# The Illusion of Understanding: Deconstructing AI Metaphors

CJ Trowbridge
cj@cjtrowbridge.com
San Francisco State University
Ethical Artificial Intelligence

## Abstract

Through a critical examination of the common metaphors used to describe artificial intelligence technologies, I argue that these metaphors, such as intelligence, rewards, punishment, attention, and learning, while widespread in AI discourse, often distort the actual capabilities, limitations, and risks of AI systems. By anthropomorphizing AI, these metaphors create unrealistic expectations and misunderstandings about the nature of AI technologies, leading to potential misapplications and safety concerns, especially in critical decision-making areas.

Through a series of vignettes, this paper deconstructs each metaphor, revealing the discrepancies between the metaphorical language and the technological realities of AI. The analysis highlights how these metaphors not only oversimplify but sometimes mislead stakeholders across various sectors, including the public, policymakers, and the AI research community.

In examining the metaphors used to describe artificial intelligence, it is crucial to understand how these metaphors construct and constrain our knowledge. Michel Foucault's concept of the "episteme," as discussed in The Order of Things, provides a valuable framework for this analysis. Foucault argues that the way knowledge is structured within a particular historical period shapes what is considered possible or true within that era (Foucault, 1966). By applying this lens to AI, we can see how the metaphors of "intelligence," "learning," and "attention" not only simplify but also limit our understanding of AI's capabilities. These metaphors create an episteme that may hinder more accurate and nuanced interpretations of AI technologies, thus affecting both public perception and policy decisions.

The goal of this paper is to foster a more nuanced and accurate understanding of AI, advocating for a shift towards terminology that reflects the true nature of these systems. (Mitchell, M. 2019) This shift is essential for the responsible development and ethical deployment of AI technologies. By promoting a clearer and more precise discourse, we aim to align AI development with societal needs and ethical standards, ensuring that AI's potential is harnessed responsibly and effectively.

# 1 Introduction

In the rapidly evolving field of artificial intelligence, the language and metaphors we use to describe AI technologies play a critical role in shaping public perception, policy decisions, and the direction of research. The title of this essay, "The Illusion of Understanding: Deconstructing AI Metaphors," reflects the central thesis: the popular metaphors currently employed in discussions about AI often impart more confusion than clarity. The risk of misrepresenting and misunderstanding what AI systems are truly capable of and how they operate is lessened by using less abstract and more accurate explanations of what is actually happening in the computational processes.

This essay aims to dissect and critique the metaphors commonly used in AI discourse—such as intelligence, rewards and punishment, attention, and learning—demonstrating how they can lead to misunderstandings and unrealistic expectations. A series of detailed vignettes explore how these terms, while convenient and appealing for their simplicity, fail to capture the true nature of AI and its capabilities. By examining these metaphors through the lens of deconstruction, as Jacques Derrida suggests in Of Grammatology, we can uncover the underlying assumptions and biases in AI discourse (Derrida, 1967). Each vignette focuses on a specific aspect of AI, unpacking the metaphorical language used to describe it and proposing a more nuanced understanding that aligns with the actual technology. As Turing (1950) noted, the question of whether machines can think is not as straightforward as it seems, and the language we use significantly impacts our understanding and expectations.

# 2 Intelligence

Intelligence in artificial intelligence is often treated as a measurable, quantifiable trait, similar to human or animal intelligence. This has led to the creation of completely arbitrary and meaningless categories in the popular discourse such as narrow, general, and super-human artificial intelligence. These categories are intended to explain that on some imaginary one-dimensional scale called Intelligence, a particular AI technology may be classified as either less, equal, or more than the average human. In reality this system of classification confuses the issue rather than providing clarity.

Gilles Deleuze and Félix Guattari's concept of rhizomatic thinking, presented in *A Thousand Plateaus*, challenges linear and hierarchical models of knowledge. They propose the rhizome as a model of knowledge that is non-hierarchical and interconnected (Deleuze & Guattari, 1980). This perspective is particularly useful for rethinking AI metaphors, which often rely on linear and hierarchical notions of "intelligence" and "learning." By adopting a rhizomatic approach, we can better appreciate the multifaceted and emergent nature of AI systems, leading to more accurate and flexible descriptions of their capabilities and limitations.

The true nature of intelligence is of course more complicated than some imaginary one-dimensional Intelligence Quotient. (Minsky, M., & Papert, S.) Ask yourself what an IQ test tells you about the intelligence of an octopus or a cat. Such classifications do not adequately capture the complexity and multi-dimensionality inherent in intelligent behavior, whether human, animal, or artificial. By examining

intelligence as a socially constructed abstraction rather than a tangible, discrete quantity, this first vignette seeks to deconstruct these commonly accepted notions.

Intelligence, particularly in the context of AI, emerges from complex interactions and cannot be accurately reduced to a single scale or continuum. (Montemayor, C. 2012) By reevaluating how intelligence is defined and discussed in AI, this vignette will highlight the importance of moving towards a more thoughtful and precise discourse, which acknowledges the abstract and emergent qualities of intelligence rather than oversimplifying it to fit within narrow definitional boundaries.

Intelligence, traditionally viewed through the lens of human cognition, is often mistaken as a concrete, measurable entity when applied to artificial intelligence. Yet, this perspective is fundamentally flawed because intelligence, by its nature, is an abstract construct—a product of social consensus rather than an inherent, quantifiable attribute. The roots of this conceptualization can be traced back to the social and historical contexts in which intelligence assessment emerged. Historically, intelligence testing began as a means to classify and stratify individuals based on perceived cognitive abilities, heavily influenced by the educational and psychological theories of the time. These tests were created with the intention of identifying and categorizing human capabilities, often reinforcing social hierarchies and biases rather than providing an objective measure of some immutable trait. When these ideas are uncritically transposed onto AI, they perpetuate a misleading representation of what AI can achieve and fundamentally misguide our understanding of artificial 'intelligence'. This discussion seeks to unpack these layers, challenging the view of intelligence as something that can be simply coded or quantified in AI systems, and highlighting the impact of societal constructs on our technological advancements.

Intelligence, particularly when considered within the context of both human cognition and artificial intelligence, emerges from a myriad of complex, high-dimensional interactions. These interactions span countless cognitive, emotional, and social domains, each contributing to the holistic construct we refer to as intelligence. This multifaceted nature makes it clear that intelligence is not a single, linear attribute that can be easily measured or distilled into a quantifiable value. In the realm of AI, this complexity is often overlooked, as the field historically attempts to model intelligence through computational methods that inherently reduce its richness to more manageable dimensions. (Mitchell, M. 2019) Such reductions fail to capture the true essence and variability of intelligent behavior, as they ignore the depth and breadth of the underlying processes. Recognizing intelligence as an emergent phenomenon that arises from these intricate and interconnected dimensions means understanding the limitations as well as the potential of artificial systems and challenging the oversimplified views that currently dominate AI research and discourse. This perspective calls for a reevaluation of how intelligence is defined and measured, urging a move away from the reductive metrics that have long been used to gauge both human and machine intelligence.

This realization is already happening in industry, where a plethora of standardized tests have emerged for challenging and scrutinizing the abilities of new models in specific, replicable, measurable ways that are easy to communicate. These include tests like MMLU (Massive Multitask Language Understanding), ARC (Abstraction and Reasoning Corpus), HellaSwag (Commonsense reasoning),

TruthfulQA (817 questions in 38 categories including health, law, finance and politics intended to illicit hallucination or confabulation), etc. The publication of any new large language model today typically includes rigorous third-party testing of these and other metrics in order to concisely and honestly communicate precise, specific, replicable measures of the capabilities of the model.

The classification of artificial intelligence into categories such as narrow, general, and superintelligence is prevalent in both academic and popular discourse. However, these categories, while useful for theoretical discussions, are fundamentally arbitrary and fail to encapsulate the true complexity of AI capabilities. These distinctions, though conceptually appealing, oversimplify the vast and nuanced landscape of AI development. They create a misleading framework that suggests a linear progression from narrow to super forms, ignoring the diverse, non-linear, and multifaceted nature of intelligence as it actually manifests in AI systems. A reliance on these categories not only skews public and professional understanding of AI but also stifles more nuanced research approaches that could lead to more profound insights into how AI operates and interacts within various contexts. Challenging these conventional categories encourages a more accurate and holistic approach to understanding artificial intelligence, one that reflects its true diversity and complexity.

Misconceptions about the nature of artificial intelligence can have profound implications for both the development of AI technologies and the expectations placed upon them. When intelligence is misunderstood as a linear, measurable trait that AI systems can possess or achieve, it leads to a host of unrealistic expectations. For instance, stakeholders may anticipate that AI systems will autonomously perform complex tasks with human-like understanding, without recognizing the inherent limitations of these systems. This misalignment can result in misguided development strategies, where resources are allocated toward achieving ill-defined goals of 'general' or 'superintelligence', potentially neglecting areas where AI could provide more immediate and practical benefits. It could also blind us to risks posed by AI which has been misclassified as less risky due to its lower "intelligence."

Furthermore, a simplistic view of intelligence promotes ethical concerns, such as the over-reliance on AI decision-making in critical areas like healthcare, legal systems, and public safety, under the false assumption that these systems operate with comprehensive understanding and impartiality. Therefore, cultivating a more nuanced appreciation of what AI can and cannot do is crucial for directing ethical AI research and development. By emphasizing a multifaceted approach to intelligence that acknowledges its emergent and context-dependent qualities, AI research can progress in a manner that is both ethically responsible and effectively aligned with realistic capabilities, thus fostering technologies that enhance human decision-making rather than erroneously attempting to replace it. Moving forward, it is crucial for both the AI research community and the public to adopt language and frameworks that avoid oversimplification and instead reflect the intricate reality of how AI technologies operate. This shift in perspective is not merely academic but essential for the ethical development and implementation of AI systems, ensuring they are used responsibly and effectively within society.

## 3 Rewards and Punishment

Metaphors like "training" and "learning" are often misused in discussions about AI to imply connections to human or animal behaviors. This misrepresentation leads to the widespread belief that AI systems can be controlled or guided through mechanisms akin to rewarding or punishing human or animal behaviors. The aim of this section is to challenge these conventional understandings and highlight the mistaken assumptions about the ability of society, organizations, engineers, and users to exert control or influence over AI models through such methods. By dissecting how these concepts are applied in AI, we seek to clarify the true nature of AI learning and optimization processes, setting the stage for a more accurate discussion about AI's capabilities and limitations.

Gradient descent is a cornerstone mathematical optimization technique in AI. It is often simplified and misrepresented as analogous to rewards and punishment. (Mnih et al., 2015) In reality, this technique involves adjusting the parameters of an AI over time in order to minimize the differences between the model's predictions and the actual data. Despite its importance, the nuances of gradient descent are frequently overlooked, leading to the misconception that AI can be guided or controlled in a manner similar to behavioral conditioning in animals. Gradient descent is fundamentally a process of mathematical calculation aimed at optimizing model performance, devoid of any sentient or volitional qualities that the terms "rewards" and "punishment" imply.

The metaphor of rewards and punishment, frequently used to describe the training processes of AI, notably misrepresents the actual mechanisms at work. This metaphor suggests that AI models respond to positive or negative stimuli, similar to how animals might learn from reinforcement. However, the reality is that gradient descent and other AI learning algorithms do not operate through sentient reactions to incentives. These processes are purely mathematical, involving the adjustment of parameters based on calculated gradients, not moral or behavioral corrections. This section argues against the anthropomorphism of AI learning processes, highlighting the inaccuracy of the rewards and punishment metaphor and its implications. A more accurate perception of AI's capabilities emphasizes that AI does not 'learn' through rewards or punishments but through iterative improvements of the statistical model based on analysis of the data over time.

Jean Baudrillard's concept of hyperreality, as articulated in *Simulacra and Simulation*, is particularly relevant when discussing the metaphors used to describe AI. Baudrillard argues that in a hyperreal world, representations and symbols come to replace and distort reality (Baudrillard, 1981). Similarly, the metaphors of "intelligence," "rewards," and "punishment" in AI discourse create hyperreal expectations about what AI systems can do. These metaphors construct a version of AI that appears more autonomous and capable than it truly is, leading to significant misunderstandings and unrealistic expectations among stakeholders and the public.

A more accurate metaphor of gradient descent would be a drop of water finding its way to a lake or ocean. This imagery illustrates the process of gradient descent as the water naturally seeks the path of least resistance, gradually flowing downhill until it reaches a body of water. Similarly, gradient descent involves the model parameters adjusting incrementally, following the path of steepest descent in the loss

landscape, to find the optimal solution or the "lowest point." This metaphor effectively communicates the essence of gradient descent as a continuous and passive optimization process, devoid of any sentient decision-making or emotional motivations. Unlike the misleading rewards and punishment analogy, the water metaphor highlights the non-volitional, systematic nature of how AI models refine their responses, providing a clearer and more factual representation of the AI optimization process. This understanding is crucial for demystifying AI operations and fostering a more accurate appreciation of how AI algorithms improve and evolve over time.

The prevailing misconceptions that imbue AI with qualities of sentient learning and decision-making is clarified by dissecting what gradient descent actually entails—an iterative, non-sentient optimization method akin to water flowing towards the lowest point. This correction is vital not only for a more accurate understanding of AI's capabilities but also for setting realistic expectations about the control and guidance we can exert over AI systems. By fostering a clearer comprehension of AI training processes, stakeholders across various sectors can make better-informed decisions, ensuring that AI technologies are used ethically and effectively in society.

## 4 Attention

The modern AI renaissance was ushered in by a now-famous academic paper entitled, "Attention is all you need." (Vaswani et al., 2017) This paper argued for a new architecture for connecting nodes inside an AI model, which it called "attention." In reality the biological concept of attention has essentially no relationship to this architectural change in AI models, but suddenly every AI was said to have attention capabilities. These mechanisms are understood to selectively concentrate on certain parts of the input data while ignoring others, ostensibly mimicking human attentional processes. However, this conventional understanding often glosses over the fundamental differences between human attention and its artificial counterpart. Unlike human attention, which involves awareness, perspective-taking, and cognitive flexibility, AI attention is purely a computational function without consciousness or understanding. (Montemayor, 2015) Exploring the philosophical and practical aspects of attention illuminates the misconceptions that lead to overestimations of AI's capabilities and advocate for more accurate and grounded discussions about the role of attention in artificial intelligence systems.

The philosophical concept of attention, deeply rooted in cognitive science and philosophy, involves not just the focus on specific stimuli but also the capacity to consider and imagine other perspectives—an essential element of human consciousness and interaction. This multifaceted understanding contrasts sharply with AI attention mechanisms, which, despite their name, do not possess true perspective-taking abilities. AI attention is designed to enhance model performance by focusing computational resources on relevant parts of data, such as key features in an image or important words in a sentence. However, it lacks the subjective awareness and the ability to engage with alternative viewpoints that characterize human attention. This discrepancy highlights a fundamental gap in AI's mimicry of human cognitive processes, where the complexity of genuine attention is reduced to a mere optimization tool within neural networks. Comparing these distinctly different concepts of

attention corrects the misunderstandings about what AI can achieve with its so-called 'attention' and emphasize the limitations inherent in current AI models.

AI attention mechanisms, often celebrated for their ability to enhance model performance, focus on specific parts of input data, such as prioritizing certain words in text or features in images. This process, while technically impressive, fundamentally lacks the depth and subtlety of human attention. Human attention involves not just the selection of sensory information but also a rich contextual understanding and the ability to shift focus in response to subtle cues and changes in environment. AI attention, on the other hand, is purely a computational tool, programmed to optimize processing efficiency and accuracy in tasks like language translation or image recognition. It does not, and cannot, replicate the conscious or cognitive processes associated with human attention, such as integrating emotional responses or adjusting based on social dynamics. This limitation is critical to understand, as it underscores that AI, despite advancements, remains a tool devoid of genuine cognitive abilities, highlighting the stark difference between optimizing data processing and achieving true cognitive comprehension.

Despite the sophisticated use of "attention" mechanisms in AI models, these systems fundamentally lack the capacity for true understanding or perspective-taking. For instance, while AI can process and respond to language inputs, it does not truly comprehend the content in the way humans do. This absence of genuine understanding is evident in tasks requiring nuanced interpersonal interactions, such as writing convincing dialogue between two characters. AI might generate text that superficially resembles conversation, but it often lacks the depth, context, and emotional subtleties that come from real understanding and perspective-taking. This limitation not only affects the quality of AI-generated content but also highlights the broader issue of AI's inability to engage with or adapt to the complexities of human thought and culture. The use of attention mechanisms may improve the efficiency of these processes, but they do not bridge the fundamental gap between data processing and genuine cognitive engagement, emphasizing the need for a clear distinction between the technical capabilities of AI and the attributes of human intelligence.

Misconceptions about the capabilities of AI attention mechanisms can lead to significant overestimations of what AI can truly achieve and misunderstandings about the nature of its functionality. By attributing human-like qualities of understanding and cognitive engagement to AI systems, stakeholders—ranging from developers to policymakers—may mistakenly assume these technologies possess greater autonomy and comprehension than they actually do. This can result in inappropriate reliance on AI for tasks that require genuine human judgment and emotional intelligence, such as in healthcare, law enforcement, and education sectors. The critical importance of accurate terminology and a clear understanding of AI's limitations cannot be overstated; it is essential for responsible AI development and communication. Clarifying these limits fosters more realistic expectations and ensures that AI is used appropriately, complementing rather than attempting to substitute the nuanced capabilities of human intelligence. This approach also improves ethical considerations in AI deployment, preventing potential harms that could arise from misaligned expectations and applications.

Donna Haraway's *A Cyborg Manifesto* offers a critical perspective on the intersection of technology, society, and identity, which is highly pertinent to the discussion of AI metaphors. Haraway argues that the cyborg, as a hybrid of machine and organism, challenges traditional boundaries and categories (Haraway, 1985). Similarly, AI technologies blur the lines between human and machine, yet the metaphors we use often reinforce outdated notions of human superiority and control. By adopting Haraway's cyborg perspective, we can develop more nuanced and inclusive metaphors that better reflect the complex interplay between AI and human cognition.

Emphasizing the computational nature of AI attention and its stark contrast to human cognitive processes highlights the necessity for precise and grounded discourse in the field of artificial intelligence. It is crucial that both the public and professionals recognize that, despite its advancements, AI does not have the ability to genuinely understand or engage with content as humans do. The popular metaphors and misleading terminologies used in AI impart more confusion than clarity. Finding that clarity is essential for setting realistic expectations, ensuring responsible development, and guiding ethical deployment of AI technologies in society.

## 5 Conclusion

These are only some of the most prevalent and inaccurate metaphors that permeate discussions about artificial intelligence. Scrutinizing the misleading use of metaphors such as intelligence, rewards, punishment, attention, and learning is vital for safely implementing or using artificial intelligence systems today. There is a vast gap between the metaphorical and the actual, highlighting the significant misunderstandings that can arise when human cognitive attributes are uncritically applied to AI systems.

Jacques Derrida's notion of deconstruction, as explored in Of Grammatology, provides a critical tool for examining the language used in AI discourse. Derrida emphasizes the importance of deconstructing the binary oppositions and hierarchical structures embedded in language (Derrida, 1967). Applying this approach to AI metaphors such as 'intelligence' and 'learning' reveals the underlying assumptions and biases that these terms carry. By deconstructing these metaphors, we can better understand the complexities and limitations of AI technologies, moving beyond simplistic and anthropomorphic representations.

The metaphors commonly employed to describe AI capabilities do not just oversimplify but often distort the reality of how AI systems operate. Equating AI learning with human learning processes, or gradient descent with rewards and punishment, leads to a profound misalignment between expectations and the reality of how AI will actually behave in the wild. Such misrepresentations can result in misplaced trust in AI's autonomy and decision-making capabilities, raising ethical concerns, especially when AI is deployed in critical areas such as healthcare, legal systems, and public safety.

By advocating for a shift towards more precise and grounded discourse, this essay underscores the importance of using accurate terminology in AI discussions. It is essential that researchers, developers, policymakers, and the public cultivate a more nuanced understanding of AI that acknowledges its computational basis and inherent limitations. (Dupuy, 2000) Only through clear and honest communication can AI development align with ethical standards and societal needs, avoiding the pitfalls of misunderstanding and misapplication. (Müller, 2020)

Moving forward, it is crucial for the AI community and society at large to engage in a more informed dialogue about what AI can and cannot do. This requires an ongoing commitment to education and transparency, fostering an environment where AI's potential can be harnessed responsibly and effectively. Dispelling the illusions cast by outdated metaphors and embracing a more accurate portrayal of AI can help stakeholders better navigate the complex landscape of technological advancement and its implications for the future.

**Works Cited**

Baudrillard, J. (1981). Simulacra and simulation. Semiotext(e).

Deleuze, G., & Guattari, F. (1980). A thousand plateaus: Capitalism and schizophrenia (B. Massumi, Trans.). University of Minnesota Press.

Derrida, J. (1967). Of grammatology. Johns Hopkins University Press.

Dupuy, J. P. (2000). The mechanization of mind. Princeton University Press.

Foucault, M. (1966). The order of things: An archaeology of the human sciences. Pantheon Books.

Haraway, D. (1985). A cyborg manifesto: Science, technology, and socialist-feminism in the late twentieth century. In Simians, cyborgs, and women: The reinvention of nature (pp. 149-181). Routledge.

Minsky, M., & Papert, S. (1969). Perceptrons. MIT Press.

Mitchell, M. (2019). Artificial intelligence: A guide for thinking humans. Farrar, Straus and Giroux.

Mnih, V., et al. (2015). Human-level control through deep reinforcement learning. Nature, 518(7540), 529-533.

Montemayor, C. (2012). Minding time: A philosophical and theoretical approach to the psychology of time. Brill.

Montemayor, C. (2015). Consciousness, attention, and conscious attention. MIT Press.

Müller, V. C. (2020). Ethics of artificial intelligence and robotics. Routledge.

Turing, A. (1950). Computing machinery and intelligence. Mind, 59(236), 433-460.

Vaswani, A., et al. (2017). Attention is all you need. Advances in Neural Information Processing Systems, 30.