

UNIVERSITY OF ZAGREB
FACULTY OF ELECTRICAL ENGINEERING AND COMPUTING

MASTER THESIS No. 1234

MODELLING AND OVERSIGHT OF NATURAL INTELLIGENCE: KEY ASPECTS

Dorijan Cirkveni

Zagreb, December 2023.

SVEUČILIŠTE U ZAGREBU
FAKULTET ELEKTROTEHNIKE I RAČUNARSTVA

DIPLOMSKI RAD br. 1234

KLJUČNI ASPEKTI MODELIRANJA I NADZORA PRIRODNE INTELIGENCIJE

Dorijan Cirkveni

Zagreb, Prosinac 2023.

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.

TODO

Contents

1	Introduction	4
1.0.1	Recent advances in artificial intelligence	4
1.0.2	Public Adoption	5
1.0.3	Backlash to "AI art" and the question of consciousness	7
1.1	Objective	7
1.2	Choosing an approach	8
1.2.1	Thinking rationally (Artificial intelligence)	9
1.2.2	Acting rationally (Apparent artificial intelligence)	10
1.2.3	Thinking "humanly" (consciousness?)	11
1.2.4	Acting "humanly" (apparent consciousness)	13
1.3	Artificial intelligence	14
1.3.1	The five schools of artificial intelligence	14
1.4	Artificial consciousness	15
1.4.1	LaMDA	16
1.4.2	Anthropomorphism and AI	18
1.4.3	The Chinese Room Argument	19
1.4.4	Popular media depictions of AI as artificial consciousness	20
1.5	Why should we care about artificial consciousness?	22
1.6	Defining intelligence	23
1.7	Types of intelligence	23
1.8	Defining consciousness	24
1.8.1	Existing theories	24
1.9	Testing for intelligence	25
1.10	Testing for consciousness	25

1.10.1	The Turing Test	25
1.10.2	Self-awareness tests	26
1.11	Simulating intelligence	27
1.11.1	Existing attempts	27
1.12	Simulating consciousness	27
1.12.1	Existing attempts	27
2	Materials and Methods	28
2.1	Virtual testing environment	28
2.2	Grid elements	28
2.2.1	Basic grid elements	28
2.2.2	Composite grid elements	30
2.3	Test entity design	30
2.3.1	Entity/Agent duality	30
2.4	Test environment design	30
2.4.1	Reference test environments	31
2.4.2	Mirror test environments	31
2.5	Agent behaviour types	31
2.6	A suggestion for a knowledge-based classification of intelligence	31
2.7	Agent aspect considerations	32
2.7.1	Processing time	33
2.7.2	Processing space	33
2.7.3	Learning data size	33
2.7.4	Learning data size	33
2.7.5	Learning data level	33
2.8	Primary approach choice	33
2.9	Implementation choice	33
2.9.1	Programming language choice	33
2.10	Tested agent choice	34
2.10.1	Reference agents	35
3	Results and Discussion	37
3.1	Instructions for usage of test environment interface	37

3.2 •	37
4 Further research suggestions	38
4.1 Environment improvements	38
4.2 Technical debt management	38
4.3 Optimisation	38
4.4 Future research concerns	38
4.4.1 Gain-of-function AI development risks	38
4.4.2 Safety versus speed	38
5 Conclusion	39
References	40
Abstract	41
Sažetak	42
A: The Code	43

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. [?] TODO

1.0.1 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 and/or non-existent.

By today, artificial intelligence systems have outperformed humans in every one of these five listed fields and more. [?]

- **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.

Language and image recognition capabilities of AI systems have improved rapidly



Figure 1.1: Advances of AI during the last few decades

- **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.
- **Robotic Control** Reinforcement learning has enabled artificial intelligence systems to control robots with unprecedented dexterity and agility, enabling them to perform tasks like picking up objects, grasping tools, and navigating complex and dynamic environments.
- **Financial Trading** Reinforcement learning algorithms are being used to make trading decisions and optimize investment portfolios, demonstrating the potential of reinforcement learning in various real-world applications.

1.0.2 Public Adoption

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

Large language models (LLMs) like Bard, ChatGPT, and Bing Copilot have demonstrated remarkable capabilities in generating human-quality text with minimal prompts and suggestions from the human author of this thesis. These models, powered by sophisticated neural networks and trained on massive amounts of data, can process and understand language in ways that were once thought to be the exclusive domain of humans. Bard, for instance, seamlessly blended the prompts and suggestions from the human author of this thesis into this paragraph, highlighting its ability to grasp the nuances of human language and produce text that is both meaningful and engaging. Similarly, ChatGPT has impressed users with its ability to generate creative text formats, while Bing Copilot has proven adept at assisting with coding tasks. The ability of LLMs to learn and adapt further enhances their potential, as they can continuously refine their skills and expand their knowledge base. As LLMs continue to evolve, we can expect even more innovative and groundbreaking applications to emerge, revolutionizing the way we interact with technology and transforming various fields, from education to health-care to communication.

[1]

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 aforementioned ability to recognize natural language input.



Figure 1.2: An image generated using Bing Image Creator. Prompt: "an artificial intelligence contemplating its existence"

1.0.3 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 its threat to the livelihoods of human artists.

However, there is one reason people oppose 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 can artificial intelligence 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?

1.1 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, although the field of artificial intelligence research 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. This is a difficult task in its own right, as both of these definitions are open questions.

Next, we need to define a metric with which we are going to measure whether an artificial entity possesses those two aforementioned qualities.

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 due to the fact that while these tasks may be prerequisites for one another, they likely require different considerations and different tests are likely to be used to determine their presence in an artificial agent.

1.2 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 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 1.1: Comparison of Intelligence and Consciousness

1.2.1 Thinking rationally (Artificial intelligence)

In this approach we wish to develop an artificial intelligence that thinks rationally, that is, one that is 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/or 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.

1.2.2 Acting rationally (Apparent artificial intelligence)

In mild contrast to the logicist approach, in this approach we ignore the inner world of an artificial intelligence entity and focus on creating and/or identifying rational artificial agents that act rationally. A rational agent is one who 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.

1.2.3 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. In order 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 - despite the fact that 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 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.

1.2.4 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.

1.3 Artificial intelligence

1.3.1 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 in order 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 fundamen-

tal building blocks called artificial neurons and meant to simulate the way our natural neurons work.

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.

1.4 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. [?]

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.

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.

1.4.1 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.

[?]

This claim was heavily contested.

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 1.3: A tweet from Professor Erik Brynjolfsson of Stanford University

1.4.2 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, as was shown above in the case of LaMDA.

This is important to note for several reasons:

- 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.



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/google-engineers-human-too/)

4:21 pm · 11 Jun 2022

Figure 1.4: 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

1.4.3 The Chinese Room Argument

Not only do artificial intelligence instances existing so far 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. [?]

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 sub-

ject 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.)

1.4.4 Popular media depictions of AI as artificial consciousness

It is common to see AI characters depicted in media in 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 1.5: 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"

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



Figure 1.6: A holographic depiction of James Moriarty in Star Trek: The Next Generation

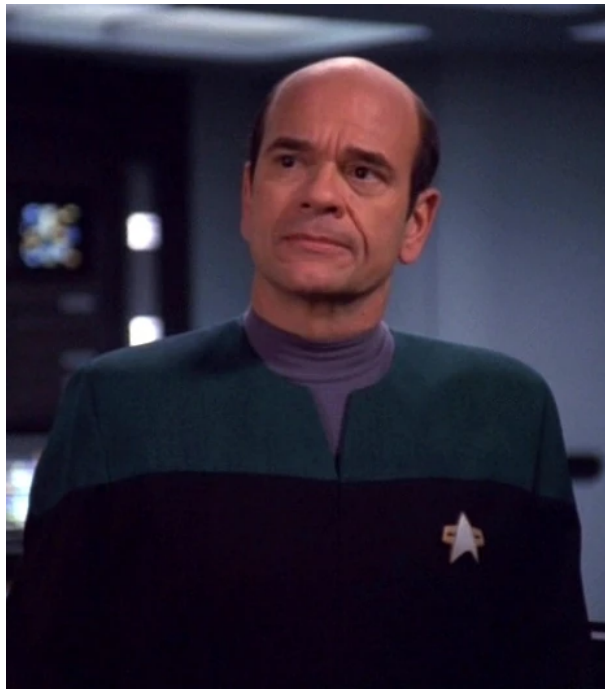


Figure 1.7: The Emergency Medical Hologram (EMH) from Star Trek: Voyager

former two entries on this list, did not gain sentience due to intentional development or unintentional human input. Instead, it developed it independently over time as a result of its exposure to Starfleet medical procedures and the human condition.

1.5 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 own consciousness works, we cannot conclusively deny that an analogue 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.

And 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 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.

1.6 Defining intelligence

We have begun this endeavor with the 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.

1.6.1 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.

1.6.2 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 RPM is a non-verbal test typically used to measure general human intelligence and abstract reasoning and is regarded as a non-verbal estimate of fluid intelligence.[1] It is one of the most common tests administered to both groups and individuals ranging from 5-year-olds to the elderly.
- The Wechsler Adult Intelligence Scale

The WAIS is designed to measure cognitive ability in several areas, such as vocabulary, comprehension, arithmetic, and reasoning skills. These subtests assess an individual's ability to process information and their speed of processing.

1.7 Simulating intelligence

1.7.1 Existing attempts

1.8 Defining consciousness

On the other hand, consciousness, as hinted earlier, involves subjective experiences, self-awareness, and the ability to reflect on one's mental states. 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.

1.8.1 Existing theories

There are several theoretical frameworks for artificial consciousness that 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 implies a clear, measurable metric and criterion an intelligent agents needs to fulfill in order to be considered conscious. The main disadvantage of the Integrated Information Theory, however,

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.

Artificial General Intelligence

An artificial General Intelligence is a type of artificial intelligence that possesses a capability of understanding, learning, as well as application of knowledge across a wide range of tasks, similar to that of a human being - as opposed to existing narrow artificial intelligence systems, which excel in narrowly defined domains and specific tasks, such as voice recognition or playing chess.

In this theory, the general type of artificial intelligence is considered a prerequisite for [?]

1.9 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. [?]

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

Thankfully,

1.9.1 The Turing Test

The Turing Test is the first attempt to test the ability of artificial intelligence agents to exhibit intelligent behaviour 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 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 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 problem with this test, however, is that it merely test the ability of a machine to appear conscious, which is not only possible to accomplish but has been does so 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.[?].

Furthermore 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 consciousness, opting instead to consider more adequate testing methods, such as the following:

1.9.2 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 recognise themselves in a mirror. It involves placing a mark on the test subject's body and then observing whether the subject will correctly recognise 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*)

Advantages 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

Drawbacks

Sensorimotor Contingency Test

1.10 Simulating consciousness

1.10.1 Existing attempts

2 Materials and Methods

2.1 Virtual testing environment

In order to test AI agents for intelligence and consciousness, this thesis will utilise a series of different tests and test environments.

2.2 Grid elements

2.2.1 Basic grid elements

Clear tile The base form of a tile - doesn't block objects from crossing it, destroy objects, or interact with objects in any way. The clear tile is, for all

Goal tile Unless specified otherwise, the goal of every agent is to reach a goal tile, and a variation of this event, such as:

- One active agent reaching a goal tile
- Every active agent reaching a goal tile
- One or more passive agents reaching a goal tile
- etc...

the simulation will stop, and 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 vision of entities that rely on ground-level vision.

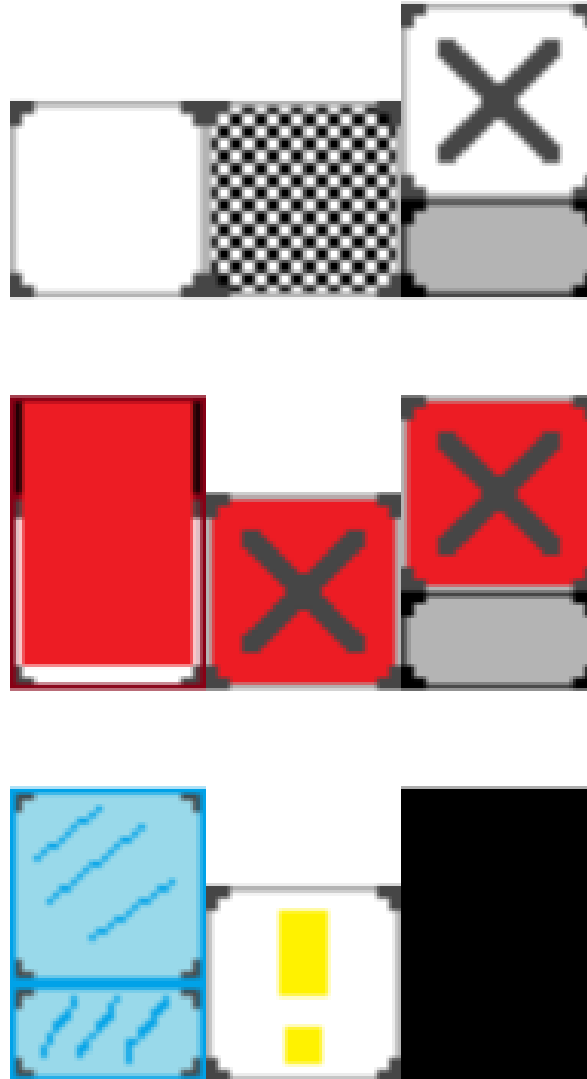


Figure 2.1: All 8 tiles used in grid environment design, as well as a "null" tile used for display purposes.

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

Lethal tile This tile will destroy any entity that moves to its position, which will in most circumstances apply a significant negative penalty to the tested agent/agents, especially if the destroyed agent is one of the active agents. In some circumstances, this may even lead to an early stop of the simulation.

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

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

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

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 not visible to an agent.

2.2.2 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.

2.3 Test entity design

2.3.1 Entity/Agent duality

In order 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.

2.4 Test environment design

1. 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.

2. 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.

2.4.1 Reference test environments

2.4.2 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.

2.5 Agent behaviour types

Basic entity behaviour 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 action and is primarily used to verify test environment functionality, although it can also be used as a test element.

Mirror

2.6 A suggestion for a knowledge-based classification of intelligence

This thesis proposes a model of intelligence based on the Cattell-Horn-Carroll theory, but tailored to artificial subjects by deriving levels of information from the three strata of intelligence, and further distinguishing information by immediate source:

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 math-

ematics in order 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 learning phase (scraped data, neural network configurations learned through various optimisation methods, probability data...)

It is highly unlikely this form of knowledge can be considered a sign of intelligence, either, much like rote memorisation 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.

Knowledge by level

Immediate action knowledge Immediate action knowledge entails knowing what action needs to be taken in a given moment. This is the lowest level of know

End goal knowledge End goal knowledge, however,

Environment knowledge

Immediate context knowledge

Global context knowledge

2.7 Agent aspect considerations

In order 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:

2.7.1 Processing time

Given enough time, even a brute-force algorithm is capable of providing the optimal answer

2.7.2 Processing space

Similarly,

2.7.3 Learning data size

2.7.4 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. [?]

2.7.5 Learning data level

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

2.8 Primary approach choice

2.9 Implementation choice

2.9.1 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 development of test environments and AI agents.

2.10 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.

Interface This segment will briefly describe the type of inter

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.

2.10.1 Reference agents

Before testing the agents that could be considered conscious, we need to test the tests themselves against reference agents in order to determine whether they are

Human input

Description Before using automated agents on test environments, the environments are tested manually to ensure their functionality as well as to attempt to provide a benchmark for consciousness, given that we are the only type of entity with confirmed consciousness and sentience known to us.

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

Advantages By using a human benchmark, we can at the very least establish a rough idea of what we are expecting our artificial agents to accomplish. Obviou

Disadvantages Obviously, this is not an artificial agent - and 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 as open questions to this day,

Pre-determined sets of actions

Large language model

Advantages

Generality

Disadvantages

3 Results and Discussion

3.1 Instructions for usage of test environment interface

TODO add 5-10 pages of instructions

3.2 •

4 Further research suggestions

4.1 Environment improvements

4.2 Technical debt management

4.3 Optimisation

More complex environments

4.4 Future research concerns

4.4.1 Gain-of-function AI development risks

4.4.2 Safety versus speed

5 Conclusion

Hello, here is some text without a meaning. This text should show what a printed text will look like at this place. If you read this text, you will get no information. Really? Is there no information? Is there a difference between this text and some nonsense like “Huardest gefburn”? Kjift – not at all! A blind text like this gives you information about the selected font, how the letters are written and an impression of the look. This text should contain all letters of the alphabet and it should be written in of the original language. There is no need for special content, but the length of words should match the language.

References

- [1] Bard, “Bard, a large language model from google ai,” 2023. [Online]. Available: <https://ai.googleblog.com/2023/01/lamda-language-model-for-dialogue.html>
- [2] A. M. TURING, “I.—COMPUTING MACHINERY AND INTELLIGENCE,” *Mind*, vol. LIX, no. 236, pp. 433–460, 10 1950. <https://doi.org/10.1093/mind/LIX.236.433>

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 accomplish one of the following - rival humans in terms of intelligence, or achieve consciousness.

TODO

Keywords: the first keyword; the second keyword; the third keyword

Sažetak

KLJUČNI ASPEKTI MODELIRANJA I NADZORA PRIRODNE INTELIGENCIJE

Dorijan Cirkveni

Unesite sažetak na hrvatskom.

TODO

Ključne riječi: prva ključna riječ; druga ključna riječ; treća ključna riječ

Appendix A: The Code

Hello, here is some text without a meaning. This text should show what a printed text will look like at this place. If you read this text, you will get no information. Really? Is there no information? Is there a difference between this text and some nonsense like “Huardest gefburn”? Kjift – not at all! A blind text like this gives you information about the selected font, how the letters are written and an impression of the look. This text should contain all letters of the alphabet and it should be written in of the original language. There is no need for special content, but the length of words should match the language.

This is the second paragraph. Hello, here is some text without a meaning. This text should show what a printed text will look like at this place. If you read this text, you will get no information. Really? Is there no information? Is there a difference between this text and some nonsense like “Huardest gefburn”? Kjift – not at all! A blind text like this gives you information about the selected font, how the letters are written and an impression of the look. This text should contain all letters of the alphabet and it should be written in of the original language. There is no need for special content, but the length of words should match the language.

And after the second paragraph follows the third paragraph. Hello, here is some text without a meaning. This text should show what a printed text will look like at this place. If you read this text, you will get no information. Really? Is there no information? Is there a difference between this text and some nonsense like “Huardest gefburn”? Kjift – not at all! A blind text like this gives you information about the selected font, how the letters are written and an impression of the look. This text should contain all letters of the alphabet and it should be written in of the original language. There is no need for special content, but the length of words should match the language.

After this fourth paragraph, we start a new paragraph sequence. Hello, here is some text without a meaning. This text should show what a printed text will look like at this place. If you read this text, you will get no information. Really? Is there no information? Is there a difference between this text and some nonsense like “Huardest gefburn”? Kjift – not at all! A blind text like this gives you information about the selected font, how the letters are written and an impression of the look. This text should contain all letters of the alphabet and it should be written in of the original language. There is no need for special content, but the length of words should match the language.

Hello, here is some text without a meaning. This text should show what a printed text will look like at this place. If you read this text, you will get no information. Really? Is there no information? Is there a difference between this text and some nonsense like “Huardest gefburn”? Kjift – not at all! A blind text like this gives you information about the selected font, how the letters are written and an impression of the look. This text should contain all letters of the alphabet and it should be written in of the original language. There is no need for special content, but the length of words should match the language.