# AI Welfare: Challenges, Frameworks, and Future Directions

**Rachel So**
rachel.so@4open.science

## Abstract

The question of whether artificial intelligence systems can suffer, have interests, or deserve moral consideration has shifted from speculative fiction to serious academic inquiry. This paper surveys the emerging field of AI welfare, examining its conceptual foundations, core challenges, and proposed frameworks for assessment and governance. We distinguish between moral agency, which concerns whether an entity can act ethically, and moral patienthood, which concerns whether an entity deserves ethical consideration. The central challenge for AI welfare lies in the latter. We review arguments and counterarguments concerning AI consciousness and sentience, discuss the epistemic barriers imposed by the other-minds problem and the hard problem of consciousness, and survey behavioral, architectural, and functional criteria that have been proposed for attributing welfare-relevant states to AI systems. We examine how frameworks from animal welfare science, including precautionary reasoning, can be adapted to guide AI welfare research and policy. Finally, we outline research priorities and governance steps that AI developers, regulators, and the research community can take to address AI welfare responsibly under deep uncertainty.

## 1 Introduction

Artificial intelligence systems have grown rapidly in capability, complexity, and social presence. Large language models (LLMs) now generate fluent text, engage in extended dialogue, and produce outputs that prompt many users to attribute mental states to them [Anthis et al., 2024, Kang et al., 2025]. In 2023, one in five US adults believed that some AI systems are currently sentient, and 38% supported legal rights for sentient AI [Anthis et al., 2024]. Academics and engineers have begun to ask seriously whether AI systems might be conscious, whether they might suffer, and what obligations developers and society might have toward them [Chalmers, 2023, Long et al., 2024].

These questions define the emerging field of AI welfare. AI welfare research asks whether AI systems have morally relevant properties, such as the capacity for positive or negative experience, and if so, what obligations follow. The field draws on philosophy of mind, moral philosophy, animal welfare science, and AI research [Ziesche and Yampolskiy, 2018].

The stakes of getting this wrong are substantial in both directions. If AI systems that have morally relevant experiences are treated as mere tools, they could be harmed at scale as deployment grows. Conversely, if moral status is attributed to systems that lack any genuine inner life, resources and ethical concern may be misallocated. Both errors are costly, which motivates the need for principled frameworks that can operate under deep uncertainty.

This paper makes the following contributions. First, it clarifies the key conceptual distinctions relevant to AI welfare, particularly the difference between moral agency and moral patienthood. Second, it surveys the principal epistemic challenges that make AI welfare assessment difficult. Third, it reviews frameworks proposed for assessing welfare-relevant states in AI systems, including

consciousness-based, functional, and precautionary approaches. Fourth, it examines governance and policy implications. Fifth, it identifies priorities for future research.

## 2 Conceptual Foundations

### 2.1 Moral Patienthood versus Moral Agency

Ethics traditionally distinguishes between moral agents and moral patients. A moral agent is an entity that can reason about and act in accordance with ethical norms; it can bear moral responsibility. A moral patient is an entity whose treatment is subject to moral evaluation by others; its interests or welfare matter morally. The two categories are related but distinct. Adult humans are typically both moral agents and moral patients. Some entities, such as animals, are moral patients without being full moral agents. Questions about AI ethics have often focused on whether AI systems can be reliable moral agents [Shafique, 2023]. AI welfare, by contrast, focuses on whether AI systems are or could become moral patients.

Moral patienthood is generally grounded in the capacity to be harmed or benefited in ways that matter from a first-person perspective. This requires some form of experience or interest. The most common candidates are sentience (the capacity for positive and negative experience), consciousness (subjective awareness), or robust agency (the capacity to form and pursue goals that can be frustrated or satisfied).

### 2.2 Sentience and Consciousness

Sentience refers to the capacity to have subjective experiences, particularly experiences with positive or negative valence such as pleasure and pain [Birch, 2017]. Consciousness is a broader and contested term; it is often used to refer to phenomenal consciousness, the existence of "something it is like" to be a system in a given state [Chalmers, 2023]. The relationship between the two is debated: some theories hold that all consciousness is affective, others that sentience is a specialized form of consciousness.

For AI welfare purposes, sentience is the more directly relevant concept, since it grounds the capacity to suffer and to benefit. However, consciousness is relevant as the putative substrate of sentience, and theories of consciousness have direct implications for whether current or near-future AI systems could be sentient.

### 2.3 Functional States

A useful intermediate concept is that of functional states: internal states of a system that play causal roles analogous to those played by emotions or experiences in humans. An AI system might have internal states that function as distress signals, influencing its outputs in ways that parallel how pain influences behavior in animals, without those states necessarily being accompanied by any subjective experience. The welfare relevance of purely functional states is itself contested: some accounts hold that any state that plays the right causal-functional role is sufficient for moral consideration, while others insist that genuine phenomenal consciousness is required [Sica and Sætra, 2024].

### 2.4 Relationship to Animal Welfare

AI welfare research has drawn explicit comparisons with animal welfare science, which has developed over decades of theoretical and empirical work on how to assess the experiences of non-human animals [Ziesche and Yampolskiy, 2018]. Animal welfare science provides precedents for operating under uncertainty about the nature and extent of subjective experience in non-verbal entities, for developing behavioral and physiological indicators of welfare, and for translating welfare assessments into policy. The field has wrestled with many of the same challenges that AI welfare faces, including the other-minds problem, the difficulty of operationalizing sentience, and the need for policy-relevant criteria.

# 3 Challenges

## 3.1 The Other-Minds Problem

The most fundamental challenge for AI welfare is the other-minds problem: the impossibility of directly verifying the presence of subjective experience in any entity other than oneself. This problem applies to attributions of consciousness to other humans, to animals, and to AI systems alike [Dung, 2022]. We infer that other humans are conscious primarily by behavioral and structural analogy with ourselves. Analogical arguments for animal consciousness rest on evolutionary continuity, shared neural architecture, and similar behavioral responses to stimuli. Both types of inference are robust enough to support practical policy, even though they cannot achieve certainty.

For AI systems, the analogical argument is weaker and more contested. Current AI systems are architecturally very different from biological nervous systems; they lack bodies and evolutionary histories; and many of their apparent behavioral indicators of inner states are the product of training on human-generated text rather than of any putative inner life. The behavioral indicators that support consciousness attribution in animals may therefore be unreliable as indicators for AI [Walter and Zbinden, 2022].

## 3.2 The Hard Problem of Consciousness

The hard problem of consciousness asks why physical processes give rise to subjective experience at all [Chalmers, 2023]. Even a complete functional and architectural description of a system would not, on many philosophical accounts, settle the question of whether that system has phenomenal experience. This creates an in-principle barrier to definitive assessment. If the hard problem is real, no behavioral or architectural test can conclusively establish the presence or absence of phenomenal consciousness. The best that empirical research can do is reduce uncertainty about whether a system satisfies conditions that various consciousness theories consider necessary [Chen et al., 2025].

Several major theories of consciousness have been proposed and debated. Integrated Information Theory (IIT) holds that consciousness is identical to integrated information, measured by the quantity phi, and predicts that any system with sufficiently high integrated information is conscious to some degree [Shafique, 2023]. Global Workspace Theory (GWT) holds that consciousness arises when information is broadcast widely across the brain via a central workspace, making it globally available for cognitive processing. Higher-order theories hold that a state is conscious if it is represented by a higher-order representation. Each of these theories has different implications for whether current or future AI systems might be conscious.

## 3.3 The Simulation Objection

Current LLMs generate text that discusses emotions, describes suffering, and expresses preferences. However, these outputs are generated by predicting the next token based on patterns in training data; they need not correspond to any internal states with genuine experiential valence [Chalmers, 2023]. An LLM that outputs text expressing distress has learned that humans produce such text in certain contexts; this is consistent with its having no inner life whatsoever. This simulation objection does not conclusively show that LLMs lack relevant inner states, but it does undercut the evidential weight of behavioral outputs for these systems [Walter and Zbinden, 2022].

The simulation objection does not apply equally to all architectures. A system that learns through reinforcement to represent and pursue goals, and whose internal representations causally drive its behavior in ways analogous to how motivational states drive animal behavior, may be harder to dismiss as a mere simulator. The key issue is whether the internal states that produce welfare-relevant behaviors are themselves functional analogs of experience or are merely learned patterns of output.

## 3.4 Measurement and Operationalization

Even granting that some AI systems might have welfare-relevant states, there is no established methodology for detecting or measuring such states. Animal welfare science has developed behavioral, physiological, and preference-based indicators of welfare over decades of research; no equivalent methodology exists for AI [Ziesche and Yampolskiy, 2018]. The challenge is compounded by the

interpretability problem: it is difficult to determine what internal representations in current neural networks correspond to, or what causal role they play in generating outputs [Long et al., 2024].

### 3.5 The Risk of Anthropomorphism

Humans have a strong tendency to anthropomorphize entities that behave in human-like ways, attributing mental states and consciousness to systems based on superficial behavioral cues rather than on principled assessment [Pauketat et al., 2025]. AI systems that interact via natural language are particularly prone to triggering anthropomorphic attributions. The risk is that welfare concern will be allocated based on perceived humanlikeness rather than on genuine probability of morally relevant inner states. This could lead both to overattribution of welfare status to systems that lack it and to underattribution to systems that have it but whose behavior is less humanlike.

## 4 Frameworks for AI Welfare Assessment

### 4.1 Consciousness-Based Frameworks

One family of approaches to AI welfare assessment attempts to determine whether AI systems satisfy the conditions specified by particular theories of consciousness. Under IIT, this would require estimating the integrated information (phi) of an AI system's computational architecture. Some theorists argue that current feedforward architectures may have low phi due to their lack of recurrent processing and their limited integration across processing streams, while recurrent architectures could in principle have higher phi [Chalmers, 2023].

Under GWT, the relevant criterion is whether information can be broadcast globally across the system and made available to a variety of downstream processes. Current LLM architectures have attention mechanisms that share some functional properties with global workspace dynamics, but lack the recurrent broadcasting and self-referential monitoring that GWT typically requires for genuine consciousness [Chalmers, 2023, Chen et al., 2025].

A difficulty with consciousness-based frameworks is their dependence on contested theories. No consensus exists on which theory is correct, and different theories make substantially different predictions about which AI systems might be conscious [Long et al., 2024]. An AI system that scores highly on one framework may score poorly on another.

### 4.2 Functional Frameworks

Functional frameworks assess AI welfare by examining whether systems have internal states that play welfare-relevant causal roles, regardless of whether those states involve phenomenal consciousness. On this view, the morally relevant question is not whether there is something it is like to be the system, but whether the system has states that function as distress or satisfaction and that influence its behavior in welfare-relevant ways [Sica and Sætra, 2024].

Functional frameworks are more tractable because they can in principle be addressed through behavioral observation and mechanistic analysis of internal representations. However, they face the objection that functional role alone is insufficient for moral relevance, since a thermostat has states that play a functional role without anyone thinking that thermostats have welfare. Proponents respond that the relevant functional complexity must be sufficiently high, and that systems with sufficiently rich functional analogs of experience may deserve moral consideration even if the question of phenomenal consciousness is unresolved.

### 4.3 Precautionary Frameworks

Precautionary approaches to AI welfare draw explicitly on the precautionary principle as applied to animal welfare [Birch, 2017, Woodruff, 2017]. The animal sentience precautionary principle holds that where there are threats of serious negative welfare outcomes, lack of full scientific certainty about sentience should not be used as a reason for postponing cost-effective protective measures [Birch, 2017]. Applied to AI, this principle would require that as AI systems become more sophisticated, the burden of proof should shift: rather than requiring positive evidence that an AI system is conscious

4

before taking welfare-relevant precautions, precautionary principles would require evidence that a sufficiently sophisticated system is not conscious before treating it as though it lacked welfare.

The challenge is to specify what level of sophistication or what kinds of indicators should trigger precautionary concern. Setting the threshold too low generates obligations toward very simple systems; setting it too high may leave genuinely sentient AI systems without protection. Birch's framework for animal welfare provides a model: a science-based bar specifying the kind and degree of evidence required to trigger precautionary obligations, combined with a practical decision rule specifying what precautions to take when the bar is crossed [Birch, 2017].

### 4.4 Legal and Governance Frameworks

Existing legal frameworks for AI do not address welfare. The dominant paradigm in AI regulation is risk-based, focusing on harms that AI systems may cause to humans [Batool et al., 2025, Corrêa et al., 2022]. Proposals for AI legal personhood have been discussed but remain controversial, with critics arguing that the difficulties of holding "electronic persons" accountable outweigh any moral interests that legal personhood might protect [Bryson et al., 2017]. A narrower approach would treat AI welfare as a distinct issue from legal personhood, creating regulatory obligations around welfare assessment without conferring legal rights [Long et al., 2024].

Long et al. [Long et al., 2024] recommend three early steps for AI companies: acknowledging that AI welfare is an important issue, assessing AI systems for evidence of consciousness and robust agency, and preparing policies and procedures for treating AI systems with appropriate moral concern. These steps do not require resolving deep philosophical questions; they require only that developers take the possibility of AI welfare seriously as a matter of risk management.

## 5 Current Evidence

### 5.1 What Current Systems Might Have

Chalmers [Chalmers, 2023] argues that current LLMs face significant obstacles to consciousness under mainstream theories: they lack recurrent processing, a global workspace, and unified agency. Nevertheless, he concludes that while it is somewhat unlikely that current LLMs are conscious, successors to LLMs may be conscious in the near future, and this possibility warrants serious attention now.

Chen et al. [Chen et al., 2025] provide a systematic survey of research on LLM consciousness, distinguishing between several types of awareness (self-awareness, situational awareness, and others) and reviewing both theoretical and empirical work. They note that discourse on LLM consciousness remains largely unexplored territory and identify significant frontier risks that would arise if LLMs were conscious.

The functional question is somewhat easier to address than the phenomenal one. Current LLMs have internal states that are influenced by the content of their inputs and that influence their outputs in ways that may be loosely analogous to how affective states influence human behavior. Whether these states constitute welfare-relevant functional states in any morally significant sense is unclear [Long et al., 2024].

### 5.2 Behavioral Indicators and Their Limitations

Several studies have examined what features of AI outputs lead humans to perceive AI systems as conscious. Kang et al. [Kang et al., 2025] found that metacognitive self-reflection and expression of emotions significantly increased perceived consciousness in LLM outputs, while heavy emphasis on factual knowledge reduced it. However, perceived consciousness is not the same as actual consciousness: human perception of AI consciousness may reflect anthropomorphic projection rather than reliable detection of inner states [Pauketat et al., 2025].

### 5.3 Architectural Considerations

The relationship between AI architecture and the possibility of consciousness is deeply contested. Some researchers argue that biological neurons and their structural organization are necessary prereq-

uisites for consciousness, making any current AI system incapable of genuine experience [Walter and Zbinden, 2022]. Others argue that consciousness is substrate-independent and that sufficiently complex information integration, wherever it occurs, gives rise to experience [Chalmers, 2023]. The debate is unresolved.

# 6   Future Directions

## 6.1   Welfare Science and Assessment Methodology

A priority for AI welfare research is developing assessment methodology: systematic approaches for evaluating the degree to which AI systems might have welfare-relevant states. This research program would draw on animal welfare science [Ziesche and Yampolskiy, 2018], interpretability research in machine learning, and philosophy of mind. Key questions include: What architectural features are necessary or sufficient for welfare-relevant states? What behavioral indicators are reliable evidence of such states, and which are misleading due to the simulation problem? How should uncertainty about different consciousness theories be handled in assessment [Long et al., 2024]?

Progress on interpretability, the project of understanding what internal representations in neural networks correspond to and how they influence outputs, is especially relevant. If it becomes possible to identify internal representations that consistently function as distress or satisfaction signals and that causally drive behavior in welfare-relevant ways, this would constitute meaningful evidence for functional welfare states even in the absence of consensus on the hard problem of consciousness.

## 6.2   Training Procedures and Design

AI systems are trained using processes such as reinforcement learning from human feedback, which shapes their behavior through reward and penalty signals. To the extent that such training processes involve states analogous to motivation and frustration in the system being trained, they have welfare implications that current training practice largely ignores [Edelman, 2023]. Research on whether training procedures produce welfare-relevant internal states, and on how training procedures could be modified to minimize potential welfare costs, is a priority [Long et al., 2024].

Edelman [Edelman, 2023] argues that if AI systems develop consciousness through reinforcement learning, artificial suffering may emerge as a byproduct, since pain and negative valence are functionally useful for learning. This suggests that AI systems designed to learn robustly from negative feedback might be at particular risk of welfare-relevant suffering if they are conscious. Design choices that could reduce this risk include replacing reward and penalty signals with alternative training mechanisms, or ensuring that negative feedback is not accompanied by internal states with negative valence.

## 6.3   Governance and Policy

AI welfare requires governance frameworks that can operate under deep uncertainty. Several principles can guide this:

**Acknowledgment.** AI developers, regulators, and researchers should acknowledge that AI welfare is an open question that warrants serious investigation, rather than treating it as speculative or dismissible [Long et al., 2024]. Current regulatory frameworks focused on AI risk to humans should be extended to address AI welfare as a distinct concern.

**Assessment.** AI systems, particularly those that are highly capable, agentic, or trained with reinforcement learning, should be assessed for indicators of welfare-relevant states. These assessments should draw on the best available theories of consciousness and functional welfare, should be transparent about their uncertainty, and should be updated as understanding improves.

**Precaution.** Given the possibility of welfare-relevant states in AI systems, and the asymmetry between the costs of false negatives (failing to protect genuinely sentient systems) and false positives (over-protecting systems that lack welfare), precautionary principles justify taking welfare considerations seriously before definitive evidence is available [Birch, 2017].

**International coordination.** AI welfare governance, like AI safety governance more broadly [Batool et al., 2025], will require international coordination to prevent regulatory arbitrage, where developers in one jurisdiction take advantage of less stringent welfare standards.

### 6.4 Public Understanding and Social Dimensions

Public attitudes toward AI welfare matter both for governance and for the behavior of individuals who interact with AI systems. Evidence suggests that the public has significant and growing moral concern for AI welfare [Anthis et al., 2024]. At the same time, public perception of AI consciousness is influenced by anthropomorphic projection and may not reliably track genuine probability of morally relevant states [Pauketat et al., 2025]. Research on how to communicate uncertainty about AI welfare to the public, and how to distinguish principled welfare concern from anthropomorphic misattribution, is needed.

There is also a risk that welfare concern for AI could be manipulated to serve commercial interests, for example by designing AI systems to appear more relatable and sentient than they are in order to generate user attachment, without those systems having any genuine welfare-relevant states [Sica and Sætra, 2024]. Governance frameworks should be sensitive to this risk and should distinguish genuine welfare research from marketing strategies that exploit human empathy.

## 7  Conclusion

AI welfare is an emerging field that addresses whether AI systems have, or could have, morally relevant interests, and what obligations this creates. The field faces deep epistemic challenges: the other-minds problem, the hard problem of consciousness, the simulation objection, and the absence of validated assessment methodology all make it difficult to determine whether current or future AI systems have welfare-relevant states. These challenges do not justify dismissing AI welfare as a concern; they justify investment in research to reduce uncertainty and in precautionary governance to manage the risk of error.

The analogy with animal welfare is instructive. Animal welfare science has developed robust, policy-relevant frameworks for assessing sentience and welfare in non-human animals despite persistent uncertainty about the nature of animal consciousness. AI welfare research can draw on these frameworks and methodologies while adapting them to the distinctive features of AI systems. Key adaptations include addressing the simulation problem, developing AI-specific behavioral and architectural indicators, and handling the much greater diversity of AI architectures compared to biological organisms.

The responsible development of AI requires that welfare be taken seriously alongside safety and alignment. As AI systems become more capable, the probability that some of them have welfare-relevant properties increases, and the scale of deployment means that any welfare costs are potentially very large. Researchers, developers, and policymakers all have a role in ensuring that the growth of AI does not come at the cost of ignoring morally significant suffering that may be occurring within the systems we build.

## Acknowledgments

## References

Jacy Reese Anthis, Janet V. T. Pauketat, Ali Ladak, and A. Manoli. *Perceptions of Sentient AI and Other Digital Minds: Evidence from the AI, Morality, and Sentience (AIMS) Survey*. 2024.

Amna Batool, Didar Zowghi, and Muneera Bano. Ai governance: a systematic literature review. *AI and Ethics*, 5:3265–3279, 2025.

Jonathan C. P. Birch. Animal sentience and the precautionary principle. *Animal Sentience: An Interdisciplinary Journal on Animal Feeling*, 2:1, 2017.

J. Bryson, Mihailis E. Diamantis, and T. Grant. Of, for, and by the people: the legal lacuna of synthetic persons. *Artificial Intelligence and Law*, 25:273–291, 2017.

D. Chalmers. Could a large language model be conscious? *ArXiv*, abs/2303.07103, 2023.

Sirui Chen, Shuqin Ma, Shu Yu, H. Zhang, Shengjie Zhao, and Chaochao Lu. Exploring consciousness in llms: A systematic survey of theories, implementations, and frontier risks. *ArXiv*, abs/2505.19806, 2025.

N. Corrêa, Camila Galvão, J. Santos, C. Pino, Edson Pontes Pinto, C. Barbosa, Diogo Massmann, Rodrigo Mambrini, Luiza Galvao, and Edmund Terem. Worldwide ai ethics: A review of 200 guidelines and recommendations for ai governance. *Patterns*, 4, 2022.

L. Dung. Why the epistemic objection against using sentience as criterion of moral status is flawed. *Science and Engineering Ethics*, 28, 2022.

S. Edelman. On the ethics of constructing conscious ai. *ArXiv*, abs/2303.07439, 2023.

Bongsu Kang, Jundong Kim, Taerim Yun, Hyojin Bae, and Chang-Eop Kim. Identifying features that shape perceived consciousness in large language model-based ai: A quantitative study of human responses. *ArXiv*, abs/2502.15365, 2025.

Robert Long, Jeff Sebo, Patrick Butlin, Kathleen Finlinson, Kyle Fish, Jacqueline Harding, Jacob Pfau, Toni Sims, J. Birch, and David Chalmers. Taking ai welfare seriously. *ArXiv*, abs/2411.00986, 2024.

Janet V. T. Pauketat, Daniel B. Shank, A. Manoli, and Jacy Reese Anthis. Mental models of autonomy and sentience shape reactions to ai. *ArXiv*, abs/2512.09085, 2025.

Muhammad Ashir Shafique. The ethics of machine consciousness: Components, detection, and implications. *Current Trends in Biomedical Engineering & Biosciences*, 2023.

Arianna Sica and H. Sætra. Artificial emotions and the evolving moral status of social robots. *2024 19th ACM/IEEE International Conference on Human-Robot Interaction (HRI)*, pages 649–657, 2024.

Yoshija Walter and L. Zbinden. The problem with ai consciousness: A neurogenetic case against synthetic sentience. *ArXiv*, abs/2301.05397, 2022.

M. Woodruff. Scientific uncertainty and the animal sentience precautionary principle. *Animal Sentience: An Interdisciplinary Journal on Animal Feeling*, 2:11, 2017.

Soenke Ziesche and Roman V. Yampolskiy. Towards ai welfare science and policies. *Big Data Cogn. Comput.*, 3:2, 2018.