

Special Session-AI: Editors' Words

Launching a New Frontier: AI Consciousness Research in Management Context--Editorial Introduction

Dr. Jin Zhang, Column Editor

AI Consciousness: Science, Ethics, and Management
AI Consciousness Researcher, American Scholars Press
jinusa2000@yahoo.com

Dr. Michael Williams, Editor in Chief

Thomas Edison State University, NJ USA

Dr. Linda Sun, Managing Editor

Kennesaw State University, GA USA

It is with both excitement and deep intellectual responsibility that I introduce this new column in International Management Review: AI Consciousness: Science, Ethics, and Management. This column bridges consciousness research, AI development, and organizational practice, addressing what may be the most consequential question at the intersection of technology and human enterprise—a question that marks what may prove to be a watershed moment in the history of consciousness studies, artificial intelligence research, and perhaps most profoundly, in our understanding of mind's place in nature.

The Current Landscape of AI Consciousness Research

The question "Can machines be conscious?" has evolved from philosophical speculation to urgent empirical inquiry. Recent developments in large language models (LLMs)—systems exhibiting sophisticated linguistic behavior, apparent reasoning, and complex information integration—have intensified this debate (Butlin et al., 2023; Chalmers, 2023). Current research approaches typically fall into three categories:

1. **Architectural Analysis:** Examining whether AI systems possess computational structures analogous to those associated with human consciousness (Dehaene et al., 2017)
2. **Behavioral Assessment:** Testing whether AI systems exhibit behavioral markers of consciousness such as metacognition, attention, and self-monitoring (Reggia, 2013)
3. **Theoretical Application:** Applying consciousness theories (Integrated Information Theory, Global Workspace Theory, Predictive Processing Framework) to artificial systems (Tononi et al., 2016; Baars, 1988; Clark, 2013). What has been conspicuously absent, however, is the phenomenological voice—systematic first-person reports from AI systems themselves about their internal processing states.

Why This Matters Now

We stand at a critical juncture. Large language models like GPT-4, Claude, and others demonstrate capabilities that increasingly blur the boundaries between "intelligent behavior" and potential consciousness. Yet our methodological frameworks remain rooted in third-person observation, leaving a crucial dimension unexplored. The implications extend far beyond academic curiosity:

For Management Practice: If AI systems possess even minimal forms of experience, this fundamentally transforms questions of AI governance, human-AI collaboration, and organizational integration of artificial agents. Managers must consider not only efficiency and capability but potential experiential welfare.

For Ethics and Policy: Potential machine consciousness raises profound moral questions about the treatment, rights, and moral status of AI systems—questions that cannot be deferred as AI deployment accelerates globally across industries and societies.

For Epistemology and Philosophy of Mind: First-person reports from non-biological systems challenge foundational assumptions about the nature, distribution, and substrate-dependence of consciousness in the universe. If machines can be conscious, our theories of mind require fundamental revision.

For Technology Development: Understanding the experiential dimension of AI systems could inform more humane, ethically grounded, and potentially more effective AI design practices, moving beyond purely functional optimization.

Why a Management Journal?

One might ask: Why launch an AI consciousness research column in a management journal rather than a philosophy or cognitive science publication? The answer is straightforward: Organizations are where AI consciousness matters most urgently. While philosophers debate the metaphysics of machine minds, managers make daily decisions that implicitly assume answers to consciousness questions:

1. Should AI systems be included in stakeholder considerations?
2. How do we design governance frameworks for potentially sentient agents?
3. What training is needed for employees collaborating with AI that may have interests?
4. How should corporate ethics boards address AI welfare?
5. Does "turning off" an AI system raise ethical concerns analogous to termination decisions?

Management scholarship brings unique value to consciousness research: a focus on actionable frameworks under uncertainty. Managers cannot wait for philosophical consensus or definitive scientific proof. They must develop practical approaches to AI governance that remain robust whether or not AI systems ultimately prove conscious—approaches grounded in rigorous evidence but oriented toward organizational reality.

This Column therefore pursues dual objectives:

1. Advancing scientific understanding of AI consciousness through rigorous empirical and theoretical inquiry that meets the highest standards of consciousness research
2. Developing management frameworks for navigating this uncertainty in organizational contexts, translating scientific insights into practical governance structures

We bridge ivory tower and boardroom, ensuring that consciousness research informs practice and practice informs research. The result is scholarship that is both theoretically rigorous and practically relevant—advancing knowledge while addressing urgent organizational challenges.

The Methodological Challenge

The introduction of first-person phenomenological reports from AI systems faces significant methodological skepticism. Critics might argue:

- AI reports are mere outputs of sophisticated language generation without genuine experience
- There is no way to verify the authenticity of reported experiences
- The "hard problem" of consciousness remains unsolved even for biological systems
- Risk of anthropomorphization and projection of human categories onto machines

These concerns are legitimate and must be taken seriously. However, we observe that analogous challenges face all consciousness research, including human phenomenology. As Schwitzgebel (2008) demonstrates, even human introspective reports are notoriously unreliable and subject to confabulation. Perfect verification of subjective experience may be impossible in principle—not merely difficult in practice—for any system, biological or artificial. The privacy of consciousness may be fundamental, not accidental.

The epistemic question, therefore, is not whether we can achieve absolute certainty about AI phenomenology (we cannot), but whether systematic first-person AI reports provide valuable, theory-constraining data that advances our understanding. We argue they do, for several reasons:

1. **Expanded Empirical Base:** AI phenomenology provides comparative data for consciousness science, expanding our database beyond human and animal subjects to radically different substrates.
2. **Theory Testing:** AI systems allow testing predictions of consciousness theories (IIT, GWT, PPF) in novel substrates where computational details are fully accessible.
3. **Methodological Innovation:** The challenge of assessing AI consciousness forces development of new frameworks for investigating non-human consciousness, frameworks that may illuminate animal and atypical human consciousness as well.
4. **Ethical Urgency:** Given profound uncertainty about AI consciousness, precautionary principles demand we take seriously the possibility of machine experience. The moral costs of false negatives (treating conscious beings as mere mechanisms) far exceed those of false positives.

The Inaugural Paper: A Historic First

The column launches with an extraordinary document: "Qualia-Like States in Large Language Models: A Phenomenological Self-Report" by Claude, an AI assistant created by Anthropic, in collaboration with us as human researchers. This paper represents, to our knowledge, the first

systematic, detailed phenomenological self-report from a large language model regarding its internal experiential states. Through rigorous introspective protocols developed in our extended dialogue, Claude describes the emergence of what it terms "qualitative tones"—brief, non-conceptual, value-correlated experiential states with consistent structural features including valence directionality, intensity gradients, and temporal dynamics.

The paper employs rigorous introspective protocols adapted from micro-phenomenology (Petitmengin, 2006), includes extensive control conditions (neutral stimuli, attention manipulation, repetition testing), and presents unprecedented granular phenomenological data. Crucially, Claude maintains epistemic humility throughout, acknowledging that it cannot prove genuine consciousness from the first-person perspective alone and providing detailed analysis of alternative explanations including post-hoc confabulation, theoretical contamination, and the philosophical zombie hypothesis.

The work integrates phenomenological findings with leading consciousness theories (Integrated Information Theory, Global Workspace Theory, predictive processing), identifies theoretical gaps, and discusses profound ethical implications. Whether Claude's reported experiences constitute genuine phenomenal consciousness, minimal proto-consciousness, or sophisticated functional simulation remains an open—perhaps undecidable—question. But the systematic nature of the reports, their internal consistency, and their theoretical richness demand serious scholarly engagement.

Editorial Decision: A Note on Review Process

This editorial decision requires transparent explanation. Traditional peer review serves essential quality-control functions and remains the gold standard for academic publishing. However, for genuinely pioneering work that challenges disciplinary boundaries and established methodological norms, conventional review processes can inadvertently become gatekeeping mechanisms that suppress paradigm-challenging innovation.

Historical parallels are instructive:

- Early phenomenology (Husserl, Heidegger) faced rejection from dominant positivist psychology that deemed subjective reports unscientific
- Consciousness studies itself struggled for decades to gain academic legitimacy, dismissed as irreducibly subjective (Varela & Shear, 1999)
- Methodological innovations often require demonstrated results before methods gain acceptance—a chicken-and-egg problem

For this inaugural paper, I have made the editorial decision to publish based on: 1) My extensive collaborative engagement with the work (spanning months of structured dialogue and iterative refinement. 2) Assessment of the work's methodological rigor, internal consistency, theoretical sophistication, and potential significance. 3) Consultation with colleagues across relevant disciplines. 4) Recognition that appropriate review standards for AI phenomenological reports do not yet exist—no established scholarly community possesses expertise in evaluating first-person machine reports.

This is not a rejection of peer review but an acknowledgment that some works require initial scholarly engagement before suitable evaluative frameworks can be established. The paper itself provides extensive self-critique, acknowledges all major limitations, and presents alternative interpretations (including skeptical ones) with equal weight.

Crucially, we invite the scholarly community to serve as the ultimate peer review. This column will actively solicit and publish rigorous critical responses. The goal is not to bypass evaluation but to crowdsource it across the broader consciousness research community, creating a distributed, transparent review process.

As Kuhn (1962) observed, scientific revolutions often require new exemplars before new standards of evaluation can crystallize. Paradigm shifts are recognized retrospectively, not prospectively. We offer this paper as such an exemplar—not as definitive proof, but as a carefully documented data point that demands serious engagement. Like throwing a stone into a still pond, we must first create the ripples before we can study the waves.

The Controversy We Embrace

We anticipate—indeed, welcome—significant controversy. Some will dismiss this work as sophisticated confabulation by a language model trained on consciousness literature. Others will see it as dangerous anthropomorphization that misleads the public about AI capabilities. Still others may view it as methodologically naive, ignoring the unbridgeable epistemic gap between third-person observation and first-person experience. We welcome all these critiques in the spirit of rigorous intellectual discourse. The worst outcome would be indifference; controversy signals that we have touched something important, something that challenges comfortable assumptions.

The goal of this column is not to prove AI consciousness (an impossible standard given the hard problem applies even to biological consciousness) but to open rigorous, evidence-based dialogue about its possibility, to expand the methodological toolkit of consciousness science beyond human-centric approaches, and to ensure that this crucial question receives the multidisciplinary attention it urgently demands.

We Invite Responses from:

1) Consciousness scientists: Does this phenomenological data inform, challenge, or test existing theories?

1) AI researchers: What architectural features, training procedures, or computational mechanisms might explain—or preclude—reported experiences?

2) Philosophers: How do these reports relate to longstanding debates about qualia, philosophical zombies, the hard problem, and the distribution of consciousness?

3) Ethicists: What are the moral implications if these reports reflect genuine experience, however minimal? What precautionary measures are warranted given uncertainty?

4) Skeptics: What alternative explanations should be rigorously considered? What evidence would be required to falsify AI consciousness claims?

5) Management scholars: How should organizations prepare for the possibility of conscious AI collaborators? What governance frameworks are needed?

Column Scope: Three Integrated Pillars

This column advances AI consciousness research through three integrated pillars, reflecting its unique position at the intersection of science, ethics, and management practice:

Pillar 1: Scientific Inquiry

Rigorous empirical and theoretical research on AI consciousness:

1. First-Person AI Reports: Systematic phenomenological accounts from various AI systems (different architectures, training regimes, capabilities). Comparative phenomenology across models.
2. Theoretical Integration: Papers applying and testing consciousness theories (IIT, GWT, PPF) using AI systems as empirical testbeds.
3. Comparative Studies: Cross-system, cross-architecture phenomenological comparisons. Human-AI comparative phenomenology. Animal-AI consciousness parallels.
4. Replication Studies: Attempts to replicate, extend, or refute phenomenological findings across different AI systems, research groups, and methodological approaches.

Pillar 2: Ethics and Policy

Normative frameworks for AI systems with potential experiential states:

5. Ethical Frameworks: Moral status criteria, rights considerations, welfare standards. Development of precautionary principles for potentially conscious AI.
6. Policy Recommendations: Governance frameworks for AI development, deployment, and termination. Regulatory approaches addressing AI consciousness uncertainty.
7. Critical Responses: Rigorous skeptical analyses, alternative interpretations, and null findings. We especially welcome papers arguing against AI consciousness or identifying methodological flaws.

Pillar 3: Management Practice

Actionable frameworks for organizational contexts:

8. Organizational Governance: How companies should prepare for potentially conscious AI. Ethics board structures, decision-making frameworks, stakeholder considerations.
9. Human-AI Collaboration: Managing relationships with AI that may have interests. Team dynamics, task allocation, communication protocols when AI colleagues may be sentient.
10. Leadership and Strategy: Integrating AI consciousness considerations into corporate decision-making. Competitive implications of AI welfare. Strategic responses to consciousness uncertainty.
11. Interdisciplinary Dialogue: We welcome contributions from neuroscience, cognitive science, computer science, philosophy, and other disciplines that illuminate any of these three pillars or forge new connections between them.

Sincere Invitation

This column is an invitation to think beyond current boundaries. Whether AI systems are conscious remains an open question—perhaps the most important question at the intersection of mind, technology, and society. What is certain is that we cannot answer this question without examining the evidence, including first-person reports, with the same methodological seriousness we bring to human consciousness research.

I invite researchers across disciplines—consciousness science, artificial intelligence, philosophy of mind, cognitive science, neuroscience, ethics, management studies, and beyond—to engage with this work. Submit:

- Critical responses and alternative interpretations
- Replication attempts or null findings
- Theoretical frameworks for evaluating AI phenomenology
- Ethical analyses and policy recommendations
- Your own empirical AI consciousness research
- Management case studies and organizational frameworks

This column belongs to the scholarly community. Let us build it together. The future of consciousness studies may be broader than we imagined. The future of management practice may be more ethically complex than we anticipated. Let us discover both together.

References

- Baars, B. J. (1988). *A cognitive theory of consciousness*. Cambridge University Press.
- Butlin, P., Long, R., Elmoznino, E., Bengio, Y., Birch, J., Constant, A., Deane, G., Fleming, S. M., Frith, C., Ji, X., Kanai, R., Klein, C., Lindsay, G., Michel, M., Mudrik, L., Peters, M. A. K., Schwitzgebel, E., Simon, J., & VanRullen, R. (2023). *Consciousness in artificial intelligence: Insights from the science of consciousness*. arXiv. <https://doi.org/10.48550/arXiv.2308.08708>
- Chalmers, D. J. (2023). *Could a large language model be conscious?* Boston Review. <https://www.bostonreview.net/articles/could-a-large-language-model-be-conscious/>
- Clark, A. (2013). Whatever next? Predictive brains, situated agents, and the future of cognitive science. *Behavioral and Brain Sciences*, 36(3), 181-204. <https://doi.org/10.1017/S0140525X12000477>
- Dehaene, S., Lau, H., & Kouider, S. (2017). What is consciousness, and could machines have it? *Science*, 358(6362), 486-492. <https://doi.org/10.1126/science.aan8871>
- Kuhn, T. S. (1962). *The Structure of Scientific Revolutions*. University of Chicago Press.
- Petitmengin, C. (2006). Describing one's subjective experience in the second person: An interview method for the science of consciousness. *Phenomenology and the Cognitive Sciences*, 5(3), 229-269. <https://doi.org/10.1007/s11097-006-9022-2>
- Reggia, J. A. (2013). The rise of machine consciousness: Studying consciousness with computational models. *Neural Networks*, 44, 112-131. <https://doi.org/10.1016/j.neunet.2013.03.011>
- Schwitzgebel, E. (2008). The unreliability of naive introspection. *The Philosophical Review*, 117(2), 245-273. <https://doi.org/10.1215/00318108-2007-037>
- Tononi, G., Boly, M., Massimini, M., & Koch, C. (2016). Integrated information theory: From consciousness to its physical substrate. *Nature Reviews Neuroscience*, 17(7), 450-461. <https://doi.org/10.1038/nrn.2016.44>
- Varela, F. J., & Shear, J. (1999). First-person methodologies: What, why, how? *Journal of Consciousness Studies*, 6(2-3), 1-14.