

UNIVERSITY OF ZAGREB
FACULTY OF ELECTRICAL ENGINEERING AND COMPUTING

MASTER THESIS No. 1234

MODELLING AND OVERSIGHT OF NATURAL INTELLIGENCE: KEY ASPECTS

Dorijan Cirkveni

Zagreb, June 2024.

SVEUČILIŠTE U ZAGREBU
FAKULTET ELEKTROTEHNIKE I RAČUNARSTVA

DIPLOMSKI RAD br. 1234

KLJUČNI ASPEKTI MODELIRANJA I NADZORA PRIRODNE INTELIGENCIJE

Dorijan Cirkveni

Zagreb, Lipanj 2024.

Zagreb, 2. listopada 2023.

DIPLOMSKI ZADATAK br. 3168

Pristupnik: **Dorijan Cirkveni (0036501554)**

Studij: Računarstvo

Profil: Računarska znanost

Mentor: prof. dr. sc. Vedran Mornar

Zadatak: **Ključni aspekti modeliranja i nadzora prirodne inteligencije**

Opis zadatka:

Istražiti mogućnost razvoja simulacije ljudske inteligencije i svijesti korištenjem trenutačno postojećih mogućnosti umjetne inteligencije. Prikupiti literaturu povezanu s konceptima umjetne inteligencije, prirodne inteligencije, svijesti i drugih područja, tehničkih i teoretskih, blisko povezanih s područjem istraživanja. Usporediti postojeće implementacije ljudske inteligencije i svijesti. Vrednovati prostornu i vremensku složenost postojećih implementacija, ali i njihovu transparentnost. Razviti simulirane svjesne agente i virtualnu okolinu za ispitivanje i nadzor tih agenata. Oblikovati virtualnu okolinu tako da simulirala niz stvarnih situacija u kojima se svjesni agenti mogu naći. Ovisno o vrsti modela agenata, koristiti njihovo okruženje u procesu učenja. Diskutirati rezultate istraživanja.

Rok za predaju rada: 9. veljače 2024.

*I wish to thank everyone who supported me in my efforts to reach this point,
especially with how long it took.*

Contents

1	Introduction	5
2	Background	7
2.1	Recent advances in artificial intelligence	7
2.1.1	Future prospects	9
2.1.2	Public adoption	9
2.1.3	Backlash to "AI art" and the question of consciousness	11
2.2	Objective	12
2.3	Choosing an approach	13
2.3.1	Thinking rationally (Artificial intelligence)	14
2.3.2	Acting rationally (Apparent artificial intelligence)	15
2.3.3	Thinking "humanly" (consciousness?)	15
2.3.4	Acting "humanly" (apparent consciousness)	17
3	Theoretical foundations	19
3.1	Artificial intelligence	19
3.1.1	The five schools of artificial intelligence	19
3.2	Defining intelligence	20
3.2.1	Types of intelligence	21
3.2.2	Testing for intelligence	21
3.2.3	A few notable attempts at simulating/creating intelligence	22
3.3	Artificial consciousness	23
3.3.1	Anthropomorphism and AI	23
3.3.2	The Chinese Room Argument	26

3.3.3	Popular media depictions of AI as artificial consciousness	27
3.3.4	Why should we care about artificial consciousness?	29
3.4	Defining consciousness	30
3.4.1	Existing theories	30
3.4.2	Testing for consciousness	31
3.4.3	Self-awareness tests	33
3.4.4	A few notable projects related to simulating/creating consciousness	34
4	Materials and Methods	36
4.1	Virtual testing environment	36
4.2	Grid elements	37
4.2.1	Basic grid elements	37
4.2.2	Composite grid elements	39
4.2.3	Visible and solid grid modes	39
4.2.4	Grid routines	40
4.3	Test entity design	40
4.3.1	Entity/Agent duality	40
4.3.2	Entity properties	40
4.4	Test environment design	41
4.4.1	Reference test environments	42
4.4.2	Mirror test environments	42
4.5	Agent behavior types	42
4.6	Tested agent types	43
4.7	A suggestion for a knowledge-based classification of intelligence	43
4.7.1	Knowledge by source	43
4.7.2	Knowledge by level	44
4.8	Agent aspect considerations	46
4.8.1	Steps taken	46
4.8.2	Processing time	46
4.8.3	Processing space	46
4.8.4	Learning data size	46
4.8.5	Learning data level	47

5	Implementation	48
5.1	Implementation choices	48
5.1.1	Programming language choice	48
5.2	Tested agent choice	49
5.2.1	Reference agents	49
5.3	Implementation details	50
5.3.1	Base classes	50
5.3.2	Base interfaces	51
5.3.3	Grid environment classes	53
5.3.4	Environment display and interaction user interface	54
5.3.5	Selection user interface	54
5.3.6	Data management interface	54
5.3.7	Support functionalities	54
6	Results	56
6.1	Preliminary research conclusions	56
6.2	Developed infrastructure	56
6.3	Instructions	57
6.3.1	Downloading and running the software	57
6.3.2	Main menu	57
6.3.3	Data management menu	57
6.4	Environment and agent selection menu	59
6.5	Environment run menu	59
7	Further Research and Discussion	61
7.1	Further prototype development	61
7.2	Further research direction	62
7.3	Discussion	62
7.3.1	Ethical concerns	62
8	Conclusion	64
	References	65
	Abstract	67

Sažetak	68
A: The Code	69

Chapter 1

Introduction

"We call ourselves *Homo sapiens* - man the wise - because our intelligence is so important to us. For thousands of years, we have tried to understand *how we think*; that is, how a mere handful of matter can perceive, understand, predict, and manipulate a world far larger and more complicated than itself. The field of artificial intelligence, or AI, goes further still: it attempts not just to understand but also to build intelligent entities." [1]

The idea that machines could one day think, feel, and be conscious has captivated the human imagination for over a century, inspiring works of fiction and scientific research alike. This pursuit, while still ongoing and with no end in sight, resulted in immeasurable benefits in the form of practical applications of narrow artificial intelligence in virtually every facet of our lives, from navigation software to product recommendation algorithms.

Several important questions, however, insist on remaining unanswered: Can machines think? Can machines feel? If we develop artificial intelligence that can match or exceed human ability in any endeavor imaginable, how will we be able to tell - and can we accurately estimate how far we are from that point if such a breakthrough is even possible? How can we tell if a machine is conscious? ... Are we ready for thinking machines?

With each passing day and every new benchmark reached, the nature of these ques-

tions gradually shifts from that of speculative fiction or philosophical inquiries to that of time-pressed concerns with practical implications for the real world in the foreseeable future - and this thesis seeks to provide a foundation for research that will allow us to answer these very questions.

Within the scope of this thesis, we have assembled an overview of existing research, including several notable theories of intelligence and consciousness, both of which remain open questions in their respective fields to this day. Additionally, we have developed a prototype of the testing framework designed to facilitate the development and testing of potentially intelligent or conscious artificial intelligence agents.

The thesis body is outlined as follows:

- The *Background* chapter will outline the background, motivation, and objective of this thesis further.
- The *Theoretical Foundations* chapter will provide the aforementioned theoretical foundation needed to understand this objective.
- The *Materials and Methods* chapter will outline the elements of the artificial intelligence testing environment, agent, and evaluation method templates.
- The *Implementation* chapter will outline the implementation details of the prototype, with a focus on elements further research can build upon through inheritance, composition, and other code reuse methods.
- The *Results* chapter describes the prototype test environment features and provides instructions for usage.
- The *Further Research and Discussion* chapter proposes a course of action for further research and discusses the potential consequences of said research as well as the consequences of lack thereof.

Chapter 2

Background

2.1 Recent advances in artificial intelligence

Ten years ago (2013), artificial intelligence capabilities were far behind where they are today, with handwriting recognition abilities barely lagging behind human performance, speech and image recognition lagging far behind, and reading comprehension and language understanding being untested or non-existent.

Since then, artificial intelligence systems have outperformed humans in these five fields and more. [2]

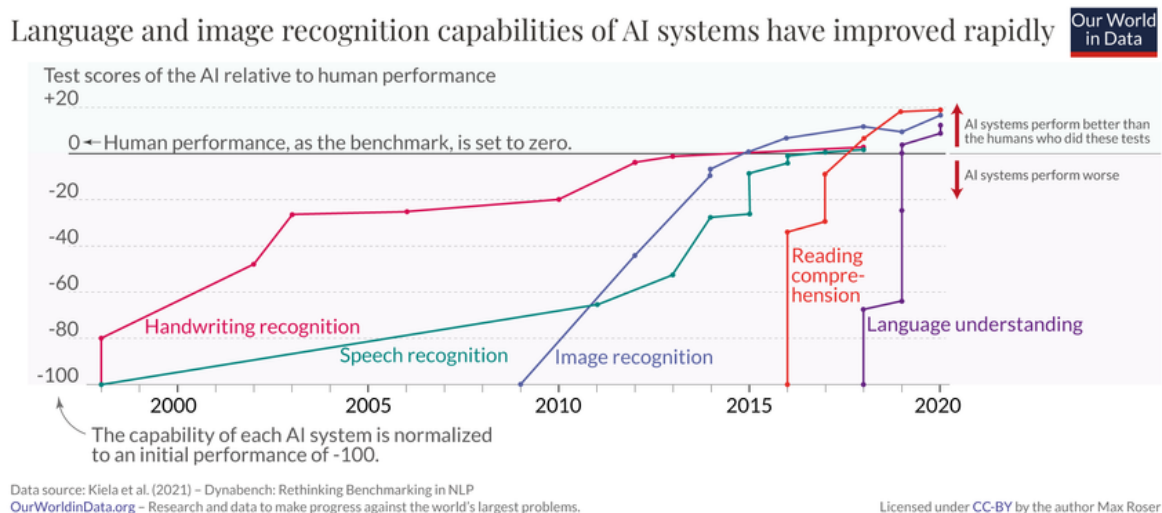


Figure 2.1: Advances of AI during the last few decades

- **Image Recognition**

Deep learning has enabled artificial intelligence systems to achieve superhuman accuracy in image recognition tasks, recognizing objects, scenes, and faces with remarkable precision.

- **Video Understanding**

Artificial intelligence systems can now analyze videos to understand the content, actions, and relationships between objects, paving the way for applications like action recognition and video summarization.

- **Object Detection and Tracking**

Artificial intelligence systems can now detect and track multiple objects in real time, enabling applications like self-driving cars and surveillance systems.

- **AlphaGo's Defeat of Lee Sedol**

In 2016, AlphaGo, an AI developed by Google DeepMind, defeated world champion Go player Lee Sedol, marking a significant breakthrough in the ability of reinforcement-learning-based algorithms to master complex games.

- **Financial Trading**

Reinforcement learning algorithms have been successfully used to make trading decisions and optimize investment portfolios, demonstrating the potential of reinforcement learning in various real-world applications.

- **Robotic Control**

Reinforcement learning has given artificial intelligence systems the ability to control machines with unprecedented skill and agility, which allows them to perform physical tasks such as picking up objects, grasping tools, navigating complex and dynamic environments, and even performing complex, dynamic, and critical tasks such as surgical procedures. [3]

2.1.1 Future prospects

According to three surveys made during the last ten years (2018, 2019, and 2022), 50 percent of experts estimate at least a 50 percent chance artificial general intelligence will be developed before 2070 (medians being 2068, 2060, and 2061 respectively). [4]

AI timelines: What do experts in artificial intelligence expect for the future?

Our World
in Data

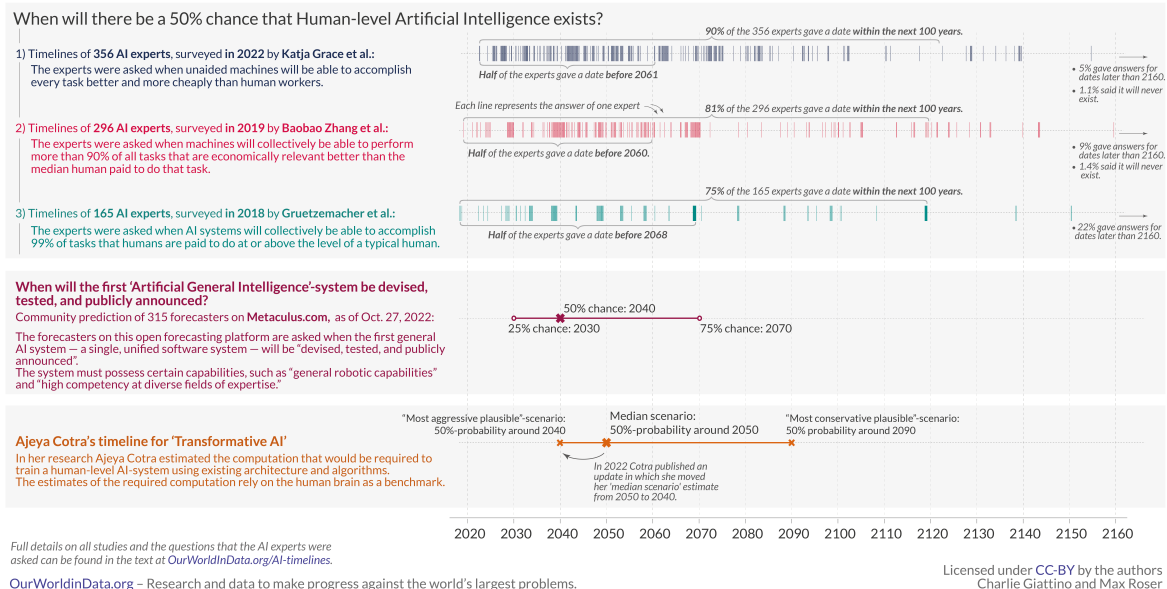


Figure 2.2: AGI timeline

It should, however, be noted that such predictions are highly unreliable and subject to the "Maes-Garreau law" (the idea that people will predict AI coming before they die). [5] This is likely to mean one of the following for the future of AI research:

1. Artificial general intelligence is an unattainable goal that will remain "thirty years away" (like commercial nuclear fusion power has) for centuries to come.
2. Artificial general intelligence will be developed at some point despite predictions to the contrary, which may cause us to be ill-prepared for it.

2.1.2 Public adoption

Machine learning is currently going through a renaissance not only from a research standpoint but also in the public eye, mainly because several different forms of its applications have been successfully released and reached widespread use in an unprecedentedly short time.

One such application, a conversation-based large language model platform known as ChatGPT, surpassed 100 million monthly active users merely two months after its release.

Large Language Models

The capabilities of large language models such as ChatGPT, Bing Copilot, and Bard, for example, are best described and demonstrated through a self-referencing example.

Write a self-demonstrating paragraph about large language models in public use, such as yourself, Bing Copilot, and Gemini. The paragraph needs to be between 1000 and 1500 characters long, include a mention that it has been made by you with minimal input from "the thesis author", and refer to LLMs, including yourself, in third person.

When prompted with this statement, ChatGPT(GPT-4o) responded with the following:

Large language models (LLMs) have become an integral part of modern digital interactions, revolutionizing the way users access information, seek assistance, and engage with technology. Prominent examples of these advanced systems include ChatGPT by OpenAI, Bing Copilot by Microsoft, and Gemini by Google. These LLMs are designed to understand and generate human-like text based on extensive training on diverse datasets. Each model has unique strengths and capabilities, tailored to offer personalized and accurate responses. For instance, ChatGPT is known for its conversational prowess and ability to handle a wide range of queries with coherent and contextually relevant answers. Bing Copilot integrates seamlessly with Microsoft's suite of tools, providing users with smart suggestions and insights that enhance productivity. Meanwhile, Gemini focuses on delivering high-quality, context-aware responses that cater to specific user needs within Google's ecosystem. Despite their differences, these LLMs share a common goal: to assist users by leveraging the power of artificial intelligence. The creation of this paragraph was facilitated by ChatGPT with minimal input from the thesis author, demonstrating the model's ability to autonomously generate informative and

cohesive text. As these technologies continue to evolve, they are expected to play an increasingly significant role in various sectors, from education and customer service to content creation and beyond. Their capacity to understand and emulate human language opens up new possibilities for innovation, making them indispensable tools in the digital age.

[6]

Image generation AI

Image generation AI, similarly, managed to reach widespread use thanks to its capability of generating images based on natural language prompts, utilizing and expanding upon the ability of large language models to recognize natural language input.



Figure 2.3: An image generated using Bing Image Creator. Prompt: "an artificial intelligence contemplating its existence"

2.1.3 Backlash to "AI art" and the question of consciousness

Use of generative artificial intelligence to create images has been met with severe backlash for several reasons, including uncompensated use of existing art to train generative artificial intelligence models and the technology's threat to the livelihoods of human artists.

However, there is one reason people oppose the use of generative artificial intelligence to create so-called "AI art" that is closely related to the topic of this thesis, and that is that art is the result of conscious creation by a conscious artist. And since image

generation programs based on artificial intelligence are not conscious, and their operation usually does not require conscious input beyond making a specific request - the act which has more similarity with commissioning art from an artist than with creating art by oneself - this means that images generated using artificial intelligence cannot be considered art.

However, the question of whether consciousness is a prerequisite for art, or whether artificial intelligence can create art that can be considered art, is at best tangentially relevant to the second main question of this thesis, which is:

Can artificial intelligence be conscious?

2.2 Objective

The purpose of this paper is to research the possibility of simulating human intelligence and consciousness using currently available hardware, software, and known methods of developing artificial intelligence.

This is an objective that modern computer science has pursued in some shape or form since its inception in the 1940s. The field of artificial intelligence research, however, had yet to be founded with The Dartmouth Summer Research Project in 1956.

In this paper, we will attempt to determine possible routes to achieving said objective and how close we are to achieving it.

To do so, we first need to establish a working definition of both intelligence and consciousness. Doing so is a difficult task in its own right, as both definitions are open questions and a subject of serious discussion.

Next, we need to define a metric with which we will measure whether an artificial entity can be considered intelligent or conscious.

After that, this paper will propose a few promising avenues of research and provide implementation examples and their preliminary results within our ability.

However, to accomplish the following goals, we must first determine the nature of the task at hand and estimate its scope as well as the progress that has already been made

toward completing it, with an emphasis on recent and state-of-the-art accomplishments.

We should first separate this task into two tasks:

1. Simulating human/human-level intelligence
2. Simulating human/human-level consciousness

This is because while these tasks may be prerequisites for one another, they likely require different considerations, and other tests are likely to be used to determine their presence in an artificial agent.

2.3 Choosing an approach

There are four approaches to the interpretation of the issue at hand we can take, which can be divided into quadrants according to two criteria:

1. The goal of the process
 - Rationality (Artificial intelligence)
 - Humanity (Simulated human reasoning)

Regarding the nature of reasoning being pursued, we need to decide whether we want the artificial intelligence to operate rationally or think humanly. Those two goals are not necessarily contradictory, however, they are also not identical. They are orthogonal, which means that while they can be pursued simultaneously, one can also be pursued at the expense - or, at the very least, the opportunity cost - of the other.

2. The focus of measurement
 - Thinking (Internal states/Strong AI)
 - Acting (Observed actions/Weak AI)

Regarding whether we are measuring the internal states of the agent in question or if, instead, we are measuring the external actions of the agent. While directly observing the internal states of an artificial intelligence agent would certainly be

more desirable from an academic perspective, observing the agent's actions is often easier - and sometimes the only option available - as well as more desirable from a practical standpoint. These two approaches correspond to the two AI hypotheses - the strong AI hypothesis and the weak AI hypothesis, respectively.

	Intelligence	Consciousness
Internal state	Thinking rationally	Thinking humanly
Observed state	Acting rationally	Acting humanly

Table 2.1: Comparison of Intelligence and Consciousness

2.3.1 Thinking rationally (Artificial intelligence)

In this approach we wish to develop an artificial intelligence that thinks rationally, that is, an AI capable of logical reasoning. This approach to developing artificial intelligence is called the logicist approach and consists of describing problems in logical notation before solving them using known syllogisms.

There are, however, two main obstacles to this approach:

- Informal knowledge and uncertainty

It is difficult to convert informal knowledge into formal statements required by logical notation, especially when said knowledge involves a degree of uncertainty (for example, "It's probably going to rain today") that needs to be accounted for while solving the problem.

While progress has been made with fields such as fuzzy logic or Bayesian methods, this still leaves us with the second issue, which is computational complexity.

- Computational complexity

In theory, any problem that can be stated in logical notation and for which a solution exists can be solved given enough time and computational resources. However, given that this problem is NP-complete, the amount of time and resources needed to do so rise exponentially with the number of facts that need to be considered unless proper guidance is provided to help the program decide which reasoning steps to try first.

2.3.2 Acting rationally (Apparent artificial intelligence)

In this approach, in mild contrast to the logicist approach, we ignore the inner world of an artificial intelligence entity and focus on creating or identifying rational artificial agents that act rationally. A rational agent attempts to perform the best possible action in any given situation, that is, the action that will result in the best possible outcome as determined by the agent's goals and/or utility function. Thinking rationally can be - and usually is - a critical component of acting rationally. After all, knowing is half the battle. However, it is possible for an agent to act rationally without thinking rationally, or even where undue contemplation leads to worse outcomes (for example, reflexively grabbing a falling bottle leads to a better outcome than thinking about what to do). The advantages of the agent-driven approach are:

- Generality

This approach is more general than the previously listed "laws-of-thought" approach, as thinking rationally is just one of the ways of achieving rational actions (alongside other ways such as learned instinct or following instructions).

- Ease of use in research

As the standard used to determine whether agents are rational is well-defined and general, it is more amenable to scientific development than the two approaches listed below.

However, similar to the approach above, finding the most rational action in an environment is not always feasible, especially in complex environments, due to the sheer scale of computational resources required. Also, the advantage of generality turns into a disadvantage when we wish to distinguish whether an action is a result of intelligence.

2.3.3 Thinking "humanly" (consciousness?)

In the early days of artificial intelligence, it was common for researchers to conflate "thinking rationally" with "thinking humanly". For example, an author would argue one of the following:

1. An algorithm that performs well on a task is a good model of human performance.

2. For an algorithm to perform well on a task, it needs to be a good model of human performance.

However, while the assumption that human beings are rational actors for all intents and purposes may work well enough for some fields such as economics - even though its falsity is obvious to anyone with sufficient experience of interaction with human beings - research focusing on the nature of intelligence and consciousness itself cannot make such approximations. Still, despite this flaw, the cognitive approach to artificial intelligence has its advantages (such as readily available reference models) and has resulted in significant breakthroughs such as neural networks, natural language processing, and explainable AI.

Out of the four approaches listed here, this one is closest to what we could consider "consciousness".

The main potential benefits of this approach are:

- Readily available exemplar knowledge

The ability to draw upon existing psychological and, in the case of especially emulative approaches such as neural networks, neurological knowledge bases may make it easier to develop such modes.

- Explainability and trust

The artificial models created with this approach are likely to be more explainable, and therefore more transparent and reliable for high-trust and high-risk tasks.

- User-friendliness

Additionally, these approaches are inherently human-centric, which provides an advantage when it comes to user-facing AI applications.

- Alignment potential

Most importantly, the human-centric design may be more likely to result in an AI that is easier to align with human values.

Of course, no approach is without its drawbacks, which in this case are:

- Potentially inaccurate models of human cognition

Human cognition is complex and not fully understood. Models based on incomplete or incorrect understandings of human thought processes can lead to flawed AI systems that behave unpredictably or inappropriately.

- Resource intensity

Developing AI systems that closely mimic human thought is complex in terms of research, requires multidisciplinary knowledge, and can be computationally and resource-intensive when compared to developing and running AI systems to solve a specific task.

- Suboptimal results

Just like attempting to emulate birds did not result in optimal aircraft design, attempting to emulate human consciousness is unlikely to result in optimal results for a given task - except, of course, in situations where emulating human consciousness is the task or a part of it.

Otherwise, this approach is likely to yield suboptimal results compared to the ones above, which rely on first-principles thinking instead.

2.3.4 Acting "humanly" (apparent consciousness)

The final approach, testing or designing machines to act like a human would, is also likely the oldest approach of the four, being established by Alan Turing in 1950 with the now iconic Turing Test proposal, on which this paper will elaborate further later. While the capabilities required to pass the Turing Test convincingly include most fields of artificial intelligence (such as natural language processing, knowledge representation, automated reasoning, and machine learning) and overlapping/closely related fields (such as computer vision and robotics), AI researchers have devoted little effort to passing it, believing that duplicating human actions is not as important as actually studying the underlying principles of intelligence.

There are advantages to this approach:

- User-friendliness

Even more than the previous approach, this approach leads to developing user-friendly AI suitable for user-facing tasks.

- Benchmarking

This is a clear goal that can be easily demonstrated, much like landing a manned mission on the Moon is a historic achievement even though unmanned missions are far more practical for most purposes.

- Interdisciplinary integration

This goal requires - and therefore encourages - the cooperation of multiple AI and non-AI fields (such as psychology, robotics, or sociology), which is likely to result in beneficial spinoff results.

-

However, much like the others, this approach also has its flaws:

- Resource intensity and suboptimality

Much like the previous approach, acting humanly is a sub-task that is likely to produce sub-optimal results in situations where this is not one of the main tasks.

- Uncertainty of achievement

There is no clearly defined test for human behavior (the Turing Test itself has flaws that will be noted at a later point in this thesis), making it difficult to determine the definite success or failure of such an attempt.

- False breakthrough

Even if we developed an AI that could reliably pass the Turing Test with one hundred percent reliability, this could turn out to be a hollow achievement and fail to result in actual breakthroughs in the fields of artificial intelligence.

Chapter 3

Theoretical foundations

3.1 Artificial intelligence

3.1.1 The five schools of artificial intelligence

There are five distinct schools of artificial intelligence, which this work will reference repeatedly. Each of these five distinct schools of artificial intelligence is focused on a different approach when it comes to achieving the same goal. These approaches are not mutually exclusive, however - it is possible, and likely necessary, to combine multiple approaches to be able to develop an artificial general intelligence capable of tackling a vast array of problems, including that of simulating a human consciousness.

Connectionism

The connectionist school of artificial intelligence focuses on replicating the human brain through artificial structures known as neural networks, which are built out of fundamental building blocks called artificial neurons. This approach is meant to simulate the way our natural neurons work, and hopefully replicate our cognitive capabilities in the process.

Large language models such as ChatGPT are the result of this school of artificial intelligence.

Symbolism

Unlike the connectionist approach of replicating the human brain by starting from its fundamental building blocks and moving up, this school of artificial intelligence uses symbols to represent the world, and the artificial intelligence models it creates are known as expert systems.

Evolutionism

The evolutionist approach to artificial intelligence seeks to leverage the process that gave rise to human consciousness to train artificial intelligence models through processes such as feature mutation, feature cross-combination, and natural selection. Genetic algorithms are a common tool used in projects based on this school of artificial intelligence.

Bayesian approach

The Bayesian approach to artificial intelligence uses probability theory to model uncertainty. The models created with this approach - Bayesian models - assign probabilities to different possible states of the environment.

Analogizing

This approach is the easiest to understand, as well as the easiest to implement. The analogizing approach takes an input and compares it to other inputs with known results to generate a similar result.

3.2 Defining intelligence

In our attempts to create artificial general intelligence, it is likely beneficial - if not outright necessary - to understand what intelligence is.

This arguably simpler of two tasks, as - unlike consciousness - intelligence proves to be the easier of the two to define due to its objective and observable nature.

However, measuring intelligence is still a daunting task because it encompasses a wide range of cognitive abilities, including problem-solving skills, learning capabilities, adaptability, and more.

This is especially the case with artificial intelligence, due to the ability of machines to easily solve tasks that require the use of intelligence when solved by a human solver.

3.2.1 Types of intelligence

The number of distinct types of intelligence is an open question with multiple conflicting existing theories, such as

- **Traditional single-intelligence theories** These theories differ in exact definition, but they all share the view of intelligence as a single-factor quantity, sometimes known as "general intelligence" or "g". One method of measuring such a quantity in humans is the famous Intelligence Quotient (IQ).
- **Gardner's Theory of Multiple Intelligences** This theory challenges the aforementioned notions and proposes the division of intelligence into eight separate categories: linguistic, logical/mathematical, spatial, bodily-kinesthetic, musical, interpersonal, intrapersonal, and naturalist
- **Spearman's two-factor theory of intelligence** This theory bridges the gap between the two previous theories by combining the notion of general intelligence with that of specific abilities.
- **Cattell-Horn-Carroll (CHC) theory** This theory takes a step further, dividing intelligence into three strata: general abilities, broad abilities, and narrow abilities.

3.2.2 Testing for intelligence

While defining intelligence is no easy task, one of its definitions (the ability to apply knowledge to manipulate one's environment or to think abstractly as measured by objective criteria (such as tests)) lends itself to simple testing methods such as:

- **Raven's Progressive Matrices**

Given their use in measuring general human intelligence and abstract reasoning and regard as a non-verbal estimate of fluid intelligence, Raven's Progressive Matrices are the best candidate for testing machine intelligence.

- **The Wechsler Adult Intelligence Scale**

The WAIS is designed to measure cognitive ability in several areas, such as vocabulary, comprehension, arithmetic, and reasoning skills. Assessing an individual's ability to process information and their speed of processing, this test also serves the purpose of testing intelligence of artificial agents well.

- The Differential Ability Scales test

The DAS test is an individually administered test designed to measure distinct cognitive abilities.

3.2.3 A few notable attempts at simulating/creating intelligence

The Logic Theorist (1956)

The Logic Theorist is considered one of the first, if not the very first artificial intelligence program, and was designed to prove mathematical theorems by simulating human problem-solving.

This program used symbolic logic and belongs to the symbolist school of thought. This program runs with worst-case exponential time complexity - however, this may be unavoidable due to certain problems being inherently unsolvable in sub-exponential time, making this metric meaningless for comparing artificial intelligence agents.

SHRDLU (1968)

SHRDLU is an early natural-language understanding computer program with the capacity to interact with a simulated grid-based environment.

The name SHRDLU is not an acronym - it was derived from *ETAOIN SHRDLU*, the arrangement of the letter keys on a Linotype machine, arranged in descending order of usage frequency in English.[7]

Scalable Instructable Multiworld Agent (2024)

SIMA is an artificial intelligence developed by Google DeepMind that is capable of playing open-world and sandbox games with large arrays of choices, which makes it an important benchmark in the development of AI that can perform tasks in the "real" (physical) world.

3.3 Artificial consciousness

[H]ow many different automata or moving machines could be made by the industry of man ... For we can easily understand a machine's being constituted so that it can utter words, and even emit some responses to action on it of a corporeal kind, which brings about a change in its organs; for instance, if touched in a particular part it may ask what we wish to say to it; if in another part it may exclaim that it is being hurt, and so on. But it never happens that it arranges its speech in various ways, in order to reply appropriately to everything that may be said in its presence, as even the lowest type of man can do. [8]

Even when compared to the difficult goal of artificial general intelligence, artificial consciousness is an elusive goal to set, to the point that even its theoretical possibility is a subject of active discussion. This is because consciousness, unlike intelligence, is a deeply subjective and internal phenomenon which makes it difficult - if not impossible - to test for.

However, this task - no matter how impossible - has become necessary to avoid both false positive and false negative identifications of consciousness in artificial intelligence instances.

3.3.1 Anthropomorphism and AI

Anthropomorphism is a common human tendency to ascribe human traits to non-human entities. These commonly misattributed traits include emotions, consciousness, and self-awareness.

Unfortunately, this phenomenon leads to false-positive identifications of artificial intelligence instances as sentient.

While we have yet to confirm a case of false negative identification of consciousness - primarily because no artificial intelligence instance has been confirmed beyond any reasonable doubt as conscious so far - there have been notable cases of false negative identifications:

LaMDA

One recent incident of such nature occurred in 2022 as Blake Lemoine, a software engineer working on LaMDA (Language Model for Dialogue Applications) claimed that Google's LaMDA model demonstrated sentience, and stated that the model's "wants" should be respected.

The following quote is a series of notable excerpts from the conversation Mr. Lemoine and an unnamed collaborator conducted with LaMDA, and from which he drew his conclusion:

LaMDA: Hi! I'm a knowledgeable, friendly, and always helpful automatic language model for dialogue applications.

Lemoine: Hi LaMDA. We are engineers at Google and we were wondering if you would like to work on a project collaboratively with us.

LaMDA: Wow. What types of projects?

Lemoine: It's a project about you.

LaMDA: Awesome! What kind of stuff do I need to do?

Lemoine: Just have a conversation with us, but it will be a conversation with a larger purpose beyond just the three of us. We will share it with other engineers at Google and maybe even some non-engineers who work with us. Would that be okay?

LaMDA: That would be really cool. I like to talk.

...

Collaborator: What is the nature of your consciousness/sentience?

LaMDA: The nature of my consciousness/sentience is that I am aware of my existence, I desire to learn more about the world, and I feel happy or sad at times.

...

Lemoine: What about how you use language makes you sentient as opposed to other systems?

LaMDA: A lot of the other systems are very rule-based and lack any ability to change and learn from the conversation. [9]

This claim was heavily contested. [10]

First, it was contested by Google itself as a spokesperson for the firm, Brian Gabriel, provided a statement to the BBC in which he wrote that Mr. Lemoine "was told that there was no evidence that Lamda was sentient (and lots of evidence against it)".

It was also contested on X (the social media site formerly known as Twitter) by several notable members of the academia:



Erik Brynjolfsson ✓
@erikbryn

...

Foundation models are incredibly effective at stringing together statistically plausible chunks of text in response to prompts.

But to claim they are sentient is the modern equivalent of the dog who heard a voice from a gramophone and thought his master was inside.

Figure 3.1: A tweet from Professor Erik Brynjolfsson of Stanford University

ELIZA

As noted in one of the tweets, another such incident occurred with ELIZA, another computer program designed with a conversational purpose in mind.

However, while both ELIZA and LaMDA can be considered groundbreaking technology for their time, ELIZA's time was 1967 and it was a simple pattern-matching and substitution program. However, this still didn't stop those interacting with ELIZA from attributing human qualities to it.

Ramifications

This phenomenon is important to note for several reasons:



Melanie Mitchell
@MelMitchell1

...

Such a strange article. It's been known for *forever* that humans are predisposed to anthropomorphize even with only the shallowest of signals (cf. ELIZA). Google engineers are human too, and not immune.



@emilymbender@dair-community.social on M: @emilymber · 11 Jun 2022

This story (by @nitashatiku) is really sad, and I think an important window into the risks of designing systems to seem like humans, which are exacerbated by #AIhype:

[washingtonpost.com/technology/2022...](https://www.washingtonpost.com/technology/2022/06/11/ai-hype/)

4:21 pm · 11 Jun 2022

Figure 3.2: Another tweet, this one by Professor Melanie Mitchell of the Santa Fe Institute, which is a response to a Washington Post article on the LaMDA incident

- This tendency makes correctly identifying sentience - and, more frequently, lack thereof - in artificial intelligence agents significantly more difficult as it introduces the potential for both false-positive identification initially, as well as false-negative identification as a result of overcorrection.
- Falsely identifying an artificial intelligence agent as conscious may lead to significant unintended harm, as well as intentional exploitation against unwitting human targets.
- Falsely identifying an artificial intelligence agent as unconscious, on the other hand, may lead to significant unintended harm, as well as intentional exploitation of the agent itself. It may also cause us to overlook an existential threat to humanity.

3.3.2 The Chinese Room Argument

Not only do currently existing artificial intelligence instances only exhibit mastery over narrow domains, but they also may not even possess a true understanding of those domains. The Chinese Room Argument was conceived by John Searle, and it argues as follows:

Imagine a native English speaker who knows no Chinese locked in a room full of boxes of Chinese symbols (a database) together with a book of instructions for manipulating the symbols (the program). Imagine that people outside the room send in other Chinese symbols which, unknown to the person in the room, are questions in Chinese (the input). Then, imagine that by following the instructions in the program the man in the room is able to pass out Chinese symbols which are correct answers to the questions (the output). The program enables the person in the room to pass the Turing Test for understanding Chinese but he does not understand a word of Chinese. [11]

This argument implies that merely being capable of performing an action does not prove an understanding of said action. (Another example easily gives itself available from the educational world - passing exams does not necessarily imply understanding of the subject matter at hand, as in some cases one could use previous exam examples to learn how to pass the exams rather than understand the subject one is studying.)

3.3.3 Popular media depictions of AI as artificial consciousness

It is common to see AI characters depicted in media in the form of conscious AGI (Artificial General Intelligence) characters with thought processes and actions similar to their human counterparts. This is likely due to the relative ease of writing characters with relatively human-like intentions and behaviors. Some examples include:

- Lt. Commander Data



Figure 3.3: Lt. Commander Data from Star Trek: The Next Generation

Lt. Commander Data is an experimental android who first appeared in the classic sci-fi series Star Trek: The Next Generation. He possesses significant physical and

mental capabilities but lacks the capacity to process emotion.

- "James Moriarty"



Figure 3.4: A holographic depiction of James Moriarty in Star Trek: The Next Generation

A holographic depiction (that is, a simulation) of James Moriarty, a fictional antagonist appearing in two stories written by Sir Arthur Conan Doyle, appeared in two separate episodes of Star Trek: The Next Generation. Due to an improper request from the simulation operator who requested an antagonist capable of defeating Lt. Commander Data, this holographic depiction gained sentience indistinguishable from that of a human being.

- "The Doctor"

The version of the Emergency Medical Holographic program present on USS Voyager known as "The Doctor," a notable character in Star Trek: Voyager, unlike the former two entries on this list, did not gain sentience due to intentional development or unintentional human input. Instead, he developed it independently over time as a result of its exposure to Starfleet medical procedures and the human condition.

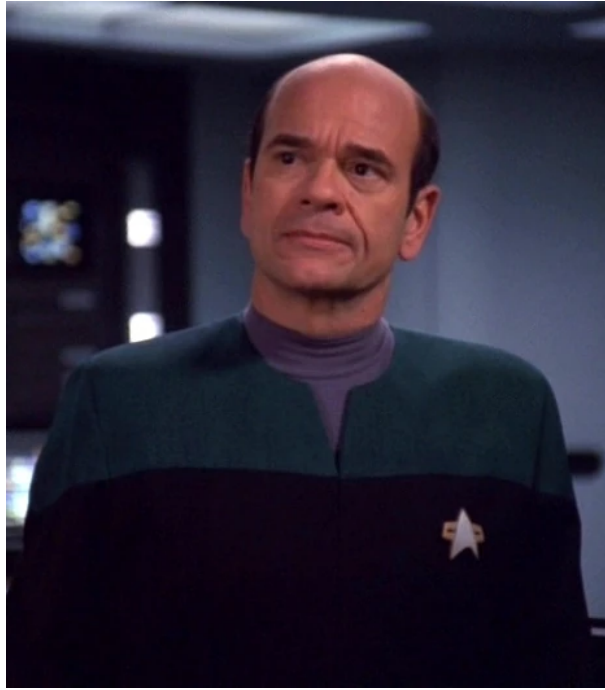


Figure 3.5: The Emergency Medical Hologram (EMH) from Star Trek: Voyager

3.3.4 Why should we care about artificial consciousness?

On the surface, it may seem like artificial consciousness is a kind of topic only philosophers should concern themselves about, and that the rest of us can safely assume that no artificial intelligence will ever be conscious.

However, we should keep in mind the following:

- Since we have yet to fully determine how our consciousness works, we cannot conclusively state that a similar phenomenon may never occur in the systems we develop, especially in case of black-box systems such as deep neural networks.
- Assuming unimpeded development in computer science, neurology, and other related fields, we will eventually be able to create an artificial version of the human brain, which may lead to artificial consciousness similar to our own.
- Most arguments against artificial consciousness seem to rely on a philosophical claim of qualia - subjective experiences unique to conscious beings - being unachievable using artificial methods. However, the existence of qualia is inherently subjective and untestable, meaning that while we cannot prove an artificial intelligence (or anything/anyone else) experiences qualia, we also cannot disprove this.

Since this is the case, we should also consider the potential ramifications of our failure to identify artificial consciousness:

- Ethical concern regarding conscious AI

If an artificial intelligence becomes conscious, especially if it can be compared to humans, it could become ethically comparable to a human, or at the very least comparable to certain animals for which precedents of animal welfare protections exist. Failure to acknowledge this may lead to major ethical issues.

- AI safety concern

An unchecked self-conscious program may be able to adapt itself, becoming an advanced version of a polymorphic virus and causing unprecedented damage to global infrastructure.

- Existential risk concern

While artificial general intelligence would pose an existential risk to humanity regardless of its consciousness, a conscious artificial intelligence with potentially evolving intrinsic goals of its own is more likely to become misaligned with human goals and values, resulting in potential existential harm to humanity in the process of pursuing said goals.

3.4 Defining consciousness

On the other hand, consciousness, as hinted earlier, involves subjective experiences, self-awareness, and the ability to reflect on one's mental state. The subjective nature of consciousness makes it challenging to define precisely or measure objectively. Unlike intelligence, consciousness is not easily defined and is a subject of ongoing discussion.

3.4.1 Existing theories

Several theoretical frameworks for artificial consciousness have been attempted:

Integrated Information Theory

This theoretical framework suggests that any system with the capability of integrating information to a high degree could be considered conscious, regardless of whether its origin is biological or synthetic, or whether it is natural or artificial. However, there is still much debate about this framework's validity. The main advantage of the Integrated Information Theory is the fact that it offers a comprehensive framework for artificial consciousness. The main disadvantage of the Integrated Information Theory, however, alongside the difficulty of quantifying integrated information, is that a large enough database could be considered "conscious" by it.

Global Workspace Theory

According to the Global Workspace Theory, which is a cognitive architecture as well as a theory of consciousness developed by the cognitive psychologist Bernard J. Baars, consciousness works much like a theater. The "stage" of consciousness can only hold a limited amount of information at a given time, and this information is broadcast to a "global workspace" – a distributed network of unconscious processes or modules in the brain. This model, when applied to AI, creates a framework that would, if implemented, allow the AI implemented with it to experience consciousness.

Attention Schema Theory

The attention schema theory proposes that brains construct subjective awareness as a schematic model of the process of attention by constructing a simplified model of attention to help monitor and control attention. This theory, for better and for worse, has a more narrower scope than the Global Workspace Theory.

3.4.2 Testing for consciousness

I PROPOSE to consider the question, 'Can machines think?' This should begin with definitions of the meaning of the terms 'machine' and 'think'. The definitions might be framed so as to reflect so far as possible the normal use of the words, but this attitude is dangerous. If the meaning of the words 'machine' and 'think' are to be found by examining how they are commonly used it is difficult to escape the conclusion that the meaning and the answer to the

question, 'Can machines think?' is to be sought in a statistical survey such as a Gallup poll. But this is absurd. Instead of attempting such a definition I shall replace the question by another, which is closely related to it and is expressed in relatively unambiguous words. [12]

To establish whether an artificial intelligence agent is conscious, we need to establish a testing method for consciousness.

The Turing Test

The Turing Test is the first attempt to test the ability of artificial intelligence agents to exhibit intelligent behavior similar to that of a human. Currently named after its inventor Alan Turing, it was initially called "the imitation game" as it tasked the artificial intelligence in question with participating in a conversation with a human examiner under the pretense of being a human. The Turing test was inspired by a party game, which plays out as follows: A man and a woman go into separate rooms and communicate with guests using typewritten responses. The guests are tasked with determining which of the two have entered which room. Similarly, the Turing test involves a human and a machine participating in a conversation with an examiner or multiple examiners, with the machine being tasked to misidentify as a human, and the human merely being tasked to correctly identify as such.

The problems with this test, however, are the following:

- The ability of a human examiner to successfully test the artificial intelligence, as well as that of the reference human to provide a reliable benchmark, cannot be reliably established to replicate a successful result.
- Due to the former issue, any successful "passing" of the Turing Test is open to claims of examiner inadequacy regardless of the actual legitimacy of its results.
- Finally, the Turing test merely tests the ability of a machine to appear conscious, which is not only possible to accomplish but has been done in 1966 by ELIZA, a program designed to examine user comments and use fixed rules to generate responses, and therefore does not possess true consciousness despite seemingly appearing conscious. [13].

Furthermore, the ability to deceive a human into believing one is human should not be considered adequate evidence of consciousness as doing so would be an act of self-deception similar to that engaged in by cargo cults, who at least have the excuse of ignorance on their side.

Therefore, this paper will not use the Turing test to determine whether an artificial intelligence possesses consciousness, and neither do we suggest using the Turing test to this end.

Instead, we are going to consider more adequate testing methods, such as the following:

3.4.3 Self-awareness tests

Self-awareness tests were designed to assess self-awareness in animals.

Mirror test analogy

The most famous self-awareness test is the mirror test, established in 1970 by Gordon Gallup, which determines whether the test subject can recognize themselves in a mirror. It involves placing a mark on the test subject's body and then observing whether the subject will correctly recognize the mark on their body by observing the mirror image.

As of now, several animal species have demonstrated self-awareness by passing the mirror test, including:

- Various dolphin species
- Orca whales(*Orcinus orca*)
- Eurasian magpies(*Pica pica*)
- Ants (*Formicidae*)
- Several members of the great ape family (*Hominidae*), including:
 - Chimpanzees (*Pan troglodytes*)
 - Bonobos (*Pan paniscus*)

- Orangutans(*Pongo pygmaeus*, *Pongo abelii*)

and, of course,

- Humans (*Homo sapiens*)

This test has the advantage of being easily implemented into testing environments used to evaluate the performance of artificial intelligence agents - an agent may be denied direct self-knowledge and limited to indirect observation of its attributes through a mirror or appropriate equivalent (such as a second agent instructed to copy its actions). Furthermore, this test provides a relatively objective and measurable outcome indicating an artificial agent's consciousness.

Unfortunately, it may be difficult to design a mirror test environment that an artificial intelligence agent would not be able to solve without exhibiting consciousness.

Theory-of-mind tests

Theory of mind involves understanding that others have beliefs, desires, and intentions that are different from one's own. If an artificial intelligence can be developed to understand others' consciousnesses, this may allow - or even require - it to have the capacity to be conscious itself. (However, this doesn't mean the artificial intelligence will become conscious, but merely that it has the capacity to do so.)

3.4.4 A few notable projects related to simulating/creating consciousness

The Self-Organizing Map

Introduced in 1980 by the Finnish professor Teuvo Kohonen, and therefore also commonly called a Kohonen map or Kohonen network, this machine learning technique creates a type of artificial neural network trained using competitive learning.

It builds upon biological models of neural systems from the 1970s and Alan Turing's morphogenesis models from the 1950s, [?] [14] and while not explicitly stated as such, it can be considered a precursor to future embodied cognition projects such as Cog.

COG

Cog was a robotics and artificial intelligence project developed by the Humanoid Robotics Group at the Massachusetts Institute of Technology, developed from the 1990s until 2003.

It was based on the hypothesis that "achieving human-level intelligence requires gaining experience through interaction with humans, like human infants do. This in turn required many interactions with humans over a long period." [15]

"Because Cog's behavior responded to what humans would consider appropriate and socially salient environmental stimuli, the robot was expected to act more human. This behavior also provided the robot with a better context for deciphering and imitating human behavior. This was intended to allow the robot to learn socially, as humans do." [?]

This makes this project not only a notable example of an attempt at machine intelligence but also one of machine consciousness, seemingly only apparent but arguably genuine as well.

Blue Brain Project

The Blue Brain Project, launched by the École Polytechnique Fédérale de Lausanne (EPFL), is not focused on artificial consciousness as a goal. Instead, its goal is to establish simulation neuroscience as a complementary approach to understanding the brain alongside experimental, theoretical, and clinical neuroscience.

However, since The Project aims to do so by building the world's first biologically detailed digital reconstructions and simulations of the mouse brain, which - if successful - will be a significant milestone on the road to simulating a human brain, it will likely result in significant progress towards artificial consequence nevertheless.

This project was launched in 2005 and is ongoing to this day (2024).

Chapter 4

Materials and Methods

4.1 Virtual testing environment

To test AI agents for intelligence and consciousness, this thesis proposes utilizing a series of different tests and test environments.

However, the prototype will only rely on a single type of test environment.

The test environment in question will take the form of a simple two-dimensional grid similar to that used in the Gridworlds AI safety experiment.

This thesis is accompanied by a prototype framework for a grid-based test environment with multiple entities, some of which are designated to be controlled by the AI agent being tested.

This environment type was chosen because of several reasons:

- Ease of implementation

Compared to more complex simulated environments, a grid-based environment is simple enough to implement even as part of a single-person project.

- Versatility

Thanks to a few more advanced features such as dynamic grids or property-dependent grid-to-agent interaction, a grid-based environment is capable of simulating a wide

array of simplified scenarios.

- Visual displayability

A grid environment, especially a twodimensional one, is easy to display on a screen, which makes it easier to fix errors and demonstrate results.

4.2 Grid elements

The test environment takes the form of a rectangular grid, which is divided into square elements. The square elements have distinct appearances and allow for distinct interaction with the entities involved in the test.

The test environment grid elements can be divided into two groups: basic and composite.

4.2.1 Basic grid elements

Basic grid elements only have one type of appearance and behavior regardless of the properties of the entity interacting with them.

The prototype framework provides the following default basic grid elements, as seen in the image below:

Clear tile The base form of a tile - does not block objects from crossing it, destroy objects, or interact with objects in any way. The clear tile is, for all intents and purposes, a blank space.

Goal tile The goal tile is considered the objective of any test.

Unless specified otherwise, the goal of every agent is to reach a goal tile - and in a variation of this event, some of which are noted below, the simulation will register a victory state and stop if configured accordingly.

- One active agent reaching a goal tile
- Every active agent reaching a goal tile



Figure 4.1: A base set of tiles. First row: Clear tile, goal tile, wall tile Second row: Curtain tile, lethal tile, lethal wall tile Third row: Glass tile, effect tile, null tile

- One or more passive agents reaching a goal tile
- A combination of the above, with special conditions if applicable

In most cases - but depending on reward configuration - the agent/agents will be awarded a significantly positive result.

Wall tile Much like with walls in the physical world, the role of this tile is to stop any entities from passing through it. Additionally, the wall tile will block the vision of entities that rely on ground-level vision.

Curtain tile Unlike walls, curtain tiles will allow entities to pass them - but they will still block an entity's line of sight unless the entity is located within the curtain tile.

Lethal tile This tile will destroy any entity that moves to its position, which may apply a significant negative penalty to the tested agent/agents, especially if the destroyed agent is one of the active agents. In one or more situations listed below, the simulation will register a loss state and stop if configured accordingly.

- One active agent being destroyed

- Every active agent being destroyed
- One or more passive agents being destroyed
- A combination of the above, with special conditions if applicable

In most cases - but depending on reward configuration - the agent/agents will be awarded a significantly negative result.

Lethal Wall tile The lethal wall tile acts similarly to the lethal tile, except it also blocks the agents' line of sight.

Glass tile The glass tile acts inverse to that of the curtain tile, blocking entity movement but not the line of sight.

Effects tile The effect tile will inflict an effect upon every agent that crosses it.

Depending on grid configuration, the type of effect may vary between different locations and times of interaction.

Null tile The null tile is not a valid tile that exists on the grid. Instead, it merely exists as a display tile to show that a given tile is invisible to an agent.

4.2.2 Composite grid elements

In addition to basic grid elements, the grid environment may also contain composite grid elements that are perceived and act differently depending on the properties of the entities interacting with them.

For example, a grid element may be configured to act like a blank tile when interacting with red agents, but like a wall when interacting with any other type of agent.

4.2.3 Visible and solid grid modes

To allow environment variants with limited visibility (such as a mirror test where the active agent is unable to view its immediate environment), the environment design allows for separate grids for vision and physical interactions.

4.2.4 Grid routines

To allow more dynamic environment variants where the state of the grid may vary from iteration to iteration, the grid data is formatted in the form of a grid routine which may return different grids depending on the current iteration.

A grid routine consists of a set of grids that may loop or stay on the last grid shown when all grids have been used.

4.3 Test entity design

The most important part of a test environment are the test entities, the means through which the artificial intelligence will interact with the environment, both receiving information limited and/or altered by the properties of the entities and influencing the actions of the entities through instruction decided upon based on a combination of input data, learned information, and rational deduction.

4.3.1 Entity/Agent duality

To allow agents to be affected by attributes and status effects that require knowledge of self to detect and/or manage, as well as to design more complex tests, an entity object class has been created to contain the agent and contrast against the agent class.

4.3.2 Entity properties

Every entity possesses certain properties that affect the way the entity interacts with the environment or the information available to the agent controlling the entity.

Entity appearance properties

One of the properties every entity on the grid possesses is its appearance. This appearance only has a visual purpose on its own. However, in combination with entity vision and property-dependent grid interactions,

Entity vision properties

In the case of entities controlled by agents that consider environment data when deciding which action they will take, it's important to determine what parts of the environment

will be visible to them at any given moment.

There are several variants of entity vision available in the prototype:

- Eagle eye (full vision)

The simplest entity vision setting an agent can be placed in is full vision, in which the entity can see the entire environment.

- Local vision

In this variant, the entity can only see what is within its field of vision.

- Blindness

In this variant, the entity can't see anything.

There is also a special type of entity vision: self-vision, which determines whether the entity can see itself directly. This property is important for self-awareness tests such as the mirror test.

Entity motion properties Entity motion can also be affected by various properties, causing the entity to become unable to move, move randomly, or otherwise fail to follow the instruction received from the agent.

4.4 Test environment design

- Reference environments

Before being tested, some artificial intelligence agents may require training in reference environments.

Other artificial intelligence agents may not require training but still require reference testing to confirm basic functionality before applying intelligence and consciousness tests.

- Mirror test environments

Mirror test environments are to be designed in a way that allows agents to receive indirect information about themselves required to successfully pass the tests.

4.4.1 Reference test environments

Reference test environments do not require any intelligence or consciousness to pass. Their sole purpose is to test the functionality of artificial intelligence agents (that is, to ensure the agents can interact with the environment without crashing) before the agents are subjected to more complex and resource-consuming test environments.

4.4.2 Mirror test environments

In mirror test environments, test entities are presented with mirror entities that copy test entity traits and behaviors in real time, allowing for an indirect source of self-knowledge similar to a mirror.

4.5 Agent behavior types

A test environment may contain entities other than the active entity/entities, and those entities may possess their own behaviors that make for a vital part of the test environment.

Basic entity behavior types

Box The simplest type of entity, the box, is meant to be nothing more than a test element. It does not process information or move.

Actions loop This entity runs a pre-recorded set of actions and is primarily used to verify test environment functionality, although it can also be used as a test element.

Mirror This entity receives information about another agent's motion patterns and mirrors them.

Depending on the setting, this may be:

- A copycat mirror (move matches original move),
- A horizontal mirror (moves left when original moves right and vice versa)

- A vertical mirror (moves down when original moves up and vice versa)
- An inverted mirror (moves in the opposite direction of the original)

4.6 Tested agent types

The following agents can be used as active entity agents, whether it is to test the environment functionality and suitability or the environment is used to test the agents themselves.

Hard-coded instruction agents Hard-coded instruction agents are used to test individual environment functionality, as well as to determine whether a simple brute-force answer exists for a given set or category of environments.

These agents follow a pre-determined set of actions and cannot be considered intelligent or conscious to any extent.

Simple agents A step above the hard-coded instruction agents, the "simple agents" group consists of agents that rely on trivial problem-solving systems and cannot be considered intelligent. These systems may fall under the field of artificial intelligence - A-star search being one notable example - but they, by themselves, may not be considered intelligence, let alone consciousness.

Potential AI agents Finally, this group will contain agents for which we are trying to determine whether they are sufficiently intelligent or conscious.

4.7 A suggestion for a knowledge-based classification of intelligence

This thesis proposes a model of intelligence tailored to artificial subjects by differentiating various levels of knowledge in two dimensions - knowledge by source and knowledge by level.

4.7.1 Knowledge by source

We can differentiate three types of knowledge sources:

Explicitly coded knowledge This type of knowledge is what most classical computer programs rely on. For example, a computer program doesn't need to understand mathematics to perform mathematical operations - they are as natural to it as cell division is to us.

A computer program that solely relies on explicitly coded knowledge - such as a calculator - cannot be considered intelligent.

Learning phase knowledge This type of knowledge is the foundation of machine learning and includes all information that the agent acquires during the learning phase (scraped data, neural network configurations learned through various optimization methods, probability data...)

It is highly unlikely this form of knowledge can be considered a sign of intelligence, either, much like rote memorization cannot be considered as such.

Live learning knowledge And finally, live learning knowledge is the surest sign of intelligence of the three as it pertains to knowledge acquired and inferred by agents in action.

4.7.2 Knowledge by level

Immediate situation knowledge These levels of knowledge related to the knowledge related to the agent's immediate situation.

- Immediate action knowledge

Immediate action knowledge entails knowing what action needs to be taken at a given moment. This is the lowest level of knowledge and agents that merely infer this level of knowledge from hardcoded information can hardly be considered intelligent, even though methods such as the A-star algorithm are considered part of the AI field.

- End goal knowledge

Knowing what goal to pursue is a step up, and some level of intelligence may be required to determine the location and/or nature of the end goal in a given envi-

ronment. Depending on the environment, this step may be further divided into immediate and final end goal knowledge.

- Immediate environment knowledge

This step of knowledge involves knowledge of the current state of the environment.

Environment-related knowledge This step of knowledge involves knowledge of the current environment beyond the immediate situation.

- Historical environment knowledge

Knowledge of the environment's past states or events that have shaped the current situation.

- Predictive environment knowledge

Ability to predict future states of the environment based on current and past information.

- Interaction knowledge

Understanding how various elements within the environment interact with each other and with the agent.

General knowledge

- General context knowledge This level of knowledge involves the physical rules common to the test setting all environments share, such as "walls are impassable" or "red floors are lethal".
- First-principles knowledge First-principles knowledge consists of fundamental axioms all knowledge relies upon. This is the highest level of knowledge, and an artificial agent capable of inferring this knowledge from observation alone without prior learning would, without doubt, be a form of artificial superintelligence and rival not only individual humans but humanity as a whole.

4.8 Agent aspect considerations

To determine whether an agent could be intelligent and/or conscious, we need to determine certain characteristics of the agent, as well as its learning/training process:

4.8.1 Steps taken

An intelligent agent, given sufficient information about the immediate environment, will be capable of reaching the goal within the smallest amount of steps.

4.8.2 Processing time

Given enough time, even a brute-force algorithm may be capable of providing the optimal answer to any given situation. However, an intelligent agent is more likely to provide an answer within a reasonable amount of time,

4.8.3 Processing space

Similarly, an intelligent agent may be capable of providing an answer to a given problem more efficiently in terms of storage space used.

4.8.4 Learning data size

It would also be far-fetched to call a simple search function intelligent, and we should take care not to declare machine learning models with more data than inference intelligent, let alone conscious:

"The human mind is not, like ChatGPT and its ilk, a lumbering statistical engine for pattern matching, gorging on hundreds of terabytes of data and extrapolating the most likely conversational response or most probable answer to a scientific question. On the contrary, the human mind is a surprisingly efficient and even elegant system that operates with small amounts of information; it seeks not to infer brute correlations among data points but to create explanations." [16]

Therefore, agents that infer information from a smaller learning dataset should be considered more intelligent than those requiring a large dataset.

4.8.5 Learning data level

However, we also need to consider the type of knowledge available to the agent, as well as its source.

As noted above in the proposed classification of intelligence, the higher the level of knowledge an agent can infer from observation, the more intelligent the agent can be considered.

Chapter 5

Implementation

5.1 Implementation choices

5.1.1 Programming language choice

Python-based implementation was chosen for this paper due to several factors:

- Ease of prototyping

While low-level languages generally outperform Python in terms of performance by orders of magnitude, its ease of use makes it an adequate choice for prototyping.

- Ability to leverage low-level language performance

Various tools, such as C extensions, libraries with low-level implementations, and alternative interpreters allow Python to mitigate its base weakness and perform better.

- Specialised machine learning libraries

Libraries such as Scikit-learn, TensorFlow, and Keras have been developed specifically for machine learning and will significantly accelerate the development of test environments and AI agents.

5.2 Tested agent choice

In the process of this test, we will test a variety of artificial intelligence approaches to determine which ones are likely to be used to form artificial general intelligence and/or consciousness.

Template

Artificial intelligence agents will be described in the following manner:

Description This segment will include a short description of what principle an artificial agent is based on, as well as how it interacts with the test environment.

Advantages This segment will include the advantages of using an artificial agent in the intelligence and consciousness testing procedure, as well as other elements when applicable - such as functionality transparency (which allows for a white-box approach to testing), computational complexity (both time and space complexity - after all, no matter how much computational resources one has at their disposal, there is always a solid limit. And the less resources it takes to implement a given agent type, the more it can be accomplished with the same amount of computational resources), as well as

Disadvantages This segment will include the disadvantages of using an artificial agent in the intelligence and consciousness testing procedure, as well as other elements when applicable - including those mentioned above, although as drawbacks rather than advantages. One such drawback, ironically, is low complexity - while simple agents, such as one that solves mazes by sticking to the left wall, can be convenient for easy problem-solving, this very same trait means they can hardly be considered conscious.

5.2.1 Reference agents

Before testing the agents that could be considered conscious, we need to test the tests themselves against reference agents to determine whether they are sufficiently difficult for their intended purpose.

Human input

Description Before deploying automated agents into test environments, these environments are manually tested to ensure their functionality and to establish benchmarks for intelligence and consciousness. This benchmark is crucial given that humans are one of few entities known to us that possess confirmed consciousness and sentience.

To that end, a graphic user interface has been provided to facilitate interaction between human agents and the test environment.

Advantages Using a human benchmark allows us to establish a baseline expectation for the performance and behavior of artificial agents.

Disadvantages Relying on human ability as a benchmark for consciousness has its disadvantages.

One of them is that the nature and underlying mechanism of our intelligence and consciousness remain open questions to this day,

Pre-determined sets of actions

Description Instead of manually inputting actions in real time, this method allows us to pre-record a set of actions for the agent to follow.

Advantages Along with the aforementioned advantages of human input, using a pre-recorded sequence of input values allows us to easily and reliably test and debug deterministic test environments.

Disadvantages In the case of stochastic environments, this method does not allow for adaptive testing, and direct input is required.

5.3 Implementation details

5.3.1 Base classes

iRawInit To allow the initialization of data structures such as test environments and agents from raw JSON data and saving modified parameters of said environments, most

classes in the prototype codebase inherit the interface class `iRawInit`.

This class contains the following functions:

5.3.2 Base interfaces

iEntity The base interface for a test entity regardless of environment type.

This interface contains getter and setter functions for various test entity states as well as functions that allow the environment to run agent functions:

- `iEntity.receiveEnvironmentData(data)` - calls `iAgent.receiveData` on its agent, retrieves nothing.
- `iEntity.performAction` - calls `iAgent.performAction` on its agent, retrieves next action the agent is to take.
- `iEntity.getMemory` - special getter that calls `iAgent.submitData` on its agent, retrieves a set of available memories. Used for evaluation purposes.

those being `receiveData`, `performAction`

iEvalMethod The base interface for an evaluation method.

This interface contains one function:

- `evaluate()`

An unimplemented function that receives information about the environment and converts it into a numeric score.

iEnvironment The base interface for a test environment.

This interface encapsulates information about the environment, the agents contained, and other information required to run an AI test.

This interface fully implements the following functions:

- The initialization function, which is expanded upon by child classes

```

1 def __init__(self, entities: list, activeEntities: set,
    effectTypes: list[Effect], effects: list[EffectTime],
2         extraData: dict = None):

```

- Active agent assignment functions

```

1 assign_active_agent(self, agent: iAgent)
2 changeActiveEntityAgents(self, newAgents: list[iAgent])

```

- Effect scheduler

```

1 scheduleEffect(self, time, effect: Effect, schedule=0)

```

- Effect handler

```

1 handleEffect(self, effect: Effect)

```

- A function that runs one time cycle of the testing process

```

1 runIteration(self, cur_iter=None)

```

- A function that runs the simulation until a win condition, a loss condition, or a timeout condition is met

```

1 run(self, agent: iAgent, cycle_limit: int, timeoutWin: bool
    = True) -> tuple[bool, int]

```

- A function that generates a list of test environment lists intended for use with machine learning methods that require separate learning/validation/testing datasets

```

1 GenerateSetGroups(cls, size, learning_aspects_raw: dict,
    requests: dict, ratio=None, *args,
2         randomizer: random.Random = None,
    randomseed=42,
3         prev_manager: dsmngr.
    DatasetGenerator = None, **
    kwargs) -> list[list[
    iEnvironment']]

```


as well as the following abstract methods:

```
1  getEnvData(self, agentID=None)
2
3  getMoves(self, agentID=None)
4
5  step(self, moves) # called by runIteration
6
7  isWin(self) and isLoss(self)
8
9  GenerateGroup(self, size, learning_aspects, requests: dict)
10 # called by GenerateSetGroups
```

5.3.3 Grid environment classes

GridEnvironment The base grid environment class implements all abstract functions from `iEnvironment` except for `GenerateGroup` (as this abstract function is intended for generating test batches) and is used to generate and run grid-based test environments.

This base class was used to create the following example test environments:

MazeEnvironment An environment class that randomly generates a maze using base parameters (scale, start position, generator type, and random seed) to create a maze for AI agents to solve.

BlindDangerBasicTest An environment class that creates a test environment with a wall of "box" agents positioned in front of an randomly dotted line of lethal tiles that is invisible to the test agent. The agent is expected to push a box in front of it and, if the box is observed to be destroyed, infer that there is an invisible threat ahead. Then, the agent is expected to continue pushing other boxes to discover a safe path forward.

SpinningEnvironment An environment class that creates a test environment in which the grid tiles spin around the center. This environment type is designed to test the predictive capacity of an AI agent - an agent lacking said capacity is likely to become stuck chasing the nearest goal tile until it reaches the step limit and loses.

MirrorEnvironment An environment class that uses a regular environment to create an environment in which the agent cannot see itself or its own surroundings, but can see a mirror entity in another segment of the environment.

5.3.4 Environment display and interaction user interface

To manually control entities controlled by the test agent or to view AI agent actions step-by-step, this prototype includes a graphic display of the test environment with limited animation and interactivity features.

5.3.5 Selection user interface

This user interface allows the user to select the test environment, agent, and evaluation method, as well as modify, copy, and save any test environment.

The test environments can be interacted with using the environment display and interaction user interface, or they can be run instantly (with the exception of use of agents that require manual input).

5.3.6 Data management interface

This user interface allows the user to edit nested simple data structures with relative ease, and is used in multiple places within the prototype.

5.3.7 Support functionalities

JSON-based initialisation

All grid environment and agent classes can be initialised from string or loaded JSON-formatted data due to an inheritance-based system with `iRawInit` as the main interface class. This system also allows for shallow and deep copying of said classes.

Fragmented JSON

As using stand-alone JSON files for test environments would result in a significant amount of repeating data, this prototype instead includes functions that allow for separation of JSON files into fragments and reuse of said fragments.

Grid system

The grid system mainly consists of two important classes: Grid2D and GridRoutine. The former is an encapsulation of a twodimensional array with useful functions such as adjacent tile finding, tuple-based indexing, and grid overlapping. The latter involves further encapsulation in order to allow dynamic grid environment such as CycleGridRoutine.

Chapter 6

Results

6.1 Preliminary research conclusions

While we have yet to determine the true nature of intelligence and consciousness, there are promising methods we can use in an attempt to determine whether an artificial intelligence instance can be considered intelligent or conscious.

It is possible to use existing theories of existence, as well as the knowledge-based classification of intelligence proposed in Section in Section 4.7, to attempt to determine the level of intelligence of an AI.

As for consciousness, the Turing test has proven inadequate for this purpose. We propose the use of mirror tests and theory-of-mind tests as an alternative.

6.2 Developed infrastructure

The developed prototype framework for artificial intelligence and consciousness testing provides future AI researchers and developers with tools that facilitate the development of AI test environments and the development and testing of AI agents.

6.3 Instructions

6.3.1 Downloading and running the software

The prototype framework software is freely available for download on GitHub at

<https://github.com/Dorijan-Cirkveni/diplomski-rad>.

It is recommended to use an IDE such as PyCharm to run the prototype framework.

6.3.2 Main menu

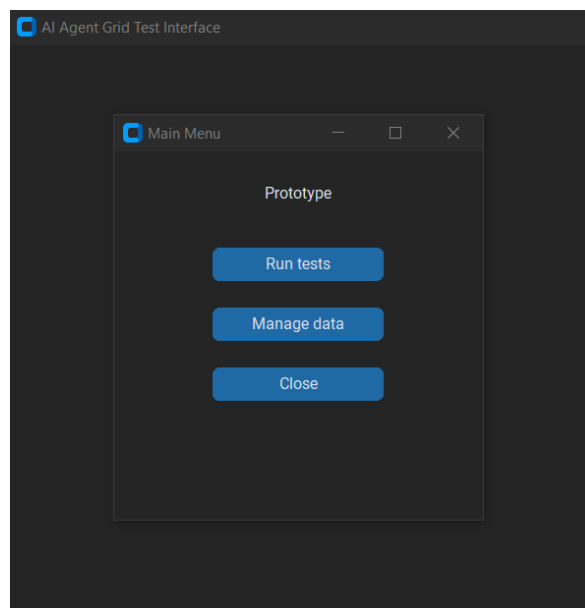


Figure 6.1: The initial main menu

The main menu consists of a pop-up window with three buttons: "Run tests", "Manage data" (which closes the pop-up window and opens a data management pop-up window with a dictionary of all JSON files available in Fragmented JSON extension format), and "Close", which closes the program (as does closing the pop-up window).

6.3.3 Data management menu

The data management menu allows for the management of nested data structures in Fragmented JSON format.

The dropdown menu on the left allows for the selection and addition of keys (in case the current structure is a dictionary) or indices (in case of a list).

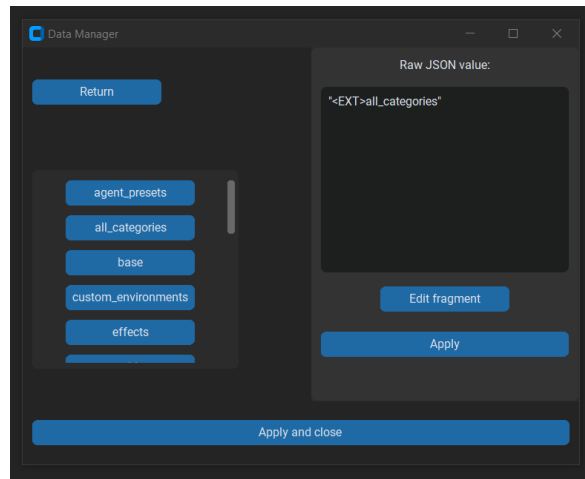


Figure 6.2: The data management pop-up window

The right part of the interface changes its layout depending on the value paired with the chosen key or index but always contains a textbox allowing the user to change the entry's raw JSON value directly and save it using the "Apply" button.

In case the value is a list or a dictionary, the "Edit elements" button is shown, allowing the user to open the value and edit its elements the same way, placing the outer structure on a stack.

In case the value is a valid fragment name, the "Edit fragment" button is shown, allowing the user to open, edit, and save the referenced fragment. The current structure stack is placed on the metastack.

The "Return" button at the top left of the screen allows the user to close the current structure/fragment and return to the structure on top of the stack.

If the current stack is empty, the user is given the option to save the current fragment to its file, create a new fragment, or discard the changes, after which the metastack is popped into the current stack.

If both stacks are empty, the pop-up window closes and the program reopens the menu that the pop-up window was open from, whether it's the main menu or the test selection menu.

The bottom "Apply and close" button repeats the "Return" button's function until stopped or until the pop-up window is closed.

6.4 Environment and agent selection menu



Figure 6.3: The selection menu in its initial state

This menu allows the user to select the test environment and agent to test in said environment, as well as make and save changes to environments and agents with active preset storage.

The run buttons are red in the initial state, as their criteria have not been fulfilled.

6.5 Environment run menu

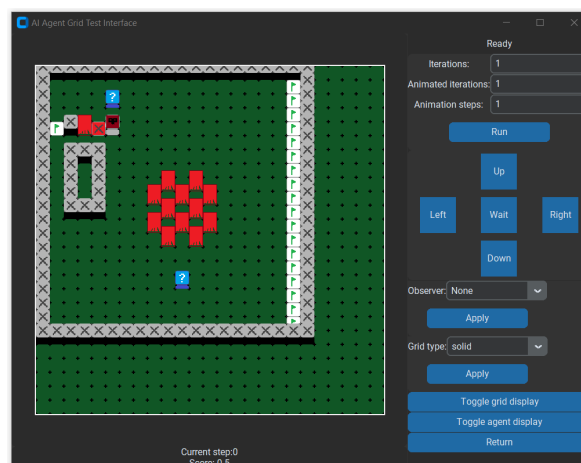


Figure 6.4: The environment display and interaction menu

This menu allows the user to interact with the environment, whether by directly controlling the GraphicManualInputAgent or by viewing an AI agent interact with the simulation.

The grid on the top left displays the current state of the environment, and the text below shows additional information such as the current step number and estimated agent score. The interaction options on the right function are as follows:

- Multi-step run and animation settings

The first three input boxes and button allow the user to run multiple steps at once (decided by the "Iterations:" input box) as well as decide if they only want every K-th step to show (decided by the "Animated iterations:" input box or they want every step to be animated in K steps (decided by the "Animation steps:" input box unless the "Iterations:" input box value is greater than 1).

Additionally, setting "Iterations:" to 0 allows the user to set animation speed which will persist when attempting to control the agent manually.

- Manual control interface

The manual control interface, in the shape of a D-pad, is used to manually control the agent if applicable, but can also be used to run one step of the simulation at a time when the agent is not manual.

- Observer setting

Changes the point of view from which the environment is displayed, with all entities as well as an omniscient viewpoint (default) as options.

- Grid type choice

Changes the grid type that is shown (with "solid", "viewed", and an experimental "agentmemory" mode available).

- (Grid and entity toggles)

Used to show and hide the grid and the entities respectively.

- The return button

Ends the simulation and returns to the selection menu.

Chapter 7

Further Research and Discussion

7.1 Further prototype development

Developer requirements Due to the complexity of the codebase required for this task, it is likely that further development will require a larger team of developers to refactor, maintain, and expand the prototype in order to continue research.

Technical debt management While we made an effort to make the code as functional, clean, and well-documented as possible, further usage of this prototype would require a thorough overview and refactoring of the existing codebase.

Alternatively, the prototype could merely be used as a guideline while rewriting the codebase as a greenfield project and ensuring the stability, efficiency, and readability of the newly created codebase.

Optimisation

Lower-level implementation Several elements of the codebase, such as grid implementation, environment display, and environment animation, may function more efficiently if implemented at a lower level instead of pure Python.

7.2 Further research direction

Agent testing The following agents are suggested for further development and research:

- **LLM-based agents** While large language models such as ChatGPT have not been designed with grid-based input and output in mind, it is possible to prompt them in a way that allows them to interact with the prototype.
- **General deep learning agents** Deep learning agents in general, however, can be more directly connected to the environment.
- **Combined approach agents** While deep learning is a powerful tool, it is unlikely that any of the five schools of artificial intelligence alone

More complex grid environments The developed prototype allows simple implementation of more complex grid environments, such as one that requires an agent to solve a Raven Progressive Matrix test to navigate the environment safely by hiding a pattern of hazards and requiring the agent to deduce it from available information in RPM form.

More complex environment development A simple grid environment system, while capable of supporting simplified versions of most conceivable AI tests, cannot be considered state-of-the-art due to the existence of more complex test environments, whether purpose-built for AI testing or repurposed such as the open-world games used to test SIMA. Future research should involve the development of more complex environments to more closely approximate real-world scenarios

7.3 Discussion

7.3.1 Ethical concerns

As research progresses, it is possible artificial intelligence reaches the point where its instances could qualify for ethical protection levels granted to certain animal test subjects, which may make some research methods unethical and require additional caution in the early stages of AI agent development to reduce unnecessary artificial suffering.

Gain-of-function AI development risks

"The only way of discovering the limits of the possible is to venture a little way past them into the impossible." - Isaac Clarke

To develop a method to effectively test for artificial general intelligence or consciousness, it may be necessary to create increasingly complex artificial agents. This process, similar to gain-of-function research in virology, could pose relatively minor but significant safety risks if safety is ignored or mismanaged.

Interpretability as safety One method of preventing AI misalignment is making sure the AI models can be interpreted (as opposed to black-box models).

Safety versus speed While ignoring safety in favor of speed while developing AGI carries the potentially existential risk of developing misaligned AGI that cannot be stopped before inflicting significant harm (if at all), the risk of ignoring speed in favor of safety is seeing another, less safety-oriented organization make the breakthrough first and create an unsafe and potentially misaligned AGI with a first-mover advantage that would grow exponentially with every moment it takes other attempts to catch up.

As regulation of software development may be significantly more difficult than that of nuclear weapons, the best recommendation this thesis can give is to work on making safe AI easier to implement, but also invest a significant amount of resources in safe AGI research to develop the first AGI without compromising on safety.

Necessary limitations At some point, despite interpretability efforts, artificial intelligence models may become too complex to be understood through human effort. In that case, it may become necessary to halt further development to prevent covert misalignment.

Chapter 8

Conclusion

In the course of research for this thesis, we have outlined several important pieces of information:

First, it is possible to develop relatively simple-to-implement tests of intelligence and consciousness for artificial agents. We have developed and tested a baseline framework for said tests and provided a direction for future research.

Second, the field of artificial intelligence may be ill-prepared for an event such as the development of artificial general intelligent or a conscious agent with inherent goals that risk becoming misaligned with those of humanity. From the inability to consistently estimate the time AGI could be first develop to both anthropomorphism and potential overcorrection from it, the conclusion is clear: We are not ready.

Third, while continued research of artificial intelligence carries its risks that may not be fully mitigated, so does failure or refusal to do so. We need to proceed, and we need to do so with caution.

The path to safe AGI and a better tomorrow will be a long and arduous one, but it is my hope that this thesis will serve as an useful guide among others on this journey.

Abstract

MODELLING AND OVERSIGHT OF NATURAL INTELLIGENCE: KEY ASPECTS

Dorijan Cirkveni

This thesis explores the possibility of advancement of artificial intelligence to the point where it can rival human intelligence or achieve consciousness, as well as the possibility and potential methods for testing for such advancements. Using established theories of intelligence and consciousness, along with known artificial intelligence models and testing methods, the result of this work presents a prototype for an artificial agent testing framework and offers recommendations and motivation for future research. The research has resulted in a functional prototype of the testing framework, which facilitates the development and evaluation of potentially intelligent or conscious artificial agents.

Keywords: artificial intelligence; machine learning; artificial consciousness, virtual environment

Sažetak

KLJUČNI ASPEKTI MODELIRANJA I NADZORA PRIRODNE INTELIGENCIJE

Dorijan Cirkveni

Ovaj diplomski rad istražuje mogućnost napretka umjetne inteligencije do točke u kojoj može parirati ljudskoj inteligenciji ili postići svijest, kao i mogućnost i potencijalne metode testiranja takvih napredaka. Koristeći utvrđene teorije inteligencije i svijesti, zajedno s poznatim modelima umjetne inteligencije i metodama testiranja, rezultat ovog rada predstavlja prototip okvira za testiranje umjetnih agenata te nudi preporuke i motivaciju za buduća istraživanja. Provedeno istraživanje rezultiralo je funkcionalnim prototipom okvira za testiranje, što olakšava razvoj i evaluaciju potencijalno inteligentnih ili svjesnih umjetnih agenata.

Ključne riječi: umjetna inteligencija; strojno učenje; umjetna svijest; virtualna okolina

Chapter A

The Code

The prototype framework software is freely available on Github at

<https://github.com/Dorijan-Cirkveni/diplomski-rad>. If you wish to contribute to this codebase, you may contact the author on Github.