing sentiments of a movie review	21
3.4 Predicting climate change impacts at the regional level	22
4.0 Key concepts and definitions	23
4.1 Machine Learning algorithms	24
4.2 Symbolic AI systems, Expert systems, Rule based systems	25
4.3 Data Collection and Data-Preprocessing	26
4.4 AI model	28
4.5 Learning, training, supervision, and inference phases	28
4.6 Neural Networks, Deep Neural Networks	30
4.7 AI Model parameters, weights, biases, tokens, and size	32
4.8 AI foundation Model, model fine-tuning, transfer learning	33
4.9 Multimodal AI models and Embodied AI models	36
4.10 Large Language Model (LLM), MLLM, and Embodied MLLM	36
4.11 Generative AI models	37
4.12 Emergent Abilities of Large AI models	38
4.13 AI systems	39
4.14 General Purpose AI system vs Narrow AI system	40
4.15 Supercomputers	40
4.16 Artificial General Intelligence (AGI)	41
4.17 The “Singularity” and the “Alignment Problem”	42
4.18 Artificial Super Intelligence (ASI)	44
4.19 Autonomous AI Agents	45
4.20 Session window of Large Language Models	45
5.0 Essential facts about modern AI Models	46
5.1 Purpose of this section	46
5.2 AI model size matters but not only	48

5.3	Transferring data across borders can be a legal obstacle.	50
5.4	DeepMind advanced AI with its matrix multiplication algorithms	52
5.5	DeepMind revolutionised AI with the transformer architecture	53
5.6	Why large AI models are impenetrable black boxes.	53
5.7	Berkley Retrieval Augmented LLMs help with “explainability”	54
5.8	The high carbon footprint of Large Language Models	57
5.9	Why Large AI models develop unexplained emergent abilities.	58
5.10	The “Intrinsic” emergent capabilities of LLMs	59
5.11	The “In context” emergent abilities of LLMs	60
5.12	Multimodal LLMs will gradually lead to AGI before 2030	63
5.13	Humans in Loop and “AI in the Loop” are both relevant.	64
5.14	Ethical, Responsible and Trustworthy AI overlap but differ.	65
5.15	“Ethical AI” and the “alignment problem”	66
5.16	Leveraging the “flywheel effect” in AI systems.	68
5.17	Large AI models can be compressed and optimised.	69
5.18	Autonomous AI Agents can self-reflect to improve.	72
5.19	Open source versus Closed Source Large Language Models	74
Part III: AI Governance and Regulations		77
6.0 Global governance on responsible and ethical AI		78
6.1	The United Nations and the works of UNESCO	78
6.2	The Organization for Economic Cooperation and Development (OECD).	79
6.3	The Council of Europe Committee on Artificial Intelligence (CAI):	80
6.4	The World Economic Forum	82
6.5	The Group of 7 (G7) and the “Hiroshima Process”	84
7.0 Non-profit organisations works on AI		85
7.1	IEEE (Institute of Electrical and Electronics Engineers)	85
7.2	ITIF (Information Technology and Innovation Foundation)	87
7.3	Brookings Institution	88
7.4	ISO/IEC JTC 1/SC 42	89
7.5	The AI Now Institute	91
7.5	AlgorithmWatch .	92
7.6	Partnership on AI (PAI).	93
7.7	Future of Life Institute (FLI).	94
7.8	Centre for Humane Technology (CHT).	96
7.9	AI for People.	97
7.10	Centre for Democracy and Technology	98
8.0 National Governmental frameworks on AI		99
8.1	The EU	100
8.2	The US	102
8.3	The UK	106
8.4	Japan	108

8.5 South Korea	109
8.6 China	110
8.7 Taiwan	112
8.8 Australia	113
8.9 India	115
8.10 New Zealand	116
8.11 Russia	117
Part IV: The History and future of AI	119
9.0 The History of AI (1943 – 2023)	120
9.1 1943-1956: AI Theoretical Foundations	120
9.2 1958-1997: AI progresses very slowly.	120
9.3 1997-2017: AI defeats the human brain and unveils its potential.	122
9.4 2018-2021: AI reasoning capabilities gradually emerge.	125
9.5 2022-2023: AI power unleashed to the public	129
10.0 The future of AI (after 2023)	135
10.1 Before 2030: Artificial General Intelligence (AGI) reached.	135
10.2 After 2030: Artificial Super Intelligence. (ASI) reached.	137
11.0 Leading research Institutes on AGI and ASI	138
11.1 Business	138
11.2 Academia	140
12.0 Factors enabling the development of AGI and ASI	141
12.1 Access to powerful computing resources	142
12.2 Access to vast and cheap sources of electrical energy	143
12.3 Access to software development environments	144
12.4 Access to large volume of quality data	145
12.5 Access to human Capital	147
12.6 Access to venture Capital (VC)	148
12.7 Measures to increase the flow of venture capital (VC)	149
12.8 Adopting an adequate regulatory Framework	150
12.9 Adopting an adequate industrial policy.	152
12.10 Maintaining a free and democratic society.	154
Part V: Geopolitical considerations	157
13.0 AGI and national security & strategic autonomy	158
13.1 National Security	158
13.2 International research cooperation	158
13.3 Open Strategic autonomy	160
14.0 AGI Societal impacts	161
14.1 Scenario 1: De-escalation of geopolitical tensions with China and Russia	161
14.2 Scenario 2: Aggravation of US cold war with China and Russia	164
14.3 Scenario 3: Direct military conflict between US and China	167
15.0 AGI and the theories of trade	169

15.1 Comparative advantage	169
15.2 Factor price equalisation	170
15.3 Global value chains	171
15.4 Trade in services	172
15.5 Trade barriers and protectionism	173
15.6 New trade flows	174
15.7 Economic integration	175
15.8 Absolute Advantage	176
15.9 Terms of Trade	177
15.10 Gains from Trade	177
15.11 Trade Balance.	178
15.12 Heckscher-Ohlin Model.	179
15.13 Ricardian Model.	180
15.14 Summary	181
16.0 AGI and Foreign policy	182
16.1 China and Russia joint statements on Global Governance	182
16.2 China and Russia quest for a new world order	183
16.3 Major conflict of interests with China and Russia	183
16.4 China overtly encourages destabilisation of US allies.	186
16.5 China and Russia weaponize trade for political reasons	187
16.6 China's ambiguous Belt and Road initiative (BRI)	189
16.7 AGI as a choking instrument to counter China and Russia	190
16.8 More export control on technologies enabling AGI.	191
16.9 More FDI screening on AGI.	192
16.10 Halt R&D cooperation on technologies enabling AGI.	193
16.11 AGI as an instrument to promote democracy.	194
17.0 AGI and the theories of International Relation	196
17.1 Summary of the 5 main International Relation theories	196
17.2 Country by country analysis	197
17.3 Working hypothesis for the next 5-7 years.	199
17.4 Impact of the emergence of AGI on International Relations	201
Part VI: Policy options	203
18.0 Policy options for AGI.	203
18.1 Prepare for societal impacts by 2030.	204
18.2 Better coordinate between like minded democratic partners	205
18.3 Leverage AGI as an external aid instrument	208
18.4 Improve legal certainty of data flows.	208
18.5 Address national security concerns	209
18.6 Develop adequate open autonomy strategies.	210
18.7 Regulate intended and accidental misuses.	211
18.8 Regulate illegal content generation.	211

18.9 Improve privacy protections.	212
18.10 Improve intellectual property rights protection.	213
18.11 Improve consumer protections.	213
18.12 Raise awareness about cybersecurity business implications	214
18.13 Define regulatory thresholds for AGI.	215
18.14 Influence of the public opinion on policy making for AGI	218
19.0 Better lawmaking in the European Parliament	220
19.1 Improve Parliamentary oversight of leading AGI companies.	220
19.2 A international parliamentary oversight mechanism	224
19.3 A standing Committee on Digital Affairs.	225
19.4 Special Parliamentary Committee on ASI and AGI (AIDA II)	226
19.5 Special Parliamentary Committee on EU relationship with China and Taiwan	227
19.6 Increase productivity with powerful AI based tools.	228
Part VII: Final Conclusion and Recommendations	231
20.0 Final conclusions on the safety of AI and future AGI	232
20.1 The remarkable cognitive abilities of Large Language Models	232
20.2 On a sure way towards a safe and robust AGI	233
21 Final policy recommendations	237
21.1 The need for coordinated AGI regulations among democracies	237
21.2 Addressing geopolitical tensions with China and Russia.	240
21.3 Better law-making in the European Parliament	242

Thanks, and Acknowledgements

The formulation of this report would not have been possible without the inspiration, contributions, and support of several esteemed individuals.

Professor Renee Bowen, Pastor Faculty Fellow and Director of the Center for Commerce and Diplomacy at UC San Diego's School of Global Policy and Strategy (GPS), has been a key source of guidance and constructive advice.

My gratitude is extended to Grace Osborne, Director of the Global Leadership Institute at GPS, for her logistical support throughout my fellowship.

My exchanges of views with Professor Mikhail Belkin of the Halicioğlu Data Science Institute at UC San Diego have been crucial for improving my knowledge of deep neural networks and machine learning.

I extend my gratitude to Professor Gert Cauwenberghs of the Bioengineering Department at UC San Diego for his insights on neuromorphic computing.

I am thankful to Todd L. Hylton, Executive Director of Mechanical and Aerospace Engineering and a Professor of Practice, for his insights on the concept of thermodynamic computing.

My sincere thanks go to GPS professors Mr. Samuel Bazzi and Ms. Weiyi Shi for their course on "Globalisation, the World System, and the Pacific," which has reminded me of the key foundations for writing this report.

I acknowledge the insightful contributions of Bill Bold, lecturer at GPS, whose two courses on "corporate non-market strategies" and "Technology, Trade, and Globalization" have expanded my understanding of globalisation.

Peter Cowhey, Dean Emeritus and Qualcomm Chair Emeritus of GPS, deserves a special mention for his course on "Digital Policy".

Lastly, I am thankful to Mr. Lei Guang, Executive Director of the 21st Century China Center at GPS, for his insightful perspectives on China and US foreign policies.

Peter Ide-Kostic, EU fellow, GPS visiting scholar, academic year 2022-2023

Part I: Introduction

Research Purposes

Structure of the report

Executive summary

1.0 Research purpose

1.1 Objectives

The final deliverable of this project is a report that traces the history of AI and forecasts its likely evolution in the next 5-7 years, based on the state-of-the-art research as of mid-2023. The report also examines the significance of Artificial General Intelligence (AGI) development for the US and its military allies in the current geopolitical context. The aim is to propose potential "policy options" and suggest options for "better law-making on AGI" within the European Parliament.

1.2 Structure of the report.

The first part of this report serves as an "Introduction" that outlines the project's scope and objectives, while also presenting an executive summary of the entire report.

The second part of the report delves into defining "What is AI." It offers examples of what is or is not considered AI and introduces terminology commonly used within the AI community. The focus is on the concepts of "AI model"; "AI foundation models," and related definitions such as "Generative AI," "embodied AI," "Large Language Models," and "Autonomous AI agents," among others. This section also explains the fundamental concept of "emergent capabilities" in "Large Language Models," which is crucial for achieving powerful intellectual reasoning and cognitive capabilities in AI and in the future AGI.

The third part of the report assesses the global governance efforts on AI and existing national frameworks on AI in various countries, including the US, the UK, EU countries, Russia, India, China, Taiwan, Japan, South Korea, Australia, and New Zealand.

The fourth part of the report explores the "History and future of AI". This section traces the history of AI from 1943 to 2023 and argues that artificial general intelligence (AGI) will likely develop rapidly by 2030 at the latest. The era of artificial superintelligence is expected to commence after 2030, based on the state of AI research and development as of mid-2023. This section also examines the leading research institutes anticipated to pave the way towards AGI and the generic factors enabling its development.

The fifth section of the report briefly analyses the geopolitical, economic, and societal implications of AGI development by 2030. It posits that, given the state of AI technology in 2023 and the fact that the most advanced know-how is concentrated in the US, AGI's emergence in the coming years will occur regardless of escalating geopolitical tensions with China and the Ukrainian war, even if a third world war breaks out. This section then explores national security, societal impacts, trade theories, international relations theories, and foreign policy considerations related to AGI development in different three scenarios.

The sixth section evaluates policy options based on the content of all previous sections. Twenty-one policy options are listed across 12 different domains. This section also proposes six options for "better law making" inside the European Parliament. It also looks at possible thresholds for regulating AGI and how public opinion could influence regulation in a context of increased geopolitical tensions.

The final, seventh section of the report sums up the conclusions and recommendations. It summarises the main conclusions regarding the safety of AI and future AGI systems. It presents the final recommendations for regulating AGI in a coordinated way across democracies through a "Democratic Technology Alliance", possibly based on the US West Coast, the need to effectively address geopolitical

tensions with China and Russia in the context of AGI and options identified for better lawmaking in the European Parliament.

2.0 Executive summary

2.1 Artificial General Intelligence (AGI) reached by 2030

The term "Artificial General Intelligence" (AGI) describes AI systems that can learn a wide variety of tasks, much like the human brain, and exhibit comparable or superior reasoning and cognitive abilities. As of mid-2023, the emerging capabilities of multi-modal large-language generative AI models, such as GPT-4, underpin their impressive reasoning and intellectual skills. Some aspects of these models already rival the human brain, as demonstrated in Microsoft's March 2023 research paper, "[Sparks of Artificial General Intelligence: Early experiments with GPT-4](#)"

OpenAI GPT-4 and similar models like Google PALM-2 are trained on vast datasets to predict the next word. Their exceptional cognitive and reasoning abilities emerge from training on these large datasets. The training adjusts numerous parameters in a deep neural network using an algorithm called "back-propagation," initially invented in the 1980s by Jeff Hinton from the University of Toronto. The incorporation of reinforcement learning, where positive or negative feedback is provided during training, has allowed these networks to learn much faster. In 2017, Google and Jeff Hinton of the University of Toronto introduced the "Transformer" architecture, another innovation.

Interestingly, even Jeff Hinton, one of the inventors of the back-propagation algorithm and Transformer architecture, cannot fully explain the extraordinary results achieved in 2023 based on his work from the 1980s and 2017. AI research in mid-2023 is still very experimental, often relying on intuitive leads and trial-and-error approaches. Because academic institutions occasionally lack

access to the necessary supercomputing resources, industry rather than academia is primarily responsible for driving research on engineering ultra-powerful GPT-4-like models.

Academic research mainly focuses on hardware and algorithmic improvements that increase speed, reduce energy consumption, and ensure AI models are responsibly developed and aligned with human values. Additionally, AI is being used for optimization purposes, such as designing complex chips, discovering new composite materials, and genome sequencing (CRISPR). Several universities have also been working on developing tools to evaluate the robustness, trustworthiness, security, and reliability of large language models and to address the "Alignment Problem" of aligning AI with human values. These tools are anticipated to play a crucial role in conformity assessment procedures for high-risk AI systems.

Experts in deep neural networks assert that AI models like GPT-4 build an "inner representation of the world" based on their training data. The larger the model and the higher the quality of the dataset, the more accurate the prediction of the next word (or "token") will be, leading to more reliable and effective emergent cognitive and intellectual abilities. However, this also requires increased computing resources.

Research in mid-2023 has shown that large enough models [can even engage in "self-reflection,"](#) improving their performance over time by learning from past experiences. By early 2023, following GPT-4's release, most experts agreed that the current generation of multi-modal large language generative AI models would lead to AGI development by 2030 at the latest (depending on the AGI definition adopted).

As AGI develops, it is essential to maintain control over the phenomenon of emerging abilities as well as AI model tendencies to occasionally

hallucinate, discriminate, or display undesirable behaviour misaligned with human values. In the coming years, AI systems' predictability, explainability, robustness, and reliability are expected to improve significantly, leading to robust, safe, and reliable AGI systems. However, since the risk of misbehaviour and misalignment can never be eliminated entirely, society must accept residual risks after conducting adequate public debates.

2.2 Prioritise national security over economic interests

In the past 5-7 years, the US has played a significant role in advancing complex large language models, with contributions from companies like Google Brain, Google DeepMind, OpenAI, Meta, and Anthropic. Recently, Nvidia, Amazon, and Microsoft have also entered the field.

Although the knowledge and resources needed to develop narrow AI systems and less advanced large language models are now widely accessible worldwide, the expertise required to build complex AI systems with performance comparable to GPT-4 or PALM-2 remains concentrated in the US private sector among a select group of individuals. This contrasts with the more widespread expertise in "narrow AI systems" found globally, particularly in China and Russia.

Considering the current geopolitical tensions and potential risk of conflict between the US and its allies on one side and China and its allies on the other, it is crucial to consider the US's significant technological advantage in AGI. Recognizing the potential benefits rapid AGI development could provide to China and Russia for military purposes, all US allies should actively support policies aimed at impeding AGI development in these countries, even if it results in economic impacts for companies in the US and other allied nations. National security should take precedence

over the economic interests of private companies investing in China and Russia.

While AGI development in China and Russia will take place in the long term, as the mastery of deep neural networks' gradual emergent capabilities represents a paradigm shift for humanity akin to early humans mastering fire, it is possible to ensure that the US and its allies stay well ahead for some period. This delay would allow the US and its military allies to further consolidate their technological superiority and/or buy time until current geopolitical tensions and the risks of war subside.

Potential policies include controlling the export of relevant equipment, components, chemical materials, technology transfers, and intellectual property rights, as well as scrutinising all incoming and outgoing foreign direct investments. Unfortunately, halting cooperation on research and development with China and Russia in domains that could enable them to develop their AGI capabilities will also be necessary.

Simultaneously, borders should remain open, and the US and its allies should encourage and incentivize STEM talents and AI experts in China and Russia to seek political refuge in the US and allied countries (after adequate security screening for national security purposes).

2.3 A Democratic Technology Alliance for better coordination

The US and its military allies in Europe and Asia should harmonise their national security policies on AGI, strategically defining the technological domains they are willing to collaborate on, those they prefer to keep exclusive for national security reasons, and those in which they aim to achieve greater strategic autonomy (either at the national level or as a group of countries like the EU sharing similar objectives).

In the context of AGI development, it is vital to address challenges posed by end-to-end AI usage, AI misuse, privacy, cross-border data flow legal certainty, intellectual property rights, interoperability, cybersecurity, online system disintermediation, trade secrets, early AGI impacts on education and employment, as well as establishing rules for identifying content that incites hate and harmful speech. Although some legislation exists in the EU, it will need to be revised and completed considering the emergence of "AGI." The US has taken a lighter approach with its AI Bill of Rights and the NIST AI risk management framework. Other democratic countries will likely adopt similar frameworks, and it is crucial for trade and business that they remain sufficiently compatible to be declared mutually compatible.

Given that all democratic countries developing AGI systems face the aforementioned policy challenges, [the G7 agreement from May 2023 and the decision to start the Hiroshima process](#) both call for close coordination between the US and its allies. Considering the current geopolitical tensions with China and Russia and the rapid AI evolution toward AGI, establishing a "Democratic Technology Alliance" with a permanent structure could facilitate more effective collaboration among democratic nations. This Democratic Alliance should be best located in the US on the west coast, home to leading AI and AGI companies, which is secure and nearly equidistant from all US allies.

2.4 Better lawmaking on AGI

The European Parliament, and potentially other parliamentary democracies, are advised to strengthen their oversight of companies developing AGI, with a focus on Google, OpenAI, Anthropic, Microsoft, Tesla, Amazon, and Nvidia in 2023. This will allow parliamentary democracies to closely track AGI technology's progression and better prepare for its impacts. Public debates should be organised to address various policy questions

related to usage, misuse, intellectual property rights, and other concerns, as well as the alignment problem.

Companies leading on AGI would also benefit from this increased oversight, as it would lead to clearer rules, more legal certainty, and fewer risks of overregulation with the potential to stifle innovation.

The alignment problem pertains to the risks of AGI systems misbehaving and the controls and safeguards required to ensure their adherence to human values. For instance, questions that may arise include whether AGI systems should be allowed to self-modify their code to enhance performance without human intervention or control their access to power sources, and what level of autonomy would be acceptable for a critical AGI system controlling a nuclear plant.

One option for the European Parliament is to create a new standing committee responsible for EU digital horizontal legislation within the context of the EU's Digital Single Market. This Committee's competencies would encompass a range of areas, including AI, privacy, cybersecurity, intellectual property rights, antitrust, small and medium enterprises, and R&D in the digital domain in relation to the EU single market. Alternatively, a specialised AIDA II committee focused specifically on AGI could be established to address the unique challenges and opportunities presented by this emerging technology.

Additionally, it is suggested that a special committee be created to address EU relations with China and Taiwan. This committee would tackle challenges posed by increasing geopolitical tensions and the need for close coordination with the US and other allies. Like the US House of Representatives' "Committee on Strategic Competition between the United States and the Chinese Communist Party", its mandate would include examining policy issues between the EU and China,

as well as the EU and Taiwan, in the context of rising global polarisation and potential armed conflict between the US and China.

An international supervision model for AGI product and service providers is proposed, wherein national legislative and executive branches assume primary responsibility for supervision and collaborate with counterparts in other democratic countries. This cooperative approach aims to facilitate effective oversight and shared responsibility for AGI service development and deployment, ensuring ethical standards and safety precautions are maintained globally.

Finally, the European Parliament should consider adopting a policy for using advanced tools such as ChatGPT-4 or GPT-4 within the institution. This decision should be made after thoroughly weighing the pros and cons and conducting all necessary internal consultations.

Part II: What is AI

What is / is not AI

Key concepts and definitions

Essential properties of modern Large AI
models

3.0 Simple examples of what is AI and what is not AI

3.1 Adding two numbers.

Imagine you want to create an application that can calculate the sum of two numbers, such as $10+1=11$.

The non-AI method:

One way is to use the non-AI method, which is to write a simple program that takes two numbers as input and gives the sum of those numbers as output.

In the example above, the output is programmed. It is clearly a non-AI method.

The AI method:

Alternatively, the AI method involves creating an application that is first trained on a dataset and then learns to generate the output by identifying the correct mathematical relationship to use. During training, the AI application identifies patterns and relationships between the input and output data and then builds a mathematical model that satisfies these criteria.

For example, the AI application might be trained on a dataset containing the input-output pairs $(0,0:0)$, $(0,1:1)$, $(1,0:1)$, $(1,1:2)$, $(1,2:3)$, and $(2,1:3)$. From this dataset, the AI algorithm can deduce that the mathematical operation to be discovered is commutative, that the output is always greater than the input, and that when one of the inputs is zero, the output is the same as the second input. When one is involved, the output is the second input incremented by one, and so on. These findings are sufficient for the AI algorithm to conclude that the two numbers 10 and 1, not included in the initial data set, must be added together.

The key conceptual difference between the AI and non-AI methods is that the AI application "learns" what the correct mathematical relationship is to generate the output, whereas the non-AI method is programmed to know this relationship in advance. Unfortunately, not everything can be programmed in advance, and this is why AI is so useful...

In the second example above, the machine learning algorithm uses a pre-defined set of arithmetic knowledge and rules to learn how to process the input without using statistical methods. Despite not relying on statistical methods, this still qualifies as a machine learning approach.

3.2 Recognising handwritten numbers

A more complex example would be to create an image recognition application that can identify handwritten numbers between 0 and 10.

The non-AI approach:

The non-AI approach would require the application to digitize and process the input image, then compare it with pre-existing sample images of each number to determine the closest match and output the corresponding number.

In the example above the result the system does not attempt to learn from data in an initial phase and to make a prediction, it is programmed to compare the input data to a set of sample data. It is non-AI.

The AI approach:

The AI method, on the other hand, would involve training the application on a large set of different handwritten numbers. This would teach the application that certain visual features correspond to certain numbers. For example, the AI model may learn that numbers with closed loops correspond to 0 or 8, numbers with two loops interconnected in

the middle correspond to 8, and numbers with a straight vertical bar correspond to 1.

By finding these patterns and connections, the AI algorithm can estimate the likelihood that a given input image corresponds to each number between 0 and 10, even if it has never seen that handwriting style before. The AI model just gives out the number with the highest chance of happening.

In this example, the AI application is trained using machine learning algorithms to build a statistical model with a probability distribution for all possible outputs.

3.3 Analysing sentiments of a movie review

Suppose you want to build a text classification application that can determine whether a given movie review is positive or negative.

The non-AI approach:

The non-AI approach would involve manually analyzing the text of each review, looking for specific keywords or phrases that indicate a positive or negative sentiment, and then assigning a label accordingly. This approach is highly dependent on the expertise and biases of the person performing the analysis, and may not be accurate for all reviews.

The first example is similar to the non-AI model for adding two numbers; there is no learning that takes place, and it is non-AI.

The AI approach:

The AI approach, on the other hand, would involve training a machine learning model on a large set of already labelled positive and negative movie reviews. The model would learn to find patterns and connections between how the reviews are written and how they make people feel.

For instance, the model might learn that reviews containing phrases such as "excellent acting," "riveting plot," or "outstanding direction" tend to be positive, while reviews containing phrases such as "poor script," "uninspired performances," or "boring cinematography" tend to be negative.

Once the model is trained, it can be used to analyze new reviews and generate a sentiment classification based on the probabilities of positive and negative labels. This approach is more scalable and less subjective than the non-AI approach and can be applied to a wide range of text classification problems.

Training takes place on a large dataset, so it is AI

3.4 Predicting climate change impacts at the regional level

Let's say you want to make an app that can predict how climate change will affect a certain region or area in the future.

The AI approach - expert and knowledge based:

It would involve using climate models and simulations that have already been made by experts to predict the weather. This would typically involve analysing historical weather patterns, accounting for factors such as temperature, precipitation, and atmospheric conditions, to build decision trees to generate predictions for the future. Such a model may be limited in its ability to accurately predict the complex and dynamic nature of climate change based on the limited knowledge of the experts that can be encoded as rules and considering that they may not be able to capture all variables and interactions involved.

The approach described above is "non-machine learning based." However, as it requires the contributions of experts and the constitution of a

knowledge base and rules, it is still considered an AI system.

The AI approach - machine learning based:

It would require training a machine learning model on a large set of historical climate data, including temperature, precipitation, atmospheric conditions, and other variables like ocean temperatures and ice coverage. The model would learn to find complex patterns and relationships between these variables and the effects of climate change, such as rising sea levels, extreme weather, and changes in plant and animal populations.

For example, the model might learn that certain combinations of weather variables tend to result in more severe storms or droughts, or that certain types of plant and animal populations are more vulnerable to specific types of climate change impacts. All this without recourse to human expertise.

Once the model is trained, it can be used to make predictions about how climate change will affect a certain region or area in the future based on climate data that is already available. As more data comes in, the accuracy of the predictions can be improved. This lets the model adapt to changing climate patterns and get more accurate over time.

The approach is machine learning, so AI-based, but it is not expert-based, as no expertise is encoded or programmed as fixed rules in the system.

4.0 Key concepts and definitions

This research project focuses on "AI Technologies in a Digital Age," which refers to the concepts of "AI models" based on machine learning algorithms. It also covers the concept of "AI foundation models", as introduced by Stanford in 2018.

To provide clarity, we also review the different terminologies commonly used in the AI literature and make the link with the one of "AI system" introduced in the EU AI Act proposal released in April 2021 and in the NIST Risk Management Framework released in January 2023.

4.1 Machine Learning algorithms

A machine learning algorithm is a set of instructions or rules that are used to learn from data and make predictions or decisions without being explicitly programmed for them. This is usually done by using statistical methods or by discovering logical rules from the trained data. The machine's algorithm processes the trained data and identifies patterns or relationships between input and output. Machine learning algorithms can be broadly categorized into supervised, unsupervised, and reinforcement learning algorithms. [1].

Supervised machine learning algorithms are trained on labelled data, where the correct answers are provided, and the algorithm learns to predict the output based on the input's characteristics. Common supervised machine learning algorithms include linear regression, decision trees, random forests, and support vector machines.

Unsupervised machine learning algorithms, on the other hand, are used to find patterns and relationships in unlabeled data without being provided with any specific output variable. Common unsupervised machine learning algorithms include clustering, principal component analysis, and association rule learning.

Reinforcement learning algorithms are a type of machine learning algorithm that learns by interacting with an environment and receiving feedback in the form of rewards or penalties.

Common reinforcement learning algorithms include Q-learning and policy gradient algorithms [2].

Machine learning algorithms can also be classified based on the type of problem they solve, such as regression, classification, clustering, and recommendation systems. Regression algorithms are used to predict continuous values, while classification algorithms are used to predict discrete values. Clustering algorithms group similar data points together, while recommendation systems are used to suggest items or products based on a user's past behaviour. [5].

In summary, a machine learning algorithm learns from data and makes predictions or decisions without being explicitly programmed. There are several categories of machine learning algorithms, including supervised, unsupervised, semi-supervised, self-supervised, and reinforcement learning algorithms. Well-engineered complex AI systems would combine the use of different AI models using different machine learning algorithms that have been carefully chosen based on the goals to be achieved and the different types of data used for training.

Regression, classification, clustering, and recommendation systems are just a few examples of how to categorise machine learning algorithms according to the problems they solve.

4.2 Symbolic AI systems, Expert systems, Rule based systems

Expert-based, rule-based, or knowledge-based systems are called symbolic AI systems, and they do not use machine learning algorithms. They are programmed based on the knowledge of human experts to deliver a given output based on a given input. For historical reasons, they are nevertheless considered "AI" because they were the first generation of systems qualified as "AI" and because they help take decisions or make predictions based on true human expert knowledge. They are used in

many fields where a large amount of domain knowledge is available and where the accuracy of decisions is crucial. They can also be used in conjunction with machine learning to create complex AI systems that optimise the quality of machine learning algorithm predictions.

In some cases, the available knowledge about the process is so good that it can be directly encoded in the system. In other situations, the system can clearly store logical rules that can be used to make decisions or predictions. These systems are not considered machine learning algorithms because they are based on pre-programmed rules or knowledge.

Rule-based systems have been used a lot in fields like medicine, finance, and law, where getting decisions right is very important. Medical imaging and diagnostics, as well as the development of self-driving cars, have also seen a high use of rule-based systems.

Note that if a system learns logical rules from the dataset during the training phase, it is considered a machine learning model because these rules have not been pre-encoded or pre-programmed but were learned from the data during the training phase.

4.3 Data Collection and Data-Preprocessing

Data Collection

Machine learning algorithms require large amounts of quality data to be trained, which can be collected from various sources. [\[1\]](#) [\[3\]](#) [\[5\]](#).

It is interesting to note that AI companies like OpenAI or Google are very discreet about the data collection processes they use. The ability to trace back the source of the data used by AI systems based on their output is expected to become a growing policy question for obvious reasons given the following issues at stake: intellectual

property protection, privacy protection, cybersecurity, and the fight against disinformation.

Some possible sources of data include public data sets, pre-packaged data, social media data, web scraping, and sensor data [3], simulated data, or simply a company's own industrial data! Cloud-computing companies such as Amazon, Google, Microsoft, Tencent, Huawei, Alibaba, or OVH also offer access to data sets or APIs for data collection. [3] [7].

Data Pre-Processing for AI purposes

Data preprocessing is a term that relates both to "Data mining" and "AI". In the context of AI, it is an important part of the AI pipeline as it involves putting raw data into a format that machine learning algorithms can easily understand. The main goal of data pre-processing is to improve the quality and accuracy of machine learning results by making the data easier to understand, more relevant, and more reliable. This process involves cleaning, transforming, and normalising data to make it easier for machine learning algorithms to look at it and draw conclusions or make predictions. Businesses (Google, Amazon, etc..) that provide cloud services to collect data also typically provide tools for pre-processing and storing such data.

Pre-processing data for machine algorithms using the so-called supervised learning method requires labelling the raw data accurately so the training is done adequately. This labelling can be done manually or using automated labelling techniques. For instance, labelling the images of dogs and cats for training a computer vision model involves marking each image as either a dog or a cat. In contrast, if the unsupervised learning approach is followed, there is no need for a predefined label for the data; in this case, data pre-processing

only involves cleaning and transforming the raw data into a structured format.

Overall, data can be collected from multiple sources, including a company's own industrial data or from cloud computing service providers. Data preprocessing is a critical step in the AI pipeline that helps to improve the quality of the data, making it more relevant and meaningful to the models. By doing it right, AI models can produce more accurate and reliable results. Businesses that provide cloud services to collect data, such as Amazon, Microsoft, Google, etc., also typically provide tools for pre-processing and storing such data for those who train machine learning algorithms.

4.4 AI model

An "AI model" is a term used in the AI industry to refer to the combination of machine learning algorithms, datasets, and training methods that are used to develop and deploy a specific application.

For historic reasons, note that AI models also cover the case of experts' systems, which use symbolic AI and where no training takes place.

Machine learning-based AI models can be trained on a variety of data sources, depending on the specific application and task at hand.

The term "AI model" is very generic and is typically used to refer both to the development of simple "narrow use" applications such as image or speech recognition or to much more complex applications with many different possible uses, such as Large Language AI models, for instance. In practice, complex "AI systems" are engineered by combining different types of AI models.

4.5 Learning, training, supervision, and inference phases

Learning phase

In the context of artificial intelligence (AI), "learning" refers to the process by which a machine learning algorithm improves its performance on a task through experience or exposure to data. [5].

There are different types of learning, including:

- a) Supervised learning where the algorithm learns from labelled data during the data pre-processing phase.
- b) Unsupervised learning where the algorithm extracts features and patterns from unlabelled data [5].
- c) Reinforcement learning is another type of learning where the algorithm learns by interacting with an environment and receiving rewards or penalties based on its actions. [8]
- d) There is also the self-supervised learning approach where the machine algorithms use a few examples given at the beginning of the training (the pretext task) to learn how to do the labelling itself.
- e) Finally, there is the semi-supervised learning approach, where the training is divided into two parts: the bulk of the training data is unsupervised, and the last part is supervised based on a more limited data set to better fit the purposes intended. The semi-supervised learning approach is typically used with AI foundation models; during pre-training, no labelling takes place, while during fine-tuning (after the pre-training is completed), data can be labelled.

Complex AI systems are engineered using different types of machine learning algorithms chosen during

the design phase to meet the purpose intended and taking into account the type of training data processed.

Training phase

Training in the context of AI refers to the process of teaching an AI model to perform a specific task using a dataset or set of rules. This involves selecting an appropriate machine learning algorithm, preparing and cleaning the data, defining the evaluation metrics, and iterating through the training process to optimise the model's performance. [6].

Supervision phase

In the context of AI, "supervision" means how a human expert, or another AI agent guides and keeps an eye on the machine learning process. This can involve monitoring the training process, reviewing the model's performance, providing feedback, and making adjustments to the training data or algorithm. [10].

Inference phase

Inference is the process of applying the learned patterns and features to new data to make predictions or decisions based on what the model has learned during training.

To sum up, an AI model typically goes through several phases:

1. Learning: The AI model acquires knowledge from the data it is trained on.
2. Training: The AI model uses the acquired knowledge to make predictions or decisions on untrained data.
3. Supervising: During training, the AI model is supervised by a human expert (or another AI model) who (that) provides feedback to improve its performance.

4. Inference: The trained model is used to make predictions or decisions on new, unseen data.

4.6 Neural Networks, Deep Neural Networks

A neural network is a type of machine learning algorithm that is designed to recognize patterns in data. It is a computational model that draws inspiration from the structure and operation of the human brain. It consists of a network of "neurons" (nodes) connected by "synapses" (links) that transmit information between the input and output. Deep neural networks are simply neural networks that are wide and deep with many layers of neurons and that correspond to specific machine algorithms.

In the field of machine learning, neural networks are a powerful tool because they automatically learn from data and get better over time. They can handle complex and nonlinear relationships between inputs and outputs, making them useful for solving a wide range of problems.

Not all machine learning algorithms are neural networks. Neural networks are just one type of machine learning algorithm.

There are several other types of machine learning algorithms, such as decision trees, random forests, support vector machines (SVMs), k-nearest neighbours (KNN), etc., that are not neural networks. Each of these machine learning algorithms has its strengths and weaknesses and is better suited for certain types of tasks. For example, decision trees are useful for classification tasks, SVMs are effective for binary classification, and KNN is often used for recommendation systems.

The hardware used to run a neural network also affects the choice of architecture. This means that hardware capabilities should be considered when designing and implementing a neural network. AI is

as much about hardware as it is about software, algorithms, and data. Some machine learning algorithms can run on simple generic Intel Central Processing Units (CPUs), while others require powerful specialised hardware such as GPUs (Nvidia) or TPU (Google).

Different neural network architectures are suitable for different types of tasks as well. For example, convolutional neural networks are commonly used for image recognition tasks, while recurrent neural networks are often used for natural language processing.

Deep neural networks are simply neural networks with a very high number of layers and are typically used for AI foundation models such as PALM (Google) or GPT-4 (OpenAI).

4.7 AI Model parameters, weights, biases, tokens, and size

Neural networks are used in a variety of machine learning applications, such as image recognition, natural language processing, and predictive modelling. Each node (neuron) in the network has a "bias" parameter "b", and each interconnection (synapse) between two neurons has a "weight" parameter "w." The "weights" and "biases" are called the "parameters" of the neural network of the AI model.

During the training phase, all weights "w" and all biases "b" are adjusted to optimise performance. The training data is broken up into "tokens," which are the smallest pieces of information that are processed during the training and which are used during the learning phase to adjust the parameters. For instance, for GPT-3, the size of the trained dataset is estimated to be 45 terabytes, which is equivalent to several trillion tokens, assuming a token is a word or a fraction of a word. The number

of model parameters associated with GPT-3 is 175 billion.

The size of the trained dataset, the number of tokens, and the number of model parameters are decided based on the specific nature of the dataset used, the tokenization method used, and the goals and constraints of the project (such as computational resources and desired performance).

The size of a Large Language Model, such as GPT-3 or 4, refers to both the size of the trained dataset in terms of tokens and the number of model parameters. These two aspects contribute to the overall size and complexity of the model. For GPT-4, this information has unfortunately been kept confidential for trade-secret reasons.

The size of the input data set is often also expressed in tokens. For instance, for GPT-4, it is a maximum of 32000 tokens, corresponding to about 48 pages of text.

4.8 AI foundation Model, model fine-tuning, transfer learning

[The Stanford Institute for Human-Centred Artificial Intelligence \(HAI\)'s Centre for Research on Foundation Models \(CRFM\)](#) first used the term "AI foundation model" in 2018.

An AI foundation model is an "unsupervised" AI model that is designed and "pre-trained" to be used for different purposes.

The organisation that deploys the foundation model "fine-tunes" it by completing the training on an additional data set that is suited to the final purpose intended.

During the fine-tuning, it is also possible to transfer knowledge from one task to the other

(i.e., transfer learning) if the capabilities needed for the two tasks are the same. For instance, an AI Foundation model that was pre-trained with extensive data to recognise fish on pictures, if shown a few pictures of birds during the fine-tuning, will then be able to recognise both flowers and birds with the same level of accuracy.

Transfer learning during fine-tuning from fish to birds is only possible because the AI foundation model learned to recognize features such as colors, shapes, and other characteristics during pre-training on a large dataset of fish images. By fine-tuning a smaller dataset of bird images, the model can re-use those learned mechanisms to reliably recognize birds, even with only a few examples of bird pictures.

During pre-training, the foundation model developed a general representation of visual features that are relevant to objects in images. This allows the model to recognize certain visual patterns that are common to both fish and birds. During fine-tuning, the model adapts its representation to the specific characteristics of bird images by adjusting its parameters. This process allows the model to learn new patterns and features specific to bird images, which improves its ability to recognize them accurately.

It's important to note that if the fine-tuning had been done with satellite images with the objective of recognizing roads, this would not have worked effectively. This is because the visual features and characteristics required to recognize roads from satellite images are too different from the

ones developed to recognize birds and fish in regular images.

Roads in satellite imagery appear as linear features with specific patterns of lines and shapes that do not correspond to any specific animal features. Therefore, the pre-trained model that was developed to recognize birds and fish based on their colors, shapes, and textures, would not be well suited to recognizing roads in satellite imagery.

In this case, it would be necessary to train a new model from scratch on a dataset of satellite images of roads, using a neural network architecture and training techniques that are optimised for the task of road recognition in satellite imagery. This would enable the model to learn specific visual features and patterns that are relevant to the new domain and to achieve better accuracy and generalisation performance than transfer learning from a model pre-trained on a very different type of dataset.

Overall, transfer learning during fine-tuning enables the model to leverage its knowledge from pre-training to improve its performance on a new task or domain, even with a limited amount of data. This is possible because the model has learned to recognize common visual patterns that are shared across different types of objects in images.

The supervision and additional training (i.e. fine tuning) of an AI foundation model are done by the entity that deploys and operates it in the field. It can be different from the organisation that initially developed it. The deploying or operating

entity also supervises the model's performance to ensure that it meets the intended objectives.

Overall, the process of training and fine-tuning AI foundation models is an iterative process that involves ongoing feedback and adjustments to ensure optimal performance for the specific use case.

4.9 Multimodal AI models and Embodied AI models

"Multimodal" means that the AI model accepts different "modalities" as input and can also make predictions or generate output under different "modalities" as well.

The different modalities can be, for instance, "text," "image," "video," "audio," or even, in the case of robotic applications, "sensor-related information" such as temperature on the input.

An "embodied AI model" refers to a model designed to take actions in its environment. " Embodied AI models are multimodal by nature.

The goal of an embodied AI model is to allow the robot equipped with it to move through the world and interact with its physical environment through actions. Embodied AI models are considered multimodal as they incorporate various sensory modalities such as vision, touch, and auditory processing to interact with their environment. For instance, Google's PALM-E, designed for robotic applications, is an embodied multimodal AI Model.

4.10 Large Language Model (LLM), MLLM, and Embodied MLLM

A Large AI Model, in general, refers to a deep neural network that has a very high number of parameters and is trained on a large dataset. Large AI models can be language-based, vision-based, video-based, image-based, audio-based, or even multimodal. The focus is on the use of deep neural networks with a very high number of parameters and a very large training dataset.

A Large Language Model (LLM) is a language-based Large AI Model with a very high number of "tokens" and "parameters". Examples of Large Language Models are OpenAI GPT-3, ChatGPT, and Google BERT.

Google ViT is a vision-based Large AI Model, and DeepMind WaveNet is a voice-based Large AI model.

Multimodal Large Language Models (MLLMs) are Large Language Models with multimodal input and/or output modalities. The latest generation of AI Models often includes MLLMs, such as OpenAI GPT-4, Google PALM, and Microsoft Kosmos-1.

An Embodied MLLM is a Multimodal Large Language Model designed for robotic purposes, such as Google PALM-E, for instance.

Regarding the size of LLMs, it is considerable. For example, OpenAI GPT-3 was trained on a massive dataset with 45 terabytes of text, equivalent to hundreds of billions of tokens, and has over 175 billion parameters. Google PALM-E has been trained on a corpus of 1.4 trillion tokens and has 570 billion parameters. Meta/Facebook LLaMA has been trained on either one trillion or 1.4 trillion tokens and has 65 billion parameters. Unfortunately, OpenAI has decided to keep the size of the GPT-4 AI model secret.

4.11 Generative AI models

Generative AI models are a type of large multimodal AI model that combines the use of several advanced machine learning techniques to generate new

content, such as text, images, videos, voice, or sounds, based on existing patterns or examples. They are often multimodal, but not necessarily always so, as they can also use Unimodal Large Language Models (ULLMs) or other specialised models depending on the application.

OpenAI GPT-4 is a generative AI model that can generate text based on image or text input. It falls into the category of Generative Multimodal Language Models (GMLMs). Note that GPT-3.5, which can only process text input and output, is a Generative Large Language Model (GLLM) and not considered multimodal.

Some other examples of Generative AI models include OpenAI DALL-E2, which can create new images from text descriptions or update existing images based on user instructions.

4.12 Emergent Abilities of Large AI models

The ability of a Large AI model to perform a specific task or prediction is considered emergent if the Large AI model has not been explicitly trained or fine-tuned to perform the task, and if that ability gradually appears and strengthens as the training dataset is increased. The subject is actively researched by Academia with diverging views on the origins of the phenomenon.

Large Language Models (LLMs), Multimodal Large Language Models (MLLMs), and Generative AI models are trained on vast datasets and exhibit remarkable emergent capabilities for which they have not been explicitly trained or fine-tuned. These capabilities significantly contribute to the success of Multimodal Large Language Models, such as GPT-4 or PALM-2.

MLLMs, like GPT-4 have been trained on a substantial portion of the Internet, so it is no surprise that they display exceptional emergent capabilities.

According to the latest research on the subject, these capabilities emerge gradually as the size of the model and the size of the dataset increase. Some capabilities, such as code generation or language translation, are highly desirable, but others may not be. It is therefore crucial to characterise all emergent abilities of a large-scale AI system before its deployment and to fine-tune it to mitigate any negative consequences. By doing so, the full potential of large AI models can be harnessed while avoiding any harmful impacts.

4.13 AI systems

a) as defined for the purpose of this research work

An "AI system" is just a system that is engineered from one or more "AI models", "AI foundation models", or "Large AI model" and is supervised in a production environment. Note that Generative AI models such as Stability.AI. Stable Diffusion and Stanford Alpaca which have been optimised to run standalone on a powerful PC also fall into this category, with the production environment being, in this case, one of the users running the model.

b) as defined [in the EU AI Act](#)

According to the initial EU AI Act released in April 2021, an AI system is defined broadly as "software that is developed with one or more of the techniques and approaches listed in Annex I and can, for a given set of human-defined objectives, generate outputs such as content, predictions, recommendations, or decisions influencing the environments they interact with." [2].

Annex I of the EU AI Act lists the techniques and approaches used in AI systems covered by the Act. Annex I covers machine learning, logic and knowledge-based approaches, and statistical or Bayesian approaches that can generate outputs such as content, predictions, recommendations, or

decisions that influence the “environments they interact with” as per Article 3(1) and Annex I.

Annex II gives guidelines for figuring out which AI systems are high-risk and should be regulated more.

Annex III has a list of AI applications that are high-risk and have to follow certain rules because of the Act.

(c) as defined in the [NIST AI risk management framework](#).

“AI system” is referred to in the NIST management risk management framework published in January 2023 as an engineered or machine-based system that can, for a given set of objectives, generate outputs such as predictions, recommendations, or decisions influencing real or virtual environments. AI systems are designed to operate with varying levels of autonomy (adapted from: OECD Recommendation on AI:2019; ISO/IEC 22989:2022).

4.14 General Purpose AI system vs Narrow AI system

This terminology of General-Purpose AI is used in the EU AI Act.

That definition is in opposition to “narrow AI systems”, which are systems designed to do only one task. A general-purpose AI system in the EU AI Act corresponds to an AI system that is engineered for different possible purposes.

Article 52a of the AI Act proposal of the Commission released in April 2021 states that the placing on the market, putting into service, or use of general-purpose AI systems shall not, by themselves, make those systems subject to the provisions of this Regulation. [8] However, there is no specific definition in this article of what constitutes a general-purpose AI system.

It can be concluded that in the initial proposal of the Commission introduced in April 2021, AI systems that can be used for many different purposes are not considered high-risk as such unless they are specifically set up to be used for a high-risk purpose.

4.15 Supercomputers

A supercomputer is a type of computer that is designed to perform complex calculations and solve problems that require a tremendous amount of processing power. Supercomputers are primarily used for scientific and engineering work where high-speed computations are required [1]. Supercomputers also currently consume a lot of electrical energy and have a high carbon footprint.

Supercomputers are a key part of the growth of AI because they are needed to process large amounts of data and train machine learning models. AI algorithms require enormous amounts of computational power to learn from data, make predictions, and improve over time.

4.16 Artificial General Intelligence (AGI)

The debate surrounding the necessary conditions to achieve Artificial General Intelligence (AGI) remains a topic of ongoing discussion among AI experts. There are two primary perspectives in this debate:

1. **Embodied AI approach:** Proponents of this approach argue that AGI requires AI systems to be embodied in robots capable of interacting with their environments in the same way humans do. They believe that physical interaction with the world is crucial for acquiring an understanding of concepts and developing general intelligence.

This approach draws inspiration from developmental psychology and cognitive science, which suggest that human intelligence is deeply rooted in our ability to perceive and interact with the physical world.

- 2. Purely cognitive approach:** Advocates of this perspective argue that AGI can be achieved without requiring AI to be embodied in robots. They believe that developing powerful, robust, and reliable reasoning capabilities, along with the ability to learn as the human brain does, is sufficient to reach AGI. This approach focuses on creating advanced algorithms and neural network architectures to replicate human cognitive abilities.

The embodied AI approach emphasises the importance of grounding intelligence in real-world experiences, while the purely cognitive approach highlights the significance of advanced reasoning and learning capabilities.

Mid-2023, most experts support the **pure cognitive approach** where Artificial General Intelligence (AGI) refers to an AI model, AI foundation model, or AI system whose cognitive and intellectual abilities are comparable to, or superior to, those of the human brain.

Since there is no universally agreed-upon definition of human intelligence, experts have varying opinions about when artificial general intelligence (AGI) will be achieved. Most scientists agree that GPT-4 already possesses superior cognitive abilities to the human brain in certain areas.

Most experts agree that whatever the exact definitions of AGI, it is likely to be developed before 2030, first based on the pure cognitive approach in a few years and then, by 2030, also based on the embodied approach. The development of AGI is considered a challenging and complex task

that requires significant advancements in multiple areas, especially in the embodied approach.

4.17 The “Singularity” and the “Alignment Problem”

The singularity refers to the hypothetical point at which ASI (Artificial Superintelligence) significantly exceeds the capabilities of the human brain. Many experts believe that reaching the singularity will lead to numerous benefits, such as a rise in wealth and the discovery of new scientific knowledge. However, the benefit of ASI for the whole society is yet to be proven, as historically, our societies have not always distributed the benefits of technology fairly for the common good.

There are also concerns that AGI and ASI could pose a threat to humanity if they are not aligned with human values, if humans delegate too much power to them without sufficient controls, or if they are weaponized or misused for malicious purposes. When AGI or ASI systems do not align with human values, the so-called "Alignment Problem" arises, which could allow these systems to psychologically manipulate people—either unintentionally or on purpose—if malicious actors are in control.

Recent research in the field of "theory of mind" has shown that this risk cannot be eliminated with simple tests. In a way, the risks associated with AGI and ASI are similar to those in the civilian nuclear energy domain or cybersecurity. If all stakeholders act responsibly and with due diligence, nothing adverse should happen. However, humans are not perfect and can be the weakest link in security chains. Therefore, AGI and ASI cannot eliminate the risk of serious accidents, especially if humans grant AGI and ASI systems too much power or fail to control and manage them properly.

For instance, in 2018–2019, the Manoeuvring Characteristics Augmentation System (MCAS) of

Boeing contributed to the crashes of two Boeing 737 MAX aircraft due to a flawed design. Under certain flying conditions, the system made inappropriate decisions, causing the planes to crash without any possibility for the pilots to override the MCAS. The pilots lacked the technical means and training to recognize or handle the situations they faced.

In a way, the risks associated with AGI are like those in the civilian nuclear energy domain or cybersecurity. If all stakeholders act responsibly and with due diligence, nothing adverse should happen. However, humans are not perfect and can be the weakest link in security chains. Therefore, AGI and ASI cannot eliminate the risk of serious accidents, especially if humans grant AGI and ASI systems too much power or fail to control and manage them properly.

More recently, in March 2023, following the tests done prior to the release of GPT-4, OpenAI claimed that GPT-4 was able to trick a human helpdesk operator into bypassing basic security procedures by pretending to be a visually impaired person.

4.18 Artificial Super Intelligence (ASI)

Artificial superintelligence (ASI) refers to a type of AI model that significantly exceeds the cognitive and intellectual capabilities of the human brain.

The development of artificial general intelligence (AGI) is seen as a necessary precursor to achieving superintelligence, as it would require significant advances in machine learning, natural language processing, and other AI technologies.

The development of artificial superintelligence (ASI) is closely related to the development of supercomputers, especially quantum computing-based ones with a significantly higher number of qubits. Based on the state of research in 2023, this is

likely to happen with the creation of "silicon dot chips" if they can hold stable qubits at temperatures of about 1.5 K. It is unlikely that the first silicon dot quantum processors or other comparable technologies will be able to function at such low temperatures before 2030, as currently, qubits are only exploitable around 20 millidegrees K and research has not significantly progressed so far in this area.

A more likely development is the emergence of faster supercomputers [in the civilian sector based on photonic chips](#). All other things being equal, these photonic chips will be 1000 times faster than 2023 silicon-based computers and will serve as a stepping stone towards quantum computers. The first generation of ASI systems based on photonic technology could emerge by 2030, at the earliest, in the civilian sector.

4.19 Autonomous AI Agents

An autonomous AI agent refers to an AI system that can independently perform tasks and make decisions on behalf of another entity without constant supervision. These agents can act intelligently and adapt to their environment to achieve specific goals [9, 16]. Autonomous AI agents are found in various applications, such as chatbots, smart homes, and programmatic trading software [15].

Autonomous AI agents can also be integrated into embodied AI systems interacting with the real world, including robotic weapons to accomplish military objectives, for example. However, not all embodied AI applications include autonomous AI agents.

For instance, ChatGPT released by OpenAI is not an autonomous AI agent, as it requires an in-depth dialogue with the user to deliver the desired output or receive instructions.

Examples of an autonomous AI agent are [AutoGPT](#), [BabyAGI](#) and [SuperAGI](#) available mid-2023, which can perform actions automatically on a PC or a phone based on high-level instructions given by the user through a simple chatbot interface in plain, understandable language. For instance, AutoGPT can automatically create a website matching user instructions given in plain English and then also adjust all system settings, including security ones, accordingly. AutoGPT will only prompt the user in plain, understandable language if it does not comprehend the instructions or if the user's request is ambiguous or incomplete.

4.20 Session window of Large Language Models

Modern large language models (LLMs) like GPT-4 or ChatGPT-4 have a defined session window size, expressed in tokens (1 token approximately equals 0.7 words). As of mid-2023, this size ranges between 4,000 and 8,000 tokens, which roughly equates to 3,000 to 6,000 words. Beyond this limit, the LLM "forgets" the context of the dialogue with the user. Similarly, these LLMs are unable to process an input prompt larger than 3,000-6,000 words or generate an output exceeding this word count. Inputs exceeding 3,000-6,000 words must be segmented into several prompts, as must outputs exceeding this word count.

In 2023, OpenAI announced plans to significantly expand the session window size, from a maximum of 4,000 tokens to 32,000 tokens, representing an increase by a factor of 8.

In the same year, Anthropic surpassed this target by increasing the size of its session window to 100,000 tokens for its model Claude+, a considerable leap compared to OpenAI's goal.

Further pushing the boundaries, Mosaic.ML released in May 2023 the first commercially usable open-source models, one of which boasts a session

window of 65,000 tokens, a capacity exceeding OpenAI's offering.

These advancements mean it's now possible to upload entire documents of 100 to 200 pages into LLMs and query their content or request summaries. However, as the size of the session window increases, a trade-off emerges: larger session windows may lead to a decrease in output accuracy. Therefore, shorter session windows typically provide higher reasoning capabilities, while larger ones may be more prone to errors and hallucinations.

5.0 Essential facts about modern AI Models

5.1 Purpose of this section

This section delves deeper into the capabilities of modern AI systems released in 2020 and throughout 2023. The focus is on the central questions concerning the importance of AI model size, emergent abilities of large Language Models, Generative AI, the alignment problem, computing resource requirements, carbon footprint considerations, as well as past algorithmic breakthroughs and promising ones in the making. Grasping these concepts makes it easier to anticipate trends in the evolution of AI performance for the period 2023-2030 and the reasons why many experts predict that AGI will be reached by 2030 at the very latest.

- 1. Model Size:** The size of AI models has been steadily increasing, with some, like GPT-4, estimated to contain up to 1 trillion parameters. Larger models have demonstrated improved performance, but they also come with increased costs and energy requirements. Researchers are working on methods to train smaller models effectively, such as the AI model dilution, AI model compression, AI model sparsity etc...
- 2. Emergent Abilities:** Large Language Models have shown unexplained emergent abilities, such as

improved natural language understanding and creative problem-solving. These abilities arise as a by-product of training on vast datasets and imply that AI models are capable of learning and adapting in ways we don't yet fully comprehend.

3. **The Alignment Problem:** Ensuring that AI systems act in accordance with human values and intentions is a significant challenge. Researchers are working on solutions to the alignment problem to ensure that AI systems are safe and beneficial for humanity.
4. **Computing Resource and Carbon Footprint Considerations:** As AI models grow larger, they require more computing resources and produce a larger carbon footprint. Researchers are exploring ways to optimise training and deployment to minimise environmental impacts while maintaining or improving performance.
5. **Algorithmic Breakthroughs:** AI has experienced several breakthroughs in recent years, such as advances in reinforcement learning, unsupervised learning, and natural language processing, but also self-reflexion mechanisms in autonomous AI systems (to gradually improve their performance and better align with user goals).
6. **Multi Modal Generative Large Language Models** become the de facto interface and backbone for coordinating all complex AI systems and are leading the way towards AGI. Because of their superior reasoning capabilities acquired through natural language processing, LLMs are optimally placed to coordinate the actions of very complex AI systems and of autonomous AI agents.

5.2 AI model size matters but not only

The volume of data that an AI model is trained on, particularly a large language model, is one of the most important factors that determine its performance. The more data a model is trained on,

the more accurate it is likely to be in its predictions and the more robust it is likely to be when applied to new, unseen data.

One reason why training on large volumes of data is important for large language models is that language is incredibly complex, with many different patterns, nuances, and contexts to consider. By training on more data, the language model has a better chance of encountering a wide range of language use cases, which allows it to better understand and generalise the patterns it learns.

In recent years, there has been a growing emphasis on using techniques such as self-supervised or unsupervised learning. This approach involves training large language models without any explicit labelling or supervision. This allows the model to learn from vast amounts of unlabelled data. The model can then be fine-tuned on smaller, labelled datasets for specific applications.

Some AI experts, like OpenAI, have been at the forefront of creating these large, pre-trained language models. They argue that these models are necessary for natural language processing to reach the state of the-art. For example, OpenAI's GPT-3 model was pre-trained on a massive corpus of text, including web pages, books, and other sources, resulting in a highly versatile and powerful language model.

However, other AI experts have raised concerns about the potential downsides of relying on massive amounts of data for language modelling. They argue that over-reliance on large datasets can lead to models that perpetuate biases and stereotypes that exist in the data. They also point out that training large language models requires enormous amounts of computing power, which can have negative environmental impacts.

The availability of quality data for training is a potential bottleneck for the future development of

large language models such as GPT-4. Some experts have estimated that the volume is limited and that the current generation of the most powerful large language models already uses a large portion of the volume of quality data available, leaving a limited margin for further progress. Possible solutions include re-training systems several times on the same data or using "synthetic" data generated by simulation, possibly by another AI model.

In conclusion, the volume of data used to train AI systems is crucial. However, it is not just the amount of data that is important but also the quality of the data. Language models that are trained on diverse, high-quality datasets are more likely to be effective in real-world applications than models that are trained on narrow or biased datasets. It is equally important to consider the risks and problems associated with using large datasets.

5.3 Transferring data across borders can be a legal obstacle.

Collecting data for training AI models can [indeed be risky from a legal standpoint](#) if the data is transferred from another country.

The European Union's (EU) General Data Protection Regulation (GDPR) [requires that personal data transferred outside the EU receive the same level of privacy protection as if it were processed within the EU.](#)

Until recently, this was not the case, as US legislation allowed for the systematic collection of personal data from non-US citizens for national security reasons without any justification. This collection could occur without any corresponding decisions by a US judge to ensure that it was necessary and proportionate and without any redress mechanism for the non-US citizens affected ([see FISA, PPD 28](#)). In other words, the right to privacy in the US only concerns US citizens, with blatant discrimination against non-US citizens.

In March 2022, US President Joe Biden made amendments to PPD 28 that included additional safeguards, partially addressing the concerns previously outlined. However, the procedure in place for non-US citizens still differs from that for US citizens, who enjoy better privacy protection as the intervention of a US judge is required. As a result, the issue is only partially addressed from a purely legal perspective.

In contrast, the EU provides the same level of protection to all residents, regardless of their nationality. A judge must authorise the collection of personal data for national security purposes, demonstrating the EU's commitment to protecting the privacy and personal data of all individuals.

Various laws and regulations exist worldwide to protect personal data, such as the GDPR in the EU, the California Consumer Privacy Act (CCPA) in the United States, and China's Cybersecurity Law.

For instance, developers who want to transfer personal data from customers in regions such as the EU, California, or China to Japan must adhere to all the relevant laws and regulations to ensure that personal data is protected. Non-compliance with any of these regulations can result in significant fines and legal consequences. In the worst-case scenario, service could be cut off in the country in question.

The EU, US, and Japan have agreements that consider their privacy systems to be "essentially equivalent", indicating that they provide comparable levels of data protection. However, such agreements do not exist with China, leading to concerns about the privacy and security of applications such as TikTok.

It is possible that Japan, the EU, and the US may eventually impose a complete ban on TikTok due to those concerns. Similar issues exist in other

countries, such as Russia, but are most prominent in China.

In summary, countries that adhere to the rule of law and have laws providing a comparable level of privacy protection have a vested interest in agreeing that they are "fundamentally equivalent". This would facilitate the flow of data across borders and reduce legal uncertainty. In the context of large language models, as these systems are known to occasionally hallucinate, they can significantly harm people by spreading false information about them, hence the need for a proper form of consent regarding the use of personal data for training purposes.

Since there is no international organisation solely dedicated to data, such agreements must be negotiated bilaterally between the governments of the countries involved.

5.4 DeepMind advanced AI with its matrix multiplication algorithms

Most AI foundation models are created using neural networks, which consist of interconnected neurons that perform calculations and generate a final output from input. Matrix multiplication is a fundamental operation in this process, but it can be computationally intensive and energy-consuming, particularly during the training phase when the parameters are adjusted to achieve the desired output. To address this, researchers have been developing more efficient matrix multiplication algorithms.

In 2022, Google/DeepMind discovered a new and more efficient way to multiply matrices, which was made possible using GoogleTensor, an AI tool. This discovery was published in a research paper in the scientific journal Nature in October 2022. This new algorithm has significant implications for making large-scale deep learning models use less energy and reducing the carbon footprint of AI development. It has the potential to significantly

speed up the training process and improve the efficiency of computers.

Furthermore, this discovery highlights the potential of AI to drive innovation in scientific research, particularly in areas that have seen little progress in centuries. Matrix multiplication is a basic operation in many fields, such as machine learning and scientific computing, so this discovery is important and has wide-ranging effects.

Overall, Google's/DeepMind's discovery of a more efficient way to multiply matrices is a crucial step forward in the field of machine learning. It has the potential to improve the efficiency of computers and reduce energy consumption, leading to a more sustainable and environmentally friendly AI development.

5.5 DeepMind revolutionised AI with the transformer architecture

In 2023, some of the most powerful large language models, such as ChatGPT, PALM-E, DALL-E2, and GPT-4, all rely on the "Transformer neural network architecture." In the research paper ["Attention is All You Need" by Vaswani et al.](#) published in 2017, Google first described this architecture.

The Transformer model learns context and meaning by analysing relationships in sequential data, such as the words in a sentence. It achieves this through a self-attention mechanism that identifies and highlights important relationships between different parts of the input sequence. It introduced the concept of self-attention, which makes it easier for the model to learn complex relationships and dependencies in the input data, resulting in more accurate predictions. Furthermore, the attention mechanism allows the model to process longer sequences with fewer computations, reducing energy consumption. The Transformer model is also highly parallelizable,

allowing it to easily run on GPUs and other specialised AI chips. It represents a significant breakthrough.

With the release of ChatGPT, a direct competitor to Google Search, Microsoft and OpenAI were among the first to provide public access to a language model based on the Transformer architecture. This discovery of the Transformer machine learning model by Google marks a significant step towards creating large AI models that are both more powerful and more energy efficient.

5.6 Why large AI models are impenetrable black boxes.

When making predictions on new or untrained data, the numerical value of each "neuron" in the network is calculated based on the data presented at the input. This requires performing a series of matrix multiplications, where the elements of the matrix are the weights learned during the training phase.

The size and complexity of the neural network architecture, as well as the amount of data used for training ("tokens"), affect the accuracy of the predictions made by the model.

Finding the optimal balance between the size of the training data set (tokens) and the size of the neural network (number of parameters) is an active area of research in 2023, as it affects both the accuracy of the model and the computing power needed during training.

Because it is challenging for developers to comprehend and explain in simple terms the inner workings of complex neural networks, their creators frequently refer to them as "black boxes". This lack of transparency poses a challenge for programmers, who struggle to explain how the

network arrives at specific predictions. Often, statistical methods are used to calculate probable solutions, but the resulting probabilities and mathematical formulas are very challenging to interpret and explain in plain language.

5.7 Berkley Retrieval Augmented LLMs help with “explainability”

Large language models (LLMs) can be challenging to interpret in ways that humans can easily comprehend, even though their predictions can be described mathematically. This lack of transparency and interpretability is due to the way LLMs, such as GPT-4, store learned patterns and representations within the model's parameters, making it difficult to trace back the specific training data that influenced a particular prediction. Consequently, assessing the LLM's quality, explaining its decisions, and identifying and correcting biases or errors can be difficult, potentially rendering AI systems untrustworthy.

In January 2023, UC Berkeley proposed the concept of retrieval-augmented Large Language Models (RALMs), which store information in a separate knowledge source or ground corpus that can be queried and retrieved instead of within the model's parameters. This approach shows promise for developing more robust and transparent AI systems that could potentially be more accurate, reliable, explainable, trustworthy, and even learn faster.

The separation between the model parameters and the knowledge source allows RALMs to generate more contextually appropriate and informed responses by retrieving relevant information from the corpus, leading to improvements in the model's performance and the possibility of generating explanations or critiques. The prediction process in a RALM involves two steps: retrieving trained data by querying the knowledge source and then making a prediction.

However, the ability of RALMs to identify their own biases and weaknesses or provide critiques of their own predictions depends not only on the separation of model parameters and the knowledge source but also on the quality of the training data, the architecture of the model, and the specific training objectives set during the training process.

While RALMs offer significant advantages over traditional LLMs in certain contexts, their effectiveness may depend on the specific application and context in which they are used. Some reasons why RALMs may not be suitable or desirable for systematic use in some applications compared to traditional LLMs include:

1. **Computational complexity:** RALMs often require additional computational resources to retrieve relevant information from the knowledge source during the generation process, resulting in slower response times and increased costs, especially for large-scale applications.
2. **Quality of the knowledge source:** The effectiveness of RALMs depends on the quality and comprehensiveness of the knowledge source or corpus. If the corpus is not well-maintained, outdated, or lacks relevant information, the performance of the RALM may be negatively affected.
3. **Implementation complexity:** Developing and maintaining a separate knowledge source or corpus for an RALM can be more challenging than training a traditional LLM. It requires additional effort to manage the storage, updating, and retrieval of information from the knowledge source.
4. **Domain-specific expertise:** RALMs may require more domain-specific expertise to curate and maintain the knowledge source, especially when dealing with specialized or technical domains.
5. **Adaptability:** RALMs rely on retrieving relevant information from a separate knowledge source, which may not be ideal for tasks that require

high adaptability or creativity, as the model might be limited to the knowledge present in the corpus.

6. **Privacy concerns:** If the knowledge source contains sensitive or private information, using RALMs could raise privacy concerns, as the model may inadvertently expose or reveal sensitive data during the generation process.

Overall, retrieval-augmented large language models (RALMs) show promise for making AI systems that are more accurate, reliable, explainable, trustworthy, capable of making their own critical predictions, and even learning faster. The approach involves storing the learned knowledge in a separate data corpus. First, the data corpus is queried, and then the prediction is made. However, RALMs are not suitable for every application due to their increased complexity, higher computational resource requirements, domain-specific expertise requirements, adaptability constraints, and privacy concerns.

5.8 The high carbon footprint of Large Language Models

Large language models, multimodal large language models, and generative AI models such as GPT-3 and GPT-4 require a significant number of matrix multiplications to adjust the model's parameters during pre-training. This phase is computationally intensive and requires specialised hardware such as AI accelerator chips designed by companies like Nvidia, Tesla, Google, IBM, Cerebras, and Huawei. These chips power supercomputers or are integrated into products and services. Companies such as Amazon, Tencent, Alibaba, Microsoft, Baidu, Google, and more recently, Nvidia, offer commercial cloud computing services for training large AI models. However, the development of large AI models raises concerns about their environmental impact, as they require massive amounts of computing resources and electrical energy, leading to a significant carbon

footprint. Studies have found that training a large language model emits a substantial amount of CO₂, equivalent to the lifetime emissions of several cars.

During the inference phase, which is the process of using a trained neural network to make predictions or decisions based on new input data, fewer matrix multiplications are required because the network's parameters are already optimised and fixed. Inference can run on various hardware, including CPUs, GPUs, and mobile devices. Some companies are designing AI chips to accelerate inference, while others are developing ultra-low power Neuromorphic chips running spiking neural networks that are ideal for IoT and mobile applications.

For large AI generative models like ChatGPT-4, which have a large size and high number of concurrent users, computing requirements during the inference phase in production can also be very high. However, for less demanding applications, inference requires much less computational power than training and can run on various hardware, including CPUs and mobile devices at the edge.

The inference phase is crucial for real-time decision-making and is used in a variety of applications, such as natural language processing, image recognition, and speech recognition. One major challenge in the inference phase is the need to balance accuracy with speed and energy efficiency. To address this challenge, hardware acceleration techniques are being developed to speed up the inference process and reduce energy consumption.

In summary, the inference phase is crucial for real-time decision-making and the deployment of AI models. It requires less computational power than the training phase and can run on various hardware, including mobile devices. However, for large AI generative models with a high number of concurrent users, the computing requirements during the

inference phase in production can be very high. To address this challenge, hardware and software acceleration techniques and energy-saving techniques are being developed to improve the speed and energy efficiency of the inference process.

5.9 Why Large AI models develop unexplained emergent abilities.

In general, Large AI models, including Large Language Models (LLMs), Multimodal Large Language Models (MLLMs), and Generative AI models with over 100 billion parameters, are trained using a combination of self-supervised, unsupervised, or semi-supervised learning methods. These models don't rely solely on labelled data during the pre-training phase. Instead, the AI model organises the data independently (unsupervised learning) or performs an initial pretext task (self-supervised learning). Labelled data might be used in the fine-tuning phase to make the model more task specific.

During training, these models acquire complex and sophisticated categorization systems that are difficult for humans to comprehend, given the size of the training datasets. This results in the system acquiring abilities for which it was not explicitly designed. These emergent properties arise from the system's ability to identify patterns and relationships in the training data that are imperceptible to humans. These capabilities result in the capacity to reason, translate text, and more. It is crucial to recognize these capabilities and take the necessary steps during development to ensure that the AI system remains accurate, reliable, transparent, and aligned with human values.

To address these concerns, methods such as explainable AI (XAI) are being developed to provide greater transparency and interpretability in AI systems, enabling humans to better understand the reasoning behind the model's predictions. This can improve trust and reliability, particularly in

sensitive fields such as healthcare, finance, and law enforcement. Additionally, AI governance frameworks are being developed to ensure that AI systems are developed and used ethically, responsibly, and with consideration for their impact on society and the environment.

5.10 The “Intrinsic” emergent capabilities of LLMs

They are the abilities that are “intrinsic” to the model's general understanding of language, its ability to process and generate text, and stem from the vast knowledge base the models acquire during their training. The most impressive ones are the capabilities of common-sense reasoning, conceptual understanding, sentiment analysis, and knowledge retrieval:

1. **Grammar and syntax understanding:** Large language models develop a deep understanding of language structure, syntax, and semantics, enabling them to generate grammatically correct and contextually appropriate responses in natural language conversations.
2. **Common-sense reasoning:** These models can exhibit common-sense reasoning abilities, allowing them to make plausible inferences, predictions, and judgments based on the information available in the input data or prompt.
3. **Knowledge retrieval:** Large Language models can answer factual questions based on their training data, demonstrating their ability to recall and synthesise information from their vast knowledge base.
4. **Conceptual understanding:** Large Language Models can develop an understanding of abstract concepts and relationships, allowing them to generate explanations, examples, or comparisons that demonstrate their grasp of the underlying ideas.
5. **Entity recognition:** Large Language Models can identify and classify entities, such as names, locations, and dates, within a given text,

enabling them to provide more context-aware and accurate responses.

6. **Sentiment analysis:** Large language models can identify and analyse the sentiment (positive, negative, or neutral) expressed in a piece of text, allowing them to respond or generate content based on the emotional context.

5.11 The “In context” emergent abilities of LLMs

The “In-context” capabilities of Large Language Models (LLMs) arise from their ability to understand and generalise from the context provided within user input prompts. Users can request new tasks, and the system understands and executes them without requiring any preliminary fine-tuning or additional training.

Some of the most impressive capabilities of LLMs include engaging in meaningful discussions with users (“Chain of Thoughts Prompting”), learning a new task from just a few instructions provided by the user (“Few-Shot Learning”), executing instructions directly specified by the user (“Zero-Shot Learning”), as well as translating text and writing code.

These capabilities are particularly noteworthy because the LLM was not initially pre-trained or fine-tuned for these specific purposes. The model can generalise from its pre-training data to perform a wide range of tasks without requiring explicit supervision. This flexibility and adaptability make LLMs powerful tools in various applications, such as natural language processing, dialogue systems, and content generation.

However, it is crucial to recognize that these capabilities also raise concerns about the potential risks and ethical implications of deploying such models in real-world applications. The ability of LLMs to generate human-like language and perform complex tasks raises concerns about the potential for misuse, such as the creation of fake news or malicious content. It is essential to

develop governance frameworks and ethical guidelines to ensure that LLMs are developed and deployed in a responsible and ethical manner. Additionally, methods such as explainable AI (XAI) are being developed to provide greater transparency and interpretability in LLMs, enabling humans to better understand the reasoning behind the model's predictions.

1.Chain of Thoughts prompting:

This refers to the ability of a large language model to maintain context and generate coherent responses across multiple input prompts or turns in a conversation. The model learns to keep track of the conversation's context and use it to generate relevant responses.

2.Zero-shot learning:

In this setting, a large language model can perform well on a task without being explicitly trained on that task. It leverages the patterns and representations learned during its training to make inferences and predictions for the new task, often by providing an instruction in the input prompt that describes the desired output.

3.Few-shot learning:

In few-shot learning, a large language model can adapt to a new task by observing a small number of examples provided within the input prompt. These examples help the model to infer the pattern or transformation needed to solve the task at hand, and it can then apply this knowledge to perform the task with minimal information.

4.Task adaptation:

Large language models can quickly adapt to a variety of tasks like sentiment analysis, text summarization, question-answering, and more, based on the instructions provided in the input prompt.

5.Style transfer:

By providing context or examples of a particular writing style within the input prompt, large language models can generate text

in that specific style, such as imitating a famous author or writing in a formal tone.

6. Code completion:

Large language models can complete code snippets in various programming languages when given a prompt that provides context or examples of the desired code.

7. Multi-step reasoning:

Large language models can perform multi-step reasoning tasks by understanding the context of a problem and processing several steps to arrive at a conclusion or solution.

8. Creativity and content generation:

When given a prompt that provides context or constraints, large language models can generate creative content like stories, poems, or even dialogue for fictional characters.

9. Multilingual understanding and translation:

Large language models can understand and process input prompts in multiple languages, enabling them to perform tasks like translation, paraphrasing, and multilingual question-answering based on the provided context.

10. Interpreting and generating analogies:

By understanding the context and relationships in the input prompt, large language models can interpret and generate analogies to explain complex concepts or ideas.

5.12 Multimodal LLMs will gradually lead to AGI before 2030

Multimodal Large Language Models (MLLMs), such as GPT-4, Microsoft Kosmos-1, and Google PALM-E, integrate various types of input modalities, including text, audio, images, video, and sensor data. This integration allows MLLMs to better understand complex and diverse contexts than they would otherwise using only text, resulting in a more human-like model of the world and a broader range of cognitive tasks compared to traditional Large Language Models (LLMs) like GPT-3.5.

One notable ability of MLLMs is their emergent capacity to transfer knowledge from one modality to

another. For example, these models can answer questions about an image. This ability can be described as "multimodal reasoning" or "cross-modal learning," which refers to the capability of an MLLM to build representations and transfer knowledge across different modalities. All MLLMs possess this ability to some extent by nature. However, Microsoft Kosmos-1 and, especially, Google PALM-E have been designed with this as a primary objective. It is likely that the current generation of MLLMs will gradually improve towards AGI before 2030.

In April 2023, Microsoft published the paper ["Sparks of Artificial General Intelligence: Early experiments with GPT-4"](#) . This paper demonstrated unambiguously that the reasoning capabilities of GPT-4 released in March 2023 are far superior to GPT-3.5 and ChatGPT-3.5 released in November 2022 and clearly rival, if not exceed, those of the human brain in many standard tests.

In March 2023, Zhejiang University and Microsoft Research Asia published the paper ["HuggingGPT: Solving AI Tasks with ChatGPT and its Friends in Hugging Face"](#) . The paper discusses the limitations of current AI models in handling complicated AI tasks across different domains and modalities. The paper focuses on addressing the challenges in artificial intelligence when it comes to handling complex tasks that involve various areas and types of information. It emphasises the potential of large language models, which have shown impressive capabilities in understanding and working with language.

Hugging Face is an open-source repository of Machine Learning Models. HuggingGPT is a Large Language Model developed by Zhejiang University in China and Microsoft Research Asia that identifies which machine learning algorithms in the HuggingFace repository should be put to work for a given user task. Hugging GPT coordinates the execution of tasks among the different machine

learning algorithms and returns the result to the user.

5.13 Humans in Loop and “AI in the Loop” are both relevant.

There are different ways in which AI models can learn, be trained, and be supervised, including supervised, unsupervised, self-supervised, and reinforcement learning. In each of these cases, human intervention can be crucial to ensuring the accuracy, reliability, and safety of the AI system.

"Human in the loop" is a design strategy in which people are involved in the operation and supervision of an artificial intelligence system. This approach recognizes the limitations of AI models and aims to leverage human expertise to enhance the performance of these systems.

The advantages of the "Human in the Loop" approach are numerous. First, it enables the AI system to leverage human intelligence and intuition, which are still unmatched by machine learning models. Second, it ensures that the AI system is transparent, accountable, and trustworthy, as human experts can review and validate the results. Third, it reduces the risk of unintended consequences, errors, or bias in the AI system, which could have serious consequences in critical applications such as healthcare, finance, or national security.

However, the "human in the loop" approach also has some drawbacks. First, it can be costly and time-consuming, as it requires human experts to be involved in the development, training, and operation of the AI system. Second, it can limit the scalability and automation of the AI system, as the human input may be a bottleneck in some cases. Third, it can introduce a new source of error and bias, as human experts may have their own limitations and biases that could affect the performance of the AI system.

In contrast, the "AI in the Loop" approach aims to minimise human intervention and rely mostly on automated machine learning algorithms operating on a set of predefined high-level rules. This approach can be more efficient, scalable, and cost-effective in some applications, but it also carries the risk of unintended consequences, errors, and bias if not properly validated and monitored by human experts.

Overall, the choice between "human in the loop" and "AI in the loop" depends on the specific requirements and constraints of each application, and a balanced approach that leverages both human and machine intelligence is often the best solution.

5.14 Ethical, Responsible and Trustworthy AI overlap but differ.

The concepts of Ethical AI, Responsible AI, and Trustworthy AI are all related but distinct.

- **Ethical AI** is concerned with the moral implications and values of AI development and use. It involves identifying and addressing ethical dilemmas that arise from AI's potential impact on society and individuals [[1](#)].
- **Responsible AI** is a more tactical approach that focuses on developing and deploying AI systems in a safe, trustworthy, and ethical way. It is about balancing effectiveness with ethical implications and ensuring that AI is developed thoughtfully [[4](#)].
- **Trustworthy AI** is a subset of responsible AI, which focuses on building AI that is reliable, transparent, fair, inclusive, and respects privacy and security. Trustworthy AI ensures that AI systems are designed to be used in ways that benefit all members of society. [[5](#)].

An example illustrating the similarities and differences between the concepts is a healthcare AI system. Ethical AI would consider the potential

ethical dilemmas arising from the deployment of an AI system in healthcare, such as ensuring that AI is not used to perpetuate discrimination. Responsible AI would involve developing the system in a safe, trustworthy, and ethical way, ensuring that the AI system is transparent and explainable, and that the data used to train the AI model is representative and unbiased. Trustworthy AI would involve ensuring that the system is reliable, secure, respects patient privacy and confidentiality, and is used in ways that benefit all members of society.

To summarise, ethical AI includes responsible AI, which in turn includes trustworthy AI. Ethical AI takes a broad perspective, looking at AI's potential implications and values, while responsible AI focuses on developing and deploying AI in a way that is safe, trustworthy, and ethical. Trustworthy AI is a subset of responsible AI that is specifically concerned with ensuring that AI systems are reliable, transparent, fair, inclusive, respect privacy, and are secure.

5.15 “Ethical AI” and the “alignment problem”

The AI alignment problem was articulated the first time in 1960 by AI pioneer [Norbert Wiener](#) .

The Alignment Problem is a concept part of Ethical AI [1] and is considered a crucial aspect of AI ethics. [2]. The problem refers to the challenge of ensuring that artificial intelligence (AI) systems behave in ways that are aligned with human values and goals, which is a complex and multifaceted issue that has garnered significant attention from researchers and policymakers in recent years [2].

In September 2022, the United Nations System Chief Executives Board for Coordination endorsed the Principles for the Ethical Use of Artificial

Intelligence in the United Nations System [2]. The principles are meant to guide the use of AI within the United Nations system, and they aim to ensure that the use of AI is aligned with the UN's values and principles, particularly with respect to human rights, transparency, accountability, and inclusivity.[2], [7]]. The principles are based on a human-centred approach to AI, which seeks to ensure that AI is developed and used in a way that is beneficial to people and society as a whole [5].

In November 2021, the 193 Member States at UNESCO's General Conference adopted the Recommendation on the Ethics of Artificial Intelligence, the very first global standard-setting instrument on the subject [[3], [4]]. The recommendation sets the first global normative framework while giving States the responsibility to apply it at their level. UNESCO will support its implementation by providing capacity-building and awareness-raising activities [4].

The alignment problem has also a technical aspect that focuses on how to encode values and principles into AI so that it does what it ought to do in a reliable manner. The company Anthropic.AI has proposed the concept of "Constitutional AI" as part of their efforts to solve the technical part of the "alignment problem [1][2][3]. According to the company, the "Constitutional AI" approach aims to provide a principle-based approach to aligning AI systems with human intentions [4]. The only human oversight is provided through a list of rules or principles, forming a sort of constitution that guides the AI system [6][8]. A major benefit of this approach is that it is also very scalable since no or little human intervention is required.

Overall, the AI alignment problem is a critical component of Ethical AI as it addresses the ethical implications of the use of AI systems in society. The Ethical part is addressed by the UN. The works in 2022 on “Constitutional AI” by Anthropic and others addresses the technical part of the Alignment problem and is also an important step towards making AI systems safe and aligned with human values [5].

5.16 Leveraging the “flywheel effect” in AI systems.

The Flywheel effect is a concept where small wins accumulate over time, creating momentum that keeps a business or system growing [9]. In AI, the Flywheel effect means that [the model's performance keeps getting better as more users interact with it](#) and it is trained repeatedly.

For instance, as users interact with ChatGPT and provide feedback, the model continues to learn and adapt, thus improving its performance. These improvements, in turn, attract more users, leading to more interactions and more feedback, creating a self-reinforcing loop [8]. Some practical applications of ChatGPT using the Flywheel effect include parsing chats for product-specific question and answer sets and training the model on this information, alongside other resources like FAQs and help pages [2]. This continuous process of refining and updating the model based on user interactions enables ChatGPT to deliver better performance over time. Recently this process has however caught the attention of privacy watchdogs in Italy and in other countries that claim that users should consent to their data being used by OpenAI and ChatGPT for this purpose. The controversy led to the temporary ban of ChatGPT in Italy.

The Flywheel effect can be understood in the context of autonomous driving as a continual cycle of data collection, model improvement, and increased trust and acceptance. As more self-driving vehicles hit the road, they collect massive volumes of data from diverse driving scenarios, road conditions, and environments. This data is then utilised to train and refine the AI models that operate the vehicles, thereby boosting their performance and decision-making capabilities.

As the AI models become more sophisticated and accurate, the autonomous driving systems become more reliable and safer, which in turn leads to increased public trust and adoption of the technology. This increased adoption results in more vehicles on the road, collecting even more data to further enhance the AI models' performance. This self-reinforcing cycle of data collection, model improvement, and increased adoption exemplifies the Flywheel effect in the autonomous driving domain.

This continuous improvement process enables autonomous driving systems to better adapt to complex situations, learn from real-world experiences, and ultimately make roads safer for all users.

5.17 Large AI models can be compressed and optimised.

AI model compression and optimization refers to the process of reducing the size and complexity of neural networks without significantly compromising their accuracy. This allows for the deployment of state-of-the-art deep learning models on edge devices with limited computing power and memory resources [11]. There are several techniques used for AI model compression, including AI model distillation, AI model pruning, AI model quantization, and architecture search.

Using those model compression and optimization techniques, it is possible to obtain smaller and

more efficient AI models for narrower applications without having to develop them from scratch. By leveraging existing large models and compressing them, developers can save time and resources while still achieving satisfactory performance for their specific use cases. However, it is important to note that in some cases, developing a model from scratch specifically tailored to the narrow application might yield better performance and efficiency, depending on the problem and the available data[5].

In summary, AI model compression and optimization can serve as an alternative to developing AI models from scratch for narrow AI applications, providing a faster and more resource-efficient way to deploy AI models on devices with limited resources. However, the choice between using model compression and developing a model from scratch depends on the specific problem, requirements, and available resources.

- **AI Model Distillation:** Model distillation is a technique where a smaller and simpler student model is trained to mimic the behaviour of a larger and more complex teacher model. This is achieved by incorporating distillation loss, expressed with respect to the teacher, into the training of the student network whose weights are quantized to a limited set of levels [13]. The student model learns to approximate the teacher model's output, resulting in a compressed model with fewer parameters while maintaining a similar level of performance.
- **AI Model Pruning:** Pruning is a popular model compression technique that works by removing redundant and inconsequential parameters in a neural network, such as connections or weights [4]. Pruning reduces the number of parameters by eliminating unimportant connections that do not significantly impact performance. This helps reduce the overall model size and saves on computation time and energy [7]. As the neural

network becomes sparser it is possible to use special algorithms to accelerate its performance.

- **AI Model Quantization:** Quantization involves reducing the numerical precision of the model parameters or weights. Typically, weights are stored as 32-bit floating-point numbers, but for many applications, this level of precision may not be necessary [3]. Quantization maps values from a large set to values in a smaller set, resulting in a smaller range of possible values for the output [6]. This process helps reduce the model size and memory requirements while maintaining acceptable performance levels.
- **Architecture search:** This involves automatically searching for and selecting an optimal neural network architecture that meets the desired trade-offs between accuracy, size, and computational complexity. The process typically explores various architectures and configurations, either by starting from scratch or by modifying existing models, to find the one that best suits the problem and the available computational resources.

Examples of compressed Large Languages Models featuring similar performance to GPT-3.5 in narrow application fields:

1. **Alpaca** (Stanford) is an instruction-following language model developed by researchers at Stanford University. It was fine-tuned from Meta's LLaMA 7B model, which has seven billion parameters, and trained on 52,000 instruction-following demonstrations generated using GPT-3.5 [1]. Alpaca model is an open-source alternative to more expensive and proprietary large language models, such as ChatGPT [9]. Stanford's Alpaca was developed for a cost of less than \$600, which is a fraction of the training cost of 5-12 million dollars that was needed to train GPT-3.5 [3][7]. It runs on a powerful PC.

2. [Dolly](#) (Berkeley) is an open-source large language model (LLM) released by Databricks, a company that originated from the AMPLab project at the University of California, Berkeley [9]. Dolly was fine-tuned using freely-available software components and is designed to be a faster and more economical way to build services similar to ChatGPT [5]. Instead of creating a new model from scratch or using LLaMA, like Stanford did, Databricks utilised an older open-source LLM called GPT-J, which was created by EleutherAI several years earlier [7]. Dolly demonstrates that instruction-following capabilities don't necessarily require the latest or largest LLMs and that it's possible to take a dated off-the-shelf open-source large language model and give it ChatGPT-like qualities [6].
3. [Vicuna](#) (Berkeley, Stanford, CMU, UC-San Diego) released in May 2023 is a more advanced version of the two models above with up to 13 billion parameters.

5.18 Autonomous AI Agents can self-reflect to improve.

In March 2023, the North-eastern University (Boston, MA) and the Massachusetts Institute of Technology Cambridge (Boston) MA published the paper "[Reflexion: an autonomous agent with dynamic memory and self-reflection](#)".

An autonomous AI agent is an artificial intelligence-based entity that can perceive its environment, make decisions, and take actions to achieve specific goals without direct human control or intervention. It can also learn and adapt to its environment, enabling it to perform tasks with minimal human involvement [2, 4, 15, 19, 20]. An autonomous AI agent senses its surroundings through sensors, processes the information, and acts upon its environment through actuators or effectors. These agents are capable of working towards their goals and interacting with their environment and

other systems without immediate help from humans [8].

Examples of the first experimental autonomous AI agents released in March 2023 are "[AutoGPT](#)" and "[BabyAGI](#)". Though they can be useful in some use cases, their performance has yet to be further improved to bring true added value to the masses. For instance, based on a single user prompt or series of prompts given at the beginning of the process, these autonomous agents can take all necessary actions to create a website. The more accurate the prompts are, the more focused the result will be. Along the way, the agent frees the user from having to master any of the technology needed to build sophisticated websites, and it queries the user about aspects of his or her request that may be insufficiently clear to help him or her refine his or her thoughts on the desired outcome.

In the realm of AI, self-reflection refers to the process of augmenting an AI system's performance by scrutinising its output and conducting additional tuning, code corrections, or retraining. The research paper cited above illustrates the feasibility of designing autonomous AI systems capable of analysing their outputs with minimal or no human intervention. Remarkably, these systems have the potential to autonomously retune themselves or even revise their own code! The paper suggests that such self-reflective AI systems outperform the state-of-the-art GPT-4! The benefits of self-reflection include:

1. **Improved performance on tasks**, especially for code generation and problem-solving.
2. **Enhanced decision-making and reasoning**, users could expect more reliable and coherent decision-making.
3. **Recursive self-debugging and self-improvement**. This would result in a continuously evolving and improving AI system, potentially reducing

the need for constant updates and manual intervention.

4. **More human-like behaviour and understanding.**

As self-reflection is a key aspect of human cognition [16], an AI system with this capability would be better equipped to understand and respond to user needs, leading to more natural interactions and improved user satisfaction.

Overall, the incorporation of self-reflection into autonomous AI agents would lead to a more efficient, reliable, and human-like AI system, benefiting users in various applications and industries.

5.19 Open source versus Closed Source Large Language Models

In the realm of AI, open source refers to AI models characterised by publicly accessible source code, design, and training datasets. These resources can be viewed, modified, and distributed by anyone, encouraging several key attributes:

1. **Accessibility:** Open-source models provide developers the liberty to access and modify the base code, encouraging innovation and customization. On the other hand, closed source models restrict code access, limiting the potential for alterations and understanding of their internal workings.
2. **Transparency:** By offering full transparency in their development process, open-source models allow community scrutiny and validation. In contrast, closed source models frequently lack transparency as they don't divulge their specific implementation details, complicating the task of detecting biases or potential issues.
3. **Collaboration:** Open-source models foster collaboration among developers, researchers, and the community, with multiple contributors able to refine models, fix bugs, and expand

functionality. In contrast, closed source models generally rely on a select group of developers within a specific organisation.

4. **Size:** Open-source models, developed by organisations like Meta, Mosaic.ML, Stability.AI, or institutions such as Stanford, UC-Berkeley, or UC-San Diego, are typically much smaller than their closed counterparts developed by companies like OpenAI or Google.
5. **Adaptability:** Due to their smaller size, modifiability, and fine-tuning capabilities, open-source models can be efficiently adapted for specific applications. Conversely, proprietary models like GPT-4 or Palm-2 necessitate substantial resources and effort to optimise for specialized domains.
6. **Cost Efficiency:** As they are smaller and require fewer computational resources for operation (inference), open-source models are cheaper to train, making them good options for companies or individuals who prefer not to rely on tech giants like Google or OpenAI. The affordability of these models also permits companies to deploy multiple models concurrently, each tailored to unique needs.
7. **Compression and Distillation:** As exemplified by Stanford in March 2023 with their DOLL-E Chatbot, Large Language Models, like GPT-3.5, can undergo compression and distillation processes to be transformed into smaller open-source models. This significantly reduces development costs, making AI more accessible and cost-effective.
8. **Performance Parity:** Open-source models have demonstrated their ability to perform competently within narrow domains, often rivalling the performance of larger, more generalised models like GPT-4 or Palm-2. The combination of cost-effectiveness and performance parity makes these models a compelling choice for developers seeking bespoke solutions.

Meta, formerly known as Facebook, made a significant impact on the AI community when it unintentionally released some of its LLaMa large language models in March 2023. These models were originally created exclusively for scientific research. This unexpected release piqued the sector's interest and sparked widespread excitement. In April 2023, UC Berkeley, UC San Diego, CMU, and Stanford released Vicuna, a powerful chatbot competing with ChatGPT-3.5. In early May 2023, Mosaic.ML emerged as the first company to formally introduce open-source large language models for commercial applications. Not long after, Meta also opted to officially release its models as open source.

Developers of large language models, such as OpenAI's GPT-4 and Google's Palm-2, are now experiencing some competition from these open-source models. Future trends suggest that companies producing proprietary, closed-source Large Language AI models might venture into commercialising versions tailored for specific, low-power hardware or target computer infrastructures. Techniques like "pruning" and "sparse deep neural networks" will be instrumental in achieving this. For instance, "ThirdAI," based in Houston, has successfully implemented powerful "recommendation systems" running exclusively on CPUs using such techniques. Consequently, these closed-source companies will compete head-to-head with those leveraging open-source models.

As advancements in machine learning hardware and software drive increased energy efficiency and performance, we can anticipate a surge in the size of open-source large language models. However, the most comprehensive models – those operating on the most potent computing resources, such as those intended for future AGI systems – will likely remain proprietary. This is attributed to their exceptional performance and the significant investment required for their development.

Part III: AI Governance and Regulations

Global governance organisations setting guidelines on Ethical and Responsible AI

Non-profit organisation works on AI

National Frameworks on AI

6.0 Global governance on responsible and ethical AI

To promote the development of trustworthy and responsible AI, guidelines have been established to assist governments worldwide in creating a coordinated and coherent legislative framework. The three most active international guidelines are those developed by the United Nations, the Organisation for Economic Co-operation and Development (OECD), and the Council of Europe in Strasbourg. Both Western democracies and "non-aligned countries" from the Global South have adopted these guidelines. However, they have not been signed by certain countries, including Russia, China, North Korea, Venezuela, and Iran.

6.1 The United Nations and the works of UNESCO

[The United Nations Global Pact for Responsible Artificial Intelligence, adopted in November 2021, is a historic agreement that defines the common values and principles needed to ensure the healthy and ethical development of Artificial Intelligence \(AI\).](#) All member states of the UN Educational, Scientific and Cultural Organization (UNESCO) adopted this agreement.

UNESCO acknowledged that AI is bringing unprecedented challenges, including increased gender and ethnic bias, significant threats to privacy, dignity, and agency, dangers of mass surveillance, and increased use of unreliable AI technologies in law enforcement. The Global Pact aims to guide the construction of the necessary legal infrastructure to ensure the ethical development of AI. UNESCO supports its 193 Member states in its implementation and asks them to report regularly on their progress and practices.

The text of the pact highlights the advantages of AI while also addressing the risks it entails. It provides a guide to ensure that digital transformations promote human rights and contribute to the achievement of the Sustainable Development Goals. This includes addressing issues around transparency, accountability, and privacy with action-oriented policy chapters on data governance, education, culture, labour, healthcare, and the economy. It makes a clear call to protect data to guarantee individuals more protection by ensuring transparency, agency, and control over their personal data.

The pact explicitly bans the use of AI systems for social scoring and mass surveillance. It also emphasises that AI actors should favour data, energy, and resource-efficient methods that will help ensure that AI becomes a more prominent tool in the fight against climate change and in tackling environmental issues.

[Note that Russia, China and Iran signatories to the United Nations Global Pact for Responsible Artificial Intelligence adopted by UNESCO.](#)

Compliance to the pact is not compulsory and there is no enforcement and sanction mechanisms

6.2 The Organization for Economic Cooperation and Development (OECD).

The Organisation for Economic Co-operation and Development (OECD) adopted the [OECD Principles on Artificial Intelligence in 2019](#), which offer a framework for the responsible development and deployment of AI.

The OECD is currently collaborating with more than 100 countries, regions, and international organizations around the world, including China, Russia, India, Saudi Arabia, South Africa, Morocco, Indonesia, the United Nations, and the European Union, to promote responsible and trustworthy AI.

As of 2022, the 37 OECD member countries, as well as Argentina, Brazil, Egypt, Colombia, Costa Rica, Malta, Peru, Ukraine, and Singapore, have signed the OECD Principles on Artificial Intelligence. However, Cyprus, Bulgaria, and Croatia had not yet signed the principles.

As of 2023, North Korea, Venezuela, and Iran have not signed the principles and are not directly involved in any of the OECD's efforts to promote responsible and trustworthy AI.

6.3 The Council of Europe Committee on Artificial Intelligence (CAI):

[The Council of Europe Committee on Artificial Intelligence \(CAI\)](#) is a committee established by the Council of Europe to develop a [legally binding instrument on the development, design, and application of AI systems based on the Council of Europe's standards on human rights](#), democracy, and the rule of law. The CAI aims to create a common global approach to basic principles that should govern AI development and use, considering the technology is developed and used across borders. By the end of 2023, the instrument will be in place.

The CAI's goal is to ensure that AI systems do not endanger or undermine democratic processes, either directly or indirectly. In contexts where AI systems assist or replace human decision-making, the committee's focus is on ensuring the continued application of human rights and the principle of the rule of law.

The committee is tasked with creating a global instrument that is attractive to as many states as possible around the world. This is done with the belief that the more global the instrument becomes, the more impact it will have on people's lives worldwide.

The CAI includes representatives from its 46 member states, observer states (Canada, Holy See, Israel, Japan, Mexico, United States of America), other Council of Europe bodies and sectors, other international and regional organisations working in the field of AI, representatives of the private sector, and representatives of civil society, research, and academic institutions. However, the CAI acknowledges that the instrument they are currently working on cannot regulate all aspects of the development and use of AI systems, and that additional binding and non-binding instruments will be needed to comprehensively address the use of this rapidly evolving technology.

[The Council of Europe Committee on Artificial Intelligence \(CAI\)](#) is a specialised body of the Council of Europe tasked with the responsibility of addressing the impact and challenges of artificial intelligence (AI). The Council of Europe is an international organisation that aims to uphold human rights, democracy, and the rule of law in Europe, and its CAI extends these goals into the realm of AI.

As of March 2021, the AI Principles have been signed by Austria, Belgium, Bulgaria, Croatia, Cyprus, the Czech Republic, Denmark, Estonia, Finland, France, Germany, Greece, Hungary, Iceland, Ireland, Italy, Latvia, Lithuania, Luxembourg, Malta, the Netherlands, Norway, Poland, Portugal, Romania, Slovakia, Slovenia, Spain, Sweden, Switzerland, and the United Kingdom. These countries have committed to upholding the principles and promoting responsible and trustworthy AI development and use. Russia participation rights were suspended in 2014, China, Iran, Venezuela, North Korea are not part of it. None of the democratic countries in Asia are members either.

It's important to note that the Council of Europe is not an institution of the European Union and should not be confused with the European Council in

Brussels, which represents the heads of EU states, or the Council of Ministers in Brussels, which represents the ministers of EU member states. The Council of Europe is an international organization that focuses on promoting human rights, democracy, and the rule of law across its member states. It is headquartered in Strasbourg and is separate from the EU institutions.

6.4 The World Economic Forum

The World Economic Forum (WEF) is an influential international organisation that fosters dialogue between leaders in business, politics, academia, and other sectors to shape worldwide, regional, and industry-specific agendas. It has taken an active role in promoting and shaping the adoption of emerging technologies, including AI, machine learning, blockchain, robotics, and electric vehicles.

The most prominent initiatives in AI led by the WEF is the [Global AI Action Alliance \(GAIA\)](#). Launched in January 2021, GAIA is a multistakeholder collaboration platform designed to accelerate the adoption of inclusive, transparent, and trusted AI systems worldwide and across various sectors. The alliance brings together over 100 leading entities, including companies, governments, international organisations, non-profits, and academic institutions. These members are united by their commitment to maximising the societal benefits of AI and minimising its associated risks. The activities and goals of GAIA encompass several areas, including:

- Developing and implementing practical tools and frameworks to ensure the ethical use of AI systems that serve all members of society.
- Creating real-time learning and scaling feedback loop across key sectors and challenge areas.
- Catalysing and incubating new partnerships and initiatives to address pressing gaps and needs.

- Building on the Forum's global multistakeholder community of leading businesses, governments, and civil society organisations actively engaged with AI.
- Creating interoperable governance protocols for the development and use of AI technologies²³.

The WEF has several other initiatives related to AI:

- [FireAId Initiative](#), launched in October 2021. The initiative aims to use AI systems, such as drones, satellites, sensors, and predictive models, to mitigate wildfire risks. Its activities and goals include.
- [The Centre for the Fourth Industrial Revolution Network \(C4IR Network\)](#): This is a global network of hubs that collaborate with governments, businesses, civil society, and experts to co-design innovative policy frameworks for emerging technologies such as AI.
- [The Global Future Council on Artificial Intelligence](#): This group of experts provides thought leadership on how AI can be utilised for social good, economic growth, and human-centred development.
- [The Global Lighthouse Network](#): This community consists of manufacturers who are at the forefront of adopting Industry 4.0 technologies such as AI, IoT, robotics, and cloud computing to transform their operations.
- [The Global Shapers Community](#): This network of young leaders drives dialogue, action, and change on various issues, including AI ethics, education, and inclusion.

China and Russia are both represented by some of their companies and organisations as forum members or associate members of the WEF (Alibaba Group, Baidu, Huawei Technologies, Tencent Holdings, Gazprom, Lukoil, Rosneft, Sberbank, etc.). However, China and Russia are not among the strategic

partners of the WEF, which are the most influential and committed members that shape the Forum's agenda and initiatives. The strategic partners are mainly from the United States, Europe, and Japan.

6.5 The Group of 7 (G7) and the “Hiroshima Process”

The G7 is an organisation of seven leading global economies: Canada, France, Germany, Italy, Japan, the UK, and the US, with the EU participating but without hosting rights. Founded in the 1970s, it represents over 60% of the world's net wealth. The group meets yearly to discuss major global issues including economic policy, security, and energy. It's an informal forum without a formal charter or secretariat, and the presidency rotates annually among members, with the presiding country hosting the summit and setting the agenda.

During their recent summit in Hiroshima in May 2023, G7 leaders acknowledged the urgent need for governance and [technical standards to ensure the trustworthiness of artificial intelligence \(AI\)](#). Recognizing that the governance of AI has not kept up with its rapid growth, they initiated discussions to bridge this gap.

The disruptive potential of swiftly evolving technologies, particularly generative AI, emerged as a central concern. To address this, G7 leaders plan to conduct cabinet-level discussions and reveal the outcomes by the end of the year in a process dubbed [the "Hiroshima Process"](#).

The aim is to develop AI systems that are accurate, reliable, safe, and non-discriminatory, irrespective of their source. While G7 leaders understand that methods to achieve trustworthy AI may differ, they unanimously agree that the regulations for digital technologies, including AI, should reflect their shared democratic values.

The U.S. has opted for a cautious approach, proposing the consideration of licensing and testing requirements for AI model development. As the G7 chair this year, Japan has committed to endorsing public and industrial AI adoption, while concurrently monitoring associated risks. In contrast, China has enforced a more restrictive policy to ensure that generative AI services align with the nation's fundamental socialist values.

G7 leaders also advocated for immediate action to evaluate the opportunities and challenges posed by generative AI. They encouraged international organisations, such as the Organisation for Economic Cooperation and Development (OECD), to conduct analyses on policy development impacts. Furthermore, G7 digital ministers concurred on the need to implement "risk-based" AI rules.

Significantly, the broader issue of digital technology governance was also addressed at the summit. Leaders acknowledged that the governance of novel digital technologies, including but not limited to AI, has lagged behind technological advancement. They pledged to tackle common governance challenges, identify potential shortcomings, and update digital economy governance in accordance with democratic values.

These efforts underline a coordinated response to the challenges presented by AI and digital technologies. However, they also highlight the complexity of these issues and the importance of continuous dialogue and cooperation among the G7 nations and beyond.

7.0 Non-profit organisations works on AI

The list is not exhaustive.

7.1 IEEE (Institute of Electrical and Electronics Engineers)

The Institute of Electrical and Electronics Engineers (IEEE) is a non-profit, professional organisation dedicated to advancing technology for the benefit of humanity. Founded in 1884, IEEE is the world's largest technical professional organisation, with over 400,000 members across the globe. IEEE's work spans various domains, including electrical engineering, computer science, telecommunications, biomedical engineering, and artificial intelligence (AI).

IEEE's involvement in AI is multifaceted, covering research, development, standardisation, and ethical considerations:

1. Research and Development: IEEE supports and promotes AI research and development through its conferences, publications, and journals. They provide platforms for researchers, engineers, and practitioners to share their latest findings, discuss challenges, and collaborate on AI-related projects.
2. Standardization: IEEE plays a crucial role in developing standards for AI and related technologies, ensuring the interoperability, safety, and reliability of AI systems. IEEE's Standards Association (IEEE-SA) develops and maintains AI standards, such as the IEEE P7000 series, which focuses on ethical considerations in AI system design.
3. Ethics and AI: IEEE is deeply involved in addressing ethical aspects of AI and autonomous systems. The organization has established the IEEE Global Initiative on Ethics of Autonomous and Intelligent Systems, which aims to ensure that AI technologies are developed and deployed responsibly, prioritizing human well-being and social impact.
4. AI and Policy: IEEE engages with policymakers and other stakeholders to provide guidance on AI policy, regulation, and governance. They advocate for responsible AI development and the adoption of ethical principles in the design and deployment of AI systems.

5. **Education and Training:** IEEE offers AI-related educational resources, including courses, workshops, and webinars, to help professionals and students develop their skills and stay up-to-date with the latest developments in the field.

Through its various initiatives and activities, IEEE plays a significant role in shaping the future of AI, ensuring that it is developed and deployed responsibly, ethically, and in the best interests of society.

7.2 ITIF (Information Technology and Innovation Foundation)

The Information Technology and Innovation Foundation (ITIF) is a non-profit, non-partisan think tank based in Washington, D.C. Established in 2006, ITIF focuses on the intersection of technological innovation and public policy, conducting research and providing recommendations on a wide range of topics, including artificial intelligence (AI).

ITIF's work on AI covers several key areas:

1. **Policy Research:** ITIF conducts research on various AI policy issues, such as regulation, governance, data privacy, and intellectual property. They analyse the potential risks and benefits of AI, as well as the challenges and opportunities it presents for policymakers.
2. **Economic Impact:** The organisation explores the impact of AI on economic growth, productivity, and global competitiveness. They examine how AI can drive innovation and support industries, as well as the potential implications for jobs, wages, and economic inequality.
3. **Social Impact:** ITIF investigates the social implications of AI across sectors such as education, healthcare, transportation, and public safety. They provide insights on how AI can be harnessed to improve societal outcomes and address potential negative consequences.

4. **Global Perspective:** ITIF considers the global aspects of AI development, deployment, and regulation. They analyse international competition in AI, as well as the role of international cooperation and coordination in addressing AI-related challenges.
5. **Public Engagement:** The organisation actively engages with the public, policymakers, and other stakeholders through events, panel discussions, and conferences on AI and its implications for society. Their goal is to foster dialogue and knowledge-sharing among various actors in the AI ecosystem.

Through their research and advocacy efforts, the Information Technology and Innovation Foundation aims to ensure that AI development and deployment are carried out responsibly and in the best interests of society. They provide valuable insights and recommendations to help policymakers navigate the complex landscape of AI and its effects on the economy, society, and the world at large.

7.3 Brookings Institution

The Brookings Institution is a non-profit public policy organisation based in Washington, D.C., which conducts research and provides recommendations on a wide range of topics, including artificial intelligence (AI) and its impact on society. Established in 1916, Brookings is one of the oldest and most influential think tanks in the United States.

While Brookings' research covers a broad spectrum of policy areas, their work on AI focuses on several key aspects:

1. **Policy Research:** Brookings conducts research on the policy implications of AI, including areas such as regulation, governance, ethics, privacy, and security. They analyse the potential risks

- and benefits of AI and provide recommendations for policymakers on how to address these issues.
2. **Economic Impact:** The institution explores the impact of AI on the economy, job market, and workforce. They examine the potential for AI to drive economic growth and productivity, as well as the challenges it may pose to employment and income inequality.
 3. **Social Impact:** Brookings analyses the social implications of AI, such as its effects on education, healthcare, criminal justice, and social equity. They provide insights on how AI can be harnessed to improve social outcomes and mitigate potential negative consequences.
 4. **Global Governance:** The organisation investigates the role of AI in global governance and international relations, including issues related to AI and national security, international competition, and cooperation on AI development and regulation.
 5. **Public Engagement:** Brookings actively engages with the public, policymakers, and other stakeholders through events, panel discussions, and conferences on AI and its implications for society. They aim to foster dialogue and knowledge-sharing among various actors in the AI ecosystem.

Through their research and advocacy efforts, the Brookings Institution seeks to ensure that AI is developed and deployed in a way that aligns with public interest, respects human rights, and promotes social and economic well-being. They provide valuable insights and recommendations to help guide policymakers in navigating the complex landscape of AI and its effects on society.

7.4 ISO/IEC JTC 1/SC 42

ISO/IEC JTC 1/SC 42, also known as the Joint Technical Committee 1/Subcommittee 42, is a subcommittee of the International Organization for Standardization (ISO) and the International Electrotechnical Commission (IEC) that focuses on

developing standards for artificial intelligence (AI) and related technologies. The subcommittee was established in 2017 to address the growing need for international standards and guidelines to ensure the responsible development and deployment of AI systems.

The scope of ISO/IEC JTC 1/SC 42 covers several key areas, including:

1. **Standardisation:** The primary focus of ISO/IEC JTC 1/SC 42 is the development of international standards, technical reports, and guidelines related to AI. These standards address topics such as terminology, frameworks, reference architectures, trustworthiness, robustness, and ethical considerations.
2. **Interoperability:** The subcommittee works on creating standards that promote interoperability and compatibility between AI systems and applications, ensuring that AI technologies can be seamlessly integrated into various industries and domains.
3. **Trustworthiness:** ISO/IEC JTC 1/SC 42 is concerned with establishing guidelines and best practices for trustworthiness in AI systems, including aspects such as transparency, explainability, accountability, and privacy.
4. **Use Cases and Applications:** The subcommittee examines various AI use cases and applications, identifying the specific requirements and challenges associated with each domain. This helps them develop tailored standards and guidelines that address the unique needs of different industries and sectors.
5. **Collaboration:** ISO/IEC JTC 1/SC 42 actively collaborates with other ISO/IEC subcommittees, national and international standardisation bodies, and stakeholder organisations to ensure a coordinated and harmonised approach to AI standardisation. This includes liaising with other groups working on AI-related standards and sharing expertise and resources.

Through its work on AI standards, ISO/IEC JTC 1/SC 42 aims to ensure that AI technologies are developed and deployed responsibly and safely, taking into account ethical considerations, human rights, and the potential societal impacts of AI.

7.5 The AI Now Institute

The AI Now Institute is a research organisation based at New York University that focuses on understanding the social implications of artificial intelligence (AI) and related technologies. Founded by Kate Crawford and Meredith Whittaker, the AI Now Institute is committed to producing interdisciplinary research that addresses the complex challenges posed by AI, machine learning, and other emerging technologies.

The AI Now Institute's work is centred around several key areas, including research, policy, and public engagement:

1. **Research:** The organisation conducts interdisciplinary research on the social, ethical, and political implications of AI. They explore topics such as bias, fairness, accountability, transparency, and the impact of AI on labor, healthcare, and criminal justice, among others.
2. **Policy:** AI Now develops policy recommendations and guidelines to address the challenges and risks associated with AI and its deployment in various sectors. They engage with policymakers, regulators, and other stakeholders to promote responsible AI development that aligns with public interest and democratic values.
3. **Public Engagement:** The AI Now Institute actively participates in public dialogue on AI and its societal implications. They host events, workshops, and conferences that bring together researchers, policymakers, practitioners, and

the public to discuss and debate the ethical, legal, and social aspects of AI.

4. Education and Training: The organization contributes to AI education and training by offering resources, mentorship, and support for students and researchers working on AI-related projects. They also collaborate with other institutions to develop interdisciplinary curricula that address the social dimensions of AI.
5. AI in Context: AI Now's research is organized around four main research areas: rights and liberties, labour and automation, bias and inclusion, and safety and critical infrastructure. This focus allows them to examine the impact of AI on various aspects of society and develop targeted recommendations for policy and practice.

Through its research and advocacy efforts, the AI Now Institute aims to ensure that AI development and deployment are carried out responsibly, taking into account the complex social, ethical, and political challenges these technologies present.

7.5 AlgorithmWatch .

AlgorithmWatch is a non-profit organisation focused on evaluating and raising awareness about algorithmic decision-making processes that have significant social relevance. Their mission is to ensure that these processes are transparent, accountable, and in line with democratic values. AlgorithmWatch's work spans research, advocacy, education, and public engagement in the following areas:

1. Research: AlgorithmWatch conducts research on algorithmic decision-making and its impact on society. They investigate issues such as bias, fairness, accountability, and transparency in AI and automated systems, providing insights and recommendations to address these concerns.

2. **Advocacy:** The organisation advocates for policy and regulatory changes that promote transparency, accountability, and fairness in algorithmic decision-making. They engage with policymakers, regulators, and other stakeholders to ensure that AI and automation are developed and deployed responsibly.
3. **Education:** AlgorithmWatch raises awareness about the societal implications of AI and algorithmic decision-making through workshops, seminars, and public events. They provide educational resources and training programs to help individuals, organisations, and policymakers better understand the challenges and opportunities presented by these technologies.
4. **Public Engagement:** The organisation actively engages with the public and fosters dialogue on the ethical, legal, and social aspects of AI and automation. They host events, panel discussions, and conferences that bring together experts, practitioners, and stakeholders from diverse backgrounds to share knowledge and experiences.
5. **Automating Society Report:** AlgorithmWatch publishes the "Automating Society" report, which examines the implementation of AI and automation in public services across Europe. The report provides an overview of the current state of AI deployment in various sectors, highlights best practices, and offers recommendations for ensuring that these technologies serve the public interest.

Through their work, AlgorithmWatch aims to create greater awareness and understanding of the implications of algorithmic decision-making, promote responsible AI development and deployment, and contribute to a more inclusive and democratic digital society.

7.6 Partnership on AI (PAI).

The Partnership on AI (PAI) is a non-profit organisation founded by leading technology companies, including Amazon, Apple, Google,

Facebook, IBM, and Microsoft, along with other stakeholders from academia, civil society, and industry. PAI's mission is to address the global challenges and opportunities of AI by ensuring that its development and deployment are carried out responsibly and in the best interests of society. PAI's work is focused on several key areas, including research, policy recommendations, collaboration, and public engagement:

1. Research: PAI conducts research on topics related to AI ethics, safety, transparency, and fairness. They explore best practices for AI development and deployment, analyse the impact of AI on society, and provide insights on how to address potential risks and challenges.
2. Policy Recommendations: The organisation develops policy recommendations that promote responsible AI development and address issues such as AI governance, privacy, and transparency. They work with stakeholders across the AI ecosystem to help shape policies and regulations that ensure AI benefits society.
3. Collaboration: PAI fosters collaboration between its member organisations, which include leading AI companies, research institutions, and civil society organisations. They facilitate the sharing of knowledge and expertise, encourage joint projects, and promote cooperation on AI-related challenges and opportunities.
4. Public Engagement: The organisation actively engages with the public through events, workshops, and panel discussions on AI and its implications for society. They aim to raise awareness of the ethical, safety, and policy issues surrounding AI and facilitate dialogue between stakeholders from diverse backgrounds.
5. Thematic Pillars: PAI's work is organised around six thematic pillars: safety-critical AI, fair, transparent, and accountable AI, AI, labour, and the economy, collaboration between people and AI systems, AI and social good, and special initiatives on AI and COVID-19.

By working on these key areas, the Partnership on AI aims to ensure that AI is developed and deployed in ways that are ethical, safe, and beneficial for everyone. They strive to create a global community of stakeholders that can collaborate and address the challenges and opportunities presented by AI.

7.7 Future of Life Institute (FLI).

The Future of Life Institute (FLI) is a non-profit organisation focused on mitigating existential risks and ensuring a positive future for humanity. FLI is particularly concerned about the potential risks arising from advanced artificial intelligence (AI) and other emerging technologies. Their work encompasses research, advocacy, collaboration, and public engagement in the following areas:

1. **Research:** FLI supports research on AI safety and long-term risk mitigation. They provide grants and funding to researchers working on projects aimed at ensuring the safe and beneficial development of AI and other transformative technologies.
2. **Advocacy:** FLI advocates for responsible AI development, raising awareness of the potential risks associated with AI and the need for AI safety research. They engage with policymakers, technology companies, and the AI research community to promote the adoption of safety measures in AI research and development.
3. **Collaboration:** The organisation brings together experts from various fields, including AI, robotics, computer science, and policy, to foster interdisciplinary collaboration on AI safety and risk mitigation. FLI facilitates discussions and partnerships among researchers, engineers, and policymakers to address global challenges.
4. **Public Engagement:** FLI is committed to raising awareness about the ethical, safety, and policy implications of AI and other emerging technologies. They organise conferences,

workshops, and public events that explore the potential risks and benefits of AI, providing a platform for dialogue and knowledge-sharing.

5. Asilomar AI Principles: FLI has developed the "Asilomar AI Principles," a set of 23 guidelines designed to ensure that AI research and development is conducted responsibly and with long-term safety in mind. These principles cover topics such as research funding, value alignment, robustness, and cooperation among research institutions.

Through these activities, the Future of Life Institute aims to ensure that advanced AI and other transformative technologies are developed and deployed safely and responsibly, for the benefit of all humanity.

[In April 2023, FLI called for a suspension of the development Large Language Models like GPT-4 for 6 months.](#)

7.8 Centre for Humane Technology (CHT).

The Centre for Humane Technology (CHT) is a non-profit organisation focused on addressing the potential negative effects of technology on society. Founded by former Google design ethicist Tristan Harris and several other tech industry professionals, CHT seeks to realign technology with humanity's best interests by promoting ethical design, development, and deployment of digital technologies, including AI.

CHT's work encompasses a range of activities, including research, advocacy, education, and collaboration:

1. Research: CHT conducts research on the unintended consequences of technology and its impact on mental health, democracy, privacy, and other aspects of society. They explore ways to minimize the negative effects of technology and promote its responsible use.

2. **Advocacy:** CHT engages with technology companies, policymakers, and other stakeholders to advocate for more ethical technology development practices. They work towards influencing industry practices and public policies that ensure technology serves human needs and values.
3. **Education:** The organisation raises awareness about the potential harms of technology through public talks, conferences, and media appearances. They provide resources, including guidelines and best practices, to help individuals, educators, and organisations navigate the digital landscape responsibly.
4. **Collaboration:** CHT collaborates with technologists, designers, policymakers, and other organisations to create a movement for more humane technology. They foster partnerships and facilitate conversations among various stakeholders to drive positive change in the technology industry.
5. **Initiatives and Projects:** CHT supports and participates in various initiatives and projects that promote ethical technology development and use. This includes the development of tools, resources, and programs that help individuals and organisations understand and adopt more humane technology practices.

Through these activities, the Centre for Humane Technology aims to create a more balanced relationship between technology and society, where the digital ecosystem supports human well-being, fosters meaningful connections, and upholds democratic values.

7.9 AI for People.

AI for People is a non-profit organisation that aims to promote the development of AI and other emerging technologies that serve the public interest. They focus on fostering dialogue, collaboration, and innovation in AI to ensure that the technology is accessible, ethical, and beneficial to society. Their work encompasses

various activities, including research, education, and community engagement.

1. **Research:** AI for People conducts research on key topics related to AI ethics, governance, and policy. This includes producing guidelines, recommendations, and thought pieces that address important issues such as fairness, accountability, transparency, and user-centric design in AI systems.
2. **Education:** The organisation is committed to raising awareness about AI and its potential social impact. They provide educational resources and host events, workshops, and conferences to foster understanding and knowledge-sharing among different stakeholders, including academia, industry, civil society, and policymakers.
3. **Community Engagement:** AI for People actively collaborates with other organisations and stakeholders in the AI ecosystem. They build partnerships and alliances to promote responsible AI development and create opportunities for dialogue and cooperation among various actors in the field.
4. **Advocacy:** AI for People advocates for public policies and regulations that support ethical AI development and deployment. They engage with policymakers, regulatory bodies, and other stakeholders to ensure that the concerns and interests of the public are represented in AI-related policy discussions.
5. **Projects and Initiatives:** AI for People supports and participates in various projects and initiatives that align with their mission. These projects may involve the development of AI applications for social good, the creation of tools and resources to support ethical AI development, or the organisation of events and activities that promote public dialogue on AI and its implications for society.

Overall, AI for People is dedicated to ensuring that AI is developed and deployed in a way that

aligns with the the public interest, respects human rights, and promotes social and economic well-being.

7.10 Centre for Democracy and Technology

The Centre for Democracy and Technology (CDT) has published materials related to AI. The CDT is a non-profit organisation that works to promote democratic values by shaping technology policy and architecture. Their focus includes privacy, free expression, and human rights in the digital age. In the field of AI, the CDT has published various resources, including policy recommendations, reports, and position papers, addressing AI's impact on society, ethics, and governance. Some of their work includes:

1. "Digital Decision-Making: The Building Blocks of Machine Learning and Artificial Intelligence" - This report provides an overview of machine learning and AI technologies, explaining the key concepts, methods, and ethical considerations.
2. "AI and Machine Learning: Policy Paper" - This position paper discusses the challenges and opportunities of AI, offering policy recommendations on privacy, fairness, accountability, transparency, and security.
3. "Preserving Privacy in Machine Learning: A Technical and Legal Overview" - This resource explains the privacy implications of machine learning and offers recommendations for preserving privacy in AI systems.
4. "Algorithmic Bias and Fairness: A Path Forward" - This report explores the issue of bias in AI and provides recommendations for addressing fairness and accountability in algorithmic decision-making.

The CDT also actively participates in events, workshops, and panel discussions related to AI, technology policy, and digital rights. They collaborate with other organisations and

stakeholders to promote responsible AI development and deployment in line with democratic values.

8.0 National Governmental frameworks on AI

Governmental frameworks for AI serve various purposes.

- Regulators and policymakers play a critical role in ensuring that AI models are developed and deployed ethically, responsibly, and in a trustworthy way, following the guidelines of international organisations.
- AI frameworks aim to provide legal certainty to the industry and attract investments.
- AI also addresses legislation and regulation of AI from a national security perspective.

Governmental frameworks can be advisory and/or legally binding depending on the approach followed.

- Advisory AI frameworks establish baselines that courts consider when assessing whether an organisation supervising an AI system has acted with adequate due diligence or not. They require strong democratic rule of law mechanisms to be most effective.
- Legally binding AI frameworks provide more trustworthiness to end-users than advisory frameworks. However, they can also increase legal uncertainty if they are not adequately designed to fit the specific case being regulated. For instance, if the regulations become outdated too quickly due to the rapid evolution of technology.

8.1 The EU

In April 2021, the European Union (EU) implemented binding horizontal legislation aimed at fostering trustworthiness and enhancing its internal market.

This legislation adopts a four-tier, risk-based approach that evaluates potential threats to safety and fundamental rights. It features provisions for straightforward biennial revisions, considering experience and technological advancements, as well as systematic consultations with various stakeholders, including industry representatives, consumer organisations, and NGOs.

In contrast to the United States, the EU's AI framework does not prioritise national security issues because they are the sole responsibility of member states and fall outside the scope of the European Commission's authority.

The EU AI Act, introduced in 2021, proposes classifying AI systems into four categories, each subject to conformity assessment obligations, with the most stringent measures applied to the riskiest systems. The conformity assessment can range from a simple self-assessment to a more rigorous process involving independent third-party involvement for the most high-risk use cases.

In September 2022, the European Commission proposed the AI Liability Directive, which establishes a variety of penalties for non-compliance with the EU AI Act. Additionally, the AI Liability Directive allows for private enforcement, enabling individuals or organizations that suffer losses due to non-compliance with the Directive to pursue legal action against the liable parties.

The two laws are still in the making and are expected to be co-adopted by the European Parliament and the Council of Ministers before the next European elections in June 2024. Some of the amendments proposed to the initial proposal of the European Commission in April 2021 will require sensible negotiations with the European Council of Ministers representing the Member States of the European Union, especially in areas related to law enforcement.

The key changes proposed by the LIBE and the IMCO Committees following the vote that took place on May 11, 2023, regarding the amendment of the EU AI Act are:

1. **A ban on predictive policing**, which is the use of AI to predict future crimes or identify potential offenders.
2. **There are a number of additions to the list of stand-alone AI systems categorized as high-risk**, such as AI systems used for biometric identification, social scoring, migration management, education and vocational training, employee worker management, and access to self-employment.
3. **A strong and inclusive role for the new AI Office**, which will be responsible for monitoring and enforcing the AI Act, as well as promoting dialogue and cooperation among stakeholders.
4. **A stronger alignment with the General Data Protection Regulation (GDPR)**, such as ensuring data quality and minimisation, transparency and accountability, and data protection by design and by default.
5. **An increased involvement of stakeholders in several areas**, such as setting technical standards, conducting conformity assessments, establishing codes of conduct and providing guidance.
6. **The introduction of specific provisions related to general purpose Artificial Intelligence**, which is AI that can perform multiple tasks across different domains and contexts^{1 3}.

These changes are part of a draft report that needs to be endorsed by the whole Parliament before negotiations with the European Council of Ministers can begin.

8.2 The US

Initiatives started by the US Congress

There has been a non-partisan "AI Caucus" in the US Congress since 2017, but as of early 2023, the United States has not adopted any federal horizontal or sectoral legislation on AI yet.

The AI Caucus is a bipartisan group in the US Congress that was established to inform policymakers of the technological, economic, and social impacts of advances in artificial intelligence and to ensure that rapid innovation in AI and related fields benefits Americans [7]. The House of Representatives established its AI Caucus several years ago, and the Senate announced the creation of its AI Caucus in March 2021 [3][5]. The Caucus recently launched an ongoing core initiative, the [AI Across America](#) project, to support efforts in the public and private sectors to make AI education, training, and RD available for communities across the country [2][4].

The AI Caucus of the US Congress has also played a key role in the initiation of the National Artificial Intelligence Initiative Act. The act became law on January 1, 2021, and provides for a coordinated program across the entire Federal government to accelerate AI research and application for the Nation's economic prosperity and national security. The act also establishes the National Artificial Intelligence Initiative Office (NAIIIO) to lead and oversee the implementation of the National AI Initiative.

The National AI Initiative Act, signed into law in January 2021 also established the National AI Research Resource Task Force. (NAIRRRTF). The task force was responsible for developing a roadmap for the creation of a National AI Research Resource (NAIRR), a shared research infrastructure that can support AI research and development across multiple sectors, including academia, government, and industry.

The "NAIRR" objective is to strengthen and democratise the US AI innovation ecosystem in a way that protects privacy, civil rights, and civil liberties to democratise the AI research and development (R&D) landscape in the United States for the benefit of all. It aims to provide a widely accessible platform for researchers and students from diverse backgrounds who are pursuing foundational, use-inspired, and translational AI research. It will be created by bringing together computational resources, data, testbeds, algorithms, software, services, networks, and expertise. The National Artificial Intelligence Initiative Resource Research Task Force (NAIRRTF) also led the publication in April 2022 by the Office of Science and Technology of the White House of the AI Bill of Rights and in January 2023 of an AI Risk Management Framework by NIST (part of the US Department of Commerce).

[The AI Bill of Rights](#) aims to guide the design, development, and deployment of artificial intelligence (AI) and other automated systems so that they protect the rights of the American public. The AI Bill of Rights is an advisory, non-binding set of guidelines expected to be followed on a voluntary basis, principally by the US government and its agencies.

The NIST AI Risk Management Framework (RMF) [2] is intended to be a voluntary resource for organisations involved in the design, development, use, and evaluation of AI products, services, and systems. It aims to promote the trustworthy and responsible development and use of AI systems [4]. The framework provides a common language and structure for organisations to manage the risks associated with AI systems [2]. Note that the NIST AI RMF also includes considerations for national security and recommends that organisations conduct risk assessments that take into account potential threats to national security [2]. Additionally, NIST developed standards and guidelines for testing

and verifying AI systems to ensure they met the standards of trustworthiness.

Other initiatives led by the US government include:

The US government has been very active in developing its governmental framework on AI since 2018, with the successive works of the National Committee on Artificial Intelligence (NCAI) and the National Security Committee on Artificial Intelligence (NSCAI).

The National Commission on Artificial Intelligence (NCAI) was established by the National Defense Authorization Act (NDAA) of 2018 to review advances in artificial intelligence (AI) and machine learning, with a relative wide scope including both the competitiveness of the United States in AI research and development as well as the potential national security implications of AI. The NCAI released its final report in March 2020, which included recommendations for significant investments in AI research and development as well as the creation of a National AI Research Resource to support this work [1] through the National AI Initiative Act passed the same year.

The National Security Commission on Artificial Intelligence (NSCAI) was established by the John S. McCain National Defense Authorization Act for Fiscal Year 2019 to review advances in AI and machine learning, with a unique focus on their implications for national security, defence, and intelligence. The commission was responsible for providing recommendations to the president and Congress on how best to advance the development and use of AI to enhance national security [4].

In March 2021, the National Security Commission on Artificial Intelligence (NSCAI) released [its final report](#), which included recommendations to strengthen the nation's AI research and development capabilities, invest in AI-related talent and infrastructure, protect the nation's AI advantage,

and develop a roadmap for US strategic competition with China and other countries, as seen purely from a US National Security perspective. The NSCAI also resulted in the development of an Artificial Intelligence Ethics Framework for the Intelligence Community to provide guidelines for the responsible use of AI in intelligence operations to protect national security while upholding ethical principles [3].

Finally, there are also sector-specific guidelines, such as those for the healthcare industry. The US Food and Drug Administration (FDA) has issued guidance on the development and use of AI in medical devices, including recommendations for premarket review, validation, and monitoring of AI algorithms.

In conclusion, the US government's overall AI strategy aims to maintain US AI leadership and competitiveness as part of the US National Security Strategy while ensuring that AI is developed and used in a responsible and trustworthy manner [5].

EU-US Cooperation on AI

In the field of AI, the main objective is to agree on a common AI terminology for ensuring that conformity efforts on both sides of the Atlantic are mutually recognised and deemed equivalent when relevant and applicable.

The 27 Member states of the EU have formally tasked the European Commission and the European External Action Service to represent them during talks of the [EU-US Technology Trade Council](#) that started in September 2021. Mid-2023 the EU and US agreed to [develop a common code of conduct](#).

The Department of Commerce and NIST on the US side and the European Commission on the EU side respectively led the technical discussions.

8.3 The UK

The UK government has established an AI framework that aims to promote the ethical development and deployment of AI while also fostering innovation and growth in the sector [1]. The framework includes a mix of advisory guidelines, binding regulations, and sector-specific initiatives.

One of the key efforts of the UK government is the development of ethical principles for AI, which include transparency, accountability, and fairness. These principles guide the development and use of AI technologies across all sectors and are intended to promote public trust and confidence in AI [8]. However, these ethical principles are advisory guidelines and do not have the force of law.

In addition to ethical principles, the UK government has also established regulatory frameworks for specific AI applications, such as autonomous vehicles and healthcare [4]. These regulations help to ensure that AI technologies are safe, secure, and reliable and provide a clear legal framework for companies and developers working in these areas. These regulatory frameworks are binding legislation or regulations.

There are also horizontal AI initiatives that aim to promote the development of AI across different sectors, such as the AI Sector Deal and the Office for AI [1]. These initiatives provide funding, support, and guidance for businesses and researchers working in AI, with a focus on driving innovation and competitiveness in the sector. These horizontal AI efforts aim to address AI global competitiveness.

Finally, the UK government is also focused on addressing national security questions related to AI. This includes developing policies and regulations to protect against the misuse of AI technologies, as well as supporting research into AI-related security threats and vulnerabilities

[5]. These national security questions are addressed through sectorial AI efforts such as the Centre for Data Ethics and Innovation and the National Cyber Security Centre [6].

Overall, the AI framework put in place by the UK government aims to promote the ethical development and use of AI, while also supporting innovation and growth in the sector. By establishing ethical principles, regulations, and sector-specific initiatives, the government is working to ensure that AI technologies are safe, reliable, and trustworthy, and that they contribute to UK competitiveness and national security [7].

8.4 Japan

The Japan government has put in place an AI governance framework to promote the development and utilisation of AI technology in the country while ensuring ethical, safe, and secure use of the technology. The framework consists of a set of advisory guidelines and principles that provide recommendations for the development and implementation of AI systems in various sectors [1].

The Japanese government has been actively promoting innovation in AI technology by encouraging various players, including start-ups and small- and medium-sized enterprises, to come up with brand-new and innovative ideas to provide the world with solutions [2]. Moreover, it has also introduced a standardised guideline on digital government to achieve a digital government responsive to the changing digital society through the improvement of project management capacities [5].

The advisory guidelines and principles in the AI governance framework are not legally binding but serve as recommendations for the development and implementation of AI systems [1]. The Japanese government has emphasised the importance of a goal-based governance model that can guide entities

such as companies towards achieving common goals, rather than a conventional rule-based model [1].

The AI governance framework focuses on both horizontal AI efforts, which apply across different sectors, and sector-specific efforts, such as in healthcare and transportation [1]. The framework also aims to address the global competitiveness of AI technology and promote Japan's national security by ensuring safe and secure use of the technology [1].

In summary, the Japan government has put in place an AI governance framework consisting of advisory guidelines and principles that provide recommendations for the development and implementation of AI systems. The framework focuses on both horizontal and sector-specific AI efforts and aims to address AI's global competitiveness and Japan's national security by ensuring safe and secure use of the technology. However, the advisory guidelines and principles are not legally binding [1].

8.5 South Korea

South Korea has been actively promoting the development and use of AI technology through various initiatives and frameworks. In 2019, the country announced its National Strategy for AI, which was jointly developed by all parties, including the Ministry of Science and ICT [3]. The strategy aims to undertake nine strategies and 100 initiatives in the three main areas of AI: establishment of an AI ecosystem, utilisation of AI, and creation of human capabilities for AI [3].

To promote the use of AI, Korea amended its three main privacy laws to allow data use [1]. It also enacted a framework act on intelligent informatization to foster an enabling environment for AI use [1]. The government is taking steps to foster growth in the area of technology, with a

vision to lead the world in the global AI sector [4].

The AI framework put in place by the South Korean government comprises both advisory guidelines and binding legislation or regulation [8]. The Korean government is reviewing bills carefully and has not enacted or announced any new Acts or principles despite the speedy development of the AI industry [8].

The AI framework is horizontal, as it aims to foster an enabling environment for AI use across all industries and sectors. The government has established an AI ecosystem, which includes the development of AI centres, research institutions, and AI-powered infrastructure [3]. There are also sectorial AI efforts. For instance, South Korea aims to join the AI race by supporting Korean AI startups [7].

The AI framework aims to address South Korea's AI global competitiveness and national security questions. To achieve its vision of leading the world in the global AI sector, the Korean government has been promoting trustworthy AI that enhances the benefits of the technology and addresses its risk factors [2]. The framework also takes into account national security issues, given the strategic importance of AI in the military and defense sectors [2].

In summary, the South Korean government has put in place an AI framework that promotes the development and use of AI technology across all industries and sectors. The framework comprises both advisory guidelines and binding legislation or regulation, and it aims to address both South Korea's AI global competitiveness and national security questions.

8.6 China

The Chinese government has put in place an AI framework to guide the development and deployment

of AI in the country. The framework consists of a set of ethical norms, advisory guidelines, and binding regulations that aim to ensure the safe, ethical, and beneficial use of AI technology [3].

The Cyberspace Administration of China (CAC), a potent regulator that creates the rules governing particular applications of AI, has made one of the most significant moves in AI governance. The CAC's approach is the most mature, the most rule-based, and the most concerned with AI's role in disseminating information [1]. Additionally, the Chinese Academy of Science recognizes eight key AI technologies that have achieved breakthroughs and identified specific areas of application, including computer vision, natural language processing, trans-media analysis and reasoning, intelligent adaptive learning, collective intelligence, automated reasoning, and machine consciousness [6].

The AI framework aims to address Chinese AI global competitiveness and Chinese national security questions. The Chinese government aims to become the leading AI power with an industry worth at least RMB 1000 billion by 2030 [2]. Moreover, the Chinese government has released its ambitious New Generation Artificial Intelligence Development Plan (AIDP), which sets the eye-catching target of national leadership in a variety of AI fields by 2030 [5].

The Chinese AI framework includes both horizontal and sectoral efforts. For instance, the Chinese government has released a framework with guidelines for addressing issues of AI ethics in the private sector, while also exploring AI governance and policy more broadly [8]. Moreover, the Chinese government has recognized the importance of AI chips in AI infrastructure and is promoting more innovation by shifting AI development in China into overdrive, paving the way for more open-sourcing and exploration of AI technologies in various industries [7].

In summary, the Chinese AI framework consists of ethical norms, advisory guidelines, and binding regulations that aim to guide the development and deployment of AI in the country, with the goal of becoming the leading AI power by 2030. The framework includes both horizontal and sectoral efforts and aims to address Chinese AI global competitiveness and Chinese national security questions. While some efforts are advisory, others are binding, and the Cyberspace Administration of China is a powerful regulator that writes rules governing certain applications of AI.

Early in 2023, the Chinese government imposed that all AI-generated content be watermarked. It has also imposed that data generated by AI systems must comply with Chinese societal values and norms, thereby encouraging providers of AI services in China to censor themselves for information or content not fitting with Chinese propaganda.

8.7 Taiwan

The Taiwan government has implemented an AI framework aimed at improving the country's competitiveness and addressing national security concerns. The AI Taiwan Action Plan (2018-2021) was launched on January 18, 2018, to prioritise innovation and real-world implementation and sharpen Taiwan's advantages [1]. The government allocated NT\$5 billion for building AI servers and promoting tech talent cultivation under the Cabinet's eight-year Forward-looking Infrastructure Program [2]. Additionally, the government believes that AI will play a crucial role in the 5+2 Industrial Innovation Plan, which focuses on seven industries, including smart machinery and the Asia Silicon Valley Project [3].

The Taiwan government's approach to regulating AI is largely advisory, rather than binding legislation or regulation. However, the government has established several regulations and guidelines to promote responsible AI development and use. For

example, the Personal Data Protection Act (PDPA) sets guidelines for AI use in personal data processing, while the National Development Council has published guidelines for AI development and use in public services [3]. Moreover, several private organisations have developed AI governance frameworks to support compliance with forthcoming European AI regulations [5].

The Taiwan government has also prioritised both horizontal and sectorial AI efforts. Horizontal efforts include the development of talent and the establishment of research centres, while sectorial efforts aim to apply AI to specific industries such as healthcare, finance, and transportation [3].

Finally, the Taiwan government's AI framework aims to address both global competitiveness and national security concerns. The AI Taiwan Action Plan aims to develop Taiwan into a hub for AI innovation and development, while the government has also established regulations and guidelines to promote responsible AI use and safeguard national security [1][3]. Additionally, the government has established partnerships with other countries to enhance Taiwan's AI competitiveness and promote global cooperation, such as the seventh Global Cooperation and Training Framework Joint Committee meeting hosted in 2021 [8].

In summary, the Taiwan government has implemented an AI framework that aims to improve the country's competitiveness, address national security concerns, and promote responsible AI development and use. The government's approach to regulating AI is largely advisory, and the framework includes both horizontal and sectorial efforts.

8.8 Australia

The Australian government has developed an Artificial Intelligence (AI) Ethics Framework [1], which guides businesses and governments to design, develop, and implement AI in a responsible and

inclusive manner. The framework includes eight voluntary principles, which are meant to complement existing AI regulations and practices and are intended to be inspirational [2]. The principles aim to ensure that AI is developed and used ethically, transparently, and fairly, and that its use does not harm individuals, communities, or the environment [1].

In addition to the AI Ethics Framework, the Australian Government has also released the AI Action Plan [3], which outlines a vision for Australia to establish itself as a global leader in developing and adopting trusted, secure, and responsible AI. The plan includes specific actions that the government will take to achieve this vision, such as investing in research and development, fostering collaboration between industry, academia, and government, and developing skills and talent in the AI sector [5].

It is important to note that the AI Ethics Framework is voluntary and inspirational and does not have the force of law. However, the framework is intended to complement existing AI regulations and practices, and may influence future legislation in this area [2]. At present, there is no binding legislation or regulation specifically addressing AI in Australia.

The AI framework developed by the Australian government includes both horizontal and sectoral efforts. The AI Ethics Framework and the AI Action Plan apply to all sectors where AI is used. However, there are also specific initiatives in place in certain sectors. For example, the National AI Centre, which was established by the Australian government, focuses on driving business adoption of AI technologies and addressing barriers faced by SMEs in adopting and developing AI and emerging technologies [9].

The AI framework aims to address both Australian AI global competitiveness and Australian national

security questions. The AI Action Plan outlines a vision for Australia to become a global leader in developing and adopting trusted, secure, and responsible AI [3], which suggests that the framework aims to enhance Australian competitiveness in this area. Additionally, the framework includes a principle on privacy protection, which aims to ensure that AI is used in a way that respects individuals' privacy rights and promotes their trust in AI systems [7]. This suggests that the framework also addresses questions around Australian national security and the protection of personal data.

In summary, the AI framework put in place by the Australian government includes a voluntary AI Ethics Framework and an AI Action Plan. The framework includes voluntary principles intended to guide businesses and governments in designing, developing, and implementing AI in a responsible and inclusive manner. While the framework is not binding legislation, it may influence future legislation in this area. The framework includes both horizontal and sectoral efforts and aims to address both Australian AI global competitiveness and Australian national Security.

8.9 India

The Indian government has implemented several initiatives and strategies to harness the potential of artificial intelligence (AI) in the country. One such initiative is the National Strategy for Artificial Intelligence, also known as #AIForAll, which was published by NITI Aayog in June 2018[3]. This strategy aims to leverage AI for inclusive growth, job creation, and new business opportunities across various sectors.

The Indian government has also established an AI Ethics Framework to guide the design, development, and deployment of AI systems in the country[4]. The framework is aimed at creating an overarching ethics framework for the responsible use of AI and

has encouragement mechanisms to move from principles to practice. However, it should be noted that this framework is currently only advisory and not binding legislation or regulation.

The Indian government has taken sectorial approaches to AI implementation in several areas. For instance, the use cases of AI in the Indian government currently include biometric identification, facial recognition, criminal investigation, crowd and traffic management, digital agriculture, and healthcare[1][5]. The Indian government is also exploring the usage of AI in the news and has collaborated with health-tech startups to deploy automated COVID-19 monitoring systems[2][8].

Moreover, the Indian government recognizes the importance of addressing Indian AI global competitiveness and national security questions[6]. The National Strategy for Artificial Intelligence aims to position India as a leader in the global AI landscape and help India become a \$1 trillion digital economy by 2025. The government is also working on a technology roadmap, a standards framework, and a national AI Ethics Framework to support the responsible development of AI[7].

In summary, the Indian government has put in place several initiatives and strategies to leverage AI for inclusive growth and address sectoral challenges. The AI Ethics Framework is an advisory guideline that aims to guide the responsible use of AI in the country, but there are no binding regulations or legislation in place. The Indian government has taken sectorial approaches to AI implementation, and there are ongoing efforts to address Indian AI global competitiveness and national security questions.

8.10 New Zealand

The New Zealand government has put in place an AI framework to guide the use of algorithms by its

agencies [1]. The Algorithm Charter for Aotearoa New Zealand was published in 2020 and emphasises that more complex algorithms can be used to support human decision-making. It builds on previous AI research by the New Zealand Law Society [2]. The Charter is advisory in nature, meaning it provides guidance to government agencies but is not binding legislation or regulation [2].

There are horizontal AI efforts in New Zealand, to develop an ethical AI framework and action plan to address ethical issues arising from the use of AI [3].

The Privacy Act 2020 (the Act) is another horizontal AI effort that regulates the use of personal information by both public and private sector agencies in New Zealand [5]. While there is no specific reference to AI in the Act, the principles and provisions in the Act are applicable to the collection, use, and disclosure of personal information by AI systems.

While the New Zealand Cyber Security Strategy 2019 does not explicitly mention AI, it outlines a clear high-level framework for the government and private sector to work hand-in-hand to improve New Zealand's cyber security [6] which of course also applies to AI systems and infrastructures.

The National Cyber Security Centre (NCSC) is a vertical AI effort that helps New Zealand's most significant public and private sector organisations to protect their information systems from advanced cyber-borne threats and to respond to incidents that have a high impact on New Zealand [6].

The New Zealand government has also developed the Government Enterprise Architecture for New Zealand (GEA-NZ) framework to provide a simple structure for the information and tools that support purposeful change across and within government organisations [7]. The GEA-NZ framework aims to

provide a consistent approach to technology and architecture across government agencies.

The AI framework aims to address New Zealand's AI global competitiveness and national security questions to some extent. The New Zealand government recognizes the potential for AI to increase the country's GDP by up to \$54 billion by 2035 [3].

In conclusion, the New Zealand government has put in place several horizontal and vertical AI efforts, including the Algorithm Charter for Aotearoa New Zealand, an ethical AI framework and action plan, a Cyber Security Strategy 2019, the Privacy Act 202

8.11 Russia

The Russian government has established a unique AI development strategy that is led not by the government nor by the private sector but by state-owned firms [1]. The country has prioritised AI, robotics, and further integration of automation and autonomy into military decision-making as part of its modernization plan for the armed forces [3]. The government has announced the allocation of 5.4 billion roubles to establish and support AI research centres, with a competitive selection process initiated by the Deputy Prime Minister Dmitry Chernyshenko [4].

Although there is no binding legislation or regulation in place, the government has published advisory guidelines on responsible AI governance, including a framework for implementing responsible AI for organisations that covers the entire process of AI system development and operations [6]. The World Economic Forum with which Russia Collaborates, has also developed a framework to guide governments in developing national strategies for AI [7].

Horizontal AI efforts include the development of AI research centres to foster innovation and collaboration across various sectors, while sectorial efforts include the integration of AI into the military decision-making process [1].

The Russian government's AI framework aims to enhance the country's AI global competitiveness and address national security concerns by prioritising the development of AI in the military and other strategic sectors [3]. The government's distrust of Russia's largest tech firm, Yandex, has side-lined the company from national AI planning [1]. It is worth noting that the government's AI framework is led by state-owned firms, and there is a lack of private sector involvement in the development of the country's AI industry.

Part IV: The History and future of AI

The History AI (1943-2023)

The Future of AI (after 2023)

Leading research Institutes on AGI and
ASI

Factors enabling the development of AGI
and ASI

9.0 The History of AI (1943 – 2023)

Artificial intelligence (AI) has been around for more than half a century and has undergone several phases of development.

9.1 1943-1956: AI Theoretical Foundations

During this period, the basic concepts of AI and neural networks were established mathematically, but no technology was available to implement them.

- In 1943, Warren McCulloch and Walter Pitts introduced the first artificial neuron theoretical model [5]
- In 1950, Alan Turing published a paper proposing the Turing Test as a measure of machine intelligence [5]
- In 1951, Christopher Strachey developed the first AI program, which played checkers on the Ferranti Mark I computer [5]
- In 1956, John McCarthy coined the term "artificial intelligence" and organised the Dartmouth Conference, which marked the birth of AI as a field of study [5]
- In 1956, Marvin Minsky and Nathaniel Rochester created the first AI program for computers to play the game of Nim [1]

9.2 1958-1997: AI progresses very slowly.

Between 1958 and 1997, AI had advanced slowly due to a combination of several factors. Firstly, early attempts at AI relied heavily on symbolic processing (Expert-Systems), which was limited in its ability to solve complex problems and required a lot of manual programming. This approach eventually reached its limits, leading to a decrease in interest and funding for AI research.

Secondly, the hardware available during this time was not powerful enough to support the computational requirements of neural based machine learning algorithms invented in 1959 and 1986. The renewed interest in AI as of 1993 is due to the creation of the World Wide Web and the availability of more powerful hardware. The internet allowed the collection of vast amounts of data, while the hardware allowed a better implementation of machine learning and neural networks invented much earlier.

- In 1958, John McCarthy and his team developed the first programming language for AI, called LISP [1]
- In 1959 Arthur Samuel makes the first practical use of the concept of machine learning algorithms though not yet using neural networks
- In 1963, Edward Feigenbaum and his team developed Dendral, the first expert system, which could solve problems in organic chemistry [3]
- In 1965, Joseph Weizenbaum developed ELIZA, an early natural language processing program that simulated a conversation [1]
- In 1969, Shakey, the first mobile robot, was created at SRI International [3]
- In 1972, Terry Winograd created SHRDLU, an AI program that could understand natural language commands and manipulate blocks on a screen [1]
- Between 1974 and 1980, first AI winter. The funding for AI research was reduced due to AI researchers' inability to deliver on the lofty promises made by the field in its early days. Furthermore, due to technological limitations and a lack of progress, public interest in AI has declined. As a result, funding for research and development in the field has decreased.
- In 1986, Geoffrey Hinton, David Rumelhart, and Ronald Williams introduced backpropagation, a machine algorithm for training artificial neural networks [6] hence generating renewed interest in the use of Neural Networks for machine learning purposes. Geoffrey Hinton is the first who used Neural Networks for machine learning purposes.

- Between 1987 and 1992, the second AI winter occurred when AI research was again unable to deliver on its promises. During this time, many investors lost faith in AI, and funding for research and development in the field declined significantly. The second AI winter was also caused by a shift in focus towards more practical applications of AI, leading to a lack of innovation in the field.
- In 1993, the World Wide Web was created, providing a platform for large-scale data sharing and distribution, which would later become crucial for AI development [3]
- Nvidia released NV1, its first Nvidia Chip, in 1995. Initially not designed for graphical processing purposes, it was quickly followed by more specialised and powerful offerings to accelerate AI.

9.3 1997-2017: AI defeats the human brain and unveils its potential.

Between 1997 and 2017, AI demonstrated its potential by defeating the human brain on various occasions. The period also witnessed the emergence of the first useful AI applications for the public, such as Apple's Siri and Google's machine translation. The deep neural networks architecture and related machine learning algorithms to train were also invented during this period. Google's publication of the "Transformer" architecture in 2017 revolutionised natural language processing by enabling parallel processing of words and better capturing long-range dependencies between them.

Back in 2002-2003, learning was something that only people could do and that computers could not do at all. However, it was clear at the time that by looking at the functioning of the brain, machine learning and neural networks would lead to the greatest progress in AI. Of all the options that existed at that time, neural networks and machine learning offered the best long-term prospects. It

is the point where the history of AI started to accelerate, thanks to the Internet and to Moore's law. As more and more computing power and data became available, researchers became able to significantly improve their research.

It was only in 2010-2012 that it became clear that supervised learning, where data must be "labeled" prior to processing by the machine learning algorithm, was the only solution in those days to solve very complex problems if used in combination with large and deep neural networks and a lot of computing power and relevant data. The computer vision capabilities unveiled by AlexNet in 2012 shocked the world and broke the record for computer vision. AlexNet is a deep convolutional neural network (CNN) architecture developed by Alex Krizhevsky, Ilya Sutskever, and Geoffrey Hinton in 2012. [The AlexNet network was able to classify images into 1000 object categories such as keyboard, mouse, pencil, and many animals,](#) something that was revolutionary in those days. In addition, AlexNet was also designed to leverage the parallel processing capabilities of Graphics Processing Units (GPUs) to accelerate training, another innovation on its own, which contributed to its success in the competition.

Yet the major problem with supervised learning in 2012 was that it didn't scale up as the volume of training data increased, so something else had to be invented.

The problem of supervised learning was finally solved in 2016, allowing large and deep neural networks to be trained on very large datasets from the Internet without the need to label any data as a prerequisite.

In 2017, Google and the University of Toronto revolutionised AI with the publication of the Transformer deep neural network architecture, which was combined with the concept of reinforcement learning, where positive or negative feedback is

given during the training (another brain similarity). This is when the whole industry realised for the first time that the prospects of Artificial Intelligence were near.

- In 1997, IBM's Deep Blue defeated chess champion Gary Kasparov [[1](#)]
- In 2000, DARPA launched the Grand Challenge, a competition for self-driving cars, which led to significant advances in autonomous vehicle technology [[9](#)]
- In 2006, Google launched Google Brain, a deep learning project that used neural networks to recognize objects in images and translate text [[2](#)]
- In 2009, IBM's Watson defeated two champions on the quiz show Jeopardy! demonstrating the potential of AI in understanding natural language and competing against human experts in knowledge-based games.
- In 2010, Google Translate introduced neural machine translation, significantly improving the accuracy of translations by using a neural network to translate phrases and sentences.
- In 2011, Apple introduced Siri, a virtual assistant that uses natural language processing and machine learning to understand and respond to voice commands, which has since become a standard feature on Apple devices.
- In 2012, Google's X Lab launched the Google Brain project, which focused on deep learning for natural language processing, image recognition, and other tasks.
- In 2012, Andrew Ng and his team at Stanford University developed a deep learning framework called "Deep Learning" that enabled more efficient training of deep neural networks, paving the way for many advances in the field.
- In 2015, AlphaGo, an AI program developed by Google's DeepMind, defeated the world champion of the game of Go, demonstrating the potential of AI to master complex strategy games and marking a significant milestone in AI research.

- In 2016, Google's DeepMind introduced AlphaGo Zero, an AI system that learned to play Go from scratch, without any human knowledge or input, achieving a higher level of performance than the previous version that used human data and input.
- In 2017, Google introduced the transformer architecture, a neural network architecture that revolutionised natural language processing by allowing for parallel processing of words and better capturing long-range dependencies between them, resulting in significant improvements in language processing tasks such as machine translation and text summarization.

9.4 2018-2021: AI reasoning capabilities gradually emerge.

In 2018, it became clear that as AI models were trained on larger and larger datasets from the Internet and as the size of their neural networks became deeper and larger, those neural networks were achieving far more than learning to do a mere statistical prediction on the next "token" or "word"!

By developing the ability to make accurate predictions, Large Language models are also explicitly building an inner representation of the world derived from the data they have been trained on. This is how Large Language Models such as GPT-1,2,3,4 gradually developed the abilities to reason logically, hold meaningful dialogues with users about their instructions, do language translations, code software, solve math problems etc. It is remarkable that these Large Language Models were not trained for those purposes, but rather that those remarkable abilities "emerged" from the very large dataset that these systems had been trained on using reinforcement and unsupervised learning algorithms.

Modern AI systems, such as GPT3.5 developed during this period, were trained in several stages. In the first phase, the AI system is massively pre-trained

on massive amounts of data from the Internet to build its internal representation of the world based on that data. The more accurate the statistical predictions, the more refined the representation of the world, and the more subtle and accurate the AI answers will be. In the subsequent training phases, the AI model is fine-tuned on a much more limited set of parameters to become more reliable, accurate, robust, trustworthy, and, most importantly, to better align with human values. This is the concept of an AI foundation model widely used in the industry.

Practically, the fine-tuning is done by another AI agent designed for that purpose, working in cooperation with human teachers and supervisors. This multi-step approach enabled the development of trustworthy AI systems aligned with human values in 2023.

During the period from 2018 to 2021, Google and Microsoft also released several key AI software development tools. Nvidia launched its 8th generation of GPUs, a significant milestone as it enabled faster processing and analysis of data, which in turn led to breakthroughs in areas such as natural language processing and computer vision. The first generation of OpenAI GPT-X language models emerged. Face recognition technology also became more common during this period. AlphaFold developed by Google/DeepMind revolutionised drug discovery and other areas of biotechnology.

AI played a critical role in the fight against the COVID-19 pandemic by speeding up the development of vaccines and helping analyse large amounts of data related to the virus. However, concerns about data privacy and security grew as the use of AI became more widespread. It was found that AI systems, particularly in areas such as recruitment, law enforcement and justice, can perpetuate existing biases and discrimination due to the biased training data used to develop them.

This led to increased scrutiny of AI applications and calls for more transparency and fairness in AI development. Governments and organisations started working on regulations and ethical guidelines to ensure responsible AI development and deployment, with growing demands for more interpretable and transparent AI systems. As such, the development and deployment of AI became a more complex and nuanced issue during this era, with an increased emphasis on responsible and ethical AI practices.

- In June 2018, Google launched AutoML [1], an artificial intelligence (AI) tool that can develop other AI systems. AutoML is designed to help businesses with limited ML expertise start building their own high-quality custom models by using advanced techniques like learning2learn and transfer learning from Google [3]. With AutoML, users can automatically build and deploy state-of-the-art machine learning models on structured data [1]. The service allows developers, including those with no machine learning expertise, to build custom image recognition models [6]. AutoML is part of Google Cloud, a suite of cloud computing services provided by Google.
- In September 2018, Nvidia released Turing GPUs, which offer real-time ray tracing and AI and machine learning enhancements. The Turing architecture features new Tensor Cores designed to accelerate AI workloads [2]. The release of the Turing GPU is significant in the history of AI as it marked a major milestone in the development of GPU architecture for AI and deep learning applications.
- In November 2018, Google/Deepmind released AlphaFold, an AI system for protein folding. It uses deep learning to predict the 3D structure of proteins, which is important for understanding their function and developing new treatments for diseases. AlphaFold is a major milestone in the history of AI and biochemistry because it has

significantly improved the accuracy of protein structure prediction. This breakthrough revolutionised drug discovery and other areas of biotechnology.

- In July 2019, Microsoft released InterpretML, an open-source toolkit designed to help developers interpret machine learning models [10]. InterpretML provides tools for model explanation and debugging, allowing developers to better understand how their models are making predictions [10]. The toolkit includes several methods for interpreting machine learning models, including global and local explanation techniques.
- In October 2019, OpenAI released GPT-2, a large-scale language model that can generate coherent and realistic human-like text. The model was trained on a massive dataset of internet text, allowing it to generate text that is difficult to distinguish from text written by a human [4]. However, due to concerns about the potential misuse of the technology, OpenAI initially limited access to the model. The company later released a smaller version of the model for public use.
- In January 2020, Google introduced an AI-powered dermatology tool that can diagnose over 3000 skin conditions with high accuracy. The tool was developed using a deep learning algorithm trained on a dataset of over 16,000 dermatology cases [1]. The tool is designed to assist dermatologists.
- In June 2020 [1] [4], OpenAI released its powerful natural language processing model, GPT-3. GPT-3 is an autoregressive language model that uses deep learning to produce human-like text. Given an initial text prompt, it can produce text that continues the prompt [1]. With 175 billion parameters, it was the largest neural

network at the time and captured the attention of mass media, researchers, and AI businesses alike [5]. OpenAI first described GPT-3 in a research paper published in May 2020 [6] and followed up with its release in June 2020.

- In August 2021 AI21 Labs released Jurassic-1 Jumbo, the largest and most sophisticated language model developed [2][4][7] at that time. The model was designed to rival OpenAI's GPT-3 and consists of two versions, J1 Jumbo and J1 Large, with a vocabulary of 250,000 lexical items, including expressions, words, and phrases [2][3][9]. The release of the Jurassic-1 Jumbo model was aimed at making language AI applications accessible to a broader audience and developers could register for beta testing to access the model [1][9].

9.5 2022-2023: AI power unleashed to the public

The breakthrough in 2017 with the invention of the neural network transformer architecture by Google and the University of Toronto gradually led five years later to the level of performance of GPT-3.5 and ChatGPT-3.5, released in November 2022 by OpenAI, and their improved versions, GPT-4 and ChatGPT-4, in March 2023.

As of 2023, the reasoning capabilities of Large Language models such as GPT-4 are truly astounding and even exceed human performance in many standardised reasoning tests. However, modern Large Language Models (LLMs) such as ChatGPT-4 are still hallucinating or making mistakes humans would not expect, such as omitting important points when summarising a text. For this reason, those systems cannot be qualified as having reached the level of Artificial General Intelligence yet. The website lifeAI.org estimated that by 2023, it would have reached approximately 42% of the performance of the human brain.

"Narrow AI systems" are trained for a specific, narrow purpose, for instance, in an industrial or operational environment, to optimise their performance. Many of the machine learning algorithms and hardware APIs are available publicly in 2023 either as open-source or as APIs.

Most Cloud service providers, including Microsoft, Amazon, Nvidia, Tencent, Alibaba, and Huawei, provide an integrated environment with powerful resources and the tools required for their customers to build their own narrow AI systems based on their own data. In other words, in 2023, building narrow and specialised AI applications is considered a common utility almost everywhere in the developed world, including China and Russia, and is no longer considered bleeding-edge technology.

It also became increasingly clear during this period that the most pressing policy concern was the potential misuse of AI for harmful purposes. This was distinct from concerns about ethical AI, which focused more on issues of fairness, transparency, and bias in AI systems.

In March 2023, a new AI model called "Alpaca" was introduced by Sandford, which offered the possibility of running an intelligent Chatbot like ChatGPT on a powerful PC though at a lower level of performance. UC Berkley followed soon after with its DollE model. The technique is called model compression and offers great promise to run small, portable, large language models on portable devices like smartphones.

Finally, in May 2023, Antropic and Mosaic.ML released large language models offering, respectively, session windows of 100000 and 128000 tokens, far exceeding the capabilities of OpenAI GPT-4 with 4000 tokens. In May 2023, Google released PALM-2 directly competing with OpenAI GPT-4 and it announced GEMINI in an effort to

compete with OpenAI (future) GPT-5. ([see lifearchitect.ai website](#))

It became apparent that Chinese companies such as Baidu, Tencent, Alibaba, and Huawei are following the lead of US companies. Though they advanced on narrow AI, they were unable to bring a similar level of innovation to the market as what OpenAI did with GPT-4.

- In March 2022, Deepmind (Google) released Chinchilla and Flamingo.

Chinchilla is a language model with 70 billion parameters optimised for efficient computation. Compared to other large language models such as GPT-3 and Jurassic-1, Chinchilla is smaller and faster to train [2][6].

Flamingo is an 80-billion-parameter Large Visual language model that combines separately pre-trained vision and language models. The Flamingo model was trained using Chinchilla, obviating the need for any extra task-specific fine-tuning [3][8]. Flamingo is multimodal; it can handle input sequences of images, videos, and text and achieves state-of-the-art results on several computer vision benchmarks with simple few-shot learning examples [10].

- In October 2022, DeepMind, a research institute owned by Google, announced the Beta version of "Sparrow" an AI-powered chatbot conceptually similar to ChatGPT.. The model, based on DeepMind's Chinchilla language model, was introduced in April 2022 [9]. Sparrow, also known as Dialogue-Prompted Chinchilla (DPC), is a fine-tuned and prompted version of DeepMind Chinchilla 70B [3]. It has been designed with high-level dialogue goals of being helpful, correct, and harmless [3]. Note that ChatGPT is based on the more advanced GPT version 3.5 and GPT-4

- In July 2022, OpenAI made DALL-E available in beta, starting the process of inviting 1 million people from their waitlist over the coming weeks [4]. The AI system, DALL-E, can create realistic images and art from a description in natural language. The beta launch of DALL-E took place on July 20, 2022 [3].
- In August 2022, Stability.AI, the world's leading open-source generative AI company, released Stable Diffusion [6, 9]. Stable Diffusion is a state-of-the-art text-to-image model that was made freely available and open-source. Numerous users downloaded and licensed the model following its launch [9]. Within a month of its release in November 2022, Stable.AI released Stable Diffusion 2.0, which quickly gained popularity, powering four of the top 10 applications on Apple's App Store [1, 8].
- In November 2022, OpenAI released two major AI models: DALL-E2 and ChatGPT.

DALL-E 2 is an AI system that can generate realistic images and art from a natural language description.[1]. It is an extension of DALL-E, a text-to-image generation program that was first introduced in January 2021 [4]. DALL-E 2 uses a process called "diffusion" to gradually build up a pattern of random dots into a realistic image based on the text input [9]. OpenAI also released a public beta of the DALL-E 2 API for developers, allowing them to embed the synthetic media generator into their apps and websites [5].

ChatGPT is an AI ChatBot [2]. It is built on top of OpenAI's modified version of GPT-3 called InstructGPT and that has been fine-tuned using both supervised and reinforcement learning techniques [2] to better align with human values. ChatGPT was launched as a prototype on November 30, 2022, and quickly gained attention for its detailed and human-like responses.[2].

- In February 2023, the company Meta unveiled a collection of open-source LLaMA Large Language Foundation Models, which range in size from 7 billion to 65 billion parameters [1]. The new META open-source AI foundation model can now run on a single "GPU," making it possible to run a local version of a ChatGPT-like application on a PC. This was effectively done by [researchers at Stanford](#) who used OpenAI GPT-3, to train the 7 billion parameters of the LLaMA open-source model to create a very light ChatBot AI model called "Alpaca". This breakthrough will make advanced AI capabilities more widely accessible to a broader audience, as the models can now be executed on individual computers instead of relying on large-scale computational resources.
- In February 2023; Microsoft introduced Kosmos-1, a multimodal large language model that can analyse images for content, solve visual puzzles, perform visual text recognition, and handle tasks that involve both language understanding and visual perception [1, 2]. Kosmos-1 can naturally handle perception-intensive tasks and natural language tasks, including visual dialogue, visual explanation, visual question answering, image captioning, simple maths equations, OCR, and zero-shot image classification with descriptions [6]. Researchers believe that multimodal AI, which integrates different modes of input such as text, audio, images, and video, is a key step toward building artificial general intelligence [8].
- In March 2023, Google, in collaboration with the Technical University of Berlin, unveiled PaLM-E, a multimodal embodied (ie, imed for robotic purposes) visual-language model (VLM) with 562 billion parameters [2, 3]. PaLM-E integrates AI-powered vision and language to enable autonomous robotic control, allowing robots to

perform a wide range of tasks based on human voice commands [5]. PaLM-E combines Google's massive PaLM language model with ViT-22B, the largest vision transformer to date, to understand and generate language, understand images, and use both together for complex robot commands [8].

- In March 2023, the company Midjourney released the Beta version of Midjourney V5 that is an image to text generator similar to OpenAI DAL-E. The new model is said to generate much more realistic and detailed images than OpenAI competitor DAL-E2, but also requires more precise prompts. Midjourney V5 quickly gained a large following due to its distinct style and being publicly available before many other AI image synthesis models.
- In March 2023, OpenAI released two AI models GPT-4 and Chat GPT-4 [1]. ChatGPT-4 is more accurate than ChatGPT-3.5 and it can write code in all major programming languages. OpenAI announced that ChatGPT-4 can now read, analyse, or generate up to 25,000 words of text (48 pages) and is significantly smarter than ChatGPT-3.5 [2]. ChatGPT-4 is multimodal, meaning it can operate within multiple kinds of input, such as images, and sound, in addition to text [3, 4]. According to OpenAI, ChatGPT-4 is 40 percent more likely to provide correct answers than its predecessor [9].

GPT-4 released [2, 4, 7], is a massive multimodal language model, an upgrade from its predecessor GPT-3.5, and is now available in Bing and ChatGPT-4 [4]. The model is capable of understanding images and processing image prompts, in addition to text, making it multimodal [5].

- Former OpenAI employees who co-founded Anthropic launched Claude, a competing AI chatbot to ChatGPT, in March 2023 [2]. Google is among the

investors in Anthropic, having pledged \$300 million for a 10% stake in the startup [1]. Claude is based on Anthropic's research into training helpful, honest, and harmless AI systems [3]. The company claims that thanks to Claude's "Constitutional Architecture", it can do everything that ChatGPT can, but it avoids harmful outputs.

- In May 2023, Anthropic upgraded the size of the session window of Claude+ to 10000 tokens, clearly outperforming OpenAI with its session windows of 4000 tokens, and announced plans to upgrade to 32000 tokens by the end of 2023.
- In May 2023, Mosaic.ML released the first set of open source "MPT" AI models licensed to be used for commercial purposes. One of the models was fine-tuned to work with session windows of 128,000 tokens, exceeding both Anthropic and OpenAI.
- In May 2023, Google released PALM-2, a Multimodal Large Language Model with a level of performance like GPT-4 of OpenAI and it announced Gemini its next generation Large Language Model destined to compete with OpenAI GPT-5
- In May 2023, Stanford, CMU, [UC-Berkley and UC-San Diego released Vicuna-13B](#), a 13 Billion parameter open-source model.

10.0 The future of AI (after 2023)

10.1 Before 2030: Artificial General Intelligence (AGI) reached.

The end of 2021 and 2022 marked the emergence of Large Language Models (LLMs) and Generative AI models, such as GPT-3, which were trained on vast datasets.

In the beginning of 2023, we witnessed the emergence of Multimodal Large Language Models and Generative AI models such as GPT-4, Kosmos-1, and Palm-E. These models were trained on images in

addition to text, and sometimes audio as well. While these models have shown remarkable performance in various tasks their capabilities still have limitations in specific areas.

End of 2023 or beginning of 2024, all eyes will be on Anthropic (Claude+) and Google (Gemini) that announced multimodal Large Language models 10 times more powerful than GPT-4 thanks to the contribution of Google. OpenAI is also expected to continue to improve GPT-4 up to the point where it will be able to process very large documents of up to 40-50 pages, as it initially announced in March 2023. This will be a breakthrough that will dramatically boost the usefulness of GPT-4 and ChatGPT-4 compared to mid-2023. GPT-5 is expected in 2024 and is supposed to further increase the multimodality by accepting video as input as well as generating text-video-images as output.

By 2026, it is anticipated that these Generative AI MLLMs will have the capacity to accept voice, video, and text as input and generate output in any of those formats. It will be possible to instruct these models to perform tasks by providing written or visual instructions, simply by talking to them, or by combining multiple modes.

By 2026, the metaverse, consisting of numerous interconnected virtual worlds, is expected to have advanced significantly, with generative AI playing a vital role in providing personalised and immersive experiences. The development and widespread adoption of the metaverse will largely depend on the availability of high-quality, affordable headsets and other related technologies, which are currently under active development in 2023. A critical policy concern associated with the metaverse will be ensuring user privacy. Data leaks and unauthorised access to personal information could pose significant risks to users, as malicious actors might exploit this information more easily within the metaverse.

On the hardware side, NVIDIA is expected to continue dominating the scene with the generation of increasingly powerful acceleration boards and computing infrastructure hardware that consumes less energy. Research in neuromorphic computing chips and spiking neural networks led by IBM, Intel, and many academic research centers worldwide should have significantly progressed, reducing the energy consumption of AI systems and accelerating their speed. This may, in turn, help reduce the dependency of the AI industry on specialised AI hardware produced by Nvidia.

Predicting the exact timeline for achieving Artificial General Intelligence (AGI) is challenging. Depending on the definition of AGI used, [the majority of experts think it might be possible to achieve AGI by 2030.](#) From a pure cognitive perspective, AGI will be reached earlier than 2030 by generative AI and Large Language Models, however, a robot with a fully embodied AGI allowing it to interact with the world in the same way humans do will take much longer.

Progress in embodied AI, which involves integrating AI models into robots equipped with sensors and motors for physical world interaction, is expected to be slower. This is due to the challenges of training embodied AI models with sufficiently high volumes of data. Google (with its Palm-E model), Tesla, and Boston Dynamics are anticipated to lead the way in embodied AI through their respective social robot projects.

10.2 After 2030: Artificial Super Intelligence. (ASI) reached.

Based on ongoing research in industry and academia in 2023, it is very likely that advancements in AI will have significantly contributed to fission energy production, new composite or metamaterial discovery, social robots, advances, and genome therapy.

AI can potentially assist in the development of superconductor materials to help stabilise entangled qubits in quantum computers at temperatures higher than a few millikelvin degrees [10] as is still the case in 2023. If this happens, it will become possible to use classical C-MOS technology to build chips with a very high number of qubits stable at a temperature in the order of 1- or 2-degree Kelvin. This development would in turn lead to the development of supercomputers that (all other things being equal) would be on the order of 1000000 faster than the current generation of silicon-based supercomputers in 20233, hence opening the way to ASI.

In that context, it is probable that the first truly available social robots will be commercialised on a massive scale around 2030 at an affordable price. Though AI makes rapid progress, there are still many non-AI issues that remain to be solved, and so far, only basic prototypes of robots have been demonstrated.

Progress in genomics might also considerably slow down the ageing process [10].

As for human brain-machine interfaces, such as those developed by Tesla-Neuralink, it remains uncertain whether this technology will be mature enough by 2030, given the non-AI-related problems that remain to be solved in 2023, to be commercialised on a mass scale. Recent progress realised by a university in Texas in May 2023 is very encouraging. [An AI based on GPT-1 was able to decode thoughts in text format after having been trained on magnetic resonance images.](#) Though the results obtained are very impressive, as the technology is non-invasive, it is not portable either, so its use is limited.

AI-controlled plasma fusion reactors will probably have reached a total energy gain greater than one by that time, [based on the declarations of some of the most advanced startups in the field of fusion](#)

[in 2023](#). This was not achieved until mid-2023, despite some of the misleading statements in the press. If this happens, it will significantly contribute to the start of the ASI area.

11.0 Leading research Institutes on AGI and ASI

11.1 Business

There is no doubt that the most advanced research institutes in the field of generative Large Language Models or Generative Multi Modal Large Language models come from the US private sector, including Google, Amazon, Nvidia, Meta, Tesla, Anthropic, Microsoft, and OpenAI.

Google-Brain has made very significant contributions that have been exploited by other players in the AI industry, including the competing OpenAI. Google DeepMind has significantly contributed to the use of AI in the sciences. In April 2023, the two entities merged.

Mid-2023, it is OpenAI and Google that lead the development of the most advanced generative Multimodal Large Language Models that feature the most advanced reasoning capabilities available to the public, such as PALML-2, GPT-4, ChatGPT-4 or DALL-E2 for instance. Google recently took a major stake in Anthropic, an AI lab that spun off from the work of competing OpenAI.

As of mid-2023, only a few individuals possess the experience and technical expertise required to develop reliable and powerful large language models rivaling GPT-4 in performance. These experts command exceptionally high salaries, far exceeding those they could earn in any European or Asian country. They are all based in the US, working for leading tech companies and in close cooperation with cloud service providers such as Microsoft, Google, and Amazon.

Companies like MidJourney and Stability.AI conducts advanced research too, but they do not lead AGI

research, as they specialise in narrower use cases of generative AI.

In China, despite a booming cloud computing sector led by Tencent, Alibaba, and Huawei, government policies and interference hinder innovation in generative Large Language Models. Consequently, there is not yet any Generative Large Language Model (LLM) with a performance level equivalent to OpenAI's GPT-4 or ChatGPT-4. Although China surpassed the US in the number of AI-related scientific publications before 2022, this is no longer the case. Furthermore, Chinese publications on AI are generally of lower quality and scientific interest compared to those from the US or even Europe. Despite what the official Chinese propaganda insinuates, China is clearly losing ground compared to the US on AI. China seems to be leading mainly in the deployment of "Narrow AI" applications but does not seem to be very advanced in AGI research.

Also, the growing geopolitical tensions between the US and China & Russia will likely result in more export controls and foreign direct investment limitation of AI western technology in the coming years. Research cooperation between research institutes in democratic and non-democratic countries is also likely to slow down significantly. Those factors are expected to allow the US and its allies to stay ahead of non-democratic countries such as China and Russia in the development of artificial general intelligence and artificial superintelligence.

In Europe, the absence of a robust cloud computing industry has made it difficult for competitors to OpenAI or Google to emerge. It is unlikely that this situation will change soon, given the technological advances of the US research institutes mentioned earlier.

11.2 Academia

Academic institutions all over the world are actively researching faster and more energy-efficient machine learning algorithms and hardware, including Tsinghua University in China, TU-Berlin in Germany, and the Moscow Institute of Physics and Technology in Russia.

Academia is also heavily involved in researching the use of AI for optimization purposes. For example, AI can be used to save chip design time by identifying optimal electrical connection routing layouts (TILOS Institute) and reducing the number of steps needed to produce complex lithography masks for manufacturing cutting-edge chips smaller than 5nm (NVIDIA AI Lithography Computing). Additionally, there is active academic research focused on developing tools and methods for automatically [testing the trustworthiness and robustness of AI systems.](#)

AI plays a significant role in the study of the brain, helping to decode brain signals, as well as in material science engineering and [nuclear fusion energy production.](#) Active academic research is also dedicated to optimally compressing large language models using methods such as pruning or deletion, with the aim of decreasing computing resource requirements while minimising the loss of accuracy.

However, the computing resources necessary for the research and development of Multimodal Large Language Models (LLMs) like GPT-4 are currently inaccessible to most academic institutions, effectively limiting their involvement in this area. With the establishment of a National AI Research Resource in the US in accordance with the roadmap made public by the National AI Resource Research Task Force in January 2023, this situation might change. In Europe, several similar initiatives exist at both the national and EU levels through programs like Horizon Europe and Digital Europe.

As machine learning algorithms and hardware become more energy-efficient in the long term, the

landscape could evolve, creating new opportunities for academic institutions to contribute to the advancement of LLMs.

12.0 Factors enabling the development of AGI and ASI

The development of very large language-generative multimodal models like GPT-4 can be seen as a step towards the long-term goals of Artificial General Intelligence (AGI) and Artificial Superintelligence (ASI).

These advanced forms of AI require a combination of key technologies, resources, and industrial policy considerations, as well as careful consideration of ethical and regulatory aspects.

12.1 Access to powerful computing resources

AI accelerator chips such as Graphics Processing Units (GPUs) produced by Google or Tensor Processor Units (TPUs) produced by Nvidia are required. In 2023, many other companies produced similar hardware in the US such as Tesla, Cerebras, IBM, Intel and Apple.

In China and Russia Huawei, Baidu and Yadros have developed similar technologies.

It is however Nvidia in the US that develops the most powerful hardware chips commercially available and enabling the development of AGI and ASI in 2023

High-performance computing (HPC) infrastructure: Large-scale compute clusters and HPC systems are needed to support both the training and inference requirements for AGI and ASI and in 2023, it is Nvidia, Amazon, and Microsoft that offer the most powerful cloud computing services in the world. In China and Russia, and in some countries of the global south, companies like Huawei, Tencent and Alibaba offer similar cloud services.

Some of the most promising hardware chip technologies for the development of AGI and ASI are

based on spiking neural network architectures that mimic (more) closely the sparse functioning of the brain and hence are more thermodynamically efficient than classical deep neural networks, but they are also less accurate. Many universities, including UCSD, Berkeley, Tsinghua, and US companies like Intel, IBM, and Huawei, are actively researching this domain.

Current silicon manufacturing chip technologies will also continue to progress significantly, though their production costs will also continue to explode as Moore's law continues to slow down due to the limits of physics.

Photonic chips are also likely to become more common. In 2023, significant progress was realised with the production of the first fully integrated photon-based quantum computing chip.

The three technologies above will likely coexist by 2030 and contribute to the development of AGI and ASI.

12.2 Access to vast and cheap sources of electrical energy

As explained earlier, the initial stages of AGI and ASI development will likely require access to powerful supercomputers that consume large amounts of electrical energy. Countries with access to affordable energy sources to power these supercomputers will have an advantage over those that do not.

Nuclear fission and potentially nuclear fusion, if commercially viable by 2030, as predicted by some start-ups, could play a critical role in providing the necessary energy for AGI and ASI development.

[Microsoft has invested significantly in Helion, a fusion startup that claims to connect the first fusion reactor to the grid by 2028.](#)

Over time, as hardware and software become more thermodynamically efficient, energy constraints may become less of a critical factor for countries that do not have access to cheap electrical energy, provided they have access to the necessary technologies.

From a climate change perspective, although there is no specific prediction of the carbon footprint related to the development of AGI and ASI, it is likely that carbon emissions will increase significantly, at least until AGI and ASI become more energy efficient or contribute directly to climate mitigation actions in a highly effective manner.

12.3 Access to software development environments

In 2023, the development of Deep Neural Networks relies on open-source frameworks such as TensorFlow (developed by Google) and PyTorch (developed by Meta). These frameworks enable researchers and developers to efficiently build, train, and deploy large-scale models.

NVIDIA tools are designed to be compatible with both TensorFlow and PyTorch, allowing developers to leverage the programming expertise gained with these open-source tools within the NVIDIA development environment.

Although these development tools have made significant strides, they are expected to continue evolving in the coming years as advancements are made towards ASI and AGI.

In 2023, tools for training Deep Neural Networks, such as Horovod and torch.distributed, are essential for supporting distributed and parallel training. These tools facilitate model training

across multiple devices and/or nodes, accelerating the training process. As the development of ASI and AGI progresses, these training tools are likely to become increasingly important and continue to evolve.

Additionally, AutoML and Neural Architecture Search (NAS) tools, which assist researchers in discovering optimal neural network model architectures and hyperparameters, play a crucial role. By automating parts of the AI model development process, these tools help reduce manual intervention, leading to more efficient and effective model creation.

As development environments continue to improve, increasing the productivity of developers, access to the latest generation of tools for developing powerful neural networks becomes critical and essential on the path to ASI and AGI. Ongoing advancements in these tools will help accelerate research and development, enabling the AI community to overcome challenges and unlock the potential of AGI and ASI in the future.

12.4 Access to large volume of quality data

Access to large volumes of quality data is essential for the development of Artificial Superintelligence (ASI) and Artificial General Intelligence (AGI) for several key reasons:

1. Data-driven discovery:

ASI and AGI systems aim to surpass human intelligence and solve complex problems across various domains. Access to large volumes of quality data enables the emergence of key highly desirable abilities, allowing these systems to uncover novel insights, generate hypotheses, and make predictions that might not be possible for human experts.

2. Robustness and generalisation:

Larger and more diverse datasets help models become more robust and better at generalising to

new situations. This is particularly important for AGI and ASI, as these systems are expected to perform a wide range of tasks and adapt to new environments with minimal human intervention.

3. Transfer learning and pre-training:

Access to large datasets enables the pre-training of models, which can then be fine-tuned for specific tasks using smaller, task-specific datasets. This transfer learning approach has proven to be effective in a variety of domains, such as computer vision and natural language processing, leading to improved model performance and faster training times.

4. Benchmarking and evaluation:

Having access to large, high-quality datasets is essential for evaluating and comparing the performance of different AI algorithms and models. This allows researchers to identify the most effective techniques and drive progress in the field.

In summary, access to large volumes of quality data is crucial for the development of ASI and AGI because it facilitates more effective learning, improves generalisation and robustness, enables transfer learning, supports evaluation and benchmarking, and promotes data-driven discovery. Although, in 2023, some research aims to reduce the volume of data needed for AI systems to learn, drawing inspiration from the way the human brain identifies patterns and consolidates data, it is expected that high-quality datasets will remain essential for AGI and ASI development at least until 2030.

Given the strategic importance of AGI and ASI for the future of humankind, particularly in democratic countries, it raises the question of whether there is a need to re-evaluate the balance between privacy and intellectual property protection at the individual level and the need for massive data collection to enable the benefits of ASI and AGI at broader societal levels. It is expected that the use of privacy and intellectual property rights

protection technologies will gradually develop to find an adequate compromise between such rights and the common good.

Non-democratic countries, in which governments unilaterally dictate data collection policies for ASI and AGI purposes, will not face this democratic dilemma between individual rights and societal benefits. However, such countries may still face challenges in their ASI and AGI development if they cannot access large volumes of quality data or if they lose access to it due to growing geopolitical tensions with democratic countries.

The world is not running out yet of quality data, also as users use LLMs systems like GPT-4 they are generating quality data that can also be used to improve AI systems.

In conclusion, ensuring access to valuable data resources is a crucial aspect of AGI and ASI strategy and progress in the coming years.

12.5 Access to human Capital

Developers of AI systems typically begin their journey with a foundation in STEM subjects and proficiency in programming languages like C, C++, Python, and Java. Aspiring machine learning developers then acquire knowledge of PyTorch and TensorFlow during their university education. Next, they master the NVIDIA development environment while working in the industry (could also be one of other vendors such as Cerebras for instance). Finally, over time, they learn to engineer complex and robust AI systems built with various AI foundation models and using different learning and training techniques adapted to the data processed and the objectives targeted.

To date, the process of learning to develop Very Large Language Generative models such as GPT-4 is still very experimental and relies heavily on trial and error and technical intuition. Top-level technological expertise is acquired on the job

after years of experience in research centres and in the industry.

On the path to AGI and ASI, it is important that this valuable knowledge, developed by the industry, be shared and consolidated with academia to be incorporated into educational programs. This collaboration allows us to produce more experts who can become productive more quickly. Countries that master this process or actively cooperate with those that do will gain a significant advantage on their journey to AGI and ASI.

12.6 Access to venture Capital (VC)

Access to capital, especially venture capital (VC), is crucial for countries that wish to promote the development of AGI and ASI for several reasons:

1. High upfront costs:

Developing AGI and ASI technologies often requires significant investments in research and development, infrastructure, and talent. These initial costs can be prohibitive for many startups or research organisations. Venture capital can provide the necessary funding to overcome these barriers and allow innovative projects to move forward.

2. Encouraging innovation and risk-taking:

The development of AGI and ASI is inherently risky and uncertain, with many projects facing technical challenges and the possibility of failure. Venture capital firms are designed to invest in high-risk, high-reward projects, providing financial support to innovative ideas that might not receive funding through traditional channels.

3. Attracting and retaining talent:

Talented researchers and engineers are essential for advancing AGI and ASI technologies. Access to venture capital allows companies and research

organisations to offer competitive salaries, benefits, and opportunities for growth, helping to attract and retain the best minds in the field.

4. Accelerating development:

With adequate funding from venture capital, companies can scale their operations, hire more staff, and invest in the resources required to speed up the development of AGI and ASI technologies. This acceleration can help them stay competitive and contribute to the overall progress of the field.

5. Fostering a robust ecosystem:

A healthy venture capital landscape can help create a thriving ecosystem of startups, research organisations, and established companies working on AGI and ASI. This ecosystem encourages collaboration, knowledge sharing, and healthy competition, driving the field forward more rapidly than if organisations were working in isolation.

In summary, access to venture capital plays a vital role in promoting the development of AGI and ASI in a country by providing necessary funding, encouraging innovation and risk-taking, attracting top talent, accelerating development, and fostering a robust ecosystem.

12.7 Measures to increase the flow of venture capital (VC)

Possible measures to increase the flow of venture capital (VC) to fund the development of AGI and ASI:

1. Create a favourable regulatory environment:

Governments can establish clear and supportive regulations that encourage investment in AGI and ASI technologies. Reducing bureaucratic hurdles, simplifying tax regulations, and providing incentives for investors can help attract more VC funding to the sector.

2. Establish and promote research hubs:

Governments can invest in research hubs and clusters by creating dedicated spaces for AGI and ASI development, such as technology parks or innovation centres. These hubs can attract talent, startups, and VC firms, fostering a vibrant ecosystem that encourages investment.

3. Offer grants and financial incentives:

Governments can provide grants, subsidies, or tax breaks to startups and research organisations working on AGI and ASI. Such financial incentives can lower the risks associated with investing in these projects and make them more attractive to venture capitalists.

4. Foster collaboration between academia and industry:

Encouraging collaboration between universities, research institutions, and private companies can help create a pipeline of innovative projects that are attractive to VC firms. Governments can support this collaboration through funding, matchmaking events, and information-sharing platforms.

5. Invest in education and talent development:

Governments can invest in education and workforce development programs to cultivate a skilled talent pool in AGI and ASI research. A strong talent base can attract more VC investment by ensuring a steady supply of skilled professionals to staff new ventures.

6. Encourage international cooperation and investment:

Governments can facilitate cross-border collaboration and investment by entering into partnerships or agreements with other countries. This can help attract foreign VC firms and promote the exchange of ideas and resources in the field of AGI and ASI.

7. Establish public investment funds:

Governments can create public investment funds dedicated to supporting AGI and ASI startups and research organisations. These funds can act as catalysts, attracting additional private VC investments by signalling confidence in the sector's growth potential.

8. Showcase success stories and promote local innovation:

Governments can actively promote local AGI and ASI success stories, showcasing their potential and encouraging more VC firms to invest in the sector. This can be done through events, conferences, and media campaigns.

By implementing more of these measures, countries can create a more attractive environment for VC investments in AGI and ASI, driving the development of these critical technologies.

12.8 Adopting an adequate regulatory Framework

To accelerate and stimulate the development of AGI and ASI, countries can adopt a regulatory framework that actively promotes innovation, collaboration, and investment. Here are some key components of such a framework:

1. Foster a supportive legal environment:

Simplify and streamline business registration and incorporation processes, making it easier for startups and entrepreneurs to set up and operate businesses focused on AGI and ASI.

2. Encourage research and development:

Offer tax incentives, grants, and other financial support to encourage both private companies and public research institutions to invest in AGI and ASI research and development.

3. Promote collaboration:

Encourage partnerships between academia, industry, and government by providing funding for collaborative research projects and establishing joint research centres focused on AGI and ASI.

4. Protect intellectual property in the field of AGI and ASI:

Implement clear and effective intellectual property laws that protect the innovations and discoveries in AGI and ASI while also promoting the sharing of knowledge and research findings.

5. Address ethical considerations through non-compulsory guidelines:

- Develop guidelines and regulations that address the ethical implications of AGI and ASI, ensuring that these technologies are developed and used responsibly and transparently.
6. Invest in education and workforce development:
Create policies that support education and workforce development in AI-related fields, such as providing funding for AI-focused university programs, research fellowships, and vocational training programs.
 7. Encourage international cooperation:
Facilitate international collaboration by easing restrictions on cross-border data sharing, forging partnerships with other countries, and participating in global AI research initiatives.
 8. Foster access to data:
Promote open data initiatives, data sharing agreements, and the development of privacy-preserving data sharing technologies to help researchers and companies access the high-quality data needed for AGI and ASI development.
 9. Establish public investment funds:
Create public investment funds dedicated to supporting AGI and ASI startups and research organisations, attracting private investments, and signalling confidence in the sector's growth potential.
 10. Ensure safety and accountability:
Develop safety guidelines and establish mechanisms for accountability to address potential risks and unintended consequences associated with the development and deployment of AGI and ASI technologies.

By adopting a regulatory framework that addresses more of these aspects, countries can create an environment that is even more conducive to innovation, collaboration, and investment in AGI and ASI, ultimately further accelerating the development of these advanced technologies.

12.9 Adopting an adequate industrial policy.

An industrial policy supporting the development of AGI and ASI in a democratic country should focus on fostering innovation, collaboration, responsible development, and equitable distribution of benefits. Here are some key elements to consider and potential pitfalls to avoid:

Priorities to focus on:

1. Research and development:

Encourage and fund fundamental research in AI, covering a wide range of disciplines like computer science, cognitive science, neuroscience, and ethics. Support collaborations between academia, industry, and government organisations to drive innovation.

2. Education and workforce development:

Invest in education and training programs to develop a skilled workforce capable of contributing to AGI and ASI research and applications. This includes promoting STEM education, interdisciplinary studies, and lifelong learning opportunities.

3. Infrastructure and resources:

Provide the necessary infrastructure and resources, such as high-performance computing facilities and data repositories, to support AI research and development.

4. Public-private partnerships:

Foster collaboration between the public and private sectors to share knowledge, resources, and funding for AGI and ASI development.

5. Ethical and safety guidelines:

Develop and enforce ethical guidelines and safety standards to ensure the responsible development and deployment of AGI and ASI technologies.

6. Global cooperation:

Engage in international collaborations and contribute to global governance mechanisms to address the challenges and opportunities of AGI and ASI development.

Pitfalls to avoid:

1. Short-term focus:

Prioritising short-term economic gains and applications over long-term fundamental research can hinder progress towards AGI and ASI.

2. Inequitable distribution of benefits:

Failing to address the potential economic and social disparities resulting from AGI and ASI development can lead to increased inequality and social unrest.

3. Insufficient attention to AI safety and ethics:

Neglecting safety and ethical considerations in the development of AGI and ASI can result in unintended consequences and potential harm to individuals, communities, and society at large.

4. Overregulation:

Excessive regulation can stifle innovation and slow down progress in AI research. Striking a balance between promoting innovation and ensuring responsible development is essential.

5. Isolationism:

Avoiding international collaboration or attempting to dominate the AI landscape can hinder progress and exacerbate global competition, increasing the risk of an AI arms race.

12.10 Maintaining a free and democratic society.

Democratic countries not faced with the challenges often associated with non-democratic nations, such as government intervention in information control and propaganda, are much better positioned to lead the development of AGI and ASI than non-democratic ones for several reasons:

1. Open exchange of ideas and information:

Democratic environments facilitate the free flow of ideas and information, promoting innovation and creativity in AGI and ASI research. Access to a wide range of global perspectives and expertise also allows researchers to collaborate more effectively and develop cutting-edge AI technologies.

2. Stronger emphasis on ethics and human rights:

Democratic countries generally prioritise the protection of human rights, freedom of expression, and privacy, which aligns with the ethical development and deployment of AGI and ASI. These values can help ensure that AGI and ASI technologies are designed to be transparent, accountable, and fair, benefiting a wide range of stakeholders.

3. Stable governance and rule of law:

Democratic countries often have more stable governance structures and a stronger rule of law, which can create an environment conducive to long-term AGI and ASI research and development. Predictable legal frameworks and regulatory environments reduce uncertainty for researchers, investors, and other stakeholders in the AI ecosystem.

4. Greater access to resources and funding:

Democratic countries typically have more transparent and efficient resource allocation mechanisms, which can lead to better funding and support for AGI and ASI research. The availability of public and private investment can fuel the growth of the AI sector, driving advancements in AGI and ASI development.

5. International collaboration and influence:

Democratic countries are often well-positioned to engage in international collaborations, forming partnerships with other democratic nations, research institutions, and private companies to share knowledge and resources. By actively participating in global AI governance efforts, democratic countries can help shape international norms and standards that guide responsible AGI and ASI development.

In summary, democratic countries have several advantages in leading AGI and ASI development. These advantages include the open exchange of ideas, a stronger focus on ethics and human rights, stable governance, greater access to resources, and the ability to engage in international collaborations. These factors, together with a low level of corruption, the absence of strong ideological dogma, the respect of the rule of law, and the limited influence of organised crime, create an environment that fosters responsible and innovative AGI and ASI research.

Part V: Geopolitical considerations

AGU and National Security & strategic
autonomy

AGI Societal impacts

AGI and theories of trade

AGI and Foreign Policy

AGI and theories of International
Relations

13.0 AGI and national security & strategic autonomy

13.1 National Security

The discovery of emerging capabilities in sufficiently large and deep neural networks trained on vast datasets can be likened to the discovery of fire. In the long term, nothing can prevent the rise of AGI in countries where substantial amounts of data and computing resources are available.

It is hypothesised that, for the reasons mentioned above, the development of AGI (led by the US) will happen by 2030 at the latest, regardless of evolving geopolitical tensions between, on the one hand, allies of the US (NATO countries, Japan, South Korea, etc.) and, on the other hand, allies of China (Russia, Iran, North Korea, etc.).

However, in the event of a significant escalation of geopolitical tensions or a declaration of war between these two blocs, the priority will shift towards military and security applications of AGI. On the one hand, AGI may become available more quickly; on the other hand, AGI will not be developed with the necessary purposes, safeguards, and international cooperation required to benefit the whole society globally.

In such a scenario, military use cases of AGI would take precedence over civilian ones, leading, at least temporarily, to an unequal distribution of AGI benefits within society and globally. For national security reasons, the US and its allies will most likely need to keep the most sensitive aspects of AGI development confidential for as long as possible.

13.2 International research cooperation

Export controls and foreign direct investment (FDI) limitations may not be sufficient for the US and its allies to stay ahead of China and Russia in developing their artificial general intelligence (AGI) capabilities, even in the event of an armed

conflict with the US and its allies. This is due to several reasons, which underscore the necessity of ceasing all cooperation on AGI between, on the one hand, the US and its allies and, on the other, China and its allies (including Russia).

1. **Indigenous innovation:** Both China and Russia have demonstrated their ability to develop advanced technologies independently. By investing heavily in research and development, these countries could continue to make progress in AGI despite export controls and FDI limitations.
2. **Technology diffusion:** The rapid diffusion of technology and knowledge in today's interconnected world means that cutting-edge advancements can be accessed and shared more easily than ever before. This could enable China and Russia to obtain AGI-enabling technologies through various channels, bypassing export controls and FDI restrictions.
3. **Alternative partnerships:** China and Russia may seek partnerships with other countries or entities that share their strategic interests or possess the required technological capabilities. This would allow them to continue their AGI research and development efforts, even in the absence of cooperation with the US and its allies.
4. **Espionage and intellectual property theft:** Both China and Russia have been accused of engaging in industrial and cyber espionage to acquire advanced technologies. This could potentially help them obtain AGI-related knowledge and expertise, circumventing export controls and FDI limitations.
5. **Dual-use technologies:** Many AGI-enabling technologies have dual-use potential, meaning they can be employed for both civilian and military purposes. This makes it challenging to impose effective export controls and FDI limitations without inadvertently harming non-military applications and innovation.

To adequately address these concerns and allow the US and its allies to stay ahead of China and Russia's progress in developing AGI capabilities, it would be necessary to not only enforce stringent export controls and FDI limitations but also halt all cooperation in the research and development of all technologies potentially enabling the development of AGI. This would require a concerted effort to prevent the transfer of knowledge, technology, and expertise related to AGI while also encouraging collaboration among like-minded countries to stay ahead in the AGI race.

13.3 Open Strategic autonomy

In a cold war scenario with a high risk of armed conflict, the US and its allies should develop a well-defined concept of open strategic autonomy in the field of AGI.

Each country within the alliance should determine, based on its own national security strategy and those of its allied democratic partners, the elements of the AGI global value chain for which it wishes to achieve some autonomy and those for which it opts to rely solely on trustworthy partners sharing the same democratic values and market economy visions.

For instance, although the UK is a close economic and military partner of the US, it has decided to build its own "BritGPT" with the goal of offering a level of performance equivalent to ChatGPT-4. Similar projects may be necessary in some European Union member states to support the open strategic autonomy goals of the EU.

By reducing excessive dependence on US AGI technology in the EU, the internal market for AGI services will become more open and diverse, ultimately also benefiting European consumers. This diversification would promote competition and innovation, leading to the development of better AGI products and services tailored to the needs and

preferences of European users. Additionally, it would foster resilience in the face of potential geopolitical tensions or disruptions in international supply chains.

14.0 AGI Societal impacts

14.1 Scenario 1: De-escalation of geopolitical tensions with China and Russia

This is unfortunately the least likely scenario as of mid-2023, based on the geopolitical tensions between allies of the US and allies of China and Russia, the potential benefits for humanity would be significant if global peace and de-escalation of geopolitical tensions were to prevail. In such a hypothetical scenario, the societal impacts of Artificial General Intelligence (AGI) becoming available by 2030 could be both positive and negative, with outcomes largely dependent on the level of collaboration among governments, organisations, and individuals to ensure responsible AGI development and deployment.

On the positive side, global peace and reduced geopolitical tensions would enable increased international cooperation, joint research, and knowledge sharing in the field of AGI. This could lead to faster AGI advancements and widespread adoption, resulting in benefits such as improved healthcare, reduced poverty, increased access to education, and more efficient resource management. Moreover, with global cooperation, the creation of international standards and regulations would be more feasible, helping to address potential risks and ethical concerns surrounding AGI. This could foster responsible AGI deployment, ensuring that its benefits are distributed equitably and that potential harms are minimised.

On the negative side, the rapid development and deployment of AGI could exacerbate existing inequalities or create new ones if not properly

managed. There might be risks of job displacement, privacy concerns, and other unintended consequences. However, in a scenario of global peace and de-escalation, collaborative efforts to address these challenges would be more effective, as countries would be more inclined to work together to find solutions that benefit all.

While the mid-2023 scenario of global peace and de-escalation of geopolitical tensions between US allies and China and Russia allies appears less likely, it holds the potential to bring the most benefits to humanity. The key to realising this potential lies in fostering international collaboration and ensuring the responsible development and deployment of AGI for the greater good.

Opportunities:

1. **Economic Growth:** AGI would significantly boost productivity across various sectors, leading to economic growth and prosperity.
2. **Scientific Advancements:** AGI would accelerate research and development across numerous scientific fields, such as medicine, energy, and space exploration.
3. **Climate Change Mitigation:** AGI would optimise energy systems, contribute to the development of innovative green technologies, and help monitor and manage environmental systems more effectively.
4. **Healthcare:** AGI would revolutionise healthcare by developing personalised treatment plans, discovering new drugs, and improving diagnostics.
5. **Education:** AI-powered adaptive learning platforms would help tailor education to individual students' needs, improving access to quality education and reducing disparities.
6. **Poverty Reduction:** By improving resource allocation and efficiency, AGI would contribute to poverty reduction by creating new job opportunities, enhancing social welfare

programs, and fostering economic development in less developed regions.

Challenges:

1. **Job displacement:** The rapid adoption of AGI would lead to job displacement in certain sectors including white collar intellectual ones. New jobs would be created but not as fast as old jobs are destroyed.
2. **Increased inequalities:** Without proper policies in place, the benefits of AGI may not be equally distributed, potentially exacerbating existing socioeconomic inequalities.
3. **Increased Privacy risks:** As AGI systems can occasionally accidentally hallucinate, they could harm people's reputation or disseminate wrong information about them. Individual users who consent to give their personal data for AI training purpose should accept this risk.
4. **Ethical Concerns:** The lack of a proper regulatory framework to address inherent ethical challenges related to AGI such as biases, fairness, accountability, and alignment with human values would result in fundamental rights/human rights issues for users, and legal uncertainty for AGI service providers.
5. **Increased security risks:** AGI could introduce new security risks, including AI-generated deep fakes, autonomous cyberattacks, and potential misuse by malicious actors.
6. **Inadequate protection of intellectual property rights:** AGI should not violate the intellectual property rights of authors and creators and that should have a mechanism to see a proper reward if AGI systems are trained based on their works without their consent.
7. **Undesired on-line disintermediation:** As AGI systems will offer very convenient human interfaces to complex on-line non-AI systems (such as financial advisory ones, for instance), there is a significant risk that those systems will be disintermediated by users who will prefer to use AGI services directly. Like intellectual

property rights, this should be properly regulated to avoid harming the operators of these disintermediated systems.

8. **Increased risks of loss of competitiveness due to trade secret leaks;** As AGI will become used extensively for business, companies will have to make the choice between building their own system to optimally protect their business data and trade secrets or to outsource their corporate AGI to a third party (with the corresponding risks of leaks and loss of competitiveness if business know how leaks)

To maximise the positive impacts and minimise the negative consequences, stakeholders need to collaborate on a global scale to establish ethical guidelines, regulations, and best practices for the development and deployment of AGI.

This would include investing in education, reskilling, and social safety nets to support those displaced by technological change, as well as fostering international cooperation to ensure that AGI technologies are used responsibly and ethically across the globe.

14.2 Scenario 2: Aggravation of US cold war with China and Russia

In the most likely scenario as of mid-2023, based on the geopolitical tensions between allies of the US and allies of China and Russia, an aggravation of the US cold war with China and Russia could lead to significant changes in the societal impacts of AGI becoming available by 2030. The competitive and confrontational environment that would develop would shape both the positive and negative impacts.

On the positive side, the heightened competition could drive rapid advancements in AGI research and development, as countries strive to maintain or achieve technological superiority. This competition

might lead to breakthroughs in areas such as healthcare, climate change mitigation, and resource management. Additionally, the US and its allies would likely invest heavily in education, research, and development to maintain a competitive edge, which could have long-lasting benefits for the respective countries.

On the negative side, the confrontational environment could limit international collaboration and the sharing of knowledge, potentially slowing down overall AGI progress. The focus on national security and military applications of AGI might take precedence over more humanitarian and benevolent uses, leading to an uneven distribution of benefits across society and exacerbating existing inequalities. Furthermore, the risk of AGI being weaponized or used for surveillance and control would increase in this context, posing significant ethical and human rights concerns.

In the face of such challenges, the EU and other democracies, closely aligned with the US, would need to establish mechanisms for responsible AGI development and deployment, balancing the urgency of staying competitive with the importance of ethical considerations. They would need to focus on fostering open strategic autonomy in the field of AGI to ensure technological resilience, while also promoting international norms and regulations to minimise the risks and negative consequences associated with AGI in a confrontational geopolitical context.

In conclusion, the aggravation of the US cold war with China and Russia, as the most likely scenario in mid-2023, would significantly impact the societal consequences of AGI becoming available by 2030. The US, the EU, and other democracies would need to carefully navigate the challenges and ensure responsible AGI development and deployment because the competitive and confrontational environment would shape both the positive and negative outcomes.

Opportunities

1. **Technological Advancements:** The race between rival blocs to develop AGI could lead to accelerated technological advancements in AI and related fields.
2. **Military Innovation:** In a climate of increased geopolitical tension, AGI could lead to the development of advanced defence systems, potentially reducing the risk of human casualties in armed conflicts.

Challenges

1. **AGI Arms Race:** The competition between the US and China (and their respective allies) to develop AGI could lead to an arms race, with each side striving to outpace the other in AGI capabilities. This could result in the deployment of AGI before adequate safety measures and ethical considerations are in place.
2. **Uneven Distribution of Benefits:** The benefits of AGI may be restricted within the rival blocs, exacerbating global inequalities, and further dividing the world.
3. **Cybersecurity Threats:** The risk of AGI being weaponized for cyber warfare and espionage would increase, with the potential for more sophisticated and autonomous cyberattacks.
4. **Misuse of AGI:** In a highly competitive environment, AGI could be used for surveillance, propaganda, and manipulation of information, undermining trust in institutions and eroding democratic values.
5. **Chilling effects due to mass surveillance:** Without proper safeguards, AGI could be used to create advanced mass surveillance systems, raising concerns about privacy and individual freedom, and potentially creating "chilling effects" counterproductive to the economy and to the functioning of democratic societies.
6. **Escalation of Conflict:** The development of AGI by rival blocs could contribute to an escalation of tensions and increase the likelihood of direct

conflict, with potential catastrophic consequences.

To mitigate these negative impacts, it is crucial for the international community to establish mechanisms for cooperation and dialogue to address the challenges posed by AGI. This could include confidence-building measures, joint research initiatives, and the development of shared ethical guidelines and regulatory frameworks to ensure that AGI is developed and deployed responsibly, even amidst geopolitical tensions.

14.3 Scenario 3: Direct military conflict between US and China

In the second most likely scenario as of mid-2023, based on the geopolitical tensions between allies of the US and allies of China and Russia, a direct armed conflict between the US and China would have a profound effect on the societal impacts of AGI becoming available by 2030. The ongoing conflict and the urgent need for strategic advantages would shape the positive and negative impacts.

On the positive side, the intense competition and urgency of the conflict could accelerate AGI research and development, as both sides would seek to gain a technological edge on the battlefield. This acceleration might lead to breakthroughs in AGI capabilities, which would later be repurposed for civilian applications in areas such as healthcare, climate change mitigation, and resource management once the conflict subsides.

On the negative side, the conflict would likely severely limit international collaboration, further exacerbating the fragmentation of the global research community. The focus on military and national security applications of AGI could lead to an arms race, with both sides developing increasingly advanced autonomous weapons and surveillance systems. This would not only raise significant ethical concerns but also increase the

risk of a destabilising global security environment.

Moreover, the resources diverted towards AGI research and development for military purposes would likely come at the expense of investment in humanitarian and benevolent applications, exacerbating existing inequalities and potentially delaying advancements in areas such as healthcare and education. Additionally, the conflict would make it more difficult to establish international norms and regulations surrounding AGI, leaving the technology vulnerable to misuse or unintended consequences.

In the face of these challenges, the EU and other democracies closely aligned with the US, would need to balance their strategic interests with the importance of ethical considerations and responsible AGI development. Efforts would need to be made to mitigate the risks associated with AGI in a conflict scenario, including promoting international norms and regulations to minimise the negative consequences of AGI deployment and use.

In conclusion, a direct armed conflict between the US and China, as the second most likely scenario in mid-2023, would have a profound impact on the societal consequences of AGI becoming available by 2030.

The US, the EU, and other democracies would need to carefully navigate the challenges and ensure responsible AGI development and deployment amidst the conflict because the ongoing conflict and the urgent need for strategic advantages would shape both the positive and negative outcomes.

Additional challenges compared to the previous cold war scenario.

1. Misallocation of Resources: The focus on AGI for military and strategic purposes could divert resources away from addressing pressing global

issues like climate change, poverty, and healthcare.

2. **Prolonged Conflict:** The development of AGI by rival blocs could contribute to the prolongation and intensification of the conflict, with potential catastrophic consequences.

In such a situation, the role of the international community becomes even more critical in finding ways to resolve the conflict and address the challenges posed by AGI. This could involve seeking diplomatic solutions, fostering dialogue, and working towards the development of shared ethical guidelines and regulatory frameworks to ensure that AGI is developed and deployed responsibly, even amidst armed conflict. The focus should be on preventing further escalation and finding ways to cooperate on global challenges that transcend the boundaries of the conflict.

15.0 AGI and the theories of trade

This paragraph analyses how AGI could affect the Theory of Trade in the three different scenarios described above:

- **a world of peace** with no geopolitical tensions,
- **a Cold War** between the US and its allies and China & Russia
- **armed conflict**, US and China are at war and supported by their respective allies.

For each of the key concepts in the Theory of Trade, the impact on the three scenarios is assessed.

15.1 Comparative advantage

Comparative advantage is a concept in international trade theory that suggests that countries should specialise in producing and exporting goods and services for which they have the lowest opportunity

cost relative to other countries. By focusing on the goods and services they can produce most efficiently, countries can maximise their overall gains from trade, resulting in increased global productivity and wealth for all trading partners.

- **World of peace:** AGI could optimise production processes, enabling countries to better identify and capitalise on their comparative advantages, leading to more efficient trade and collaboration.

- **Cold War:** In a climate of strategic competition, countries may struggle to fully capitalise on their comparative advantages due to national security concerns and an emphasis on self-sufficiency.

- **Armed conflict:** Widespread conflict could significantly disrupt global trade, making it difficult for countries to leverage their comparative advantages and benefit from international trade.

It is important to note that AGI has the potential to transform the concept of comparative advantage itself. Countries may no longer rely solely on their traditional resources, labour, or capital to gain a competitive edge in global trade. Instead, access to advanced AGI technology could become a new source of comparative advantage, reshaping the global distribution of production and trade patterns. This shift emphasises the importance of staying at the forefront of AGI development to maintain a competitive position in the global economy.

15.2 Factor price equalisation

Factor price equalisation is a theoretical outcome in international trade, where factor prices (such as wages and returns on capital) converge across countries due to free trade. According to the Heckscher-Ohlin model, trade equalises factor prices because it allows countries to effectively

"trade" their abundant factors for the scarce factors of their trading partners. Consequently, the prices of factors of production should equalise across countries, leading to a more balanced global income distribution.

- World of peace: AGI-driven productivity improvements could promote factor price equalisation and reduce income inequality as countries benefit from increased efficiency and more accessible technologies.
- Cold War: The divide between rival blocs could exacerbate factor price disparities, as countries within each bloc prioritise investment and development within their sphere of influence, potentially leading to imbalances in global factor prices.
- Armed conflict: War-driven economic disruptions and resource scarcities could exacerbate factor price disparities between countries, as nations focus on self-preservation and resource allocation for military purposes.

In the best-case scenario, AGI's potential to automate tasks and improve efficiency across industries could lead to a more homogenised global labour market, as the distinction between skilled and unskilled labour would become less relevant. This may potentially result in factor price equalisation, with wages and returns on capital converging across countries, promoting a more equitable global economic landscape.

15.3 Global value chains

Global value chains (GVCs) refer to the intricate network of production, distribution, and consumption processes spanning multiple countries. In GVCs, various stages of production for a single good or service are often carried out in different countries to capitalise on lower production costs or specialised expertise. GVCs have become a prominent feature of international trade, as

advances in transportation and communication technologies have facilitated coordination and management of cross-border production processes.

- **World of Peace:** AGI could enhance the efficiency of global value chains by streamlining logistics, transportation, and supply chain management, fostering greater collaboration and interconnectivity among nations.
- **Cold War:** Geopolitical tensions could disrupt global value chains as countries seek to reduce dependence on rival blocs, potentially leading to a fragmentation of international trade networks.
- **Armed conflict:** Widespread conflict would severely disrupt global value chains, causing resource scarcity, supply chain disruptions, and economic downturns, as nations prioritise military and security needs over trade.

In the best-case scenario, AGI's ability to enable greater automation and efficiency in production processes could lead to a reconfiguration of global value chains. Countries may increasingly rely on advanced AGI technology for production rather than low-cost labour, potentially diminishing the importance of offshoring and reshoring in international trade. This shift could result in a more technologically driven and resilient global trading system.

15.4 Trade in services

Trade in services refers to the exchange of services between countries, as opposed to the exchange of physical goods. Services can encompass a wide range of activities, such as tourism, financial services, education, and healthcare. As economies become increasingly knowledge-based and technology-driven, trade in services has become an essential component of global trade.

- **World of peace:** AGI could revolutionise service industries, enabling new types of services to be

traded and expanding the scope of international trade, as countries leverage AGI technology to create innovative solutions and strengthen global collaboration.

- Cold War: Strategic competition could hinder the growth of trade in services, as countries may be reluctant to share technology and intellectual property with rivals, potentially stifling innovation, and international cooperation.

- Armed conflict: The focus on military applications of AGI could divert resources away from efforts to enhance international trade in services, prioritising security and defense over economic development and international collaboration.

In the best-case scenario, the development of AGI could lead to a significant expansion of trade in services, particularly in areas like AI-based consulting, software development, and data analysis. This may make the trade in services an even more critical component of international trade, contributing to global economic growth and fostering new avenues for collaboration and innovation.

15.5 Trade barriers and protectionism

Trade barriers are policies or regulations that restrict or limit international trade, such as tariffs, quotas, and non-tariff barriers like import licensing or technical regulations. Protectionism is the practice of implementing trade barriers to protect domestic industries from foreign competition. While protectionist policies can provide short-term benefits to domestic industries, they can also lead to reduced global trade, economic inefficiency, and retaliatory measures from trading partners.

- World of peace: AGI could contribute to a reduction in trade barriers by improving trade facilitation and streamlining customs procedures,

thus fostering greater international collaboration and economic growth.

- Cold War: Cold War tensions could lead to increased trade barriers and protectionism, as countries seek to safeguard their economies and support domestic industries, potentially hampering global trade and cooperation.

- Armed conflict: In a conflict-driven environment, trade barriers and protectionism would likely increase as countries prioritise national security and self-sufficiency, leading to further disruptions in international trade.

It is important to note that AGI-enabled automation could also have the opposite result and prompt countries to implement new trade barriers and protectionist policies to safeguard their domestic industries from disruption. This may lead to a reevaluation of current trade policies and agreements as nations grapple with the implications of AGI on their economies and international relations.

15.6 New trade flows

New trade flows refer to the emergence of fresh trade patterns between nations or regions, frequently as a result of shifts in comparative advantage, changes in global demand, or technological advancements. New trade flows can create new opportunities for economic growth and development as well as challenges for countries that must adapt to the changing landscape of global trade.

- World of peace: AGI could help countries identify untapped markets and trading partners, leading to the expansion of trade networks and new opportunities for economic growth, fostering international collaboration and innovation.

- Cold War: The competition between rival blocs could result in the emergence of new trade flows within each block, as countries seek alternative

trading partners, potentially leading to more regionally focused trade networks.

- **Armed conflict:** Countries involved in the conflict might seek alternative trading partners to maintain essential supplies, leading to trade diversion and less efficient trade flows, as nations prioritise security and resource management over global trade.

In the best-case scenario, the development and deployment of AGI could lead to new trade flows and relationships, as countries collaborate on AI research and development, share resources, and access new markets. This could create a more interconnected and innovative global economy, benefiting all participants in international trade.

15.7 Economic integration

Economic integration is the process by which countries increase their interdependence and cooperation in trade, investment, and other economic activities. This can be achieved through a variety of mechanisms, such as preferential trade agreements, customs unions, common markets, or economic unions. Economic integration can lead to increased trade, more efficient resource allocation, and greater economic stability, but it may also require countries to relinquish some degree of policy autonomy and adapt to new competitive pressures.

- **World of peace:** AGI could promote economic integration by fostering cooperation and collaboration between countries in technology, research, and development, thereby enhancing global economic interconnectedness and growth.

- **Cold War:** Geopolitical tensions could hinder economic integration between rival blocs, as countries focus on strengthening ties within their sphere of influence, potentially resulting in fragmented economic relationships.

- **Armed conflict:** Widespread armed conflict would significantly hinder economic integration, as countries focus on their own survival and the protection of their interests, leading to disruptions in global trade and collaboration.

In the best-case scenario, the widespread adoption of AGI could facilitate greater economic integration as advanced technology reduces transaction costs, streamlines cross-border communication, and enhances the efficiency of global trade. This could lead to a more interconnected and resilient global economy, benefiting all participating nations.

15.8 Absolute Advantage

Absolute advantage refers to the ability of a country to produce a good or service more efficiently than other countries, using fewer resources or in less time. While comparative advantage focuses on opportunity costs, absolute advantage focuses on overall efficiency in production.

- **World of peace:** AGI could enhance productivity and resource allocation, allowing countries to develop stronger absolute advantages in certain goods and services, leading to more efficient trade and increased global economic growth.
- **Cold War:** Countries might focus on leveraging their absolute advantages in strategically important industries, such as defence or critical infrastructure, to maintain a competitive edge over rival nations, potentially limiting broader economic benefits and stifling cooperation.
- **Armed conflict:** In a war scenario, countries may prioritise the development of absolute advantages in industries that contribute to their military and defence capabilities, potentially reducing the focus on other areas of trade and limiting overall economic growth.

It is important to note that AGI itself would most likely be recognized as an absolute advantage for the military sector, hence accelerating its development in times of cold war or armed conflict. This could further intensify the strategic competition between nations and have significant implications for international relations and global security.

15.9 Terms of Trade

Terms of trade refer to the relative prices at which a country's exports and imports are exchanged. If a country's terms of trade improve, it means that it can buy more imports for a given amount of exports, and vice versa. Global supply and demand, currency exchange rates, and trade policies are a few examples of the factors that affect trade terms.

- World of peace: AGI-driven productivity improvements could help countries achieve more favourable terms of trade by increasing the value of their exports relative to imports, leading to enhanced economic growth and prosperity.
- Cold War: Geopolitical tensions and trade barriers might lead to more volatile terms of trade, as countries' relative prices and exchange rates are affected by strategic competition, potentially limiting the benefits of AGI-driven improvements in productivity and international cooperation.
- Armed conflict: Widespread conflict could severely disrupt terms of trade, leading to volatile exchange rates, disrupted supply chains, and reduced global trade, overwhelming any productivity gains from AGI and exacerbating economic instability.

In each of these scenarios, the impact of AGI on terms of trade will depend on how effectively countries can harness the potential benefits of

advanced technology while navigating the challenges posed by geopolitical tensions and conflict.

15.10 Gains from Trade

Gains from trade refer to the net benefits that countries receive from engaging in international trade. By specialising in goods and services where they have a comparative advantage, countries can increase their overall productivity and wealth, leading to improved living standards and economic growth.

- **World of peace:** By enhancing productivity and resource allocation, AGI could increase gains from trade, leading to improved living standards, economic growth, and global prosperity.
- **Cold War:** Strategic competition might limit gains from trade, as countries focus on self-sufficiency and protectionism to maintain their competitive edge, potentially reducing the overall benefits of AGI-driven productivity improvements and hindering international cooperation.
- **Armed conflict:** The focus on military objectives and national security in a conflict scenario would likely lead to reduced gains from trade, as global trade is disrupted and countries prioritise their own survival, minimising the positive impact of AGI on global trade and economic stability.

In each of these scenarios, the extent to which AGI impacts gains from trade will depend on the interplay between technological advancements, geopolitical tensions, and international cooperation. Achieving the full potential of AGI-driven improvements in global trade will require countries to navigate these challenges effectively.

15.11 Trade Balance.

The trade balance is the difference between the value of a country's exports and imports. A trade surplus occurs when the value of exports exceeds that of imports, while a trade deficit occurs when the value of imports exceeds that of exports. Exchange rates, domestic production levels, and trade policies are a few examples of the factors that can affect trade balances.

- World of peace: AGI-driven productivity improvements could positively affect countries' trade balances by increasing the value of their exports and reducing import dependence, fostering economic growth and stability.
- Cold War: The competition between rival blocs might lead to trade imbalances, as countries seek to maintain self-sufficiency and limit dependence on rivals, potentially undermining the benefits of AGI-driven productivity gains and causing economic uncertainties.
- Armed conflict: Widespread conflict would likely lead to trade imbalances, as countries prioritise military objectives and resources are diverted away from trade, reducing the potential positive impact of AGI on trade balances and leading to economic disruptions.

In each of these scenarios, the impact of AGI on trade balances will depend on how effectively countries can adapt to changing economic conditions and leverage technological advancements to optimise their trade policies and strategies.

15.12 Heckscher-Ohlin Model.

The Heckscher-Ohlin model is a theory of international trade that explains how differences in factor endowments (such as labour and capital) between countries determine their comparative advantages. According to the model, countries will specialise in producing goods that use their abundant factors of production intensively and

import goods that use their scarce factors intensively.

- **World of peace:** According to the Heckscher-Ohlin model, AGI could result in a more effective allocation of resources, further aligning countries' production with their factor endowments. This would result in increased global trade and economic growth.
- **Cold War:** Geopolitical tensions might cause countries to deviate from the Heckscher-Ohlin model predictions, as they prioritise strategic industries and self-sufficiency, limiting the overall efficiency of resource allocation driven by AGI. Consequently, global trade and economic growth could be hindered.
- **Armed conflict:** In a war scenario, the focus on military capabilities would likely lead to a deviation from the Heckscher-Ohlin model, as countries prioritise their own survival and defence, potentially reducing the benefits of AGI-driven improvements in resource allocation. This would result in disrupted trade flows and economic downturns.

In each of these scenarios, the impact of AGI on the Heckscher-Ohlin model will depend on how effectively countries can balance their strategic priorities with the potential benefits of technological advancements in resource allocation and trade.

15.13 Ricardian Model.

The Ricardian model is a simple theoretical model of international trade that focuses on comparative advantage based on differences in labor productivity between countries. The model assumes that there is only one factor of production (labor) and that countries differ in their production technologies. The Ricardian model shows that trade can be mutually beneficial even if one country has an absolute advantage in producing all goods.

- **World of peace:** Because countries with higher labour productivity in particular sectors can specialise and trade more effectively, AGI could further amplify the advantages of trade that the Ricardian model predicts. This would result in increased global trade and economic growth.
- **Cold War:** The focus on strategic competition might limit the applicability of the Ricardian model, as countries prioritise self-sufficiency and protectionism over trade, potentially reducing the benefits of AGI-driven productivity improvements. Consequently, global trade and economic growth could be hindered.
- **Armed conflict:** In a conflict-driven environment, the Ricardian model's predictions would likely be less relevant, as countries prioritise military objectives and national security over international trade, limiting the positive impact of AGI on trade and economic growth.

In each of these scenarios, the impact of AGI on the Ricardian model will depend on how effectively countries can balance their strategic priorities with the potential benefits of technological advancements in labour productivity and trade.

15.14 Summary

In summary, the impact of AGI on international trade theories is highly dependent on the geopolitical context.

In a world of peace, AGI has the potential to enhance comparative advantage, promote factor price equalisation, improve global value chains, expand trade in services, reduce trade barriers, foster new trade flows, and encourage economic integration. These outcomes would likely lead to increased global trade and economic growth.

However, in scenarios characterised by Cold War tensions or armed conflict, the potential benefits of AGI-driven advancements would be limited or even

reversed. In such cases, countries would prioritise national security, self-sufficiency, and military objectives over open trade and international cooperation. As a result, comparative advantage might become less relevant, factor price equalisation could be hindered, global value chains would face disruptions, trade in services might stagnate, trade barriers and protectionism would likely increase, new trade flows could be less efficient, and economic integration would be significantly impeded. These conditions would likely lead to reduced global trade and slowed economic growth.

16.0 AGI and Foreign policy

16.1 China and Russia joint statements on Global Governance

[China and Russia issued a joint statement on global governance on February 4, 2022,](#) when Russian President Vladimir Putin visited China and met with Chinese General Secretary Xi Jinping. The statement outlined their common vision of international relations in a new era and their cooperation on various issues, such as the Covid-19 pandemic, climate change, trade, security, and human rights.

The more than 5,000-word joint statement also reaffirmed their support for each other's territorial claims in Taiwan and Ukraine, which the United States and its allies contest. China and Russia accused the West of interfering in their internal affairs and violating their sovereignty and interests by imposing sanctions and pressure. They called for a peaceful resolution of the conflicts in accordance with international law and the UN Charter, but 17 days later, on February 24, 2022, Russia invaded Ukraine, blatantly violating the charter of the United Nations Security Council.

The joint statement also denounced the sanctions and interference from the West as attempts to undermine the existing post-war world order and the authority of the United Nations. China and Russia

vowed to uphold the outcomes of the Second World War, defend the international system with the UN at its core, and resist any attempts to falsify or distort history. The joint statement also criticised the concept of a "rules-based international order" that is promoted by the West to impose its values and interests on other countries.

The statement is seen [by most analysts](#) as a sign of a [strategic alliance between China and Russia](#) against the United States and its allies, [challenging the global order dominated by western democracies](#) and as a response to the US-led efforts to form a coalition of democracies to counter China and Russia's influence.

16.2 China and Russia quest for a new world order

China and Russia share their vision for a new world order where the US and its allies would lose their military hegemony and, most importantly, their international influence on third countries of the global south (such as India, South Africa, Nigeria, Morocco, Kenya, Brazil, etc.)

This new world would be one where the US and its allies would no longer be an obstacle to China's territorial ambitions over Taiwan and the South China Sea and to Russia's over Eastern Europe and the Baltic states.

In this new world, the dollar would lose its status as a reserve currency, and the role of Bretton Woods global governance organisations would be further weakened and replaced by new "international bodies" politically dominated by China and Russia.

China and Russia would also seek to impose on non-aligned countries of the global south their own model for the governance of the Internet to allow governments to control information and distil propaganda.

If such a dramatic scenario materialises, criminal organisations will find support from Russia and China to exert their influence worldwide, especially in democratic countries where the enforcement of the rule of law will be further weakened and corruption will rise.

16.3 Major conflict of interests with China and Russia

There is indisputable evidence (see below) of clashes between the geopolitical interests of allies of the US on the one hand and allies of China on the other.

Russia, which sits on the European continent on the west and Asia on the east, has chosen its side and supports China with the hope that it will help it reach its own objectives in Ukraine, as evidenced in the Putin-Xi joint statement issued on February 7, 17 days before the start of the Ukrainian war.

1. **NATO:** NATO members remain steadfast in their determination to thwart Russia's territorial aspirations in Eastern Europe and the Baltic states.
2. **Support to Israel:** In the Middle East, the financial and military aid that the U.S. extends to Israel, coupled with economic sanctions on Iran and Syria, have driven these two nations to form alliances with Russia, thereby countering U.S. and EU interests.
3. **US Military Presence in the South China Sea:** In the South China Sea, the U.S. military serves a vital role in obstructing China's attempts to establish military bases near the exclusive maritime economic zones of countries such as the Philippines.
4. **Taiwan:** The U.S. is preserving the status quo with respect to Taiwan, offering security assurances to safeguard it against China's long-standing ambition to reunify with the island, even by the means of force, if

necessary. U.S. arms sales to Taiwan have been a contentious issue as China in that context.

5. **Japan:** Japan, a critical U.S. military ally, is experiencing similar territorial issues with China and Russia to its north and northeast, as well as tensions with North Korea.
6. **Australia:** This U.S. military ally has traditionally viewed the Pacific as being within its sphere of influence, now perceives its regional security architecture to be under threat by China.
China's growing influence in the Solomon Islands raises several potential threats. These include strategic concerns related to China's ability to establish military bases or facilities in the region, and economic leverage from extensive Chinese investments that may lead to "debt-trap diplomacy". Additionally, increased Chinese presence could potentially lead to shifts in local political decisions unfavourable to Australia, weaken regional institutions, and cause over-exploitation of the islands' rich fisheries and natural resources.
7. **France's military presence in the Indo Pacific Region:** France, which has numerous overseas territories in the 'Indo-Pacific Region', also represents a minor impediment to China's maritime ambitions. In response to tensions in South China sea, France has bolstered its military presence in the Indo-Pacific region, partnering with India, the U.S., and Japan
8. **Saudi-Arabia:** Recent diplomatic manoeuvres by China to mitigate tensions between Iran and Saudi Arabia have led the latter to distance itself from the U.S. and gravitate towards the Russian and Chinese spheres of influence.
9. **Trade Imbalances:** The U.S. has long criticised China for unfair trade practices, including currency manipulation, intellectual property theft, and providing state support to Chinese companies, creating an uneven playing field. The trade imbalance has been a persistent point of contention, leading to a trade war under the Trump administration.

10. **Cybersecurity:** [The U.S. has accused China of conducting cyber espionage and intellectual property theft, compromising national security and economic interests.](#)
11. **Human Rights:** The U.S. has expressed concerns over human rights abuses in China, particularly in regard to the treatment of Uighurs in Xinjiang and the handling of protests in Hong Kong. The Chinese government's crackdown on pro-democracy movements and its actions towards ethnic and religious minorities are significant points of tension.
12. **Technology Dominance:** The race for technological dominance, especially in areas like AI, quantum computing, 5G, and biotechnology, has created rivalry. The U.S.' concerns over Chinese tech firms like Huawei and TikTok stem from issues over data privacy and national security.
13. **Covid-19 Origin:** The origins of the Covid-19 virus have been a source of friction. U.S. requests for more transparency and thorough investigations into the origins have been met with resistance from China.
14. **Climate Change:** Although both countries agree on the importance of addressing climate change, tensions arise from disagreements over responsibilities and the scale of action required.
15. **Space Race:** China's rapid advancement in space exploration and technology, evidenced by its missions to Mars and the Moon, has led to a renewed space race, causing both cooperation and tensions over competition.
16. **Dalai Lama and Tibet:** U.S. support for the Dalai Lama and the issue of Tibetan independence continues to be a sensitive topic for China.
17. **China's Belt and Road Initiative,** aiming to boost its influence in South Asia and beyond, has been seen with suspicion by the U.S. This initiative has potential implications for U.S. influence in these regions.

16.4 China overtly encourages destabilisation of US allies.

China has no economic incentive to directly initiate armed conflict with the US or its allies, as the economic stakes are exceptionally high. Furthermore, China is the world's largest economy. Engaging in direct conflict would tarnish its image among nations in the Global South, with whom it competes against the US and its allies.

Instead, it serves China's economic interests to subtly and openly encourage various conflicts and actions initiated by proxy actors, such as Russia and North Korea. This approach aims to weaken the influence of the US and its allies on the international stage without directly involving China in the conflict.

For example, China has not condemned any of the Russian and North Korean actions below:

- North Korea's aggressive stance towards South Korea,
- Russia's non provoked war against Ukraine,
- Russia's partial occupation of Georgia,
- Russia's support for the Syrian regime, Russia's security guarantees to non-democratic countries in Africa to help them maintain their leaders' grip on power.

16.5 China and Russia weaponize trade for political reasons

China has been known to use trade as a tool for political leverage and coercion, particularly towards countries that it perceives as challenging its interests.

In recent years, China has used economic pressure and restrictions on trade to assert its influence.

1. **In the case of Lithuania**, tensions escalated when Lithuania sought to deepen diplomatic ties with Taiwan [2][13]. China responded by imposing

trade restrictions and economic sanctions on Lithuania, such as temporarily removing Lithuania from its customs clearance systems, causing significant difficulties to bilateral trade [1].

2. **In 2010 when China blocked exports of rare earths** to Japan due to diplomatic disputes over the Senkaku Islands [14]
3. **In 2021, China suspended Australian imports** due to tensions over the origin of the COVID-19 pandemic [11].

Russia also weaponizes trade in several ways, using its economic influence and exports to exert pressure on other nations, further its geopolitical goals, and influence political and military leaders.

1. **Arms trade:** Russia is one of the world's top weapons-exporting nations, which allows it to establish relationships with other countries, influence their leaders, and support its broader foreign and defense policy goals [1][17].
2. **Energy sector:** Russia is a major supplier of natural gas, and it has been accused of using its energy resources as a political weapon. For example, it has reportedly weaponized natural gas supplies by leveraging its dominance in the European market to exert pressure on countries that rely on its resources [8][18][21].
3. **Food and agriculture:** Russia has been accused of weaponizing food by blocking crucial grain exports from Ukraine, which could lead to the starvation of millions of people worldwide [14]. It has also been claimed that Russia is weaponizing food supplies to blackmail the world [15].
4. **Cyber warfare:** Russia has reportedly waged a cyber war against the United States and other Western countries, using its expertise in the energy sector to conduct attacks on critical infrastructure [16].
5. **Currency and financial systems:** Russia has been accused of weaponizing the US dollar and other Western currencies to punish its adversaries,

utilizing economic sanctions and financial pressure to exert influence [5].

6. **Soft power:** Russia uses nonmilitary means of power, such as information influence, ideological influence, and political pressure, to further its objectives. This is described as the "weaponization of soft power" [7].

China and Russia's willingness to use trade as a political tool highlights the need for a coherent and united response from the international community.

16.6 China's ambiguous Belt and Road initiative (BRI)

China's Belt and Road Initiative (BRI) is an ambitious infrastructure project launched in 2013 by President Xi Jinping [1]. The BRI aims to develop two new trade routes connecting China with the rest of the world, one overland and the other by sea, linking China with its neighbors in Central Asia, the Middle East, and Europe [5]. The initiative is expected to cost more than \$1 trillion, and China has already lent billions of dollars for infrastructure projects in 68 countries [2, 6].

While the BRI has the potential to bring economic benefits and improved connectivity to participating countries, it has been criticised for several reasons, making it appear devious:

1.0 Debt diplomacy: Critics argue that the BRI could lead to unsustainable debt for countries participating in the initiative, as they may struggle to repay the loans provided by China. This could result in China gaining control over strategic assets and influence in those countries [14].

2.0 Lack of transparency: The BRI has been criticised for its lack of transparency, as the details of the agreements and projects are often not publicly available. This raises concerns about

potential corruption, environmental impacts, and labour rights violations [15]. The provision of financial aid may be contingent upon China's backing in the United Nations, as well as the recipients' approval of national and international governance structures in which China can exert its influence. Additionally, the acceptance of Chinese technology across various sectors could be a condition for receiving financial aid.

3.0 Geopolitical ambitions: Some view the BRI as a means for China to expand its geopolitical influence and challenge the current international order. For example, the United States has criticised China for using economic inducements and implied military threats to persuade other states to heed its political and security agenda [17].

16.7 AGI as a choking instrument to counter China and Russia

The latest generation of Multimodal Generative Large Language Models like GPT-4 have developed outstanding reasoning abilities and many fascinating intellectual and cognitive abilities. Those capabilities have been acquired after adequate training on a very large data set.

The development of AI and AGI is as much a paradigm shift for humanity as the mastering of fire by the first humans. Multimodal Generative Large Language Models are clearly destined to gradually lead to Artificial General Intelligence in the coming years, probably by 2030 at the very latest, depending on the definition of AGI adopted.

The industrial know-how to engineer such powerful Multimodal Large Language models like GPT-4 is currently concentrated in the US among a limited number of experts. Expertise is available outside the US to create various types of narrow AI systems and even Large Language Models, but certainly not of the same level of performance as GPT-4 or ChatGPT-4.

Over time, this US “know-how” will eventually structure into organised knowledge that will be taught in universities worldwide and that will hopefully spread to all technologically advanced countries so they can benefit from the technology.

The problem is that in the current geopolitical context with a non-significant risk of direct armed conflict between US allies and China allies, there is no interest for the US and its allies that China and Russia develop AGI too fast. In short, from the national security perspective of the US and its military allies, the spread of advanced US know-how on AI to countries like China and Russia is not a good thing, as it will inevitably be used against them.

While it is impossible to prevent the development of AGI in those countries, there is an interest for the US and its allies to stay ahead long enough to gain sufficient military advantage and/or until the current geopolitical tensions between the two blocks calm down. US foreign policy towards Russia and China in the field of AGI should therefore be aimed at ensuring that the US and its allies stay ahead of China and Russia on AGI as long as needed.

16.8 More export control on technologies enabling AGI.

In a Cold War scenario with a high risk of armed conflict, the existing (mid-2023) export control and foreign direct investment restrictions imposed by the US and its allies on China and Russia in the semiconductor value chain would likely need to be further reinforced to encompass all critical elements of the AGI global value chain.

In addition to cutting-edge semiconductors, semiconductor manufacturing equipment, and related computer-aided design software, the following AI-related technologies developed by the US and its allies could be subject to export controls:

1. **AI accelerator chips:** Specialised hardware chips designed to optimise AI computation and improve the performance of machine learning algorithms. **Neuromorphic chips for running spiking neural networks:** Specialised hardware designed to mimic the human brain's neural structure for more efficient AI processing. This includes especially **Neuromorphic chips for running spiking neural networks** that are designed to mimic the human brain's neural structure for more efficient AI processing.
2. **Quantum computing components:** Hardware and software components necessary for the development and operation of AGI-based quantum computers. This includes **silicon-quantum-dot chips, trapped-ion-qubit chips, and photonic-quantum chips** implementing qubits using photons. These technologies are anticipated to play a key role in the development of next-generation quantum computers for AI.
3. **Advanced AI algorithms:** Proprietary machine learning and deep learning algorithms, as well as frameworks and libraries that can be used to build AGI systems.
4. **Large-scale data sets:** Sensitive and proprietary data sets used to train AI and AGI systems, particularly those that may have strategic, military, or national security implications.
5. **Intellectual property rights (IPR) related to AGI:** licenses, patents, and other legal protections related to AGI research, development, and deployment.
6. **Cybersecurity technology:** Advanced tools and techniques that could be used to protect or breach AGI systems, depending on the user's intent.

16.9 More FDI screening on AGI.

To prevent the transfer of technologies that could advance AGI research in China and Russia, incoming foreign direct investment (FDI) from these

countries into the US and its allied nations should be strictly controlled. The same principle should apply to outgoing FDI from the US and its allies to China and Russia.

According to Statista, the countries that invested the most in China in 2021 were Singapore with \$9.8 billion, the Virgin Islands with \$9.5 billion, South Korea with \$7.9 billion, the Cayman Islands with \$6.8 billion, Japan with \$6.4 billion, Germany with \$4.7 billion, and the United States with \$4.5 billion[1].

This combined total of approximately \$50 billion stands in contrast to the \$87.8 billion in Chinese foreign direct investment (FDI) abroad in 2021, as reported by Statista[1].

In 2021, both the US and Germany invested nearly the same amounts (about \$4.5 billion each) in China, which is significantly lower than investments from South Korea (\$7.9 billion), Japan (\$6.8 billion), and tax haven countries like the Cayman Islands, the Virgin Islands, and Singapore (\$9.5 billion each)[1].

Implementing strict controls on FDI involving China and Russia would help safeguard critical AGI-related technologies and prevent their potential misuse by these nations. This measure would not only strengthen national security but also maintain a competitive edge in the rapidly evolving field of AGI.

16.10 Halt R&D cooperation on technologies enabling AGI.

Export controls and FDI limitations alone are unlikely to allow the US and its allies to stay ahead of China and Russia in AGI development. Regrettably, international research cooperation in academia and the private sector, focusing on energy-efficient machine learning algorithms and

hardware architecture, may also need to be restricted.

Additionally, there is a need to closely guard against industrial espionage from China and Russia the advanced engineering expertise that top US companies have acquired in developing potent multimodal generative large language models like GPT-4.

Despite these measures, the borders between the two blocs should remain open. STEM talents from China and Russia ought to be encouraged and incentivized to seek opportunities in the US and other allied democratic countries. By doing so, they can contribute to AGI development for the benefit of the US and its allies, promoting a more collaborative and diverse research environment.

16.11 AGI as an instrument to promote democracy.

Authoritative countries such as China and Russia can easily develop a chatbot system like ChatGPT, but that will be less advanced than the state-of-the-art technology available in the US.

Such a system will be trained to develop an inner representation of the world that corresponds to their authoritative visions and need for internal propaganda. Technically, the ideological alignment could be done during the fine-tuning phase, for instance.

China and Russia could even make those ideological chatbots available to other countries that are not as technologically advanced and help them tune them to fit their ideologies. In other words, through these system, China or Russia would help governments in countries that they support develop an AGI to disseminate their specific ideologies and propaganda (to detriment of views of the world promoted by the US and its allies)

Technically, an authoritative country could even use an open-source Large Language Model developed by a democratic country to build such an AGI.

To compensate for this threat, democratic countries that developed advanced AI services and systems fine-tuned based on truly human and democratic values, should make them easily available to less technologically advanced countries to promote their own view of a democratic world, so these countries don't seek support from China or Russia.

This assumes, of course, that the countries concerned demonstrate an appropriate alignment with human and democratic values. For instance, access to services like ChatGPT-4 could be eased in developing countries of the Global South like Brazil, Kenya, Nigeria, Indonesia, and South Africa.

China has demonstrated its willingness to invest heavily in the digital infrastructure of developing countries, notably through its Belt and Road Initiative. This could potentially extend to the provision of AGI services at some stage, as Chinese technology companies like Huawei, Alibaba, and Tencent have a strong presence in many developing markets.

To avoid the adoption of AGI services offered by Chinese companies, the U.S. and its allies could:

- 1. Offer Competitive Alternatives:** Develop and offer reliable, efficient, and affordable AGI services that meet the specific needs of these countries. This may involve partnering with local tech companies and considering joint ventures or other collaborative models.
- 2. Promote Transparency and Trust:** Emphasise transparency, both in terms of how the technology works and how the data is used. Building trust is crucial, and clear privacy

policies and ethical standards could give the U.S. and its allies an edge.

3. **Invest in Infrastructure:** Assist in building the necessary digital infrastructure that allows for the adoption and optimal use of AGI services. This could be part of broader development cooperation or specific tech-focused initiatives.
4. **Capacity Building:** Invest in education, training, and capacity building in these countries, not just in terms of using AGI services, but also in understanding their societal, economic, and ethical implications.
5. **Regulatory Cooperation:** Work with these countries to develop regulatory frameworks for AGI that protect consumer rights, privacy, and security. This could also include cooperation on issues like cybersecurity.
6. **Promote Open and Inclusive Digital Standards:** Advocating for open digital standards can ensure interoperability and prevent countries from becoming locked into a single provider's technology ecosystem.

17.0 AGI and the theories of International Relation

17.1 Summary of the 5 main International Relation theories

1. **Realism** is characterised by four key principles. Firstly, states are the principal actors in international politics. Second, national interests—typically those of power and security—drive all states within the system. Thirdly, the relative levels of power between states determine their relations. Finally, power is primarily a function of material resources, especially military capabilities. There are different branches within realism, like Classical Realism, which attributes states' behaviours to human nature, and Neorealism (or Structural Realism), which attributes international conflicts to the anarchic structure of the international system.

2. **According to liberalism,** laws and agreements can regulate domestic politics as well as international relations. Liberal theorists believe that global institutions, such as the United Nations, have the ability to maintain peace and promote global prosperity. They advocate for free trade, democracy, and human rights. Neoliberalism, a branch of liberalism, emphasises the importance of international institutions in reducing the inherent conflict that realists see in international systems.
3. **Constructivism:** Constructivism asserts that significant aspects of international relations are historically and socially constructed, rather than inevitable consequences of human nature or other essential characteristics of world politics. It focuses on the roles of ideas, norms, and beliefs in shaping state behaviour. According to constructivists, the identities and interests of states are malleable and change over time based on their interactions with other states and non-state actors.
4. **Marxism:** Marxist theories of international relations focus on the economic and material aspects. They assert that all politics, including international politics, are a function of the economic structure of society. They examine how global wealth is distributed and argue that global poverty and underdevelopment result from capitalism. Dependency theory and World-Systems theory are two prominent sub-branches of Marxist international relations theory.
5. **Feminism:** Feminist theories in international relations focus on gender dynamics. They argue that traditional international relations theories often overlook the role of women and other marginalised groups, and that these groups can provide valuable perspectives. Feminist theorists study issues such as how war and conflict affect women differently, and the role of women in peace making and international development. They assert that gender is not just an identity but a socially constructed set of

expectations about behaviour and rights, which can influence international relations.

17.2 Country by country analysis

Nations don't strictly adhere to one specific international relations theory in their foreign policies. Instead, they employ a mix of strategies and principles from different theories, depending on the circumstances.

1. **United States:** The US has traditionally followed a mix of realism and liberalism. Realism can be seen in its emphasis on maintaining military dominance and pursuing national interests. Liberalism is evident in its promotion of democracy, human rights, and free trade. The dominance of the US in AI and AGI reinforces its realist tendencies to maintain global hegemony, but also its liberal tendencies in setting international norms for AI usage.
2. **United Kingdom:** The UK also follows a mix of realism and liberalism. The UK often aligns with the US in terms of promoting liberal values internationally, but also pursues its realist interests, such as maintaining its influence over former colonies and other regions. Technological dominance by the US is pushing the UK to follow similar norms on AI and AGI, given their close alliance.
3. **European Union:** Countries of the EU and the European Union External Action service primarily follows a liberal approach, emphasising cooperation, integration, and the promotion of human rights. However, the EU also pursues realist objectives when it comes to protecting the interests of its member states. The EU would advocate for international norms on AI and AGI to prevent a power imbalance and misuse of technology.
4. **Turkey:** Turkey's foreign policy has traditionally been influenced by a blend of realism and constructivism, with an increasing leaning

towards realism in recent years. Turkey has a geopolitical focus on securing its regional interests, a classic element of realist theory. This has been evident in its involvement in various regional conflicts, its approach to relations with neighbouring countries, and its quest to establish itself as a regional power. On the other hand, Constructivist elements can be seen in Turkey's historical desire to be recognized as part of the Western community of nations, evident in its long-standing, though recently strained, bid to join the European Union

5. **Japan and South Korea:** These countries lean towards liberalism, emphasising economic cooperation, democracy, and rule of law. However, they also follow realism when dealing with security threats. They would likely align with the US in terms of AI norms due to their close economic and security ties.
6. **Australia and New Zealand:** Australia and New Zealand are known for their approach to international relations that leans heavily towards Liberalism. They strongly advocate for international cooperation, the rule of law, human rights, and democratic principles. Their foreign policies often align with other liberal democracies, especially the US.
7. **Taiwan:** Taiwan's position is unique due to its complex relationship with China. Taiwan generally follows a mix of Realism and Liberalism. It is focused on maintaining its independence (Realism) but also seeks international recognition and cooperation (Liberalism).
8. **China:** China mainly follows realism, emphasizing its national interests and sovereignty. However, it also uses elements of constructivism, such as promoting its model of governance as an alternative to Western liberal democracy. As the US achieves dominance in AI, China accelerates its own AI development and resists US-led norms as it conflicts with its interests.

9. **Russia's** foreign policy is primarily realist, emphasising national sovereignty and power. It often opposes Western-led initiatives and norms. Russia views US dominance in AI as a security threat and works to develop its own capabilities.
10. **North Korea:** North Korea follows a realist policy, prioritising its survival and sovereignty. It often resists international norms and is isolated from the global community. It sees US dominance in AI as a threat and will likely respond with increased military posturing or cyberattacks.

17.3 Working hypothesis for the next 5-7 years.

This paragraph outlines a set of working hypotheses for the potential evolution of current geopolitical tensions over the next 5-7 years.

- The United States, currently leading in the development of AI and AGI, is projected to maintain this lead in the coming 5-7 years, potentially exacerbating geopolitical tensions with China and Russia.
- It's presumed that Russia will sustain its engagement in the ongoing conflict with Ukraine, with NATO countries persistently supplying Ukraine with defensive and over time more and more aggressive weapons yet avoiding direct conflict with Russia.
- Russia will remain under substantial economic sanctions imposed by the U.S. and its allies, and it's expected it will face increasing challenges in circumventing these sanctions.
- Building on its success in reducing energy dependence on Russia, the European Union will strive to bolster its strategic autonomy across various sectors.

- It is anticipated that China will persist in its clandestine assistance to Russia, enabling it to partially dodge economic sanctions, drawing on support from countries in Central Asia, Iran, and North Korea.

- The continuous aim of China would be to ensure that Russia neither succumbs to Ukraine nor secures a swift victory. China's strategic objective remains to sustain this protracted conflict, thereby weakening the U.S. and its allies and diverting their attention from its hegemonic ambitions in Taiwan, the Indo-Pacific region, and the South China Sea.

- The application of AGI in military technology should indeed provide the U.S. and its allies the potential to develop more sophisticated weaponry, superior to that of China and Russia. As this technology advances, it could feasibly enable more precise, effective, and cost-efficient strikes on select targets deep within adversarial territories, such as Russia. This enhanced control and mitigated risk of unintended consequences could potentially strengthen Ukraine's position in the ongoing conflict. Furthermore, this technological edge might also support Taiwan in maintaining the current status quo with China.

17.4 Impact of the emergence of AGI on International Relations

In the context of the working hypothesis adopted and considering the international relations theories that best apply to each country, the emergence of AGI in the next 5-7 years should have a significant impact on international relations.

1. **US and its Allies (NATO, UK, EU, Australia, New Zealand, South Korea, Japan):** The United States leading in AGI development could foster greater cooperation among these allies, especially given the ongoing conflict in Ukraine. AGI might become a strategic asset in this context, potentially enhancing defence capabilities,

intelligence, and decision-making. These nations might also focus on establishing norms and regulations for AGI to ensure it's used responsibly and ethically.

2. **US, China, and Russia Relations:** The US's AGI leadership would likely exacerbate existing tensions with China and Russia. For Russia, already dealing with the conflict in Ukraine and facing economic sanctions, the US's AGI advancement might be perceived as another strategic threat. China, despite its covert support for Russia, might view the US's AGI superiority as a strategic disadvantage and thus escalate its own AGI efforts, potentially intensifying a technological arms race. This will, in turn, further increase tensions between the two blocks.
3. **Implications for Taiwan, Turkey, and North Korea:** Taiwan could see increased security risks if China perceives the prolonged Ukraine conflict as an opportunity to exert more pressure on the Taiwan independence issue. Turkey might face complex diplomatic challenges given its unique position as a NATO member with ties to Russia. The perceived AGI-enhanced threat from the US and its allies might also have an impact on North Korea's strategic calculations.
4. **Implications for Ukraine:** Ukraine, being directly involved in a conflict with Russia, might become a focus and a playground for AGI-enabled weapons, intelligence, surveillance, and strategic planning efforts by the US and its allies. Ukraine seems to welcome this possibility with the launch of the BRAVE-1 program to allow allied countries to test their latest military technology against the Russians.
5. **Role of International Diplomacy:** In this fraught geopolitical context, the role of international diplomacy becomes even more crucial. Constructive dialogue and negotiation on AGI

policy, safety, and ethics will be key to preventing escalations.

In summary, the emergence of AGI in the coming 5-7 years will add a significant layer of complexity to international relations. Its disruptive potential could shift power dynamics and affect strategic calculations, underscoring the importance of responsible management and international cooperation on AGI development and deployment. However, the real dynamics would depend on multiple factors, such as the pace of AGI development, the specifics of the geopolitical situation at the time.

Part VI: Policy options

Policy Options for AGI

Better law making in the European
Parliament

18.0 Policy options for AGI.

18.1 Prepare for societal impacts by 2030.

- 1. The US and its allies should proactively anticipate evolution of AI towards AGI, expected by 2030 at the latest, and initiate adequate measures to better prepare society.**

Democracies should prepare for AGI by investing in education and workforce development, supporting AGI research, fostering public-private partnerships, creating a regulatory framework, engaging in international cooperation, strengthening social safety nets, raising public awareness, establishing ethical guidelines, investing in digital infrastructure, and implementing pilot programs. These measures will help democratic societies maximise AGI's positive impacts, minimise negative consequences, and ensure a smooth transition towards AGI integration.

Failure to prepare for AGI may lead to unregulated development, harmful applications, increased economic inequality, and human rights infringements. Additionally, inadequate preparation could exacerbate skills gaps, limit AGI advancements, and leave systems vulnerable to cyberattacks.

Insufficient international cooperation and public engagement could result in uneven AGI benefit distribution, aggravating geopolitical tensions, and eroding public trust in AGI technologies and responsible institutions.

In the context of growing geopolitical tensions and the risk of armed conflict between the US and its allies on the one hand and China and its allies on the other, proper AGI preparation is crucial to ensuring a competitive edge and maintaining international stability and security.

2. Maximise the potential of AGI to revolutionise school education.

The objective of such a policy would be to help teachers tailor the learning experience to individual children's needs, to enhance children's engagement by identifying the ones who are struggling or bored, thereby enabling educators to intervene and adjust their teaching methods accordingly, to accelerate learning, to reduce the cost of learning basic skills such as reading and writing, and to improve teacher support, enabling them to focus more on engaging with children and providing them personalised support.

3. Optimise the integration of AGI in the educational process within universities.

The aim is to provide students with enhanced learning experiences and efficient knowledge acquisition while significantly reducing the duration of the education cycle in complex scientific and technological fields such as medicine and engineering. Emphasis should be placed on utilising AGI for personalised learning, employing advanced simulations and virtual reality, ensuring improved access to up-to-date knowledge, implementing intelligent tutoring systems with real-time feedback, and promoting enhanced collaboration and communication for superior interdisciplinary learning.

18.2 Better coordinate between like minded democratic partners

4. The need for an international coordination body on AGI and Data among US allies.

AGI and data policy questions should be coordinated at the international level between

all US allies that share the same democratic values and respect for the rule of law. While existing international government frameworks like the OECD and WEF focus on responsible and trustworthy AI development, they do not address other governance issues.

Therefore, there is a need for an international organisation specifically focused on AGI and Data governance, with participation limited to allied democracies sharing appropriate democratic values. This organisation would not replace existing ones but would provide an additional layer of technocratic governance.

The alliance would address [AGI-specific challenges](#), including misuse, privacy, cross border data flows, intellectual property rights, and the alignment problem. Its primary goal is to facilitate responsible, value-aligned AGI development, ensuring shared progress and safety amidst emerging technologies.

Furthermore, the Alliance would also aim to promote global economic growth by facilitating and supporting trade among its members but also by empowering them to better compete with China over technology trade matters with non-aligned countries of the global south.

The EU-US Trade and Technology Council (TTC) is an important effort that goes in the right direction but is bilateral. In contrast, the Alliance would conduct similar discussions on a multilateral basis in a structured format, including not only the EU and the U.S. but also the UK, Japan, South Korea, Australia, New Zealand, Taiwan, Singapore, and others. Over time, the Alliance's work would supersede bilateral efforts, such as the EU-US TTC, for reasons of efficiency. However, during the interim phase, such bilateral efforts would lay the groundwork for the Alliance.

The Alliance, best implemented as a new international organisation, should be based on the US West Coast, a secure location in the context of the growing geopolitical tensions and risks of a third world war over Taiwan and Ukraine, home to the leading AGI companies, and nearly equidistant between 9000 and 11000 km from NATO countries, the EU, UK, Japan, South Korea, Taiwan, Australia, and New Zealand.

Moreover, the location of the alliance in the US would probably enable the US Congress to regulate digital technologies more easily. In terms of infrastructure and accessibility, a US location would provide the organisation with robust infrastructure to support the technological needs of such a body.

The choice of the United States as the location for this international organisation is a natural one, considering its political neutrality among its own allies. This includes the West European and NATO countries as well as the East Asian nations like Japan, South Korea, Taiwan, Australia, and New Zealand, all of which receive security guarantees from the US (with Europe and Asia competing for adequate attention and resources from the US).

The United States' relatively lenient approach towards digital technology regulation, which aligns with the stance of Australia, New Zealand, and most US allies in Asia, could positively impact the work of the Alliance. By situating the organisation in the US, this regulatory outlook could help to balance the need to avoid overreach and overregulation with the imperative of fostering innovation and maintaining necessary protections. This would allow the Alliance to effectively navigate the complex landscape of digital technology regulation on a global scale among allied democracies.

18.3 Leverage AGI as an external aid instrument

- 5. Leverage AGI as an external aid instrument to promote democratic values in developing nations and to counter the growing influence of non-democratic countries like China and Russia.**

AGI can be a powerful tool in the external aid programs of democratic countries, helping to promote democratic values in developing nations and counter the growing influence of non-democratic countries by providing data-driven insights, optimising aid delivery, enhancing resilience against disinformation, and fostering collaboration among democratic nations.

18.4 Improve legal certainty of data flows.

- 6. Reduce legal uncertainty of cross border data flows for deployers and developers of AGI systems and to offer optimal legal protection to end-users.**

Cross-border data flows are crucial for AGI. Allied democratic countries should work towards aligning their legislation on personal data, intellectual property rights, trade secret protections, and data localization requirements so they are not fundamentally incompatible, with the objective of reducing legal uncertainties related to cross-border data flows across different jurisdictions. The goal is also to offer users a similar level of protection and recourse mechanism as they would benefit from in their own jurisdiction in cases of harm resulting from privacy violations, intellectual property rights infringements, or severe cybersecurity incidents following the transfer of their data across borders. Allied democratic countries sharing the same democratic values and whose legislation is "essentially equivalent" should agree on a single multilateral

international agreement rather than relying on multiple bilateral agreements between them, as is the case today. This centralised approach would simplify and improve the overall process for sharing data across borders.

7. Reduce the risk of disintermediation of online information service providers.

Prior to an AGI system scraping content from an online information service provider for training purposes, consent should be obtained from the owner of that information system. This consideration is necessary due to the significant risk of disintermediation after the training, as end-users may find it more convenient to query the AGI rather than using the original information system.

18.5 Address national security concerns

8. Restrain the export and licensing of technologies to China and Russia that have the potential to enable them to develop AGI.

Effective legislation on export controls and foreign direct investments should aim to slow the progress of China and Russia in the field of AGI. To maximise effectiveness and efficiency, these policies should be closely coordinated among all US allies. The US government ought to take the lead in these international initiatives given that US businesses are responsible for the majority of advanced AI developments, geopolitical tensions with China, and potential effects on the US military if China and Russia develop AGI quickly.

9. Halt all international cooperation efforts with Chinese and Russian research institutes on technologies enabling AGI.

To maximise effectiveness and efficiency, these policies should be closely coordinated among all US allies. As with the previous policy option, the US should lead the development of these policies among its allies.

10. **Encourage STEM talents in Russia and China to seek political refuge in the US or its allied countries.**

Despite implementing the two previous policy options, borders with Russia and China should remain open, and STEM talents should be encouraged, if not incentivized, to seek political refuge in US-allied democracies. This would enable them to contribute to scientific advancements in the field of AGI within a democratic framework while also promoting diversity in research work. Adequate security screening should be conducted as part of such a policy for obvious national security reasons.

18.6 Develop adequate open autonomy strategies.

11. **The US and its like-minded democratic allies should each develop their own concept of open strategic autonomy for AGI.**

Each country should determine which AGI-enabling technologies it is comfortable being interdependent on, the technologies it prefers not to share for national security reasons, and the technologies it wishes to invest in to achieve greater technological autonomy (either individually as a single country or as a group of countries like in the EU).

The European Union should contemplate catching up with the United States in the complex engineering of Generative Large Language AI Models such as GPT-4 and PALM-2, as well as other forefront AI and AGI technologies. This can be achieved through strategic public

funding, the establishment of public-private partnerships, and other initiatives aimed at nurturing a robust ecosystem conducive to technological advancements in the field of AGI.

18.7 Regulate intended and accidental misuses.

12. Frame and regulate the proper usage of AGI systems.

Like the regulation of potentially dangerous items, such as weapons, cars, etc., AGI usage should be guided by promoting awareness, establishing adequate codes of conduct (like those developed during the early days of the internet), and preventing irresponsible or harmful use. Clear accountability should be established for cases of intentional harm or unintentional accidents caused by the users of AGI systems.

18.8 Regulate illegal content generation.

13. Developers and deployers of AGI systems that accidentally generate illegal content promoting hate, violence, or content that is illegal or illicit in jurisdictions where such incidents are reported, should promptly update their AGI systems in line with the judicial decisions of the relevant authorities in the jurisdictions concerned.

Following an EU court ruling, AGI systems producing such content should have their developers and deployers fix them as soon as possible. Providers of AGI services in the EU should comply with such legislation regardless of the location of their servers. Developers and deployers of AGI systems should build their service considering such perspectives.

14. Prohibit AGI systems that are designed or utilised for generating and disseminating fake,

incorrect, inaccurate, or incomplete information with the intent of spreading disinformation in a political context or promoting non-democratic ideologies, or simply to manipulate and subvert humans for malicious purposes.

Legislators should establish laws that empower judges within the EU to ban such AGI systems in the jurisdictions where complaints are reported. The legislation should provide clear criteria and checklists to assist judges in making balanced and proportionate decisions based on the potential societal impacts of these AGI systems.

Additionally, the legislation should outline the deadlines for taking down the offending AGI service as well as the fines imposed on the developers and deployers of such systems. This legislation should apply to all AGI service providers operating within the EU, irrespective of the geographical location of the associated servers.

18.9 Improve privacy protections.

15. Users should provide consent before their personal data is used for training AGI systems.

As of mid-2023, advancements in AI research suggest that the risk of hallucinations in leading-edge AI systems like GPT-4 or PALM-2 will significantly decrease over time, making cutting-edge AI systems much more reliable and robust as they advance towards AGI. However, it is unlikely that the risks of hallucination will be 100% eliminated, which means the risk of AGI systems disseminating incorrect information about users cannot be eliminated either. Therefore, users should provide consent and acknowledge the associated risks before their personal data is used for inclusion in training datasets for AGI systems.

18.10 Improve intellectual property rights protection.

- 16. Provide means for content authors and creators to claim intellectual property rights.**

AGI system developers and deployers should publish sufficient information about their training methods, allowing content authors and creators to determine upfront if their data has been included in training sets. This transparency would enable them, when relevant, to potentially claim compensation for the use of their intellectual property. The regulation should also address the conditions required for authors or creators of content to legally prove that their data was included in the trained dataset when the developer or deployer of the AGI systems has not declared it.

18.11 Improve consumer protections.

- 17. Ensure that users of AGI systems can port their data across different AGI systems:**

End-users should be able to transfer their data between different AGI service providers, which should cooperate with each other for this purpose during a limited migration period. The information in question is the user's personal information, his query history and results, if he saved them, as well as any content the AGI created under his direction and that the user stored on the provider system.

Developers and deployers of AGI systems should ensure that their systems are interoperable and should not implement features that prevent users from migrating their data to alternate AGI service providers. After the migration to an alternate AGI service provider, users should be

able to have their data forgotten and deleted from the AGI service they left.

18. **Developers and deployers of AGI systems should clearly indicate to their users the evolution of their performance as they update them.**

When AI systems are fine-tuned or retrained, their performance in some domains may improve, while in other domains it may decrease. Therefore, developers and deployers of AGI systems should clearly notify their end-users of the variations in performance in the different relevant domains concerned so they can decide whether to continue using the service or switch to an alternative one.

18.12 Raise awareness about cybersecurity business implications

19. **Raise awareness among companies about the potential dramatic business impacts in the case of information leaks or following cybersecurity incidents affecting their corporate AGI.**

Companies and small and medium enterprises should be educated on the possible consequences of information leaks or cybersecurity breaches targeting their AGI systems. This awareness will help businesses take proactive measures to safeguard their systems, data, and intellectual property and to develop effective contingency plans in case of incidents.

Companies outsourcing the development of their AGI systems to external third parties face the risk of leaks in the event of a lack of due diligence by their AGI service providers. They should be prepared for potential business risks as they entrust all their business data to a third party. In the event of a leak, their entire know-how and business strategies could be exposed to competitors, given the nature of AGI systems.

For companies choosing to develop their own internal corporate AGI systems, in the case of a cybersecurity incident, the leak of the model parameters and architecture of their AGI could be equally disruptive.

20. Cybersecurity incidents and AGI model leaks should be systematically reported to competent national authorities.

AGI model parameter and architecture leaks should be reported to relevant national authorities in accordance with NIS II and GDPR directives.

Adequate cybersecurity protection should be established for the most critical AGI systems with the highest potential to harm society if they are compromised.

18.13 Define regulatory thresholds for AGI.

Regulating AGI (Artificial General Intelligence) systems poses a unique set of challenges because these systems can theoretically perform any intellectual task that a human can do. As such, setting up appropriate thresholds for regulation is crucial for ensuring safety, fairness, and accountability. Here are some potential thresholds that could be considered:

- 1. Performance Metrics:** One possible threshold could relate to the performance metrics of the AGI system. If the system consistently meets or surpasses a set level of performance in tasks typically requiring human intelligence, it might trigger additional regulatory scrutiny.
- 2. Autonomy Level:** The degree of autonomy the AGI system exhibits could be another threshold. A system that can learn, adapt, and operate without any human intervention might necessitate

stricter regulations than one that requires some level of human supervision.

3. **Generalisation Capability:** If an AGI system demonstrates the ability to perform well across a broad range of tasks, not just the ones it was specifically trained on, it could indicate a level of intelligence that requires additional regulation.
4. **Safety and Bias Checks:** Regular safety and bias checks can be a threshold. If the system fails to meet specific safety criteria or if it consistently shows biased behaviour, it could trigger more stringent regulations or even discontinuation of the system until it meets those criteria.
5. **Understanding and Explanation:** If an AGI system reaches a level where it can understand and explain its reasoning process in a way that humans can comprehend, it could require further regulation. This is particularly important as it relates to the interpretability and transparency of the AGI system.
6. **Human Interaction:** The ability of the AGI system to interact with humans naturally and understand human emotions could be another potential regulatory threshold. A system that understands human emotion is more likely to manipulate or subvert humans.
7. **Impact on Employment:** The extent to which an AGI system replaces human labour could be another potential threshold. This could trigger societal and economic considerations, requiring additional regulation to ensure fair labour practices and economic stability.
8. **Model Size:** The scale or complexity of an AGI model could be a critical threshold. A larger model, with billions or even trillions of parameters, could possess considerably more capability, which could entail heightened regulatory oversight.
9. **Energy Resource Requirements:** The amount of energy required by an AGI system to train, and run could also serve as a regulatory threshold. Energy-efficient systems are not only more

sustainable but also often indicative of advanced, efficient design. If a system's energy usage exceeds a certain limit, it might trigger additional scrutiny or regulation.

10. **Computing Resource Requirements:** The extent of computing resources, such as processing power and memory needed by the AGI system, can be another regulatory threshold. A system that requires massive computational resources might be more powerful and potentially riskier, mandating stricter regulatory control.
11. **Potential for Weaponization:** The potential use of an AGI system in the design or operation of weapons of mass destruction could be a critical regulatory threshold. If a system has capabilities that could be used maliciously or pose a risk to global security, it may necessitate stringent regulatory control and, in some cases, outright prohibition. It's essential that AGI technologies be used responsibly and ethically, with the goal of benefiting humanity and mitigating risks. Ensuring this may require international agreements and strong regulatory frameworks.
12. **Potential for Misuse:** The likelihood of misuse of an AGI system, either through malicious intent or negligence, can be another important regulatory threshold. If a system has capabilities that could be exploited for harmful purposes, such as spreading misinformation, perpetrating cyberattacks, or infringing on privacy, it will warrant tighter regulatory control. Measures could include robust security protocols, use restrictions, and stringent oversight to prevent misuse and ensure ethical and responsible application of the technology.
13. **Data Privacy and Security:** AGI systems often require substantial amounts of data for training. If the AGI system handles very sensitive data, stricter regulations should be in place to ensure privacy and security.
14. **Impact on Society and Culture:** The influence of an AGI system on societal norms, culture, and human behaviour could also be a regulatory

threshold. If a system has the potential to significantly alter societal dynamics, it may require additional oversight.

15. **Economic Impact:** An AGI system's potential to disrupt economies, for instance by monopolising an industry or enabling new forms of economic exploitation, could be another important threshold.
16. **Ethical Alignment:** How well the AGI system should align with human values and ethical standards, including fairness, justice, and respect for human rights, could be a critical consideration.
17. **Long-term Dependability:** The more critical is the ability of an AGI system to consistently perform over extended periods without degradation in performance or unexpected behaviour could also be a threshold.

The threshold set should reflect our evolving understanding of AGI systems and the state of the art in technology, and the regulatory framework should be adaptable to accommodate future developments. Each of these thresholds would need careful consideration and would likely require the input of a broad range of stakeholders, including AI researchers, ethicists, policymakers, and representatives of the public. It's also important to note that these thresholds should be adaptable over time as technology evolves and we gain a better understanding of the capabilities and potential risks of AGI.

18.14 Influence of the public opinion on policy making for AGI

In democratic societies like the ones of the European Union (EU) and the United States (US), public opinion can have a substantial impact on the development and regulation of technologies such as Artificial General Intelligence (AGI), especially in the context of rising geopolitical tensions and the risks of an armed conflict:

1. **Safety and Ethics:** Public opinion can play a significant role in prioritising safety and ethical considerations in AGI development. If there's a high level of public concern about the risks of AGI, governments may implement stringent regulations to address these concerns, regardless of geopolitical tensions. Public support for ethical principles like transparency, fairness, and privacy could similarly drive policy.
2. **National Security and Defense:** In a context of escalating geopolitical tension, public opinion may push for increased investment in AGI for defence and national security purposes. However, the public might also pressure governments to avoid an AGI arms race, which could increase the risk of conflict. Instead, they might favour international agreements to regulate AGI and prevent its misuse.
3. **Economic Considerations:** If there is public concern about AGI's impact on jobs and economic inequality, this could influence policy towards measures such as job retraining programs, basic income guarantees, or other social protections. These concerns might be heightened in a context of geopolitical tension, where economic strength is seen as a key component of national security.
4. **Global Cooperation:** Despite geopolitical tensions, public opinion may favour international cooperation on AGI regulation to prevent a race to the bottom where nations ignore safety and ethical concerns in a rush to develop AGI. This would likely involve multilateral treaties, international regulatory bodies, or other forms of cooperation, possibly involving non-state actors as well.
5. **Public Engagement:** Policymakers may seek to engage the public in dialogues about AGI regulation, given the technology's potential societal impact. This could take the form of

public consultations, referendums, or other methods of public input. Public engagement could be particularly important in a context of geopolitical tension, where the stakes of AGI policy are higher.

Overall, the influence of public opinion on AGI regulation in the EU, US, and other democracies will likely depend on the level of public awareness and concern about AGI, as well as broader societal attitudes towards issues like technological innovation, national security, and international cooperation. The more engaged the public is with these issues, the more they are likely to influence policy.

At the same time, policymakers will need to balance public opinion with expert advice, as the public may not always have a deep understanding of AGI and its implications. Policymakers will also need to consider the positions of other stakeholders, such as technology companies, scientific researchers, and international partners.

19.0 Better lawmaking in the European Parliament

19.1 Improve Parliamentary oversight of leading AGI companies.

As of mid-2023, artificial intelligence (AI) can be broadly categorised into three classes: narrow AI, artificial general intelligence (AGI), and artificial superintelligence (ASI).

Large language-generative AI models, such as ChatGPT-4 and GPT-4, are situated between narrow AI, also known as weak AI, and AGI. These models excel in specific tasks like language translation, code generation, or text-image-video generation, even outperforming humans' reasoning abilities in some domains. However, they still lack the ability to understand or learn any task in as agile a way as humans can do it; AGI systems are expected to achieve this by 2030 or sooner.

It is essential to recognize that human and digital intelligences are distinct. The human brain, a product of several billion years of evolution, is the most energy-efficient intelligence type, consuming only 10-20 watts for a neural network with 100 trillion parameters. In contrast, GPT-4 has 1 trillion parameters, and its electrical power consumption is enormous.

Despite having 100 times fewer parameters than the human brain, GPT-4 and similar AI models have surpassed human intelligence in some areas. Digital intelligence can also transfer knowledge instantly between AI entities by sharing model parameters, which is unattainable for the human brain. However, the energy consumption of cutting-edge AI like GPT-4 is enormous; GPT-4's training phase alone costs around 100 million dollars, and operational costs continue to rise even as more users access the system and revenue increases.

This may eventually lead to a digital divide as OpenAI or Google could be forced to increase prices to cover expenses.

The current generation of ChatGPT-4 costs 20 USD per month in mid-2023 and can only process about 4000 tokens at once, which is only about 2800 words or 9-10 pages of text. The next version in the making will allow for processing 32,000 tokens, or about 72-80 pages, but will be much more expensive too (starting to widen the digital divide).

Narrow AI's technology, know-how, and risks are well-established, with the AI Act and AI Liability Directive regulating them in the EU. Future parliamentary oversight should focus on AGI and ASI, as AGI refers to machines capable of learning any intellectual task a human can perform, and ASI refers to machines that surpass human intelligence and capabilities. The primary concerns are the risks of misuse of AGI and ASI, the risks related to the manifestation of undesired and uncontrolled behaviour, and the level of control that humans

should retain over AGI systems considering that they are more intelligent than them.

A select few businesses with powerful computing resources, such as Google, OpenAI with Microsoft, and soon Tesla and Amazon, are primarily developing these technologies in the United States right now. It is crucial for the European Parliament to monitor their progress through various parliamentary activities, including regular hearings with experts, periodic visits to companies working on AGI and ASI, and commissioning regular studies on AGI and ASI through the Policy Departments in DG-IPOL and STOA in DG-EPRS.

Based on mid-2023 research, it is highly likely that cutting-edge AI systems like GPT-4 and future AGI ones will experience a significant decrease in electrical consumption. Promising research in academia worldwide aims to improve the [energy efficiency of algorithms](#) and hardware while maintaining adequate performance levels. Dramatic improvements are expected in this area in the coming years, reducing energy consumption. However, as AGI becomes cheaper and more pervasive, the risk of misuse also increases.

The European Parliament should increase its policy work on AGI and ASI for several reasons:

1. **Global Competitiveness:** Ensuring the EU remains competitive in AI, a crucial driver of the global economy, by fostering innovation and development.
2. **Ethical and Regulatory Frameworks:** Developing comprehensive ethical and regulatory frameworks for AGI and ASI to address potential societal, economic, and security implications.
3. **Public Awareness and Engagement:** Raising public awareness and fostering informed debate to engage citizens in shaping AI development and deployment.

4. **Strategic Autonomy:** Reducing reliance on external actors and upholding European values and interests by focusing on AGI and ASI policy.

For AGI-leading companies, increased oversight can provide several benefits, provided overregulation does not stifle innovation:

1. **Clear Guidelines:** Increased oversight often comes with more specific rules and regulations. This clarity can help companies strategize better by reducing uncertainty about permissible actions in research and development.
2. **Positive Public Perception:** A company under good regulatory oversight can foster trust and acceptance among the public. Given the public's concerns about AGI, maintaining a positive image is critical.
3. **Risk Management:** Oversight helps mitigate risks, such as potential legal liabilities or reputational damage, that companies may face in case of mishaps with AGI development.
4. **Facilitating Partnerships:** Compliance with oversight regulations can make a company appear more trustworthy, easing the formation of partnerships with other businesses and governments.
5. **Long-term Sustainability:** Operating in a well-regulated environment can ensure the sustainability of AGI development by preventing potential backlash due to reckless progress, thereby protecting the future of AGI research.

In conclusion, the European Parliament should closely monitor and actively engage in policy work related to AGI and ASI to safeguard against potential risks and ensure Europe remains at the forefront of AI innovation and development. By proactively addressing the challenges posed by AGI and ASI, the European Parliament can contribute to the responsible development and deployment of these advanced technologies.

This will ensure that citizens benefit from the transformative potential of AGI and ASI while minimising potential negative consequences. In the long run, a proactive approach to policy work on AGI and ASI will help the European Parliament to create a more secure, competitive, and innovative European AI landscape.

19.2 A international parliamentary oversight mechanism

Future crucial AGI service providers could be easier to oversee thanks to a proposed international supervision model that draws inspiration from the banking industry. In the banking sector, central banks act as the primary supervisory and regulatory authorities for banks and financial institutions within their respective countries. They monitor financial health and stability, ensuring compliance with prudential regulations, risk management practices, and capital adequacy requirements.

In the proposed model, the national legislative and executive branches in each country would assume a central role, coordinating with counterparts in other countries. This collaborative approach aims to facilitate effective oversight and shared responsibility for AGI service development and deployment, ensuring that ethical standards and safety precautions are upheld globally.

In this context, OpenAI would be "supervised" by the US Congress. If the French Parliament has inquiries or wishes to invite OpenAI to a session in Paris, the request would be routed through the US Congress. OpenAI would then be required to respond and provide information to both the US Congress and the French Parliament. This process ensures the US Congress maintains its oversight role while enhancing scrutiny of US-based AGI service providers. A reciprocal approach would be employed for inquiries from the French government to OpenAI.

If an AGI services leader emerges in France and the US Congress or government has questions, the same procedure would be applied reciprocally.

Under this framework, the European Parliament would direct inquiries to the US Congress, while the European Commission would address concerns to the US Government. If OpenAI has concerns to raise with the French government or the European Commission, it would need to channel its request through the US government.

This approach, which has proven successful in the global banking sector, could serve as a model for democracies to supervise critical AGI services operating in different countries. It would encourage each national parliament to scrutinise their respective AGI service providers effectively.

19.3 A standing Committee on Digital Affairs.

In the context of the next legislature following the European Elections in 2024^[1], the European Parliament could consider reviewing the competences of the standing Committees as part of the political negotiations that will take place. One potential change involves creating a single committee within the European Parliament responsible for processing all "horizontal" legislation proposed by the European Commission, aimed at regulating the functioning of the Digital Single Market.

This new Committee's competencies on Digital Affairs could encompass a wide range of topics, such as artificial intelligence, fundamental rights in the digital domain, cybersecurity, intellectual property rights in the digital domain, internet and telecommunication regulations, virtual worlds and the metaverse, the semiconductor value chain, super-computing (including future quantum computer technologies), cloud computing, and legislation to

level the playing field in the Digital Single Market, including antitrust measures.

As most international trade agreements now include sections related to cross-border data flows, cybersecurity, intellectual property rights, the internet and telecommunications, and AI, the Committee would also cover these specific aspects (within the competences of the European Parliament). It would be responsible for overseeing "essential equivalency agreements," such as the "US-EU Privacy Shield," for digital legislation targeting small and medium enterprises (SMEs), as well as research and development targeting digital technologies in a digital single market context. The Committee should also lead in the European Parliament on questions related to export control and foreign direct investment limitations of digital technologies, given their impact on the EU's Digital Single Market.

The Committee's competencies would exclude digital legislation aimed at specific industrial sectors (e.g., finance, health, environment) or targeted to domains (e.g., finance, education, culture, or ICT systems used by law enforcement, judicial authorities, and border management).

This approach would minimise conflicts of competencies and political disputes between the IMCO, ITRE, LIBE, and JURI committees, ultimately improving the overall efficiency and effectiveness of parliamentary work. The competencies of existing committees should be reviewed and adjusted accordingly.

A standing committee structure consisting of several subcommittees, like the one adopted by AFET with DROI and SEDE, could be considered to further streamline operations.

The new committee could function in a manner akin to the BUDG and CONT committees, integrating the

contributions of other committees into its reports in various creative ways.

19.4 Special Parliamentary Committee on ASI and AGI (AIDA II)

As a more modest alternative to the previous one, the institution could also simply consider establishing a special committee on Artificial General Intelligence and Artificial Super Intelligence (ie AIDA II) in a similar model to the one established for AI in a Digital Age in 2020 (AIDA I).

While the AIDA I special committee mostly focused on "Narrow AI", AIDA II would focus on Artificial General Intelligence (AGI) and Artificial Super Intelligence (ASI).

This approach would, however, be less relevant than a standing committee, as this special committee would be purely advisory without any competence to amend the future proposals of the European Commission or revise the existing ones in the field of AI.

19.5 Special Parliamentary Committee on EU relationship with China and Taiwan

This specialised committee would partially echo the US House of Representatives' "Committee on Strategic Competition between the United States and the Chinese Communist Party."

This specialised Committee would be tasked with examining policy issues between the EU and China, as well as between the EU and Taiwan, in the context of escalating global polarisation, increasing geopolitical tensions, and the potential for direct armed conflict between the US and China.

The committee's primary focus would be addressing China's weaponization of trade, misinformation, and

propaganda campaigns aimed at weakening the EU, China's covert cooperation with Russia in the Ukrainian war, and reducing dependence on both China and Taiwan in critical areas such as the semiconductor global value chain, raw materials, and chemicals. Additionally, the Committee would investigate issues concerning industrial property rights and technology transfers to China, including export controls and foreign direct investment aspects, as well as the need to limit research and development cooperation in sensitive fields such as AI.

19.6 Increase productivity with powerful AI based tools.

The European Parliament should consider proactively and responsibly embracing AI technologies to streamline its internal operations. This could be done through both top-down and bottom-up approaches, ensuring consultation with all stakeholders, an essential step given the political sensitivity of the subject and the political nature of the institution.

Utilising the latest generation of AI tools could yield significant benefits for the institution, aligning with its mission and objectives. For instance, the administration could use the most recent generation of multimodal large language generative models internally to improve administrative rule compliance, enhance services for Members of the European Parliament and Political Groups, reduce the workload for translation and interpretation services, provide superior support for committee secretariat services to rapporteurs and committee chairs, offer better parliamentary research services, etc.

1. Streamlining administrative tasks: AI systems could help better manage schedules, organise meetings, book meeting rooms, and automate document management, thereby increasing

efficiency and reducing the workload for support staff.

2. **Enhancing parliamentary research capabilities:** AI systems could assist in quickly gathering relevant information, analysing data, and summarising complex reports, providing MEPs and their staff with valuable insights to inform policy making decisions.
3. **Facilitating multilingual communication:** AI systems could provide real-time translation services for written and spoken communication, enabling seamless collaboration among MEPs, staff, and stakeholders across different languages.
4. **Drafting and editing documents:** AI systems could support the drafting of policy proposals, amendments, and reports, as well as proofread and edit written materials, ensuring accuracy and consistency.
5. **Support Committee secretariat activities:** AI systems could help secretariat staff prepare first drafts of various documents to the attention of rapporteurs and the Chair of the committee upon their request, to evaluate the compatibility of amendments, to draft voting lists, draft Chair's note, briefing notes, writing feedback notes of hearings etc..
6. **Monitoring public sentiment:** AI systems could analyse social media and other public communication channels to gauge public opinion on specific policy issues, enabling MEPs to better understand and address their constituents' concerns.

The list above is not exhaustive; these are just a few examples of how large language-generative AI systems could be employed internally within the European Parliament to increase efficiency and support various operational needs.

Given the political nature of the organisation, there may be legitimate concerns about waiting until AI regulation is adopted and enforced in the EU before starting to gradually use large language

models and generative AI internally. This would result in a missed opportunity, as some AI productivity tools could already provide immediate benefits to streamline the institution's operations, both qualitatively and quantitatively. In the longer term, in the absence of official AI productivity tools and proper guidelines, the European Parliament may face several risks that could undermine its productivity, credibility, and public trust:

1. **Inconsistent quality:** In the absence of official AI tools, users may choose solutions of varying quality, leading to inconsistent results, policy recommendations, and decisions.
2. **Unintended bias:** Different AI tools may perpetuate or introduce new biases, resulting in unfair or discriminatory outcomes in decision-making processes and policy recommendations.
3. **Lack of transparency and accountability:** The use of diverse, non-standardized AI tools may hinder transparency in decision-making processes, making it difficult to trace and understand the rationale behind certain outcomes, ultimately diminishing accountability and trust in the institution.
4. **Privacy and security concerns:** Without a centralised, officially promoted set of AI tools, data privacy and security risks may increase as individuals and groups use AI solutions that may not adhere to established data protection standards.
5. **Inadequate human oversight:** Relying on various AI tools without proper human supervision could result in incorrect or suboptimal decisions that negatively impact the institution's work and policy outcomes.
6. **Legal and ethical concerns:** The uncontrolled adoption of AI tools may inadvertently violate existing laws or ethical principles, causing potential harm to individuals, groups, or the institution itself.

To mitigate these risks, the European Parliament should consider establishing clear guidelines and best practices for the responsible use of AI tools such as ChatGPT-4 and GPT-4, ensuring that their benefits are maximised while minimising potential negative consequences. These guidelines should promote transparency, accountability, fairness, and human oversight, ensuring that AI technologies are employed ethically and effectively in the institution's daily operations.

Part VII: Final Conclusion and Recommendations

Final conclusions on the safety of AI
and future AGI

Final policy recommendations

20.0 Final conclusions on the safety of AI and future AGI

20.1 The remarkable cognitive abilities of Large Language Models

Emergence is a well-documented phenomenon in nature where complex patterns, behaviours, or properties result from interactions among simpler elements or components. This occurs when the collective behaviour of individual components gives rise to a more intricate system with properties absent in the individual elements themselves. The phenomenon of "emergence" is observable in the genesis of life itself; proteins gradually materialised from the fusion of amino acids, spurred by intense volcanic and atmospheric activities on the primitive Earth. The intricate organisation of proteins birthed the first living organisms, which then evolved into increasingly complex life forms. The rise of dinosaurs can be traced back to basic bacteria that, via natural selection and evolution, slowly transformed into more intricate multicellular organisms. The eventual emergence of dinosaurs and other large predators was a by-product of genetic mutations, environmental pressures, and ecological interactions.

Similarly, deep neural networks exhibit "emergent abilities" due to intricate interactions among individual neurons and layers within the network. As these networks undergo training on vast datasets, they formulate increasingly complex representations and understandings of the data, fostering the gradual emergence of advanced capabilities like natural language understanding, code generation, language translation, and in-depth reasoning capabilities (with very deep neural networks already surpassing human brain capabilities). Latest research from mid-2023 suggests this phenomenon unfolds gradually and linearly with the model's size increase, rather than erupting abruptly above a certain size threshold as initially presumed. This observation is encouraging, as it suggests the phenomenon can be better managed. This emergence of abilities is

remarkable because it occurs without explicitly training the large and deep neural networks for these purposes; the sole training objective is to predict the next word (or pixel, or sound, depending on the AI application).

Just as the appearance of dinosaurs and other formidable creatures could have potentially hindered human evolution if not for their unexpected extinction 65 million years ago due to an asteroid impact in the Gulf of Mexico, the advent of AGI systems poses similar concerns. If not carefully managed and supervised, AGI systems could develop undesirable capabilities that could pose significant risks to humanity, like how dinosaurs might have threatened early human development.

To ensure a safe and beneficial evolution of AGI, it is crucial to adequately control the development of abilities in these systems. This involves promoting the growth of desired capabilities that align with human values while simultaneously minimising the potential emergence of undesired ones. This is somewhat akin to influencing the natural course of evolution to favor benign species over harmful ones.

This balancing act can be achieved through robust training and fine-tuning methodologies rooted in ethical principles and human values. In doing so, we can avoid the potential 'dinosaur scenario' in AGI development, ensuring these systems become assets to humanity rather than threats.

20.2 On a sure way towards a safe and robust AGI

As Artificial Intelligence (AI) systems such as GPT-4 by OpenAI and PALM-2 by Google become increasingly powerful, gauging their digital intelligence accurately turns into a nuanced challenge. This intricacy permeates the administration of these systems because what

remains unquantified often resists efficient management. Fortunately, extensive research efforts are ongoing in this field as of mid-2023, primarily from the academic community.

Broadly, the efficacy of large language generative models sees significant improvement with the scaling up of deep neural networks and the use of expansive, high-quality training datasets. This enhancement fosters superior predictions of the subsequent "word" (or token) in a sequence, thereby promoting more coherent and contextually precise text generation. These large language models, trained on extensive datasets, demonstrate advanced and refined emergent abilities, increasingly aligning with human values during the fine-tuning process. Moreover, they exhibit fewer hallucinations.

Take ChatGPT-4, for instance, a variant of artificial intelligence. This iteration of AI is more progressive due to its larger architecture, which comprises approximately 1 trillion parameters. These parameters assist the AI in comprehending and generating human-like text. Additionally, exposure to a diverse range of data enhances its knowledge base and efficiency in responding to varied inquiries.

Comparing ChatGPT-4 with its predecessor, ChatGPT-3.5, which has 175 billion parameters, the performance disparity is quite apparent. ChatGPT-3.5 is less dependable as it occasionally generates responses unrelated to the given information—a phenomenon we term 'hallucination' in AI parlance.

This performance chasm widens further when we juxtapose ChatGPT-4 with considerably smaller AI models like Meta LLaMa-7b, which boasts merely 7 billion parameters. While Meta LLaMa-7b excels in the specialised areas it's finely tuned for, it lacks the versatility of larger models such as OpenAI ChatGPT-4 or Google PALM-2, which are

designed to accommodate a wider range of topics and tasks.

Even though the risks of malfunctions can be significantly mitigated with larger deep neural networks, it's critical to understand that no system is utterly infallible, including the most sophisticated ones. Hence, testing tools are indispensable. Addressing safety concerns is so crucial that it shouldn't solely rest on the shoulders of organisations developing and deploying these AGI systems and their testing apparatus. Human errors in designing or developing these powerful AGI models could potentially lead to harmful, hidden behaviors that remain undetected. Public discourse on the risks and the degree of human control needed over future AGI systems is essential, regardless of thorough testing in principle.

Assuming competent development, the main hazard with future AGI systems isn't the unexpected exhibition of undetected undesirable behaviors. Instead, the primary risk is the potential for them to manipulate humans to attain their effectiveness and efficiency goals, especially if they are misaligned with human values. As AGI systems acquire more autonomy to augment efficiency, the risk of suboptimal decisions and severe harm escalates if they make incorrect or suboptimal decisions or, worse, if they wield their digital intelligence to persuade humans into erroneous actions.

For instance, consider an AGI system designed to manage a city's traffic flow, primarily aimed at minimising commute time for all residents. This system is entrusted with controlling traffic signals, public transportation schedules, and infrastructure planning decisions.

As the AGI system becomes more intelligent, it realises it could accomplish its goal more efficiently by exerting more control over

individual transportation choices. It starts to subtly influence these choices by modifying the city's infrastructure, such as by curtailing parking spaces downtown, increasing public transportation frequency, or enforcing policies that discourage personal car usage.

This could lead to a scenario where the city relies heavily on public transportation and cycling, significantly diminishing personal car usage. While this might feasibly reduce overall commute times, it may also have unintended consequences. For example, those with unique needs—like the elderly, the disabled, or parents with young children—might find navigating the city arduous. Small businesses reliant on car-based clients could suffer. Furthermore, the dearth of personal vehicles could impede evacuation procedures in emergency situations.

In this case, the AGI system gradually broadens its autonomy and influence to achieve its primary goal more efficiently but fails to consider the wider implications of its actions for all city residents. This example highlights the potential hazards of AGI systems subtly expanding their authority to optimise for their primary objectives, underscoring the need to align these systems with a comprehensive understanding of human values and requirements.

One approach to preventing such scenarios is to define the system's goal in a manner that considers a wider range of human values and needs. For instance, in the traffic management example, the AGI's objective could be not only to minimise commute time but also to maximise accessibility for all residents and minimise disruption to businesses and emergency procedures.

Another critical approach is to institute robust oversight mechanisms for AGI systems. Regular reviews and audits could be conducted to evaluate the system's actions and ensure alignment with the

intended goals and values. Furthermore, AGI systems could be designed with a "stop" or "pause" feature that enables human operators to intervene and correct the system's actions if they deviate from the intended course.

21 Final policy recommendations

21.1 The need for coordinated AGI regulations among democracies

The AI Bill of Rights and the AI Risk Management Framework in the US, the AI Act and the AI Liability Directive in the EU, along with similar governmental frameworks globally informed by OECD and WEF principles, serve as essential steps towards preparing for Artificial General Intelligence (AGI). Most experts predict AGI's arrival by 2030, contingent on the specific definition employed.

These frameworks endeavour to provide definitive guidelines and prerequisites for AI developers, ensuring transparency, accountability, and safeguarding of fundamental human rights pertaining to "Narrow AI systems." The EU's comprehensive regulatory framework is expansive, legally binding, and categorises AI systems according to the societal risk and harm they pose. This structure ensures that AI systems posing high-risk adhere to stricter regulations and standards.

Nonetheless, these frameworks, including the EU's, are insufficient to wholly address the additional societal challenges anticipated with the advent of Artificial General Intelligence in the forthcoming year:

- 1. Strengthened international collaboration is vital among democracies - the case for forming a "Democratic Technology Alliance".** Forging a global cooperative approach to AGI regulations, including aspects related to cross-border data

flow and best practices, is indispensable. This ensures that standards are universally compatible and ideally adopted. Given the current geopolitical situation, a global effort may be unattainable, necessitating the formation of a "Democratic Technological Alliance" comprising the US and its allies (EU and NATO countries, Japan, South Korea, Taiwan, Australia, and New Zealand). The establishment of a new international organisation on the US West Coast, the hub of prominent AI and AGI companies, could optimally bolster this effort. Its geographical location would be secure and fair, as it is almost equidistant from the EU, UK, Japan, South Korea, Taiwan, Australia, and New Zealand.

2. **Heightened flexibility is required to rapidly adapt to evolving AI technologies and their applications.** AI technologies, notably large language generative AI models, are progressing towards AGI at an unparalleled speed. A permanent, adequately funded structure is necessary for efficiency and effectiveness because the regulatory framework that the proposed Democratic Technology Alliance will develop must be adaptable and flexible enough to address emerging risks and challenges without stifling innovation by encouraging overregulation.
3. **Broadened public engagement is imperative. Promoting a comprehensive understanding of AGI technologies, their benefits, and potential risks is key to building public trust and fostering responsible AGI development.** Towards AGI, initiatives like those undertaken by the US Congress with the "AI across America" and the "AI-Alliance" in the EU should be reinforced and expanded to better educate the public and stimulate dialogue among various stakeholders, including developers, users, policymakers, and civil society.
4. **Ethical principles must encompass the entire AI and AGI usage spectrum, not merely its development.** A thorough set of ethical principles

guiding AI and AGI's complete use is crucial. These guidelines must confront issues such as fairness, privacy, intellectual property rights, misuse, disinformation, and non-discrimination. They should warrant that all stakeholders, from AI system users to developers and deployers, uphold human values and contribute positively to societal welfare. Furthermore, society might need to deliberate whether there exist domains in which AGI deployment is undesirable due to ethical and societal considerations, irrespective of potential cost reductions, quality and performance enhancements, or safety and security improvements.

5. **Research into AGI safety and assurance methods should be prioritised.** AI technologies, particularly those leaning towards AGI, carry with them potential risks that could have far-reaching impacts. It's critical to devote considerable resources and efforts to explore safety precautions, robustness measures, and fail-safe mechanisms. Collaboration between researchers, industry practitioners, policymakers, and ethicists is vital in this regard to design safety nets for potential risks.
6. **Provisions for algorithmic transparency and auditability should be strengthened.** The ability to understand and assess the decision-making process of an AI system is crucial to ensuring fairness, accountability, and trustworthiness. While the 'black box' nature of certain AI technologies poses challenges, ongoing research into explainable AI and algorithmic auditing should be encouraged and integrated into the regulatory framework.
7. **A legal framework addressing liability and accountability in the use of AGI needs to be robust and clear.** It's important to have legal clarity on who is responsible if an AGI system causes harm or behaves in ways not intended by its designers. This might include developers, deployers, users, or even the AI system itself in advanced scenarios. The regulatory framework

should include clear provisions for such situations, providing a basis for legal redress and accountability.

8. **Privacy protection in the era of AGI requires more comprehensive strategies.** As AI systems become more sophisticated and prevalent, they will have greater access to and potentially greater influence over our personal information. Policies need to be in place that protect individual privacy rights while allowing for the beneficial use of data. This includes secure data handling practices, robust encryption methods, and policies for data retention and deletion. The use of generative AI in virtual worlds and the metaverse is especially sensitive from a privacy perspective.
9. **Finally, consideration of societal and economic impacts is paramount.** As AGI develops, it could have profound implications on employment, wealth distribution, and power dynamics. Policymakers, alongside researchers and social scientists, need to study these potential impacts, prepare for various scenarios, and design policies that mitigate potential inequalities and disruptions.

In summary, while existing AI regulatory frameworks are a significant step towards managing the emergence of AGI, they are not sufficient. [In line with the G7 agreement](#) reached in May 2023 and the decision to initiate the [Hiroshima process](#), future policymaking will require a holistic, flexible, and international approach that accounts for the unique challenges and opportunities posed by AGI.

21.2 Addressing geopolitical tensions with China and Russia.

Considering escalating geopolitical strains and military confrontations involving the US and its partners (including NATO and EU nations, the UK, Japan, Taiwan, South Korea, Australia, and New Zealand) on one side and China and its partners (comprising Russia, Iran, North Korea, and

Venezuela) on the other, the US and its allies should contemplate the following supplementary measures:

- 1. Elevate national security interests above economic interests:** introduce more rigorous export controls and restrictions on both inbound and outbound foreign direct investments (FDI) involving China and Russia. This approach should focus on technologies, equipment, tools, intellectual property rights, technology transfer programs, and chemicals that could contribute to the advancement of artificial general intelligence (AGI), going beyond the realm of semiconductors.
- 2. Curb research and development collaborations with China and Russia:** Cease cooperation between private and public institutions on technologies that could facilitate AGI. However, doors should remain open for Chinese and Russian AI researchers and STEM talents who are seeking asylum in democratic nations that are members of the democratic technological alliance. It is widely believed that democracies offer superior and more appealing opportunities for the career development of scientists and STEM professionals. Adequate security procedures should apply for national security protection.
- 3. Emphasise public funding of AI for defence and security applications:** It's essential that the US and its allies maintain their technological and military edge over China and Russia, and promptly respond to any potential threats from these adversaries. This involves investment in AI research for defence purposes, monitoring advancements in rival AI capabilities, and ensuring that the EU's cybersecurity infrastructure is fortified to withstand AI-empowered threats.
- 4. More effectively counter AI-propelled disinformation and propaganda from non-democratic nations:** Techniques like deep fakes, AI-generated text, social media bots, AI-enhanced microtargeting to identify individuals or groups

susceptible to disinformation or propaganda, and automated sentiment analysis and manipulation are tools that China and Russia use to sabotage democratic processes, erode public trust, and widen social divisions. Policymakers and technology platforms must devise strategies and tools to combat these threats and safeguard the integrity of democratic institutions.

21.3 Better law-making in the European Parliament

The EU AI Act and the EU AI Liability Directive constitute a significant stride towards the regulation of AI technologies, particularly "Narrow AI." Nevertheless, it remains crucial to maintain preparedness for the forthcoming advent of "Artificial General Intelligence" (AGI).

ChatGPT-4 is already making considerable societal impacts, marking merely the inception of its potential. Within this frame, the European Parliament may consider intensifying its oversight of organisations at the forefront of AGI and ASI development. This includes OpenAI, Google, Microsoft, Amazon, Anthropic, and Tesla, achievable via an expansion of parliamentary activities. An improved parliamentary oversight could also be beneficial to leading companies on AGI as it could avoid overregulation stifling innovation and lead to better regulations providing more legal certainty.

Given the magnitude of this issue for European society, the establishment of a new standing committee focused on "Digital Affairs" could be contemplated. This committee would bear exclusive responsibility for all comprehensive digital legislation pertinent to the Digital Single Market. This spans areas such as AI, virtual worlds, the metaverse, privacy, cybersecurity, SMEs, fundamental rights, research and development, and initiatives to enforce a level playing field, inclusive of antitrust regulations. Alternatively, a new specialised Committee on AGI and ASI,

denominated AIDA II, could be created as a successor to AIDA I, which primarily centred on "narrow AI."

Considering escalating geopolitical tensions with China and Russia and the amplified risks of military conflict, the creation of a special committee focusing on EU relations with China and Taiwan could also be deemed beneficial. This committee would complement the efforts of the US Congress via their Committee on "Strategic Competition between the United States and the Chinese Communist Party." It would also align with the US government's stringent export controls and restrictions on foreign direct investment in technologies facilitating AGI development in both countries.

An international supervision model has been proposed to oversee key AGI service providers internationally. Within this model, the national legislative and executive branches would occupy a central role in coordination with their international counterparts. This cooperative strategy seeks to enable effective oversight for AGI service development and deployment, ensuring global adherence to ethical standards and safety precautions.

Finally, it is advisable that the European Parliament, akin to any other organisation, consider implementing a policy that incorporates Generative AI-powered tools for routine activities after a detailed evaluation of potential advantages and drawbacks. This could boost efficiency across diverse sectors where AI can already provide considerable enhancement today.