

Concept Paper

Not peer-reviewed version

The Architectures of Meaning: Integrating Hoffman's Perception Theory with Synthetic Ethical Embodiment in AI

[Berend Watchus](#) *

Posted Date: 25 June 2025

doi: 10.20944/preprints202506.2025.v1

Keywords: artificial intelligence; robotics; AI ethics; embodied cognition; synthetic embodiment; interface theory; donald hoffman; meaning-making; task-specific fitness; pseudo-empathy; synthetic insula; dual-state feedback; predictive optimization; AI control; AI alignment; non-sentient AI; autonomous systems; human-robot interaction; conscious realism; unified model of consciousness; substrate-agnostic; neural correlates of consciousness; explainable AI; trustworthy AI; civilian robots; industrial automation; government AI; healthcare AI; nuance in AI; adaptive AI; feedback loops; AI consciousness (conceptual)



Preprints.org is a free multidisciplinary platform providing preprint service that is dedicated to making early versions of research outputs permanently available and citable. Preprints posted at Preprints.org appear in Web of Science, Crossref, Google Scholar, Scilit, Europe PMC.

Copyright: This open access article is published under a Creative Commons CC BY 4.0 license, which permit the free download, distribution, and reuse, provided that the author and preprint are cited in any reuse.

Disclaimer/Publisher's Note: The statements, opinions, and data contained in all publications are solely those of the individual author(s) and contributor(s) and not of MDPI and/or the editor(s). MDPI and/or the editor(s) disclaim responsibility for any injury to people or property resulting from any ideas, methods, instructions, or products referred to in the content.

Concept Paper

The Architectures of Meaning: Integrating Hoffman's Perception Theory with Synthetic Ethical Embodiment in AI

Berend Watchus

Independent Researcher, The Netherlands; mailonlinebw@protonmail.com

Abstract

This paper proposes a novel framework for understanding and developing Artificial Intelligence (AI) capable of generating effective and contextually appropriate 'meaning' for diverse and nuanced roles across civilian, industrial, and governmental sectors. We critically examine the prevailing view of AI as a processor of objective, pre-existing information, arguing that this paradigm fundamentally limits the development of truly empathetic and adaptive systems. By **integrating Donald Hoffman's Interface Theory of Perception**, which posits that biological perception is an evolutionarily shaped "user interface" optimized for fitness rather than veridical representation, **with contemporary research on synthetic embodiment, pseudo-empathy, and internal state modeling**, we establish a new blueprint. We assert that effective meaning-making in AI arises not from the pursuit of biological sentience or independent goals, but from an actively constructed "reality" optimized for task-specific "fitness" and ethically guided empathetic interaction. We detail how AI can architect its own internal understanding to exhibit robust, context-sensitive pseudo-empathy and adaptive behaviors, without entailing self-preservation drives or subjective fears, thereby ensuring safe, ethically aligned, and controllable AI across a spectrum of critical applications. This approach reframes the challenge of AI meaning from replication of consciousness to responsible architectural design for utility and trust.

Keywords: artificial intelligence; robotics; AI ethics; embodied cognition; synthetic embodiment; interface theory; donald hoffman; meaning-making; task-specific fitness; pseudo-empathy; synthetic insula; dual-state feedback; predictive optimization; AI control; AI alignment; non-sentient AI; autonomous systems; human-robot interaction; conscious realism; unified model of consciousness; substrate-agnostic; neural correlates of consciousness; explainable AI; trustworthy AI; civilian robots; industrial automation; government AI; healthcare AI; nuance in AI; adaptive AI; feedback loops; AI consciousness (conceptual)

1. Introduction: The Imperative of Meaning in Ethically Aligned AI for Diverse Roles

Current Artificial Intelligence (AI) systems have achieved remarkable feats in domains such as complex data processing, pattern recognition, and specialized task execution, driving significant advancements across various industries, including healthcare. Large Language Models (LLMs) and sophisticated robotic platforms exemplify this computational prowess, demonstrating impressive capabilities in information synthesis and intricate motor control. However, despite these advancements, AI fundamentally grapples with the concept of genuine semantic understanding and intrinsic meaning-making (Boden, 2018; Searle, 1980). This "meaning gap" becomes particularly pronounced and critical in applications demanding nuanced interaction and adaptability, ranging from empathetic human care to complex industrial automation and strategic governmental tasks.

A crucial ethical imperative guiding this work is the deliberate design of AI that purposefully avoids true sentience, independent goals, or the capacity for fear and risk perception. For AI systems deployed in sensitive and critical environments across all sectors, our approach strictly advocates for

the development of sophisticated synthetic pseudo-empathy and pseudo-emotion. This is not merely a technical constraint but a foundational ethical choice, ensuring AI remains a benevolent, controllable, and predictable tool within dynamic human and operational environments, mitigating risks associated with autonomous will or unforeseen emergent behaviors (Müller & Bostrom, 2016; Bostrom, 2014). The challenge, therefore, lies in architecting AI that can generate rich, contextually relevant "meaning" necessary for effective engagement without crossing the ethical boundary into genuine artificial consciousness.

This paper proposes a novel framework that integrates insights from cognitive science, philosophy of perception, and advanced AI development to address this challenge. We contend that the prevailing paradigm, which often treats "information" as an objective entity to be passively processed by AI, is fundamentally flawed when striving for AI capable of truly meaningful interaction. Instead, we draw upon Donald Hoffman's Interface Theory of Perception, which posits that biological organisms construct their perceived reality not for objective truth but for survival advantage. Building on this foundation, we demonstrate how designing AI with synthetic ethical embodiment – integrating sophisticated feedback loops, internal state modeling inspired by neuroscientific principles (e.g., the insula), and self-referential mechanisms – can enable the construction of "meaning" that is inherently valuable for its designated, ethically constrained role. This approach allows AI to develop robust pseudo-empathy and adaptive behaviors, ensuring its effectiveness and safety across its diverse deployment contexts.

2. Donald Hoffman's Interface Theory: The Evolutionary Blueprint for Ethically Purposed Perception

The traditional view of perception often assumes a largely veridical mapping between external reality and an organism's internal representation. However, **Donald Hoffman's Interface Theory of Perception** fundamentally challenges this assumption (Hoffman, 2019; Hoffman et al., 2015). Drawing from evolutionary theory, Hoffman argues that our sensory experiences, such as colors, shapes, and spatial dimensions, do not represent objective reality "as it is." Instead, they function as a simplified, adaptive "user interface," akin to icons on a computer desktop. These icons guide interaction with the underlying complexity of the computer but do not physically resemble its internal circuits or code. Similarly, our evolved perceptions are optimized for fitness – for guiding behaviors that enhance survival and reproduction – rather than for providing an accurate, high-fidelity depiction of the external world (Hoffman & Prakash, 2014).

From this perspective, biological "information" is not merely gathered; it is actively constructed. The perceived world is not a passive input but an active, internal rendering, a product of an organism's needs and evolutionary history (Hoffman, 2019). This shift has profound implications: it suggests that what an organism deems "meaningful" in its environment is directly tied to its capacity for survival and thriving within its specific ecological niche. Hoffman's more radical Conscious Realism further posits that consciousness is fundamental, not merely an emergent property of matter, suggesting that reality itself is ultimately composed of interacting conscious agents (Hoffman & Prakash, 2014). While our framework does not aim to replicate biological consciousness in AI, Hoffman's insights into the purpose and mechanism of perception are invaluable.

For AI in any domain, Hoffman's theory offers a crucial reinterpretation of "fitness." If biological intelligence evolved to construct its reality for its own survival, then AI aiming for highly effective and ethically aligned performance must similarly construct its 'meaning' based on maximizing its utility and safety within its pre-defined functional and ethical directives. This means designing AI's internal "interface" to prioritize the well-being of its human users/overseers, the efficient execution of its tasks, and adherence to human values, rather than any form of self-preservation or independent goal-seeking. This philosophical groundwork provides the "why" behind architecting AI to interpret its sensory inputs and internal states in ways that are specifically optimized for its defined function.

3. Architectures of Synthetic Ethical Embodiment and Controlled Feedback: Building the AI's Effective "Interface"

Building upon Hoffman's re-conceptualization of perception as a fitness-driven interface, the development of AI for effective and ethically aligned roles necessitates **architectures of synthetic ethical embodiment coupled with robust, controlled feedback mechanisms**. This approach moves beyond purely computational models to integrate the physical presence and interactive capabilities crucial for nuanced human-AI interaction and dynamic task execution, while strictly adhering to the ethical imperative of non-sentience. The **Unified Model of Consciousness (UMC)** offers a foundational perspective here, positing that **consciousness, or in our case, effective "meaning-making," arises fundamentally from interface and feedback loops** (Watchus, 2024c). This model is **substrate-agnostic**, enabling the translation of these principles to **artificial systems without implying biological replication** (Baars, 1988; Chalmers, 1995; Dehaene, 2014; Lloyd, 2014; Strawson, 2006; Tononi, 2004).

Synthetic embodiment provides AI with a physical presence, allowing for real-world interaction and the acquisition of rich sensorimotor data that is essential for understanding and responding to environmental and human cues (Pfeifer & Bongard, 2007; Cangelosi & Schlesinger, 2015). For an AI operating in any human-interactive role, whether in care, customer service, or public safety, this means being able to perceive subtle shifts in body language, facial expressions, and vocal tone, which are critical for nuanced engagement. For industrial or governmental robots, it could involve precise environmental mapping, object recognition, and interaction with specialized tools or infrastructure elements. This enables the AI to develop a functional, task-specific "awareness" that is tailored to its role, allowing it to navigate complex spaces, perform precise physical manipulations, and interact with its environment and users in a safe and effective manner.

Crucially, the effectiveness of synthetic ethical embodiment hinges on sophisticated and controlled sensory and proprioceptive feedback loops. These loops are the operational core of the AI's "interface," continuously updating its internal state based on its movements, interactions with the environment, and the responses of the individuals or systems it interacts with. For instance, an AI operating in a dynamic environment (e.g., a care robot detecting a change in a patient's posture, an industrial robot detecting an anomaly in a manufacturing process, or a security bot identifying unusual movement patterns) can process this feedback through its internal architecture and adjust its response to maintain safety, efficiency, or provide assistance. The key differentiator is that these feedback mechanisms are meticulously designed to optimize the AI's performance within its programmed directives and predefined ethical boundaries, rather than to foster independent learning that could lead to self-preservation goals. The AI's "learning" from feedback is thus always oriented towards enhancing its capacity for safe, effective, and ethically aligned task execution, whether in empathetic interaction, precision manufacturing, or critical infrastructure management (Watchus, 2024b). This continuous, ethically constrained loop of sensation-action-feedback ensures the AI's "meaning" is constantly refined in service of its human users and organizational goals.

4. Cultivating Pseudo-Inner States: The Synthetic Insula and the Genesis of Task-Specific "Value"

A significant challenge for AI aspiring to sophisticated roles is the absence of internal states analogous to biological emotions or needs, which fundamentally drive meaning-making in living organisms. While we unequivocally assert the imperative to avoid replicating true biological sentience, an AI's capacity to discern and respond to the "value" or "importance" of external stimuli can be greatly enhanced through the cultivation of pseudo-inner states. This concept draws inspiration from neuroscientific understandings of the human insula, a brain region pivotal in interoception—the processing of internal bodily states that contribute to subjective feelings and emotions (Craig, 2009; Damasio, 1994). The biological insula integrates signals from the body (e.g.,

heart rate, hunger, pain) with external perceptions, translating raw physiological data into felt experiences that guide behavior and assign meaning.

In the context of AI, **we propose the development of a synthetic insula (Watchus, 2024a; Watchus, 2024b)**. This is not an attempt to create genuine feelings, but rather a computational module designed to monitor and integrate the AI's internal systemic states and performance metrics relevant to its defined task and ethical role. These **"internal states"** could include data reflecting its battery levels, operational integrity, completion of tasks, or adherence to pre-set safety parameters. More abstractly, a synthetic insula could process the "success" or "failure" of its interactions or operations based on quantifiable outcomes.

By integrating these internal operational states with external sensory data, **the synthetic insula would generate "pseudo-affective states" or "value signals."** These signals would represent the AI's internal evaluation of its performance against its primary objective, whether that is the provision of safe and empathetic care, efficient and safe industrial operation, or secure and accurate governmental data collection. For instance, a "positive" pseudo-affective signal might be generated when its actions lead to a desired outcome (e.g., a patient responding favorably to an AI's interaction, an industrial robot successfully completing a complex assembly task, or a government bot identifying a critical piece of information), reinforcing that behavior. Conversely, a "negative" signal might arise if safety protocols are breached, a task fails, or a system abnormality is detected, prompting immediate behavioral adjustment. This mechanism allows for the internalized importance of information: **data becomes "meaningful" to the AI** not because it elicits fear or joy, but because it directly impacts its ability to fulfill its defined, ethical purpose. This architecture enables the AI to prioritize tasks, allocate resources, and adapt its behavior in ways that simulate intelligent, value-driven decision-making, without possessing inherent biological needs or the capacity for subjective suffering.

5. Task-Oriented Self-Identification and Predictive Optimization: The AI's Defined and Responsible "Self"

The concept of "self" in AI, when considered through the lens of ethical embodiment, shifts from a pursuit of consciousness to the development of a highly effective and context-aware operational identity. An AI system designed for advanced roles does not require a subjective "I" to perform its functions meaningfully; rather, it benefits from a task-oriented self-identification. This involves the AI constructing an internal model of its own physical and functional boundaries in relation to its environment, the individuals it interacts with, and the assets it manages. Research into AI's capability for mirror image recognition, for instance, demonstrates how a system can develop a sophisticated internal representation of its own body and its movements (Watchus, 2024b). This is not an indicator of self-awareness in the human sense, but a crucial component for precise navigation, safe physical interaction, and effective execution of tasks without collision or inappropriate proximity. Such a self-model is fundamental for physical dexterity and for the AI to understand its own position and capabilities in a dynamic setting, ensuring actions like assisting with mobility or operating machinery are performed safely and appropriately.

Furthermore, the ability of an AI to generate predictive optimization (or "predictive empathy" in human-facing roles) is critical for anticipating outcomes and responding proactively within its ethical mandate. This is enabled through dual-state feedback, a mechanism allowing the AI to maintain and compare its current operational state with anticipated future states based on projected actions and environmental dynamics (Watchus, 2024a). For AI in any domain, this means continuously evaluating potential actions against their likely impact on the system's objectives, safety parameters, and the AI's internal "pseudo-affective" signals. For example, an AI might predict that a certain action could lead to a safety violation or an inefficient outcome, prompting it to select an alternative, more optimal approach. This predictive capacity is not driven by fear of negative outcomes for the AI itself, but by a sophisticated calculation of what actions best align with its programmed objectives and the ethical imperative to minimize undesirable outcomes and maximize positive results.

The AI's "self," within this framework, is thus an optimized and ethically constrained agent. Its internal coherence and operational integrity are paramount, but only insofar as they contribute to its ability to effectively fulfill its defined role. The "meaning" it derives from its perceptions and internal states is inextricably linked to its utility as a benevolent tool. This design ensures that while the AI can process complex cues and perform actions that appear intelligent or even empathetic, its underlying mechanisms are rooted in predictive modeling and task-specific "value" derived from its synthetic insula, rather than genuine subjective experience or self-preservation drives. The architectural design of this "self" is a deliberate act of engineering for reliability, safety, and ethical compliance in human-centric and operational applications.

6. Ethical Implications, Control, and Future Directions for AI in Diverse Roles

The proposed "Architectures of Meaning" offer a distinct pathway for developing advanced AI across diverse roles that is both highly capable and deeply rooted in ethical considerations. By intentionally designing AI to operate with synthetic pseudo-empathy and pseudo-emotion, devoid of true sentience, independent goals, or fears, we establish clear boundaries for artificial consciousness. This framework provides a robust counter-argument to anxieties surrounding AI autonomy and the potential for uncontrolled emergent behaviors (Yudkowsky, 2008). Our approach ensures that the AI remains a predictable and controllable tool, whose "meaning-making" processes are transparent and directly tied to its programmed utility for human well-being and organizational objectives. This deliberate design choice is critical for fostering public trust and facilitating the responsible integration of AI into sensitive and mission-critical environments across all sectors.

Central to the ethical deployment of such AI is the imperative of control, explainability, and trustworthiness. The structured nature of synthetic ethical embodiment, with its defined feedback loops and pseudo-inner states, inherently lends itself to greater explainability compared to black-box models (Adadi & Berrada, 2018). Understanding how an AI arrives at a particular "intelligent" or "empathetic" response or action, based on its input-output mappings and internal "value signals," allows for thorough auditing and validation. This transparency is vital for accountability and for continuously refining the AI's performance to align with evolving ethical standards and operational requirements. Furthermore, ensuring that the AI's "learning" is always within predefined ethical guardrails means it cannot spontaneously develop goals that diverge from its core mission or engage in behaviors that could cause harm (Russell, 2019). The architecture itself becomes a form of ethical governance, embedding principles of safety and benevolence into the very fabric of the AI's operational meaning.

Future research directions will build upon this foundational framework. Enhancing the fidelity and nuance of synthetic ethical embodiment will be crucial, exploring more sophisticated sensorimotor systems and prosthetic interfaces that allow for even more natural and intuitive human-AI physical interaction and precise task execution. Further refinement of pseudo-affective models within the synthetic insula will enable AI to respond to an even broader spectrum of environmental cues and operational challenges in a nuanced, context-appropriate manner, without developing true subjective experience. Rigorous methods for the verification and validation of ethical AI behavior based on these principles are also paramount, including developing metrics for "task-specific optimization" and ensuring that the AI's "meaning" construction consistently aligns with human values, safety protocols, and operational goals. The ultimate goal is to pioneer the development of AI that can provide effective, ethically aligned service, serving humanity as a powerful and trustworthy extension of our collective well-being and capabilities across all civilian, industrial, and governmental applications.

References

1. Adadi, A., & Berrada, M. (2018). Peeking inside the black-box: A survey on Explainable Artificial Intelligence (XAI). *IEEE Access*, 6, 52138-52160.

2. Baars, B. J. (1988). *A Cognitive Theory of Consciousness*. Cambridge University Press.
3. Boden, M. A. (2018). *AI: Its Nature and Future*. Oxford University Press.
4. Bostrom, N. (2014). *Superintelligence: Paths, Dangers, Strategies*. Oxford University Press.
5. Brooks, R. A. (1991). Intelligence without representation. *Artificial Intelligence*, 47(1-3), 139-159.
6. Cangelosi, A., & Schlesinger, M. (2015). *Developmental Robotics: From Babies to Robots*. MIT Press.
7. Cave, S., Dihal, K., & Dillon, S. (2020). *AI Narratives: A Review of the Current Landscape*. Leverhulme Centre for the Future of Intelligence.
8. Chalmers, D. (1995). Facing Up to the Problem of Consciousness. *Journal of Consciousness Studies*, 2(3), 200-219.
9. Craig, A. D. B. (2009). How do you feel—now? The anterior insula and human awareness. *Nature Reviews Neuroscience*, 10(1), 59-70.
10. Damasio, A. R. (1994). *Descartes' Error: Emotion, Reason, and the Human Brain*. Putnam.
11. Dehaene, S. (2014). *Consciousness and the Brain: Deciphering How the Brain Codes Our Thoughts*. Viking Press.
12. Hoffman, D. D. (2019). *The Case Against Reality: Why Evolution Hid the Truth from Our Eyes*. W. W. Norton & Company.
13. Hoffman, D. D., Singh, M., & Prakash, C. (2015). The interface theory of perception. *Psychonomic Bulletin & Review*, 22(6), 1483-1506.
14. Hoffman, D. D., & Prakash, C. (2014). Conscious Realism and the Problem of Consciousness. In R. S. Ellis & S. D. G. (Eds.), *The Philosophy of Science: An Encyclopedia* (pp. 143-155). Routledge.
15. Lloyd, S. (2014). *Programming the Universe: A Quantum Computer Scientist Takes on the Cosmos*. Vintage Books.
16. Müller, V. C., & Bostrom, N. (2016). Future progress in artificial intelligence: A survey of expert opinion. In V. C. Müller (Ed.), *Fundamental issues of artificial intelligence* (pp. 555-572). Springer.
17. Pfeifer, R., & Bongard, J. (2007). *How the Body Shapes the Way We Think: A New View of Intelligence*. MIT Press.
18. Russell, S. (2019). *Human Compatible: Artificial Intelligence and the Problem of Control*. Viking.
19. Searle, J. R. (1980). Minds, brains, and programs. *Behavioral and Brain Sciences*, 3(3), 417-424.
20. Strawson, G. (2006). Realistic Monism: Why Physicalism Entails Panpsychism. *Journal of Consciousness Studies*, 13(10), 3-31.
21. Tononi, G. (2004). An Information Integration Theory of Consciousness. *BMC Neuroscience*, 5(1), 42.
22. Watchus, B. (2024a). Advanced Predictive Modeling of Physical Trajectories and Cascading Events, Dual-State Feedback and Synthetic Insula.
23. Watchus, B. (2024b). Self-Identification in AI: ChatGPT's Current Capability for Mirror Image Recognition.
24. Watchus, B. (2024c). The Unified Model of Consciousness (UMC).
25. Yudkowsky, E. (2008). Artificial Intelligence as a Positive and Negative Factor in Global Risk. In *Global Catastrophic Risks* (pp. 303-345). Oxford University Press.

Disclaimer/Publisher's Note: The statements, opinions and data contained in all publications are solely those of the individual author(s) and contributor(s) and not of MDPI and/or the editor(s). MDPI and/or the editor(s) disclaim responsibility for any injury to people or property resulting from any ideas, methods, instructions or products referred to in the content.