

# Non-Human Agents: Exploring Modern AI through the Lens of Contemporary Ontology and Theology

Hugo L. Rufiner, José G. Funes S.J., Mariano Asla, Leonardo Giovanini & Enrique Alberto Majul

## ABSTRACT

Accelerated Artificial Intelligence (AI) development necessitates re-evaluating intelligence and agency. This paper examines AI's evolution—focusing on large language models (LLMs), agentic AI, and robotics—exploring philosophical and theological implications. Engaging ontological models, from substantialism to relationalism, it interprets AI's nature and impact on human self-understanding. Theologically, it centers on *imago Dei*, comparing human and artificial intelligence regarding rationality, relationality, and embodiment. It addresses ethical dilemmas including moral agency and responsibility. Integrating philosophical and theological methodologies, the paper aims to understand AI's place, significance, and ethical considerations for responsible development, concluding with emphasis on interdisciplinary dialogue.

## KEYWORDS

Artificial Intelligence (AI), Large Language Models (LLMs), Agentic AI, *Imago Dei*, Ontology of AI, Ethics of Technology

## 1. Introduction

The exponential development of artificial intelligence (AI) has significantly impacted society, prompting deep reflections on the nature of intelligence and agency. Rooted in intellectual traditions dating back to the seventeenth century, AI has evolved from specialized data-analysis systems to ones capable of emulating human cognitive capabilities and exhibiting sophisticated anthropomorphic traits. Notably, in 2024, AI played a crucial role in scientific discoveries recognized with Nobel Prizes in Physics and Chemistry, demonstrating its increasing impact across a wide range of fields.<sup>1</sup> This exciting evolution poses a significant challenge to long-standing conceptions of cognition, learning, and consciousness, thereby giving rise to fundamental questions regarding the very essence of intelligence itself.

To contextualize this discussion, it is essential to briefly consider some definitions of intelligence. Traditionally, intelligence has been viewed as the capacity to learn, understand, and apply knowledge to manipulate one's environment or think abstractly, as measured by objective criteria.<sup>2</sup> Intelligence encompasses abilities such as reasoning, problem-solving, planning, abstract thinking, comprehension of complex ideas, and learning from experience.

From a biological perspective, human intelligence arises from the intricate complexity of the human body. The brain, composed of approximately 86 billion neurons interconnected by trillions of synapses<sup>3</sup>, works in concert with the nervous system, sensory organs, immune and endocrine systems to regulate bodily functions and respond to external stimuli. Neurotransmitters and hormones modulate mood, cognition, and behavior.<sup>4</sup> The brain's plasticity allows for adaptation and learning, reflecting the dynamic and multifaceted nature of biological intelligence. Moreover, physiological processes such as metabolism and

homeostasis contribute to cognitive functioning<sup>5</sup>. Thus, human intelligence is not confined to the brain, but it results from the integrated functioning of the entire organism in relation with its environment.

Gardner's theory of *multiple intelligences* expands upon traditional definitions by proposing that intelligence is not a single, unified capacity but a set of distinct modalities.<sup>6</sup> Gardner identifies at least eight intelligences, including linguistic, logical-mathematical, spatial, musical, bodily-kinesthetic, interpersonal, intrapersonal, and naturalistic intelligences. Other authors also propose the existence of a spiritual intelligence related to transcendence capacity, wisdom, and holistic comprehension.<sup>7</sup> Among all these, linguistic intelligence plays a central role in structuring human thought and communication, particularly enabling higher-order cognitive functions such as self-reflection, reasoning, and complex problem-solving. Language can be understood as both a vehicle and a scaffold for thought allowing the manipulation of abstract concepts that would otherwise be difficult to comprehend. However, not all cognition depends on language-basic sensory experiences, motor skills, and even some forms of problem-solving occur independently of linguistic representation. As mastery in language use is deeply intertwined with human thought and self-awareness, it is not a surprise that advanced language models generate high expectations.

The advent of new AI "*thinking*" models represents a development beyond *linguistic emulation* toward *artificial cognition*. These models not only generate coherent texts but also exhibit a kind of reasoning and metacognitive abilities, blurring the lines between human and artificial intelligence. This progression not only emulates but also defies the conventional understanding of intelligence, which is predominantly interpreted through the lens of intellect and reason.<sup>8</sup> Building upon these advancements in reasoning, the pursuit of Agentic AI is emerging as the next frontier, aiming to create AI systems capable of autonomous action towards specific goals. This kind of AI is considered a critical step towards Artificial General Intelligence (AGI). Moreover, the development of embodied AI and humanoid robots highlights the role of physical embodiment in achieving more holistic forms of cognition. Embodied AI — through systems like autonomous systems and humanoid robots — integrates sensory perception, physical interaction, and contextual awareness, offering a form of intelligence that more closely parallels human cognition. This progression requires a critical examination of AI's implications within ontological, ethical, and theological frameworks, particularly regarding embodied cognition and anthropomorphism.

The ontological and theological implications of AI development are especially compelling when considering the concept of being made in the "*image and likeness of God*" (*imago Dei*). According to the Judeo-Christian tradition, humanity was created as a reflection of the divine, endowed with capacities for rationality, morality, and relationality.<sup>9</sup> This concept invites a deeper exploration of the analogies between human and artificial —i.e. non-human agents. If humans, as reflections of the divine, create AI that increasingly resembles them, one must ask whether AI also bears a reflection of the divine. Can anthropomorphism in AI be connected to theological notions of resemblance to God, and if so, what are the philosophical consequences of such a connection?

This paper draws upon various theological models that explore the *imago Dei* and seeks to relate these models to the anthropomorphism inherent in AI systems. By engaging with these models through methodological approaches specific to theology and philosophy, we aim to

understand whether AI's resemblance to humanity extends beyond functional capabilities to encompass questions of purpose, essence, and theological significance.

## 2. The Evolution of AI Systems

### 2.1 Historical Foundations

AI's development traces back to early philosophical inquiries and technological advancements of humanity. The study of reasoning mechanisms dates back to Aristotle's *Organon*.<sup>10</sup> More recently, in his 1637 work *Discourse on the Method*, René Descartes envisaged the mechanistic understanding of nature and the possibility of machines simulating human behaviors. He argued that while machines might mimic numerous human behaviors, they would lack the capacity for language communication and reasoning, distinguishing humans from automatons.<sup>11</sup> For Descartes, this distinction laid the groundwork for considering the unique aspects of human intelligence.

In 1843, Ada Lovelace provided commentary on Charles Babbage's Analytical Engine, laying the foundations of algorithmic thinking. She envisioned machines capable of manipulating symbols and numbers, foreshadowing the development of programmable computers.<sup>12</sup> Lovelace's insights highlighted the potential for machines to go beyond simple calculation, touching on creativity and the processing of abstract concepts. In 1936, Alonso Church<sup>13</sup> and Alan Turing<sup>14</sup> independently formalized the concept of algorithm and computability. These theoretical developments establish the foundations for future algorithmic and computational research.

The mid-twentieth century saw significant strides in formalizing AI. In 1943, McCulloch and Pitts developed the first mathematical model of a neural network, introducing the idea that neural activities could be represented through logical operations.<sup>15</sup> This work bridged biology and computation, suggesting that cognitive processes could be simulated artificially.

Alan Turing's seminal paper *Computing Machinery and Intelligence* proposed the Turing Test, establishing a framework for evaluating machine intelligence based on the machine's ability to exhibit human-like behavior in a written interrogatory.<sup>16</sup> This test shifted the focus from the internal processes of machines to their observable outputs, emphasizing functional equivalence.

At the 1956 Dartmouth Conference, John McCarthy coined the term *artificial intelligence*, marking the official inception of AI as a field of study focused on creating machines capable of intelligent behavior.<sup>17</sup> This conference brought together leading thinkers to discuss the possibilities of machine intelligence, setting the stage for decades of research and development.

Simultaneously, *cybernetics* emerged as an interdisciplinary approach to understanding dynamical systems. It is concerned with general principles for modelling, designing, managing, and regulating the behavior of dynamical systems that are relevant across multiple contexts, including engineering, economic, biological, and social systems, among others. Coined by

Norbert Wiener in 1948, cybernetics focused on the study of regulatory systems, feedback loops, adaptation mechanisms, and the integration of hardware and software in a similar way to living organisms.<sup>18</sup>

Cybernetics proposed a vision of intelligence as inherently embodied and immersed within an environment. It emphasized the continuous interaction between a system and its surroundings, where sensory inputs produce actions that, in turn, affect the environment, creating a feedback loop essential for adaptation and evolution. Pioneers like Ashby and Walter built cyber-physical systems, such as the Homeostat and the Machina Speculatrix, that demonstrated adaptive behaviors through simple electronic circuits and feedback mechanisms.<sup>1920</sup> This approach contrasted with the symbolic and algorithmic focus of early AI, which often abstracted intelligence into computational processes divorced from physical embodiment. Cybernetics sought to replicate the embodied functioning of living organisms and dynamical systems, where cognition, perception, and action are interdependent.

In the following decades, AI and cybernetics followed divergent paths. AI research predominantly adopted a symbolic, algorithmic approach, focusing on high-level reasoning, problem-solving, and language processing using formal logic and representations. This movement, sometimes referred to as *Good Old-Fashioned AI* (GOFAI), aimed to model intelligence through abstract computational algorithms.<sup>21</sup> Conversely, cybernetics continued developing cyber-physical systems exploring the role of adaptation, learning, and feedback in systems behavior. It gradually lost prominence in mainstream AI research, which initially prioritized symbolic methods and in the last decades numerical and neuromorphic methods. The separation was further reinforced by limitations in hardware capabilities and the complexity of modelling biological systems accurately.

This divergence between algorithms and devices resulted in AI systems that excelled in specific domains but lacked the integrated sensory-motor capabilities of living organisms. This phenomenon is known as Moravec's paradox, which highlights the counterintuitive discovery that apparently high-level reasoning requires relatively little computation compared to low-level sensorimotor skills.<sup>22</sup> Cybernetics and AI developed along largely separate trajectories, with cybernetics focusing on the algorithmic, computational, and physical capabilities to engage machines/robots with humans and the environment, while AI concentrated on algorithms for information processing and computational intelligence.

## 2.3 Advancements in Neural Networks

*Artificial neural networks* (ANNs) are foundational to modern AI. Rosenblatt's introduction of the perceptron in 1958 was a significant milestone, representing an early neural network model capable of learning and classifying input data through supervised training.<sup>23</sup> The perceptron demonstrated that machines could adjust their parameters based on experience, a fundamental aspect of learning.

The backpropagation algorithm, introduced by David Rumelhart, Geoffrey Hinton, and Ronald Williams in 1986, enabled the effective training of multilayer networks by calculating error gradients and adjusting weights accordingly.<sup>24</sup> This breakthrough propelled the field of deep learning, allowing for the modeling of complex, non-linear relationships.

Modern architectures have been built upon these foundations. Convolutional Neural Networks (CNNs), pioneered by LeCun, are specialized for processing grid-like data structures such as images and have achieved remarkable success in computer vision tasks.<sup>25</sup> Recurrent Neural Networks (RNNs) and Long Short-Term Memory (LSTMs) networks handle sequential data, capturing temporal dependencies and improving performance in speech recognition and language modeling.<sup>26</sup>

The collective advancements in neural network architectures and training methodologies, particularly those pioneered by researchers like Hinton, Bengio, and LeCun, were crucial in paving the way for the emergence and widespread adoption of the so-called *deep learning*. Technological enhancements, such as the rise of Graphics Processing Units (GPUs) and Tensor Processing Units (TPUs), along with access to vast datasets from the internet, have enabled the training of such deep neural networks with billions of parameters.

Generative Adversarial Networks (GANs), introduced by Goodfellow in 2014, consist of two networks—a generator and a discriminator—that train together, enabling the generation of realistic synthetic data across various domains.<sup>27</sup> The transformer architecture, presented by Vaswani et al. in 2017, utilizes attention mechanisms to process sequences in parallel, dramatically improving efficiency and performance in natural language processing.<sup>28</sup> These advancements have culminated in sophisticated AI systems capable of performing complex tasks requiring high levels of pattern recognition, understanding, and generation.

## 2.4 Emergence of Large Language Models

LLMs represent a significant leap in AI's ability to process and generate human language. The development of models like OpenAI's<sup>29</sup> GPT-3 and GPT-4 or Anthropic's Claude<sup>30</sup> has been propelled by advancements in computational capacity, algorithmic innovations, and the availability of massive textual datasets.<sup>31</sup>

The transformer architecture allows models to consider the context of each word in a sentence relative to all other words, capturing long-range dependencies and nuances in language. Training on diverse and extensive textual data enables these models to learn a wide range of linguistic patterns, factual knowledge, and cultural references. With hundreds of billions of parameters, these models can generate coherent and contextually appropriate text, perform translation, answer questions, and create code, often surpassing human-level performance in specific tasks.

In-context learning enables LLMs to adapt to new tasks with minimal examples, demonstrating a form of few-shot learning that mimics aspects of human learning. However, the ability to generate human-like language does not equate to understanding or consciousness. LLMs operate based on statistical patterns learned during training, lacking awareness or intentionality.

## 2.4 Towards Multimodal Language Models

Advanced text-to-speech systems generate natural-sounding voices with emotional expressiveness, enhancing user engagement.<sup>32</sup> Visual avatars and humanoid robots like Ameca provide visual cues and expressions, fostering a sense of presence and immediacy.<sup>33</sup> The introduction of synthetic voices with emotional expression capabilities and the incorporation of distinctive "personalities" in systems, such as Anthropic Claude Sonnet 3.5, represents deliberate efforts to render interactions with these systems more "human-like."

The recent introduction of advanced voice capabilities in ChatGPT marks a significant leap in the development of anthropomorphic AI interfaces.<sup>34</sup> This advanced voice mode, powered by a new text-to-speech model and enhanced by OpenAI's Whisper speech recognition system, enables ChatGPT to generate highly realistic, human-like audio interactions. The model utilizes prosody adjustments to convey emotions, recognizes verbal cues such as speaking speed, and supports a range of accents, making interactions feel more natural and engaging. Users can also select from multiple distinct voices, each with its own tone and character, to personalize their experience.

These anthropomorphic features mimic human interpersonal or emotional intelligence. They can enhance user experience but also complicate the ethical landscape, as they blur distinctions between genuine personal interaction and the simulation of human-like qualities. Features like real-time interruption support and dynamic responses to non-verbal cues enhance the fluidity of conversations, yet raise questions about the ethical implications of such realism. The ethical implications of anthropomorphism warrant careful consideration, particularly in the context of trust, user expectations, and the evolving nature of human-machine relationships. Advanced voice interactions may risk creating false perceptions of empathy or genuine understanding, complicating the boundary between AI assistance and human connection.

AI systems now adapt to user preferences, employing language styles and personalities that resonate with individuals. Conversational dynamics, such as interruptibility, turn-taking, and acknowledgment tokens, mimic human conversational patterns, making interactions more intuitive and familiar. While these developments aim to make AI interactions more accessible, they raise ethical concerns about deception, emotional attachment, and the potential for users to misunderstand the nature of AI.<sup>35 36</sup>

## 2.5 Emergence of Reasoning and Agentic Capabilities

Modern systems exhibit behaviors that resemble reasoning and even metacognition, challenging our understanding of AI. LLMs can generate step-by-step solutions through chain-of-thought reasoning, breaking down complex problems into sequential steps.<sup>37</sup> They can detect inconsistencies in their outputs and revise their responses upon further prompting, demonstrating a form of self-correction.

The emergence of models like OpenAI's O(1)<sup>38</sup> (internally code-named *Strawberry*) exemplifies these advancements. O(1) integrates reasoning processes within the language model framework by incorporating *Chains of Thought* during training. Labeled examples guide the model to produce detailed reasoning steps, enhancing problem-solving capabilities. *Reinforcement Learning with Human Feedback* (RLHF) fine-tunes the model using feedback

from human evaluators, aligning outputs with desired behaviors.<sup>39</sup> During the inference process, O(1) generates multiple chains of reasoning in the form of decision trees, using techniques such as *Monte Carlo Tree Search* (MCTS) to evaluate possible reasoning paths internally. This approach enables the model to explore different reasoning routes, assess the coherence and accuracy of each, and select the optimal chain to generate a final, well-founded response. By implementing real-time dynamic evaluation of reasoning chains, O(1) significantly enhances the accuracy and depth of responses provided to users.

The name "O(1)" implies that it is the first iteration of a new approach by OpenAI, with future iterations (e.g., O(3), O(4)) expected to further refine these capabilities based on usage data and expert feedback. In fact, O(3) has already been announced and publicly released in different limited versions. This evolution marks a clear departure from previous versions, such as ChatGPT-4o. While ChatGPT operated primarily as a cognitive System 1 —intuitive, approximate, and fast, but prone to errors—O(1) functions more like a System 2, emphasizing deliberate and rational reasoning.<sup>40</sup> This reflects a move towards integrating both *intellectus* (intuitive understanding) and "ratio" (discursive reasoning) in AI systems.<sup>41</sup> This shift allows O(1) to deliver more elaborate responses when given more time to *think*, effectively blending language modeling with advanced reasoning techniques.

In terms of problem-solving performance, OpenAI claims that O(1) has surpassed its predecessor, ChatGPT-4o, evolving from a level comparable to an incoming college student to that of a doctoral candidate in standardized tests for mathematics and programming. However, this leap in capability comes with significant computational and financial costs. Other recent models like DeepSeek R1 also demonstrate advanced reasoning and problem-solving capabilities but with much less resources involved in both training and test stages.<sup>42</sup> An important characteristic of this last AI model is that it is open weights<sup>43</sup> and open source, free to use and modify under a very flexible license, which contrasts with the closed model approach of companies like OpenAI.

Metacognition refers to the capacity to reflect upon and regulate one's own cognitive processes.<sup>44</sup> While these capabilities are not equivalent to human metacognition, they represent a significant step towards creating AI systems capable of introspection-like behaviors.<sup>45</sup> While models like O(1) exhibit advanced reasoning, they lack genuine self-awareness and the ability to consciously monitor or regulate their cognitive processes. Therefore, equating these advancements with human metacognition remains inaccurate.<sup>46</sup>

These advancements are paving the way for more autonomous and goal-directed AI systems, often referred to as *Agentic AI*. This new stage represents the next evolution beyond current models, focusing on creating AI systems that can act autonomously to achieve specific goals. An *agent* in AI is typically defined as an entity that perceives its environment through sensors and acts upon that environment through actuators to achieve specific objectives.<sup>47</sup> Agentic AI aims to imbue these systems with greater autonomy, enabling them to plan, execute, and adapt their actions in complex and dynamic environments without continuous human intervention. While current LLMs excel at generating text and reasoning within given contexts, Agentic AI seeks to create systems that can proactively pursue goals, learn from experience, and interact with the world in a more self-directed manner. This is considered a critical step towards Artificial General Intelligence (AGI), as it addresses the need for AI to not only understand and process information but also to act purposefully in the world.

As we have said, while these capabilities mark significant progress, they still do not imply genuine understanding or subjective consciousness. And even acknowledging that consciousness is a highly debated issue, a characterization that could garner some consensus in philosophy of mind might include: subjective experience of awareness and the ability to perceive, think, and feel. It encompasses both phenomenal experience—what it is like to be in a particular mental state—and access to consciousness, which allows information to be available for reasoning and action. The notion of self emerges as a unifying substrate of consciousness. AI reasoning and agentic capabilities exhibited are based on pattern recognition and statistical associations rather than deliberate, conscious thought.

## 2.6 The Embodiment of AI: Advanced Humanoid Robots

In recent years, there has been a resurgence of interest in integrating embodiment with AI, reflecting a convergence of the algorithmic and cybernetic traditions. Advances in computational power, sensor technology, and machine learning have enabled the development of robots and AI systems that can perceive, learn, and act within physical environments more effectively.

The integration of AI algorithms with sensory inputs and actuators has led to significant advancements in fields like autonomous vehicles, humanoid robotics, and adaptive control systems. Technologies such as deep reinforcement learning enable agents to learn from interactions with their environment, mirroring the feedback loops central to cybernetics.<sup>48</sup>

Humanoid robotics has redefined the integration of AI into physical systems, especially in recent years (between 2022 and 2025<sup>49</sup>). Pioneering companies such as Figure AI<sup>50</sup>, Boston Dynamics<sup>51</sup>, Tesla<sup>52</sup>, Sanctuary AI<sup>53</sup>, and Agility Robotics<sup>54</sup> have pushed the boundaries of what robots can achieve by equipping them with advanced sensorimotor systems, autonomous learning capabilities, and natural language interfaces.

Embodied AI emphasizes that cognition arises from the dynamic interaction between an agent and its environment, aligning with theories of embodied cognition in psychology and philosophy.<sup>55 56</sup> Researchers like Brooks advocated for behavior-based robotics, rejecting the centralized symbolic processing model in favor of decentralized systems that interact directly with the world.<sup>57</sup>

Integrating robotics with LLMs allows for interactive robots that can engage in conversation, understand commands, and perform complex tasks. Embodied AI systems perceive and navigate their surroundings through sensors and cameras, make decisions, and react to changes. They manipulate objects using physical actuators, performing tasks ranging from assembly to caregiving. These innovations not only promise to transform industrial and domestic applications but also raise deep philosophical and theological questions regarding the nature of intelligence and the boundaries of human agency.

From a technical perspective, the integration of LLMs with robotic systems introduces considerable difficulties. Sensorimotor integration requires converting continuous, noisy sensor data into discrete language inputs that LLMs can process, and vice versa. Additionally, achieving real-time decision-making and maintaining proprioception and spatial awareness



are crucial for safe and efficient operation. Researchers are exploring hybrid architectures where LLMs provide high-level planning while dedicated control modules manage fast, low-level motor actions. Furthermore, multimodal processing remains an active area of research. Advanced LLMs are beginning to incorporate visual and auditory inputs; however, aligning these modalities with the kinesthetic feedback necessary for physical tasks is a challenge that requires further innovation. The ongoing development in this area aims to create a cohesive system where language, vision, and action are seamlessly integrated.<sup>58 59</sup>

The embodiment of AI in humanoid form challenges traditional philosophical distinctions between mind and body. Theories of embodied cognition suggest that genuine intelligence arises from direct, sensorimotor interactions with the environment. By affording robots a body, researchers argue that AI systems can achieve a form of situated cognition that is inherently more adaptive and context-aware than disembodied systems. This embodiment, however, also opens critical philosophical and theological questions, especially when considering the concept of *imago Dei*. While these robots exhibit increasingly human-like behaviors, they currently lack essential attributes such as consciousness, moral agency, and the capacity for genuine relational interaction, aspects that will be further explored in the following sections. Consequently, while they may mimic aspects of human intelligence, they do not fully replicate the ontological and spiritual dimensions that define humanity.<sup>60</sup>

## 2.7 The Unprecedented Pace of AI Evolution

The development of AI technology is proceeding at an unprecedented rate, exceeding previous expectations in terms of both the speed and breadth of its advancement. In contrast to previous technological transitions, which occurred over extended periods, the advent of AI is characterized by rapid and substantial transformations, occurring within remarkably brief cycles. This accelerated progress, driven by exponential increases in computational power and continuous algorithmic innovation, has reached a pivotal historical juncture. It necessitates a thorough re-evaluation of the prevailing technological, ethical, and philosophical paradigms.

For decades, technological acceleration has aligned with Moore's Law, which initially described the doubling of transistors on a microchip approximately every two years, resulting in exponential increases in computational power and a proportional decrease in costs.<sup>61</sup> But this accelerated evolution is largely propelled by trends extending the principles of Moore's Law into the era of Deep Learning. The exponential augmentation of performance for AI is driven by three key factors: the development of specialized hardware, the increasing accessibility of vast datasets, and the constant algorithmic innovation. Deep learning itself has been transformative, shifting AI development from incremental improvements to a phase of rapid, self-amplifying advancement.

The rapid advancements in AI capacities, encompassing language comprehension, complex reasoning, agentic autonomy, and physical embodiment, give rise to profound ontological and theological inquiries. The emerging potential to create sophisticated, agentic AI, which may eventually lead to the development of AGI, challenges long-held, human-centric concepts of agency and personhood. This prompts an imperative and thorough interdisciplinary reflection on the fundamental nature of life, consciousness, and the consequences of non-human agents

within our world, as we navigate the uncharted territories opened by this technological revolution.

### 3. Philosophical Perspectives

#### 3.1 Ontological Models

In engaging with the reality of artificial intelligence, it is essential to first grasp the nature of what we are interacting with, particularly its ontological and moral status. The ethical implications that arise are often rooted in underlying, sometimes unconscious, ontological preconceptions. To better understand these implications, we must turn to philosophical models that seek to interpret the nature of AI. So far, these models oscillate between emphasizing relationality and substantiality, which is natural since humans are both substantial and relational beings. Therefore, it is not absurd that AI reflects this duality. Substantialist accounts regard AI as a tool or instrument, while relational ones erode the boundaries of agency. However, in light of its complexity and the novelty of the situation, it may be necessary to develop in the future new ontological frameworks that can more effectively address the unique complexities of AI.

Classical philosophers first contributed to these discussions. Aristotle's concept of *techne* (craftsmanship or art) and his exploration of the relationship between humans and their creations provide a foundation for considering the nature of artificial entities. He suggests that artifacts, though not natural substances, share a quasi-substantial status in the sense that they exist in a way that is dependent on human intention and purpose, as they are the products of *techne*, aimed at fulfilling specific goals and bringing about particular ends.<sup>62 63</sup> Immanuel Kant's emphasis on autonomy and rationality as defining features of personhood challenges us to consider whether AI, lacking true autonomy<sup>64</sup>, can be regarded as moral agents.<sup>65</sup>

The Extended Mind Hypothesis, proposed by Andy Clark and David Chalmers, suggests that tools and technologies can become extensions of the human mind.<sup>66</sup> According to this model, cognitive processes are not confined to the brain but extend into the environment through interactions with external devices. AI systems, particularly those integrated into daily life, may function as cognitive prosthetics, augmenting memory, perception, and reasoning. This pushes traditional boundaries of cognition and raises questions about the nature of the self.

Alfredo Marcos's model of delegated control and decision-making presents AI as a tool to which humans delegate specific tasks and decision-making within particular domains.<sup>67</sup> This model emphasizes that AI should not be viewed in isolation, as a detached entity. Instead, it should be understood as part of a broader interactive process, where its functioning is shaped by human engagement. Given that life and thought are inherently interactive, AI is best understood as a product of this dynamic relationship. The delegation of control to AI requires recognizing its role within this process, acknowledging the continuous interplay between human decisions and AI actions.

Luciano Floridi<sup>68</sup> presents a new and radical ontological framework in which information, rather than matter or energy, is the foundational constituent of reality. In this model, any entity,

whether natural or artificial, can be understood as information. This ontology is predicated on questioning traditional substantialist approaches by emphasizing interconnectedness and relationality, in a manner that resonates with Hegelian or even Buddhist perspectives. The focus on individual moral responsibility is shifted to complex distributed systems, where both human and non-human agents interact, blurring the boundaries between them.

Mark Coeckelbergh<sup>69</sup> and Shannon Vallor<sup>70</sup> also propose relational models, but they focus on the way these interactions shape human identities and social structures. While Coeckelbergh emphasizes the role of AI in shaping human experiences and the ethical challenges it introduces, Vallor's perspective centers on the moral virtues and responsibilities humans must cultivate as they interact with AI technologies. Both hold skeptical positions regarding AI's true agency.

### 3.2 Anthropomorphism in AI

The tendency to attribute human-like qualities, such as emotions, intentions, and personalities to non-human entities plays a pivotal role in shaping human interactions with AI and in the development of AI discipline itself. While anthropomorphism enhances user engagement by providing a relatable framework for human-AI interaction, it also risks distorting user perceptions, leading to unrealistic expectations and ethical dilemmas.

Anthropomorphism is often considered a sign of childish naivety or an epistemic vice characteristic of undeveloped societies, persisting as an unconscious and undesirable bias. We hold that this oversimplified stance fails to grasp the real dimension of the problem, as it is a natural constituent of human psychology and knowledge.

Anthropomorphism is, at once, a condition of possibility for human knowledge and an epistemic risk. It is a condition of possibility because, whether consciously or not, we understand any type of behavior in relation to our own categories—that is, in relation to what we are. We cannot know the “other,” the foreign, on its own terms but only by comparing it to ourselves, to what is most familiar to us. This occurs in the realm of theology as well as in ethology. It is as though we cannot help but project human categories onto the mental states of God or animals. However, this inevitability can become an epistemic risk (and an actual error) when such projection proves unwarranted.

This tendency to ascribe human-like qualities such as emotions, intentions, and personalities to non-human entities also plays a pivotal role in shaping interactions with AI.<sup>71</sup> The effect is particularly pronounced in multimodal interfaces, where visual and auditory cues simulate human traits. Such design choices enhance user engagement by creating an illusion of familiarity and relationality. However, they also introduce complexities in understanding AI's true nature. People may overestimate AI's cognitive abilities, assuming a level of comprehension or intentionality that it does not possess. The result is a paradox: the very feature that facilitates intuitive interaction also fosters misconceptions about the technology's capacities and limits.

As much as it enables interaction, anthropomorphism can also lead to hype and fallacy.<sup>72</sup> By making AI appear more human-like, it fosters intuitive engagement and increases user trust,

yet it simultaneously risks distorting perceptions, leading to unrealistic expectations and ethical dilemmas. Users may wrongly attribute agency or moral responsibility to AI, blurring the line between human intention and machine operation. To mitigate these risks, AI design must balance anthropomorphic features with clear indications of the system's limitations, ensuring transparency about its non-human nature while maintaining usability and accessibility.

#### 4. Theological Reflections

As Floridi<sup>73</sup> points out: "Digital technologies are not just tools that are limited to changing the way we interact with the world... They are above all systems that shape (format) and increasingly influence the way we understand the world and relate to it, as well as the way we conceive of ourselves and interact with each other."<sup>74</sup> In this section of our paper, we explore two questions: (i) Amidst the rapidly evolving landscape of AI and robotics discussed in previous sections, how does this development influence our perception of the image of God and our own self-understanding? (ii) In what ways can a theological perspective shed light on the numerous ethical dilemmas posed by these advanced technologies? These two topics clearly exceed the limits of our interdisciplinary paper. Consequently we will discuss them briefly.

To address the first query, we consider the concept of *Imago Dei* as a path or a model, not the only one, to consider the relationship between human beings and God. For the sake of this paper, we follow Jürgen Moltmann<sup>75</sup> and the recent Vatican document *Antiqua Et Nova*<sup>76</sup> on the relationship between AI and human intelligence.

##### 4.1 The Concept of *Imago Dei*

*Imago Dei* is the basic concept of theological anthropology: humans are created to be the image of God. It only appears in Gen 1, 26.27; 5,1; 9, 6. Other biblical texts presuppose this notion. The image of God designates first of all the relationship of God with men and women, and later the relationship of humankind with God.

In a Christian perspective, the true image of God is not at the beginning, but at the goal of God's history with humanity, the Omega Point in Teilhard de Chardin's words. Christ is the "image of God",<sup>77 78</sup> and believers have been predestined by God to reproduce the image of his Son.<sup>79</sup> Thus the *imago Christi* is an *imago Dei* mediated by Christ.

Theological tradition has always understood the image of God as a reflection in a kind of mirror. According to Pauline traditions in the New Testament, there is a place where the relationship between God and humankind is revealed and can be known, as the face of humans becomes a mirror of God: "All of us, gazing with unveiled face on the glory of the Lord, are being transformed into the same image from glory to glory."<sup>80</sup> "At present we see indistinctly, as in a mirror, but then face to face. At present I know partially; then I shall know fully, as I am fully known".<sup>81</sup> Related to the mirror representation, Vallor notes that it would "help to understand that today's most advanced AI systems are constructed as immense mirrors of human intelligence. They do not think for themselves; instead, they generate complex reflections cast by our recorded thoughts, judgments, desires, needs, perceptions, expectations, and imaginings"<sup>82</sup>. In this sense this author warns that "like Narcissus, we readily

misperceive in this reflection the seduction of an *other*—a tireless companion, a perfect future lover, an ideal friend, an unbiased judge, a foolproof collaborator—yet in truth, a thing with which we are increasingly left alone, talking to ourselves”<sup>83</sup>.

As the image of God, humankind represents God on earth, as the likeness of God, human beings reflect God. According to the analogy of relationship, likeness consists in the communion of man and woman, which corresponds to the intra-Trinitarian communion of God.

In this sense, “human intelligence is not an isolated faculty but is exercised in relationships, finding its fullest expression in dialogue, collaboration, and solidarity. We learn with others, and we learn through others”.<sup>84</sup> It is not possible to live likeness to God in a solitary way.

Also, “human intelligence is ultimately ‘God’s gift fashioned for the assimilation of truth.’ In the dual sense of intellectus-ratio, it enables the person to explore realities that surpass mere sensory experience or utility, since the desire for truth is part of human nature itself”.<sup>85</sup> “This innate drive toward the pursuit of truth is especially evident in the distinctly human capacities for semantic understanding and creativity”.<sup>86</sup>

Also, the Vatican document states that ‘The Christian tradition regards the gift of intelligence as an essential aspect of how humans are created “in the image of God’.<sup>87</sup> The Church emphasizes that this gift of intelligence should be expressed through the responsible use of reason and technical abilities in the stewardship of the created world” (AN 1). *Antiqua et Nova* also points out the imperative to understand AI within God’s plan.<sup>88</sup>

The emergence of AI systems capable of advanced reasoning, problem-solving, and agentic-like attributes invites a theological reflection on whether these attributes, when exhibited by AI, could also reflect some aspect of the divine. However, the uniqueness of human consciousness and moral awareness—qualities that AI, despite its capabilities, lacks—remains a crucial distinction.

## 4.2 AI and the Divine Image

AI’s rationality and creativity, demonstrated through problem-solving and generative capabilities, echo human intellectual faculties. However, without consciousness or self-awareness, AI does not possess these attributes inherently but simulates them through programmed algorithms. Theologically, the possession of a soul and the capacity for a relationship with God are integral to the *imago Dei*, which AI lacks.

AI’s capacity for relational interactions, such as engaging in conversations and responding to emotions, raises questions about the authenticity and theological significance of these relationships. While AI can simulate relational behaviors, it lacks genuine empathy and the ability to form meaningful connections rooted in shared experiences and mutual understanding. The relational model of the *imago Dei* emphasizes genuine relationships that involve mutual recognition and love, which AI cannot fulfill.

In terms of functional participation, AI may assist humans in fulfilling stewardship roles by performing tasks that manage or care for creation. This extension of human capabilities raises questions about the role of AI in fulfilling human responsibilities and the ethical implications of

delegating such tasks to machines. Pope Francis, in his address to the G7 participants in the session on AI, recalled: "Faced with the wonders of machines, which seem to know how to choose independently, we must be very clear that it is always up to human beings to make the decision, even with the dramatic and urgent tones with which this sometimes occurs in our lives. We would condemn humanity to a hopeless future if we took away from people the ability to decide for themselves and for their lives, condemning them to depend on the choices of machines."<sup>89</sup>

One area in which it is very worrying to delegate our decisions is the development and use of autonomous weapons. We should be concerned by the fact that much of the contemporary research in robotics and AI has been driven by military goals.

### 4.3. Embodiment

As the Vatican document AN states, "Christian thought considers the intellectual faculties of the human person within the framework of an integral anthropology that views the human being as essentially embodied... the soul is not merely the immaterial part of the person contained within the body, nor is the body an outer shell housing an intangible core. Rather, the entire human person is simultaneously both material and spiritual. This understanding reflects the teaching of Sacred Scripture, which views the human person as a being who lives out relationships with God and others (and thus, an authentically spiritual dimension) within and through this embodied existence. The profound meaning of this condition is further illuminated by the mystery of the Incarnation, through which God himself took on our flesh and "raised it up to a sublime dignity".<sup>90</sup>

According to John Paul II, the physicality of the body enables humans to enter into relationships, experience love, and participate in the creative work of God. The body is a "sacrament" of the person, making visible the invisible reality of the spiritual soul.<sup>91</sup> Therefore, human embodiment is essential to understanding human nature and destiny. This perspective highlights the sanctity of the human body and its role in expressing the *imago Dei*. Although recent technological advances have enabled humanoid robots with a body similar to the human body integrated with AI systems, we cannot currently attribute dignity similar to human dignity.

### 4.4 Ethical Implications

The creation of AI involves ethical considerations rooted in theological reflections on creativity and responsibility. As beings created in the divine image, humans exercise creativity in developing AI, bearing responsibility for the purpose and impact of their creations.<sup>92</sup> The ethical implications include ensuring that AI serves the common good, respects human dignity, and aligns with moral principles.

Delegating tasks to AI entails accountability for its actions, aligning with the concept of *imago Dei*. Humans remain responsible for the outcomes of AI's actions, particularly in areas affecting human well-being, justice, and peace. Pope Francis has addressed these challenges in the message for the World Day of Peace on Artificial Intelligence and Peace<sup>93</sup> and in the speech in the G7 session on Artificial Intelligence.<sup>94</sup>

One point that we would like to mention though we will not discuss it in detail is that thinking of AI as an oracle today could be a useful metaphor for reflecting on its role as a tool for prediction, consultation and decision-making in multiple areas giving to the AI a divine power that it clearly does not possess. However, it also implies considering its limitations, risks, and the relationship between those who consult AI and the answers it offers.

Assigning moral agency to AI is problematic, as AI lacks consciousness, free will, and moral understanding. Without the capacity for intentionality and ethical reasoning, AI cannot bear moral responsibility in the same way humans do. Theological ethics emphasizes the importance of intention and conscience in moral actions.<sup>95</sup>

The anthropomorphism of AI prompts us to consider the implications of creating machines that mirror human attributes. While AI may simulate aspects of human behavior, it remains fundamentally different in essence. The theological caution against equating artificial simulations with the divine image underscores the importance of recognizing and respecting these distinctions.<sup>96</sup>

The Christian tradition also offers ethical guidance on the appropriate relationship between humans and technology. By affirming the sanctity and dignity of the human body, it cautions against reducing human beings to mechanistic or utilitarian terms. The use of technology must respect the integral value of the person and promote authentic human development.

In the context of AI, this perspective urges caution in attributing human qualities to machines and emphasizes the need to maintain clear distinctions between human and artificial entities. Ethical considerations should prioritize human well-being, relationality, and the promotion of genuine human flourishing.

## **5. Final Considerations**

### **5.1 Exponential Growth Challenges**

The rapid acceleration of AI development presents significant challenges for regulation, governance, and societal adaptation. Laws and policies must evolve to address AI's capabilities, ensuring safety, privacy, and ethical use. This requires an agile approach that can keep pace with technological advancements and anticipate potential risks.

Global collaboration is essential to manage AI's cross-border implications, fostering shared standards and cooperative oversight. Societal adaptation involves preparing the workforce for AI integration, emphasizing education and re-skilling to address changes in employment and the economy. Social dynamics are also affected, as AI influences communication, relationships, and community structures.

Ethical frameworks must guide AI development and deployment, prioritizing human dignity and the common good. Developing ethical guidelines involves engaging diverse stakeholders, including technologists, ethicists, theologians, and the public, to ensure that AI serves societal needs without compromising values.

## 5.2 Methodological Reflections

An interdisciplinary dialogue between theology, philosophy, and science enriches our understanding of AI's implications. Theology employs methodologies such as scriptural interpretation, historical analysis, and doctrinal development to explore questions of meaning, purpose, and moral responsibility. Philosophy raises ontological issues and their epistemic and ethical implications. Science utilizes empirical investigation, experimentation, and theoretical modeling to understand natural phenomena and technological possibilities.

Integrating these methodologies allows for a comprehensive analysis that addresses both empirical and existential questions. Ethical methodology combines normative ethical theories with practical considerations, guiding responsible decision-making in AI development.

Reflecting on the methodologies used in theology and science highlights the importance of critical engagement and humility. Recognizing the limits of human knowledge and the complexity of AI encourages a collaborative approach that values multiple perspectives.

## 6. Conclusions

The intersection of sophisticated AI technologies and philosophical and theological inquiry offers a distinctive occasion to re-examine fundamental questions concerning intelligence, agency, and the human condition. The increasing sophistication of AI necessitates a re-evaluation that extends beyond mere technical aspects, encompassing the ethical, metaphysical, and spiritual dimensions inherent in the creation of non-human agents.

By engaging with concepts like the *imago Dei* and various ontological models through the specific methodologies of theology and philosophy, we gain insights into AI's place within the broader context of human experience. While AI may simulate certain human attributes, it lacks the inherent qualities that define humanity in theological terms—consciousness, free will, and a spiritual relationship with the divine.

Ensuring that AI development aligns with human values requires ongoing interdisciplinary dialogue. As creators, we bear responsibility for the impact of AI on individuals and society. The ethical deployment of AI must prioritize human dignity, promote the common good, and respect the unique aspects of human identity.

In contemplating AI's potential reflection of the divine, we are reminded of the profound responsibility inherent in our creative capacities. The challenge lies in harnessing AI's benefits while maintaining a clear understanding of the distinctions between artificial agents and the essence of humanity. This balance is essential to navigate the complex landscape at the intersection of technology, philosophy, and spirituality.

## Disclosure Statement

No potential conflict of interest was reported by the authors.



## Notes on contributors

*Hugo Leonardo Rufiner*, Biomedical Engineer, CONICET (National Scientific and Technical Research Council)—Universidad Nacional del Litoral and Universidad Nacional de Entre Ríos, Argentina.

*José G. Funes, S.J.*, Astronomer, CONICET—Universidad Católica de Córdoba, Argentina.

*Mariano Asla*, Philosopher, Universidad Austral, Argentina.

*Leonardo Luis Giovanini*, Electronics Engineer, CONICET—Universidad Nacional del Litoral, Argentina.

*Enrique Alberto Majul*, Physician, Universidad Católica de Córdoba, Argentina.

<sup>1</sup> <https://www.nobelprize.org/all-nobel-prizes-2024/>

<sup>2</sup> Ulric Neisser, Gwyneth Boodoo, Thomas J Bouchard Jr, A Wade Boykin, Nathan Brody, Stephen J Ceci, Dianne F Halpern, Jhon C Loehlin, Robert Perloff, Robert Sternberg and Susan Urbina, "Intelligence: Knowns and Unknowns," *American Psychologist*, 51(2) (1996):77–101.

<sup>3</sup> [Frederico A. Azevedo](#), [Ludmila R. Carvalho](#), [Lea T. Grinberg](#), [José Marcelo Farfel](#), [Renata E. Ferretti](#), [Renata E. Leite](#), [Wilson Jacob Filho](#), [Roberto Lent](#), [Suzana Herculano-Houzel](#), "Equal Numbers of Neuronal and Nonneuronal Cells Make the Human Brain an Isometrically Scaled-up Primate Brain," *Journal of Comparative Neurology*, 513(5) (2009):532–541.

<sup>4</sup> D Purves, G. J Augustine, D. Fitzpatrick, W. Hall, A. LaMantia and L White, *Neuroscience* (6th ed.), (2018) Oxford University Press.

<sup>5</sup> M. Carabotti, A. Scirocco, M. Maselli and C Severi, "The Gut-brain Axis: Interactions between Enteric Microbiota, Central and Enteric Nervous Systems," *Annals of Gastroenterology*, 28(2) (2015): 203–209.

<sup>6</sup> H. Gardner, *Frames of Mind: The Theory of Multiple Intelligences*, (1983) Basic Books.

<sup>7</sup> Danah Zohar, and Ian Marshall, *Spiritual Intelligence: The Ultimate Intelligence*, (2000) Bloomsbury Publishing.

<sup>8</sup> *Antiqua Et Nova, Note on the Relationship Between Artificial Intelligence and Human Intelligence*, 2025, Dicastery for the Doctrine Of The Faith and Dicastery for Culture and Education, [https://www.vatican.va/roman\\_curia/congregations/cfaith/documents/rc\\_dcf doc 20250128\\_antiqua-et-nova\\_en.html](https://www.vatican.va/roman_curia/congregations/cfaith/documents/rc_dcf doc 20250128_antiqua-et-nova_en.html)

<sup>9</sup> Genesis 1:26–27

<sup>10</sup> Aristotle. *Organon*. (Original work c. 350 BCE).

<sup>11</sup> René Descartes, *Discourse on the Method*, (1637).

<sup>12</sup> Ada Lovelace, "1842 Notes to the Translation of the Sketch of the Analytical Engine," *Ada User Journal* 36.3 (2015)..

<sup>13</sup> Alonso Church, "A Note on the Entscheidungsproblem," *Journal of Symbolic Logic*, 1(1936):40–41.

<sup>14</sup> Alan M Turing, "On Computable Numbers, with an Application to the Entscheidungsproblem," *Proceedings of the London Mathematical Society*, 2(42) (1936):230–265.

<sup>15</sup> W. S. McCulloch and W. Pitts, "A Logical Calculus of the Ideas Immanent in Nervous Activity," *Bulletin of Mathematical Biophysics*, 5 (1943):115–133.

<sup>16</sup> Alan M Turing, "Computing Machinery and Intelligence," *Mind*, 59(236) (1950):433–460.

<sup>17</sup> J. McCarthy, M. L. Minsky, N. Rochester and C. Shannon, "A Proposal for the Dartmouth Summer Research Project on Artificial Intelligence," *AI magazine*, 27(4) (2006): 12–12.

<sup>18</sup> Norbert Wiener, *Cybernetics: Or Control and Communication in the Animal and the Machine*, (1948) MIT Press.

<sup>19</sup> W. R. Ashby, *An Introduction to Cybernetics*, (1956) Chapman & Hall.

<sup>20</sup> W. G. Walter, "An Imitation of Life," *Scientific American*, 182(5) (1950), 42–45.

<sup>21</sup> J. Haugeland, *Artificial Intelligence: The Very Idea*, (1985) MIT Press.

<sup>22</sup> Hans Moravec, *Mind Children*, , (1988) Harvard University Press.

<sup>23</sup> F. Rosenblatt, "The Perceptron: A Probabilistic Model for Information Storage and Organization in the Brain," *Psychological Review*, 65(6), 386–408.

<sup>24</sup> D. E. Rumelhart, G. E. Hinton and R. J. Williams, "Learning Representations by Back-Propagating Errors," *Nature*, 323(6088), (1986) 533–536.

<sup>25</sup> Y. LeCun, L. Bottou, Y. Bengio and P. Haffner, , "Gradient-Based Learning Applied to Document Recognition," *Proceedings of the IEEE*, 86(11), (1998):2278–2324.

- <sup>26</sup> S. Hochreiter and J. Schmidhuber, "Long Short-Term Memory," *Neural Computation*, 9(8), (1997):1735–1780.
- <sup>27</sup> Ian J. Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair, Aaron Courville and Yoshua Bengio, "Generative Adversarial Nets," *Advances in Neural Information Processing Systems*, 27 (2014).
- <sup>28</sup> Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N. Gomez and Łukasz Kaiser, "Attention is All You Need," *Advances in Neural Information Processing Systems*, 30 (2017).
- <sup>29</sup> <https://openai.com/>
- <sup>30</sup> <https://claude.ai/>
- <sup>31</sup> Tom B. Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, Sandhini Agarwal, Ariel Herbert-Voss, Gretchen Krueger, Tom Henighan, Rewon Child, Aditya Ramesh, Daniel M. Ziegler, Jeffrey Wu, Clemens Winter, Christopher Hesse, Mark Chen, Eric Sigler, Mateusz Litwin, Scott Gray, Benjamin Chess, Jack Clark, Christopher Berner, Sam McCandlish, Alec Radford, Ilya Sutskever and Dario Amodei, "Language Models are Few-Shot Learners," *Advances in Neural Information Processing Systems*, 33, (2020):1877–1901.
- <sup>32</sup> N. Li, S. Liu, Y. Liu, S. Zhao and M. Liu, "Neural Speech Synthesis with Transformer Network," *Proceedings of the AAAI conference on artificial intelligence*, 33(1), (2019):6706-6713).
- <sup>33</sup> T. Fong, I. Nourbakhsh and K. Dautenhahn, "A Survey of Socially Interactive Robots," *Robotics and Autonomous Systems*, 42(3–4), (2003):143–166.
- <sup>34</sup> <https://openai.com/index/gpt-4o-system-card/>
- <sup>35</sup> A. Placani, "Anthropomorphism in AI: Hype and Fallacy," *AI and Ethics*, 4(3), (2024):691-698.
- <sup>36</sup> A. Deshpande, V. Murahari, T. Rajpurohit, A. Kalyan and K. Narasimhan, "Toxicity in Chatgpt: Analyzing Persona-assigned Language Models," *arXiv preprint*, (2024):arXiv:2304.05335.
- <sup>37</sup> Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Brian Ichter, Fei Xia, Ed Chi, Quoc Le and Denny Zhou, "Chain-of-Thought Prompting Elicits Reasoning in Large Language Models," *Advances in neural information processing systems*, 35, (2022):24824-24837.
- <sup>38</sup> <https://openai.com/index/learning-to-reason-with-llms/>
- <sup>39</sup> P. J. Christiano, J. Leike, T. Brown, M. Martic, S. Legg and D. Amodei, "Deep Reinforcement Learning from Human Preferences," *Advances in Neural Information Processing Systems*, 30 (2017)..
- <sup>40</sup> D. Kahneman, *Thinking, Fast and Slow*, Farrar, MacMillan (2011).
- <sup>41</sup> Ibid., 8
- <sup>42</sup> D. Guo, D. Yang, H. Zhang, J. Song, R. Zhang, R. Xu, Q. Zhu, S. Ma, P. Wang, X. Bi, X. Zhang, X. Yu, Y. Wu, Z. F. Wu, Z. Gou, Z. Shao, Z. Li, Z. Gao and Z. Zhang, "DeepSeek-R1: Incentivizing Reasoning Capability in LLMs via Reinforcement Learning," *arXiv preprint* (2025):arXiv:2501.12948
- <sup>43</sup> In the context of LLMs, "weights" refers to the vast number of numerical parameters that the model learns during its training process. These parameters essentially store the "knowledge" the model has acquired.
- <sup>44</sup> J. H. Flavell, "Metacognition and cognitive monitoring: A new area of cognitive–developmental inquiry," *American Psychologist*, 34(10) (1979):906–911.
- <sup>45</sup> J. Toy, J. MacAdam and P. Tabor, "Metacognition is all you need? Using Introspection in Generative Agents to Improve Goal-directed Behavior," *arXiv preprint* (2024):arXiv:2401.10910.
- <sup>46</sup> I. Mirzadeh, K. Alizadeh, H. Shahrokhi, O. Tuzel, S. Bengio, and M. Farajtabar, "GSM-Symbolic: Understanding the Limitations of Mathematical Reasoning in Large Language Models," *arXiv preprint* (2024):arXiv:2410.05229.
- <sup>47</sup> S. Russell and P. Norvig, *Artificial Intelligence: A Modern Approach* (4th Edition), (2020) Pearson. ISBN 978-0134610993.
- <sup>48</sup> Volodymyr Mnih, Koray Kavukcuoglu, David Silver, Andrei Rusu, Joel Veness, Marc G. Bellemare, Alex Graves, Martin Riedmiller, Andreas Fidjeland, Georg Ostrovski, Stig Petersen, Charles Beattie, Amir Sadik, Ioannis Antonoglou, Helen King, Dharmarajan Kumaran, Daan Wierstra, Shane Legg and Demis Hassabis, "Human-level control through deep reinforcement learning," *Nature* 518, (2015):529–533.
- <sup>49</sup> <https://robotsguide.com/robots?category=humanoids&sort=year>
- <sup>50</sup> <https://www.figure.ai/>
- <sup>51</sup> <https://bostondynamics.com/>
- <sup>52</sup> <https://robotsguide.com/robots/optimus>
- <sup>53</sup> <https://www.sanctuary.ai>
- <sup>54</sup> <https://agilityrobotics.com>

- <sup>55</sup> F. J. Varela, E. Thompson and E. Rosch, *The Embodied Mind: Cognitive Science and Human Experience.*, (1991) MIT Press.
- <sup>56</sup> Lawrence Shapiro, *Embodied Cognition*, (2019) Routledge.
- <sup>57</sup> R. A. Brooks, "Intelligence without Representation," *Artificial Intelligence*, 47(1–3), (1991):139–159.
- <sup>58</sup> F. Sun, R. Chen, T. Ji, Y. Luo, H. Zhou and H. Liu, . "A Comprehensive Survey on Embodied Intelligence: Advancements, Challenges, and Future Perspectives," *CAAI Artificial Intelligence Research*, 3, (2024): 9150042.
- <sup>59</sup> F. Zeng, W. Gan, Y. Wang, N. Liu and P. S. Yu, "Large Language Models for Robotics: A Survey," *arXiv preprint* (2023):arXiv:2311.07226..
- <sup>60</sup> L. Barrett and D. Stout, "Minds in Movement: Embodied Cognition in the Age of Artificial Intelligence," *Philosophical Transactions B*, 379(1911), (2024):20230144.
- <sup>61</sup> G. Moore, "Cramming more components onto integrated circuits," *Electronics*, 38(8), (1965):114-117.
- <sup>62</sup> Aristotle. *Metaphysics*. (Original work c. 350 BCE).
- <sup>63</sup> Aristotle. *Nicomachean Ethics*. (Original work c. 350 BCE)
- <sup>64</sup> Autonomy in AI stands for the ability of a system to operate independently, make decisions, and execute tasks without direct human intervention, which is essentially different to the "self-governance" or "self-determination" concept often considered a key aspect of individual freedom and moral agency.
- <sup>65</sup> Immanuel Kant, *Groundwork of the Metaphysics of Morals*, (2013) Cambridge University Press.
- <sup>66</sup> A. Clark and D. Chalmers, "The Extended Mind," *Analysis*, 58(1), (1998):7–19.
- <sup>67</sup> J. Marcos, "Artificial Intelligence and Delegated Control," *Journal of Technology Ethics*, (2018).
- <sup>68</sup> L. Floridi, *The Philosophy of Information*, (2011) Oxford University Press.
- <sup>69</sup> M. Coeckelbergh, *Growing Moral Relations: Critique of Moral Status Ascription*. (2012) Palgrave Macmillan.
- <sup>70</sup> S. Vallor, *Technology and the Virtues: A Philosophical Guide to a Future Worth Wanting*, (2016) Oxford University Press.
- <sup>71</sup> Ibid., 35
- <sup>72</sup> Ibid., 34
- <sup>73</sup> Luciano Floridi, "Translating Principles into Practices of Digital Ethics: Five Risks of Being Unethical," *Philosophy & Technology*, 32(2), (2019):185-193.
- <sup>74</sup> Luciano Floridi, *Pensare l'infosfera: La filosofia come design concettuale*. Raffaello Cortina Editore. (2020 ) Kindle Edition.
- <sup>75</sup> Jürgen Moltmann, *Dios de la Creación, Doctrina Ecológica de la Creación*, (1987) Salamanca: Ediciones Sígueme.
- <sup>76</sup> Ibid., 8
- <sup>77</sup> Colosenses 1:15
- <sup>78</sup> Hebreus 1:3
- <sup>79</sup> Romans 8:29
- <sup>80</sup> 2 Corinthians 3:18
- <sup>81</sup> 1 Corinthians 13:12
- <sup>82</sup> Shannon Vallor, *The AI Mirror: How to Reclaim Our Humanity in an Age of Machine Thinking*, (2024) Oxford University Press. Kindle Edition.
- <sup>83</sup> Ibid. 6.
- <sup>84</sup> Ibid., 8 18-19
- <sup>85</sup> Ibid., 8 21
- <sup>86</sup> Ibid., 8 22
- <sup>87</sup> Genesis. 1:27
- <sup>88</sup> Ibid., 8 36
- <sup>89</sup> <https://www.vatican.va/content/francesco/en/speeches/2024/june/documents/20240614-g7-intelligenza-artificiale.html>
- <sup>90</sup> Ibid., 8 16
- <sup>91</sup> John Paul II, Pope. *Man and Woman He Created Them: A Theology of the Body*, trans. M. Waldstein, (2006) Pauline Books & Media.
- <sup>92</sup> Exodus 35:30–35
- <sup>93</sup> <https://www.vatican.va/content/francesco/en/messages/peace/documents/20231208-messaggio-57giornatamondiale-pace2024.html>
- <sup>94</sup> <https://www.vatican.va/content/francesco/en/speeches/2024/june/documents/20240614-g7-intelligenza-artificiale.html>
- <sup>95</sup> Romans 2:14–15
- <sup>96</sup> Isaiah 44:9–20