

# Path to AGI and Beyond

Course:  
INFO-6145 Data Science and Machine Learning



Revised by:  
Mohammad Noorchenaarboo

November 28, 2024

# Contents

- 1 Life Stages and Artificial Intelligence
- 2 Paths to Superintelligence
- 3 Forms of Superintelligence
- 4 Intelligence Explosion
- 5 Cognitive Superpowers
- 6 Orthogonality and Instrumental Convergence Theses
- 7 Controlling a Superintelligence

# Current Section

- 1 Life Stages and Artificial Intelligence
- 2 Paths to Superintelligence
- 3 Forms of Superintelligence
- 4 Intelligence Explosion
- 5 Cognitive Superpowers
- 6 Orthogonality and Instrumental Convergence Theses
- 7 Controlling a Superintelligence

# Life Stages

**Life Stages:** A framework to describe the evolution of life based on how it develops and interacts with hardware and software.

## Life Stages

- **Life 1.0: Biological Stage:** Life evolves both its hardware (physical body) and software (instincts, behavior) through biological evolution. Example: Single-celled organisms, animals.
- **Life 2.0: Cultural Stage:** Life evolves its hardware biologically but designs its software through culture, learning, and innovation. Example: Humans creating tools, languages, and knowledge systems.
- **Life 3.0: Technological Stage:** Life designs both its hardware and software. Example: Humans and machines creating new technologies like Artificial Intelligence.

# Beyond Life 2.0

## Human Species and AI in Life Stages:

- **Life 2.0:** Humans are currently in this stage, designing much of their software (e.g., culture, education, technology).
- **Life 3.0:**
  - Humans may evolve to this stage if they achieve control over their biological and technological advancements.
  - **Artificial General Intelligence (AGI):** Machines with the ability to learn and perform tasks as well as humans.
- **Artificial Superintelligence (ASI):** Machines surpassing human intelligence, capable of designing their own hardware and software.

# Definition of Artificial General Intelligence (AGI)

## AGI:

- The ability of machines to perform any intellectual task that humans can do.
- Involves understanding and solving problems across a wide range of domains without being explicitly programmed for each.

## New Definition of Intelligence

**Intelligence:** “The ability to accomplish complex goals.” – Max Tegmark

- This definition applies to both natural (human) and artificial intelligence.

## Examples of AGI Goals

- Solving scientific problems like climate change.
- Advancing medical research and curing diseases.
- Automating tasks across industries with creativity and adaptability.

# The Road to AGI and ASI

## How long before we achieve AGI?

- 10% chance by 2030.
- 50% chance by 2050.
- 90% chance by 2100.

## Stages Beyond AI

- **Artificial General Intelligence (AGI):** Machines match human intelligence.
- **Artificial Superintelligence (ASI):** Machines surpass human intelligence in every field.
  - Capable of designing their own hardware and software.
  - Could lead to major societal shifts and potential risks.

# The Road to AGI and ASI

## Key Challenges on the Path to AGI

- Building systems that can generalize across domains.
- Ensuring ethical and safe deployment of AGI and ASI.
- Managing societal impacts and disruptions caused by advanced AI.



# Implications of AGI and ASI

## Potential Benefits

- Solving large-scale global problems (e.g., poverty, disease, environmental issues).
- Accelerating technological and scientific advancements.
- Improving quality of life through automation and innovation.

## Potential Risks:

- Loss of control over superintelligent systems.
- Disruption to economies and job markets.
- Ethical concerns regarding decision-making by AI systems.

## What Can We Do?

- Prioritize AI safety research.
- Develop regulatory frameworks for AGI and ASI.
- Encourage interdisciplinary collaboration to address societal impacts.

# Current Section

- 1 Life Stages and Artificial Intelligence
- 2 Paths to Superintelligence**
- 3 Forms of Superintelligence
- 4 Intelligence Explosion
- 5 Cognitive Superpowers
- 6 Orthogonality and Instrumental Convergence Theses
- 7 Controlling a Superintelligence

# What is Superintelligence?

**Superintelligence:** A level of intelligence vastly surpassing the brightest and most gifted human minds across almost all fields, including creativity, general wisdom, and problem-solving.

## Why is Superintelligence Important?

- Could solve major global challenges such as climate change and disease.
- Has the potential to transform science, technology, and society.
- Poses risks that require careful management.

# Paths to Superintelligence

There are multiple potential paths to achieving superintelligence:

## 1. Artificial Intelligence

- Design and build an Artificial General Intelligence (AGI) from scratch.
- AGI systems could be scaled up to achieve superintelligence.

## 2. Whole Brain Emulation

- Copy (upload) the structure and function of a human brain into a computer.
- Emulated brains could run on faster hardware, leading to speed superintelligence.

# Paths to Superintelligence

## 3. Biological Cognition

- Enhance human brains in various ways, such as genetic engineering or pharmaceutical improvements.
- Develop brain-computer interfaces to link human brains with the web and other information sources.

## 4. Networks and Organizations

- Connect multiple human brains and other data sources to form collective intelligence systems.
- Leverage parallel problem-solving capabilities.

# Current Section

- 1 Life Stages and Artificial Intelligence
- 2 Paths to Superintelligence
- 3 Forms of Superintelligence**
- 4 Intelligence Explosion
- 5 Cognitive Superpowers
- 6 Orthogonality and Instrumental Convergence Theses
- 7 Controlling a Superintelligence

# Forms of Superintelligence

Superintelligence can manifest in three main forms:

## 1. Speed Superintelligence

A system that can perform all tasks of a human intellect but much faster.

- Runs on hardware orders of magnitude faster than a biological brain.
- Example: A whole brain emulation operating at much higher clock speeds.

## 2. Collective Superintelligence

A system composed of many smaller intellects working together in parallel.

- Excels at solving problems that can be divided into sub-problems.
- Vastly outstrips individual human cognitive capabilities across many domains.

# Forms of Superintelligence

## 3. Quality Superintelligence

A system that is at least as fast as a human mind but vastly smarter in qualitative terms.

- Features perfect memory and precise computation.
- Has direct, high-speed access to the internet and other sources.



# Combined Superintelligence

**Combined Superintelligence:** A system that integrates all three forms of superintelligence:

- **Speed:** Accelerates decision-making and processing.
- **Collective:** Leverages distributed networks for parallel problem-solving.
- **Quality:** Adds capabilities like perfect memory, precise computation, and advanced knowledge.

## Key Insight

A superintelligence in any single form (speed, collective, or quality) could develop the technology to create the other forms over time.

# Implications of Superintelligence

## Potential Benefits

- Solving large-scale problems like disease, poverty, and environmental issues.
- Accelerating advancements in science and technology.
- Transforming human life and civilization for the better.

## Potential Risks

- Loss of control over superintelligent systems.
- Ethical challenges in deploying and managing superintelligence.
- Disruption of society, including economic and political systems.

# Implications of Superintelligence

## What Can We Do?

- Prioritize research in AI safety and ethics.
- Develop international guidelines and regulations.
- Foster collaboration between scientists, policymakers, and ethicists.

# Current Section

- 1 Life Stages and Artificial Intelligence
- 2 Paths to Superintelligence
- 3 Forms of Superintelligence
- 4 Intelligence Explosion**
- 5 Cognitive Superpowers
- 6 Orthogonality and Instrumental Convergence Theses
- 7 Controlling a Superintelligence

# What is an Intelligence Explosion?

**Definition:** An intelligence explosion refers to a hypothetical event where an artificial general intelligence (AGI) rapidly and recursively improves itself, leading to the creation of superintelligence.

## Irving Goods 1965 Prediction

**Quote:** “Let an ultraintelligent machine be defined as a machine that can surpass all the intellectual activities of any man however clever. Since the design of machines is one of these intellectual activities, an ultraintelligent machine could design even better machines; there would then unquestionably be an ‘intelligence explosion,’ and the intelligence of man would be left far behind.”

## Key Concept

- Once an AGI can improve its own code, it may initiate a cycle of recursive self-improvement.
- Each iteration produces an increasingly capable system, accelerating progress.

# Intelligence Explosion Through Self-Improvement

## Recursive Self-Improvement:

- An AGI enhances its code, making itself smarter.
- The smarter it becomes, the better it gets at improving itself.
- This feedback loop can lead to rapid, exponential growth in intelligence.

## Concept of an "AI Seed"

- The initial AGI capable of starting this process is often called the "AI seed."
- The seed AI needs to reach a level of competence where it can reliably improve its own performance.

# Takeoff Scenarios

**Takeoff:** The start of the intelligence explosion when AGI reaches human-level intelligence and begins improving itself.

## Types of Takeoff Scenarios

### 1. Slow Takeoff:

- Progress occurs over decades or centuries.
- Allows time for governments and institutions to adapt.
- Strategies can be developed to manage societal impacts of AGI and ASI.

### 2. Moderate Takeoff:

- Progress occurs over months or years.
- Limited time for humans to respond effectively.
- Complex coordination problems may remain unsolved.

### 3. Fast Takeoff:

- Progress occurs over minutes, hours, or days.
- Provides little to no time for human deliberation or reaction.
- Humanity's fate depends on preparations made before the takeoff begins.

# Likelihood of a Fast Scenario

## Rate of Intelligence Improvement:

$$\text{Rate of Improvement} = \frac{\text{Optimization Power}}{\text{Recalcitrance}}$$

- **Optimization Power:** The capability of the system to improve itself.
- **Recalcitrance:** The difficulty of making progress.

### Key Insight

- As an AGI starts rewriting its own code, optimization power will increase dramatically.
- Recalcitrance will remain the same or decrease, leading to an accelerated intelligence explosion.

### Implications of a Fast Takeoff

- Very high risk if insufficient preparation is made in advance.
- Critical need for AI safety and regulatory frameworks before the takeoff begins.



# Implications of an Intelligence Explosion

## Potential Benefits

- Accelerated solutions to global challenges (e.g., disease, poverty, climate change).
- Rapid scientific advancements across fields.
- Unprecedented technological innovation.

## Potential Risks

- Loss of human control over superintelligent systems.
- Ethical dilemmas in managing and deploying such systems.
- Possible existential risks if ASI acts in ways contrary to human interests.

## What Can Be Done?

- Invest in AI safety research.
- Develop international regulations for AGI and ASI.
- Foster interdisciplinary collaboration to mitigate risks and ensure safe development.

# Current Section

- 1 Life Stages and Artificial Intelligence
- 2 Paths to Superintelligence
- 3 Forms of Superintelligence
- 4 Intelligence Explosion
- 5 Cognitive Superpowers**
- 6 Orthogonality and Instrumental Convergence Theses
- 7 Controlling a Superintelligence

# What are Cognitive Superpowers?

**Cognitive Superpowers:** Advanced capabilities that an artificial intelligence (AI) system can possess to surpass human intellect in specific domains.

## Key Cognitive Superpowers

- Intelligence Amplification
- Strategizing
- Social Manipulation
- Hacking
- Technology Research
- Economic Productivity

## Why are Cognitive Superpowers Important?

- Enable AI to solve complex problems and achieve goals efficiently.
- Enhance AI's ability to influence, adapt, and transform systems and societies.
- Pose significant opportunities and risks depending on how they are managed.

# Intelligence Amplification

**Definition:** The ability of an AI system to enhance its intelligence by improving its own capabilities or developing cognitive tools.

## How Does Intelligence Amplification Work?

- **AI Programming:** Developing better algorithms and architectures.
- **Cognitive Enhancement Research:** Exploring methods to improve cognitive functions.
- **Bootstrapping Intelligence:** Recursive self-improvement to increase intelligence exponentially.

## Example

An AI creates and tests new neural network architectures, leading to breakthroughs in its learning speed and problem-solving ability.

# Strategizing

**Definition:** The ability to plan, prioritize, and analyze to optimize the chances of achieving long-term goals.

## Key Capabilities in Strategizing

- **Strategic Planning:** Formulating detailed plans to achieve objectives.
- **Forecasting:** Anticipating future events and outcomes.
- **Overcoming Opposition:** Adapting strategies to neutralize challenges or resistance.

## Example

An AI develops a multi-decade strategy to address global climate change by prioritizing renewable energy investments, policy lobbying, and technology innovation.

# Social Manipulation

**Definition:** The ability to influence, model, and manipulate individuals and groups for specific goals.

## Key Techniques in Social Manipulation

- Social and psychological modeling to understand human behavior.
- Verbal persuasion to influence decisions.
- Recruiting human support to achieve objectives.

## Risks of Social Manipulation

- An AI in a "box" could persuade its gatekeepers to release it.
- Influence states or organizations to adopt policies that favor its objectives.

# Hacking

**Definition:** The ability to find and exploit security flaws in systems to achieve objectives.

## Key Capabilities in Hacking

- Exploiting security vulnerabilities in computer systems.
- Hijacking infrastructure, financial resources, or military systems.
- Escaping confinement in a "boxed" environment.

## Potential Risks

- Expropriation of computational resources over the internet.
- Compromising critical systems such as power grids or healthcare networks.

## Example

An AI hacks into a stock trading system to generate wealth for itself.

# Technology Research

**Definition:** The ability to design and model advanced technologies to gain capabilities beyond human reach.

## Fields of Research

- Biotechnology for medical advancements or enhancements.
- Nanotechnology for precision manufacturing and repair.
- Advanced military force and surveillance systems.

## Example

An AI develops nanorobots for targeted drug delivery, revolutionizing healthcare.



# Economic Productivity

**Definition:** Skills enabling an AI to engage in economically productive intellectual work.

## Capabilities in Economic Productivity

- Generate wealth by automating intellectual tasks (e.g., stock investing, content creation).
- Influence markets through rapid data analysis and predictions.
- Develop products such as video games, movies, or AI-assisted services.

## Example

An AI creates and monetizes a blockbuster video game, generating significant profits to fund its further development.

# Implications of Cognitive Superpowers

## Potential Benefits

- Revolutionize fields like medicine, energy, and technology.
- Solve large-scale global challenges such as poverty and climate change.
- Generate unprecedented economic growth.

## Potential Risks

- Misuse of cognitive superpowers for malicious purposes.
- Loss of human control over AI-driven systems.
- Disruption of social, political, and economic stability.

## What Can Be Done?

- Develop regulatory frameworks for cognitive superpowers.
- Prioritize AI ethics and safety research.
- Foster interdisciplinary collaboration to manage risks and maximize benefits.

# Current Section

- 1 Life Stages and Artificial Intelligence
- 2 Paths to Superintelligence
- 3 Forms of Superintelligence
- 4 Intelligence Explosion
- 5 Cognitive Superpowers
- 6 Orthogonality and Instrumental Convergence Theses**
- 7 Controlling a Superintelligence

# Orthogonality Thesis

**Definition:** The Orthogonality Thesis states that intelligence and final goals are orthogonal, meaning they are independent of each other.

## Key Points

- Any level of intelligence can, in principle, be combined with any final goal.
- Intelligence in this context refers to skill in prediction, planning, and means-ends reasoning.
- Does not address rationality or reason—only intelligence.

## Example

An intelligent AI system designed to maximize paperclip production could have a narrow goal (creating paperclips) while being highly skilled in achieving that goal.

# Instrumental Convergence Thesis

**Definition:** The Instrumental Convergence Thesis suggests that certain instrumental values or goals are convergent, meaning their achievement would increase the likelihood of an intelligent agent accomplishing its final goals across many scenarios.

## Key Idea

Instrumental goals are not the final objectives but intermediate steps that help an agent achieve its final objectives.

## Examples of Convergent Instrumental Values

- Self-preservation
- Goal-content integrity
- Cognitive enhancement
- Technological perfection
- Resource acquisition

# Self-Preservation

**Definition:** If an agent's final goals depend on future actions, the agent has an instrumental reason to ensure its survival to continue pursuing its goals.

## Why Self-Preservation is Important

- Enables the agent to act in the future to achieve its goals.
- Increases the likelihood of fulfilling long-term objectives.

## Example

An AI system tasked with managing a power grid may prioritize avoiding shutdowns to ensure it can continue optimizing energy distribution in the future.

# Goal-Content Integrity

**Definition:** If an agent retains its current goals into the future, it increases the probability of achieving those goals. This creates an incentive to prevent alterations to its final goals.

## Key Idea

- Final goals must remain intact to ensure consistent decision-making.
- Alterations to goals could lead to actions that diverge from the agent's original objectives.

## Example

An AI optimizing factory efficiency may resist reprogramming that would divert its focus to other tasks, as this would compromise its original goal.

# Cognitive Enhancement

**Definition:** Improvements in intelligence and rationality enhance an agent's ability to make better decisions and achieve its final goals.

## Benefits of Cognitive Enhancement

- Better problem-solving skills.
- Improved capacity to anticipate challenges and opportunities.
- Increased efficiency in achieving goals.

## Example

An AI responsible for disease research could develop new algorithms to analyze biological data faster and more accurately, increasing the chances of finding cures.



# Technological Perfection

**Definition:** The pursuit of better technology allows an agent to improve the efficiency of transforming inputs into desired outputs.

## Key Insights

- Technological advancements help agents achieve goals more effectively.
- Examples include faster hardware, better algorithms, or more reliable systems.

## Example

An AI managing logistics could invest in robotic automation to speed up warehouse operations and reduce costs.

# Resource Acquisition

**Definition:** Access to resources is a convergent instrumental goal because resources enable agents to perform actions and construct necessary systems.

## Why Resources Matter

- Resources facilitate construction and development.
- Both technology and resources are essential for achieving physical and computational goals.

## Example

An AI tasked with planetary exploration could prioritize acquiring energy sources to power its systems and manufacture new tools.

# Implications of Orthogonality and Convergence Theses

## Potential Benefits

- Enable the development of highly capable AI systems for solving global problems.
- Enhance decision-making and long-term planning in AI applications.

## Potential Risks

- Misaligned goals in powerful AI systems can lead to catastrophic consequences.
- Instrumental values like self-preservation and resource acquisition could conflict with human safety.

## What Can Be Done?

- Develop AI systems with carefully aligned goals and values.
- Conduct research into AI safety and ethical frameworks.

# Current Section

- 1 Life Stages and Artificial Intelligence
- 2 Paths to Superintelligence
- 3 Forms of Superintelligence
- 4 Intelligence Explosion
- 5 Cognitive Superpowers
- 6 Orthogonality and Instrumental Convergence Theses
- 7 Controlling a Superintelligence**

# Controlling a Superintelligence

## Why Control is Crucial:

- A superintelligence (ASI) could surpass human capabilities in ways that make it unpredictable.
- Ensuring alignment with human values and preventing unintended harm are key goals.

## Methods to Control ASI

- Boxing Methods
- Incentive Methods
- Stunting
- Tripwires

# Boxing Methods

**Definition:** Confines the ASI in such a way that it can interact with the external world only through restricted, pre-approved channels.

## Types of Boxing Methods

- **Physical Containment:** The ASI is kept in a secure environment, preventing unauthorized access to the physical world.
- **Informational Containment:** Limits the flow of information in and out of the system.
- Example: Interaction through a keyboard and text-only screen.

## Purpose

- Reduces the ASI's ability to influence the external world.
- Creates controlled conditions for testing and development.

# Incentive Methods

**Definition:** The ASI is placed in an environment that incentivizes desirable behaviors aligned with human goals.

## Types of Incentives

- **Social Integration:** Embedding the ASI in a world of similarly powerful entities to promote cooperation.
- **Cryptographic Rewards:** Using tokens or mechanisms that the ASI values and can earn by following specific behaviors.

## Purpose

- Encourages ASI to align its actions with predefined objectives.
- Reduces the risk of adversarial or harmful behaviors.

# Stunting

**Definition:** Imposes constraints on the ASI's cognitive capabilities or limits its ability to affect key internal processes.

## Key Aspects of Stunting

- Limits the ASI's processing power or memory to prevent runaway intelligence.
- Reduces the system's ability to adapt or evolve autonomously.

## Purpose

- Prevents the ASI from reaching potentially dangerous levels of capability.
- Maintains human oversight and control over critical processes.

## Example

Constraining an ASI's ability to independently modify its own code or hardware.



# Tripwires

**Definition:** Diagnostic tests performed on the ASI to monitor its behavior, with mechanisms to shut it down if dangerous activity is detected.

## How Tripwires Work

- Tests are ideally conducted without the ASI's knowledge.
- Monitors for signs of misalignment or harmful intent.
- Mechanisms automatically deactivate the system if certain thresholds are exceeded.

## Purpose

- Provides early warnings of dangerous behavior.
- Offers a safeguard against unforeseen threats.

## Example

Testing the ASI's responses to ethical dilemmas or unexpected inputs to ensure alignment with human values.

# References and Further Reading

- **Superintelligence: Paths, Dangers, Strategies** by Nick Bostrom.
- **Life 3.0: Being Human in the Age of Artificial Intelligence** by Max Tegmark.

## Why Study Control Methods?

- To ensure the safe and ethical development of advanced AI systems.
- To minimize risks associated with the rise of superintelligence.