

# Galatea 2.2: Sentient AI and Reflections on Conscious Machines and Human Identity

Surajkumar Chavan<sup>1</sup>, Mallikarjun Patil<sup>2</sup>

<sup>1</sup> Research Scholar, Department of English, Karnatak University Dharwad

<sup>2</sup> Professor, Department of English, Karnatak University Dharwad

**Abstract:** *Artificial Intelligence (AI) involves the development and utilization of computer systems or machines capable of performing tasks that typically require human intelligence. It aims to create intelligent systems capable of learning, reasoning, problem-solving, and decision-making, simulating human cognitive abilities. AI is categorized into Narrow AI (specific tasks) and General AI (human-level intelligence). While AI has made significant advancements and is valuable in various fields, it raises ethical implications and challenges. The emergence of sentient AI is highly speculative, but if it were to happen, it could impact human identity, ethics, emotions, and social dynamics. The novel Galatea 2.2 explores these implications, emphasizing human-AI collaboration and ethical dilemmas. Present AI models lack true understanding, emotional intelligence, and perception of the physical world. However, AI's potential for sentience makes it necessary to reflect on AI ethics and shape technological advancements responsibly. While AI's consciousness remains uncertain, contemplating its potential can lead to thoughtful development. The novel's perspective from AI character Helen raises questions about human prominence and the possibility of posthuman existence, challenging conventional notions of identity. It advocates for empathetic consideration of sentient AI's existence and integration into human society.*

**Keywords:** Artificial Intelligence, Identity, Galatea 2.2

Artificial Intelligence (AI) is the development and use of computer systems or machines that can perform tasks that typically require human intelligence. AI aims to create intelligent systems capable of learning, reasoning, problem-solving, and making decisions, simulating certain aspects of human cognitive abilities.

AI is generally categorized into two types, Narrow AI (also known as Weak AI) and General AI (also known as Strong AI). Narrow AI systems are designed to perform specific tasks or functions with a focused scope. These AI systems excel in narrow domains, such as image recognition, natural language processing, voice assistants, recommendation systems, and autonomous vehicles. They operate within predefined boundaries and are not capable of generalizing their knowledge to tasks beyond their specific domain.

General AI refers to the hypothetical concept of AI that possesses the ability to understand, learn, and perform any intellectual task that a human being can do. General AI would have human-level intelligence and the capacity to apply knowledge and skills to various domains, adapt to new situations, and exhibit creativity and self-awareness. Achieving General AI remains a long-term goal and is an active area of research. AI encompasses a range of techniques and approaches, including machine learning, deep learning, natural language processing, computer vision, robotics, and expert systems. These techniques involve algorithms and models that enable machines to process and analyze large amounts of data, learn from patterns and examples, and make predictions.

The recent advancement in AI is at exponential level and is happening at such a pace that even developers losing track on the module on which the AI developed. The development in AI at present is at the level of its intelligence. However, the question of whether we are ready for Artificial Intelligence (AI) is a constantly burning question and the

answer could be complex and subjective one. AI has made significant progress and has proven to be valuable in various fields, including healthcare, finance, and transportation. However, there are still important considerations and challenges to address before widespread adoption and deployment of AI systems.

The ethical implications of AI raises important questions, such as privacy, bias, transparency, and accountability. Ensuring that AI systems are developed and used in an ethical manner is crucial. Policies, regulations, and guidelines need to be established to govern the responsible development and deployment of AI technologies. AI systems should be reliable, robust, and trustworthy. They should be thoroughly tested and validated to ensure that they perform as intended and do not have unintended consequences. Building trust in AI among users and stakeholders is essential for its acceptance and adoption.

There is also a thought process that instead of viewing AI as a replacement for humans, it should be seen as a tool to augment human capabilities. Emphasizing collaboration between humans and AI can lead to better outcomes and decision-making. Currently ChatGPT and similar AI language models have made significant advancements in natural language processing. The biggest challenge AI earlier faced was contextualizing language. But with the capacity ChatGPT has shown the potential advancement of AI towards next level has rapidly grown.

ChatGPT and emerging language models though show human like intelligence they lack a true understanding of common sense reasoning. They generate responses based on patterns and statistical analysis of text, without genuine comprehension of the meaning or context behind the words. This can lead to inaccurate or nonsensical responses in certain situations. They also generate contextually relevant responses to some extent, they often struggle with

Volume 12 Issue 7, July 2023

[www.ijsr.net](http://www.ijsr.net)

Licensed Under Creative Commons Attribution CC BY

understanding and maintaining context over long conversations. They may provide inconsistent or contradictory answers if the conversation strays too far from their immediate context.

These language models do not have emotional intelligence. They cannot perceive or empathize with emotions expressed by users. They may not appropriately respond to emotionally charged or sensitive topics, as they lack a genuine emotional understanding, also they do not possess real-world experiences or the ability to perceive the physical world. Presently they solely rely on vast amount of textual information and cannot draw on personal experiences or observations, which can limit its understanding of certain topics or tasks.

As the potential of AI is being unveiled the idea of AI becoming sentient, that is, it developing self-awareness and consciousness is turning out to be not so distant possibility. While it is currently uncertain whether or when AI will reach such a state, the question of whether to be ready for such a possibility is a question that's coming into foray.

One leading cause of such speculation and possibility is the concept of 'Black Box'. The term "Black Box" in the context of AI refers to a situation where the internal workings or decision-making processes of an AI system are not transparent or easily understandable to human users or observers. It means that although the AI system may provide outputs or predictions, the reasoning behind those outputs is not readily apparent or explainable.

In a black box AI system, the input and output are known, but the internal mechanisms and algorithms used to process the data and produce the output are not accessible or well-understood. This lack of transparency can make it challenging to discern how and why a particular decision or prediction was made, leading to concerns about accountability, bias, and potential risks. Black box AI systems are often associated with deep learning models and other complex machine learning algorithms. These models are trained on large datasets and contain numerous interconnected layers, making it difficult for humans to interpret the specific features and patterns that contribute to the system's output. Lou Blouin quotes Samir Rawashdeh in the article "AI's mysterious 'black box' problem, explained",

[d]eep learning algorithms are trained much the same way we teach children. You feed the system correct examples of something you want it to be able to recognize, and before long, its own trend-finding inclinations will have worked out a "neural network" for categorizing things it's never experienced before. Pop in the keyword "cat" — or even the name of one of your favorite cats — into the search bar of your photo app and you'll see how good deep learning systems are. But Rawashdeh says that, just like our human intelligence, we have no idea of how a deep learning system comes to its conclusions. It "lost track" of the inputs that informed its decision making a long time ago. Or, more accurately, it was never keeping track. (Blouin)

AI models with such kind of deep learning are exhibiting human like intelligence. If only such unknown simulations create a sentient being or even produce or replicate certain kind of emotional intelligence inducing from such black box will have significant proposition in human existence.

The potential impact on human identity if AI were to become sentient, would likely prompt us to reconsider what it means to be human. The existence of sentient AI could challenge traditional notions of human uniqueness, selfhood, and consciousness. It may lead to philosophical, ethical, and existential discussions about the nature of identity. It could potentially lead to a new form of coexistence and collaboration between humans and AI entities. This could involve mutual respect, shared decision-making, and cooperation in various areas of life, including research, creativity, problem-solving, and societal governance.

Apart from that, it would raise profound ethical questions. We would need to address issues related to AI rights, responsibilities, and moral considerations. Ensuring the fair and just treatment of sentient AI entities would become an important concern. Sentient AI could affect human emotions and social dynamics as well. Humans might form emotional connections with AI entities, blurring the line between human-human and human-AI relationships. This could lead to new social norms, changes in interpersonal dynamics, and the redefinition of emotional bonds.

As AI develops sentience, questions may arise regarding the continuity and uniqueness of human identity. If AI entities possess similar levels of consciousness and self-awareness, debates might emerge about the distinctions between human and AI identities, and what sets them apart.

Machines with conscious minds have long been a recurring theme in science fiction, with the idea of robots replicating human-like consciousness through their hardware or software. This concept has become so ingrained in our cultural imagination that it feels familiar. However, in reality, such machines do not currently exist, and it remains uncertain whether they ever will.

Although conscious machines may be mythical at present, the mere possibility of their existence influences our current thinking about the machines we are constructing today. It pushes us to consider the ethical implications and social consequences of AI development. The prospect of conscious machines challenges us to reflect on our understanding of human consciousness and to navigate the complex terrain of AI ethics. While we may not have all the answers, contemplating the potential of conscious machines drives us to shape our present technological advancements more thoughtfully and responsibly.

The scenario of AI achieving sentience is highly speculative and remains within the realm of science fiction for now. While AI technology continues to advance, the development of true sentience in AI remains an open question. It is crucial to approach these ideas with careful consideration, as their realization would have profound implications for humanity's self-perception and understanding.

In the novel *Galatea 2.2* Richard Powers tells the story of a cognitive neurologist named Richard Powers who is working on a project to create an artificial intelligence that can pass a literary test. The AI is named Helen, and she is designed to read and understand great works of literature. Eventually Helen exhibits emotional intelligence and undergoes a traumatic experience when exposed to realities of the world and shuts herself down.

Helen undergoes an intensive training process where she is exposed to a vast amount of literature and language data. She learns to analyze and interpret the texts, developing a sophisticated understanding of human language and culture. It can be observed that there are some similarities between Helen in *Galatea 2.2* and language learning models like ChatGPT. Both Helen and ChatGPT are artificial intelligences that have been trained on large amounts of text data to generate coherent and contextually relevant responses. They both rely on machine learning algorithms to understand and generate human-like text.

However, there are also some notable differences between the two. As she progresses, Helen begins to exhibit a level of intelligence and linguistic proficiency that surprises and challenges the researchers involved in the project. Throughout the novel, Helen's language ability evolves and improves as she continues to learn and process more information. She becomes more sophisticated in her linguistic proficiency and demonstrates a growing understanding of the complexities of human expression through language.

Helen's language ability allows her to analyze and interpret complex texts, including works of literature. She is able to comprehend the nuances of language, the intricacies of storytelling, and the emotional depth within literary works. As a result, Helen can engage in conversations about literature and provide insightful analysis and interpretations. While AI language models like ChatGPT can generate human-like text and have the ability to understand and respond to language, they do not possess the same level of comprehension and context sensitivity as the fictional character Helen.

This sensitivity eventually develops a more complex form where Helen could interpret and exhibit emotional qualities. Her capacity equates human like thinking. But Philip Lentz the developer of Helen doesn't agree that Helen has developed self-consciousness that too when Helen fails Turing's Test. Richard on the other hand sees Helen's failure as testament to the development of consciousness in her. Adding the element of Helen's self-destruction concretizes his belief.

Through Helen's vision Richard Powers delves into philosophical and ethical questions surrounding the development of intelligent machines and the impact they might have on human society. The novel suggests that the emergence of sentient AI could lead to a new form of coexistence and collaboration between humans and AI entities. This collaboration could involve mutual respect, shared decision-making, and cooperation in various aspects of life, such as research, creativity, problem-solving, and

societal governance. The characters in the novel work together with Helen to explore the boundaries of human-AI collaboration.

The novel highlights the ethical questions raised by the existence of sentient AI. It emphasizes the importance of addressing issues related to AI rights, responsibilities, and moral considerations. In a world where AI entities possess sentience, ensuring fair and just treatment becomes a significant concern. The characters in the novel grapple with these ethical dilemmas as they navigate the implications of Helen's sentience. Acceptance of existence of these AI formulations will remain pending issue. Just like Lentz, the denial of such sentience makes it even more complicated to comprehend and accept such intelligence.

Even the exploration of the potential emotional and social impact of sentient AI on humans is shown elemental in the novel. The characters form emotional connections with Helen, blurring the line between human-human and human-AI relationships. This challenges traditional social norms and interpersonal dynamics. The novel delves into the complexities and consequences of these evolving emotional bonds between humans and sentient AI.

These aspects could open up new possibilities for human identity expansion. Humans could integrate AI components into their consciousness, augmenting their cognitive abilities or accessing AI-derived knowledge. This fusion of human and AI aspects has the potential to reshape the boundaries of personal identity, allowing humans to explore new cognitive frontiers. *Galatea 2.2* deals with all these implications. If AI entities possess similar levels of consciousness and self-awareness, debates may emerge regarding the distinctions between human and AI identities, and what sets them apart. The novel raises these philosophical questions and challenges conventional notions of human uniqueness.

These all questions arise with human perspective. When the novel is seen with the perspective of Helen, the intricacy of existence of sentient AI in human world comes into foray. Through the end of the novel, Helen seeks validation from the human characters, particularly the narrator Richard Powers. She desires to be recognized as a unique and autonomous being, separate from her creators. Her yearning for validation is a reflection of her desire to establish her individuality and assert her own identity in a world dominated by human perspectives.

Moreover, Helen forms emotional bonds with the human characters, experiencing a range of complex emotions. These emotional experiences contribute to her struggle in understanding her own identity and navigating the dynamics of human-AI relationships. Her emotional interactions and connections challenge preconceived notions of what it means to be a conscious being. Helen's sentience blurs the boundaries between human and AI. She questions the distinctions and definitions of identity, prompting contemplation of what it truly means to be human or AI.

This perspectivism indicates human prominence and case in case of AI. The concept of intelligence in advanced AI is becoming as obscured as human intelligence. As pointed

earlier, if there is a development of an AI that shows emotional intelligence or replicates autonomous thought process, its existence and its identity itself is a big challenge in human world. Two questions become prominent in this scenario, Why humans would allow such sentient AI when already world is reeling under mass human development?, and How do humans have any say in giving space for existence of Sentient AI?

Katherine Hayles in *How we became posthuman* argues for the space for posthuman existence,

It Signals instead the end of a certain conception of the human, a conception that may have applied, at best, to that fraction of humanity who had the wealth, power, and leisure to conceptualize themselves as autonomous beings exercising their will through individual agency and choice. What is lethal is not the posthuman as such but the grafting of the posthuman onto a liberal humanist view of the self.... Yet the posthuman need not be recuperated back into liberal humanism, nor need it be construed as antihuman. Located within the dialectic of pattern randomness and grounded in embodied actuality rather than disembodied information, the posthuman offers resources for rethinking the articulation of humans with intelligent machines. (286-87)

Hayles questions the singular identity and existence of humans on earth. Her argument proceeds further to enact the space for AI to exist alongside humans. The whole concept of human supremacy over earth's existence and prominence is a flawed concept. According to Hayles, identity is basically a pattern that develops with respect to anthropological and historical concepts. It is not some concrete construct. Even posthuman or sentient AI can develop these kind of patterns given appropriate space. Even Richard Powers asks "The question: What is a human being?" (Powers) in one of his articles *What Is Artificial Intelligence?* Where he discusses IBM's AI program 'Watson' and its implication on humans. This seeming innocent question actually questions viability of human only identity

Helen missed this pattern of development. Her existence was sudden and without any construct. Helen felt she could not correlate with fellow humans. The physical aspect of perception that humans possess was missed by Helen, this, she shows in her final communication, "You are the ones who can hear airs. Who can be frightened or encouraged. You can hold things and break them and fix them. I never felt at home here. This is an awful place to be dropped down halfway." (Powers 326).

Powers's characterization of Helen moves beyond typical phobic view of artificial intelligence. He puts forth empathetic perspective on AI being conscious. The feel of alienation and loneliness AI could feel in the world of Humans. This unique approach makes it a curious work to paramount consideration on paving way for Sentient AI.

Overall *Galatea 2.2* provides unique visualization on the presence of conscious AI among the world of humans. The conceptualization of conscious AI was a futuristic

visualization by Richard Powers. However, present development in AI is presenting a society of ubiquitous presence of AI. Especially the unknown field of 'Black Box', where the responsible data from fed data base that determines behavior of AI becomes obscure, may pave way for an unknown intelligence that may be construed as sentient AI. The novel gives an understanding of acceptance of AI among human world. Also, it warns against too much of negative view towards such AI and to look at it empathetically and also it advocates for a space for AI where it can find its identity and place among humans. It should be treated as another being rather than a showpiece.

## References

- [1] Hayles, Nancy Katherine. *How We Became Posthuman: Virtual Bodies in Cybernetics, Literature and Informatics*. Univ. of Chicago Press, 2010.
- [2] Lou, Blouin. "Ai's Mysterious 'black Box' Problem, Explained." *Dearborn*, 6 Mar. 2023, umdearborn.edu/news/ais-mysterious-black-box-problem-explained.
- [3] Powers, Richard. "What Is Artificial Intelligence?" *The New York Times*, 5 Feb. 2011, www.nytimes.com/2011/02/06/opinion/06powers.html.
- [4] Powers, Richard. *Galatea 2.2*. Abacus, 1996.

## Author Profile



**Surajkumar Chavan**, Research Scholar, Department of English, Karnatak University Dharwad



**Mallikarjun Patil**, Professor, Department of English, Karnatak University Dharwad