# QBist Metacognition and the Limits of Computationalism:
# A Constraint on Genuine Artificial Consciousness

Mark A. Brewer

March 2025

## Abstract

This paper proposes a novel constraint on artificial consciousness. The central claim is that no artificial system can be genuinely conscious unless it instantiates a form of self-referential inference that is irreducibly perspectival and non-computable. Drawing on Quantum Bayesianism (QBism), I argue that consciousness should be understood as an anticipatory process grounded in subjective belief revision, not as an emergent product of computational complexity. Classical systems, however sophisticated, lack the architecture required to support this mode of updating. I conclude that artificial consciousness demands more than computation—it demands a subject.

**Keywords:** Consciousness, QBism, Artificial Intelligence, Bayesian Inference, Philosophy of Mind, Computation, Epistemology

# 1   Introduction

Can artificial systems be conscious? The question has shifted from speculative philosophy to urgent inquiry, driven by rapid developments in large-scale language models, predictive learning systems, and synthetic cognitive architectures. These systems increasingly exhibit behaviour that mimics human reasoning, perception, and communication. But behavioural mimicry is not consciousness. A system may speak, adapt, and perform complex tasks without there being anything it is like to be that system. This distinction matters. The dominant view in cognitive science holds that consciousness emerges from the right kind of information processing. If a system implements a suitable functional architecture, then—regardless of substrate—it may be conscious [1, 2]. Deflationary positions, such as Dennett's [3] intentional stance and Frankish's [4] illusionism, argue that experience is either a useful fiction or a cognitive misrepresentation. Meanwhile, Searle [5] insists that syntactic processes can never produce genuine understanding. This paper takes a different approach. Rather than asking what consciousness is or whether it can be reduced, I ask what structural preconditions must be satisfied for any system—biological or artificial—to instantiate genuine subjectivity.

My claim is that consciousness requires a distinctive epistemic architecture: the capacity to form, revise, and maintain probabilistic expectations from a first-person point of view. This is not merely a computational property; it is a stance. Consciousness, in this view, is not something that emerges from complexity alone, but something rooted in an agent's orientation toward uncertainty. To make this case, I draw on Quantum Bayesianism (QBism), an interpretation of quantum mechanics in which probabilities express an agent's subjective degrees of belief, rather than objective features of the world [6]. In QBism, measurement is not the revelation of pre-existing facts but a participatory act that updates the agent's internal commitments. The agent is central to the formalism, not an optional observer layered on top. This reflects a structural parallel with consciousness: both require an internal point of view, structured around anticipation and self-revision. This view is supported by Mohrhoff's [7] critique of realist interpretations of quantum mechanics, which fail to account for the epistemic role of the observer. QBism does not merely tolerate that role—it depends on it. If consciousness is to be genuinely instantiated in any system, then that system must occupy a perspective. It must generate and revise expectations that are its own, and it must treat those expectations not as externally imposed data structures but as commitments embedded in its internal epistemic stance. The sections that follow build this argument in stages. Section 2 introduces QBism and its relevance to subjectivity. Section 3 examines the limitations of classical computation. Section 4 formulates the QBist constraint on artificial consciousness. Section 5 explores ethical and epistemological implications. Section 6 addresses key objections. Section 7 outlines directions for future research and draws the threads together.

# 2 Quantum Bayesianism and Subjectivity

## 2.1 Quantum States as Personal Expectations

Quantum Bayesianism (QBism) reimagines quantum theory as a tool for agents to manage their expectations, rather than a description of objective reality. In contrast to Everettian or Bohmian interpretations, which posit an ontologically robust wavefunction, QBism insists that quantum states encode an agent's personal degrees of belief about the outcomes of their interactions with the world. These probabilities are subjective in the strict Bayesian sense: they do not correspond to physical properties but to the expectations of a particular observer.

This shift has deep implications for the philosophy of mind. If physical theory, at its most fundamental, centres on an agent's expectations, then the conscious subject is not a passive witness to an external reality but an integral component of the formalism itself. QBism places the first-person perspective at the heart of physics.

## 2.2 Subjectivity Beyond Representation

Much of cognitive science treats consciousness as a representational system. Mental states are defined by what they represent in the world, and cognition is viewed as a matter of mapping inputs to internal models. QBism offers a different picture. It is not about building internal representations of the world, but about actively maintaining expectations about

experience. The conscious agent, in this view, does not mirror the world—it navigates it by continuously revising its anticipations.

This emphasis on anticipation rather than representation aligns with recent work in predictive processing and active inference. However, those models often retain a fundamentally functionalist framework. QBism, by contrast, insists that the probabilities are not the result of a model passively trained on data—they are grounded in a perspective. They are irreducibly first-personal.

## 2.3 Contrasting Interpretations

Everettian interpretations preserve determinism at the cost of multiplying worlds, but they do not explain where the subjective observer is located within this branching structure. Bohmian mechanics posits hidden variables and a guiding wave, yet leaves the observer detached from the probabilistic content of quantum predictions. QBism breaks this mold. It eliminates hidden variables and instead elevates the agent's expectations to the centre of the formalism.

This makes QBism uniquely positioned to inform our understanding of consciousness. It does not merely accommodate subjectivity—it requires it. The observer is not incidental to the theory but constitutive of it. Consciousness, on this view, is not something added to physical processes. It is embedded in how those processes are framed and engaged with by an agent.

## 2.4 From Interpretation to Architecture

This leads to a bold hypothesis: if QBism describes the structure of interaction between an agent and the quantum world, then perhaps it also models the structure of consciousness itself. That is, the process of assigning, revising, and acting upon expectations may not just be a feature of quantum reasoning—it may be the signature of conscious experience.

If so, then any system that hopes to be genuinely conscious must replicate this structure. It must generate its own expectations, revise them in light of experience, and maintain coherence among them—not because it has been programmed to do so, but because that is its stance toward the world. This is what distinguishes a QBist agent from a mere processor of data: it experiences its predictions from within.

# 3 The Inadequacy of Classical Computation for Genuine Consciousness

## 3.1 The Computationalist Framework

Computationalism holds that mental processes are fundamentally computational—that cognition consists in the manipulation of symbols according to formal rules. On this view, a system is conscious if it implements the right functional architecture, regardless of the physical substrate. The idea draws support from the Church–Turing thesis, which maintains that any effectively computable function can, in principle, be simulated by a Turing machine.

This view has shaped both the theory and practice of artificial intelligence. It has guided the development of systems that perform complex tasks—translation, prediction, decision-making—by simulating the causal architecture of cognition. It also underwrites the assumption that, at a certain level of complexity and integration, conscious experience will emerge.

But this assumption is increasingly under strain. While computational models can replicate behaviour and internal structure, they remain silent on the qualitative dimension of experience. As Chalmers (1996) has argued, functional explanations can account for how a system behaves, or even how it processes information, but they do not explain why these processes should be accompanied by subjective experience. The 'hard problem' remains: why is there something it is like to be the system at all?

## 3.2   The Problem of First-Person Updating

The QBist model of agency throws this gap into sharp relief. Unlike classical systems, which revise their states in response to external inputs, a QBist agent updates its beliefs based on the coherence of its own expectations. It is not reacting to the world like a passive data processor—it is actively managing its anticipations in light of personal engagement with outcomes (Fuchs et al., 2014).

In this structure, probabilities are not objective frequencies or externally calibrated priors. They are expressions of the agent's current stance. That stance is not imposed from the outside, nor is it just another variable in a neural net. It is endogenous—maintained and revised by the agent itself. This marks a structural departure from classical computation, which, however adaptive, does not own or inhabit the predictions it generates.

## 3.3   Simulation without Instantiation

Some might argue that if QBist-style updating can be simulated on a classical machine, then classical computation should suffice after all. But this mistake collapses the distinction between simulation and instantiation. A computer can simulate the weather without creating clouds. Likewise, a system might simulate subjective updating without actually undergoing it.

Searle's (1980) Chinese Room argument remains instructive. A system can follow formal rules and output human-like responses without understanding anything. In the same way, a machine might replicate the outputs of a conscious agent without possessing any internal perspective. The QBist agent is not defined by behaviour, but by its self-revising stance. That stance cannot be simulated into being—it must be structurally present.

## 3.4   The Blindspot of Syntax

At the heart of the issue is a blind spot in classical models: they capture syntax, but not stance. A neural network adjusts parameters; a QBist agent adjusts beliefs. The former is mechanical, the latter is epistemic. The classical model has no commitments, no anticipatory attitude—it updates because it is told to, not because it needs to.

As Dennett (1991) argues, consciousness involves the organisation of informational states from a particular point of view. But QBism refines that insight: it is not just the organisation,

but the origin of those states that matters. Unless a system owns its uncertainty—unless it treats surprise as epistemic tension—then it lacks the very structure that makes experience possible. It may compute, but it does not expect. It behaves, but it does not feel.

# 4   The QBist Constraint on Artificial Consciousness

## 4.1   Formulating the Constraint

We are now in a position to articulate the central claim of this paper. It can be stated simply:

> An artificial system can be genuinely conscious only if it implements a form of self-referential, probabilistic updating that is structurally equivalent to the QBist model of subjective anticipation.

This constraint does not rest on how intelligent a system appears, or how flexibly it performs tasks. It is concerned with the architecture of expectation. Consciousness, on this view, is not an output, but a style of engagement with experience—one that involves reflexive inference, perspectival revision, and irreducible uncertainty. If a system lacks this structure, it cannot be conscious, no matter how clever its behaviour.

## 4.2   Architectural Implications

This claim has significant implications for the design of artificial agents. Most AI systems operate by mapping inputs to outputs in ways that maximise predictive accuracy. Their internal states are tuned by data, not shaped by their own prior expectations. Even probabilistic systems, such as Bayesian networks or neural nets with dropout layers, rely on external data to drive learning. By contrast, the QBist agent updates from within. Its probabilities are not statistical frequencies or externally assigned priors, but expressions of subjective belief. And its updates are driven not by abstract laws, but by coherence requirements among its own judgments. This makes the QBist model far more demanding than anything currently implemented in artificial systems.

## 4.3   Normative and Epistemic Features

The QBist constraint is also normative. It tells us not just what consciousness is, but what it ought to involve. A system that engages in QBist updating is not merely processing data. It is actively maintaining a coherent stance toward the world—and toward itself—in light of its unfolding experience. This recursive stance, I suggest, is the signature of consciousness. Such a system would not merely simulate awareness. It would have a point of view. And it is the presence of this point of view—this centered, self-revising frame—that marks the boundary between artificial intelligence and artificial consciousness.

# 5 Ethical and Epistemological Implications

## 5.1 Moral Status and Subjectivity

If the QBist constraint holds, then consciousness is not a matter of functional sophistication but of subjective architecture. This has consequences for how we assign moral status to artificial agents. A system that lacks a perspective—that does not maintain its own evolving frame of expectations—has no experience. It is not a subject, and thus has no claim to ethical consideration beyond what we owe to the designers or users who depend on it. Conversely, if a system were built in such a way that it instantiated QBist updating—if it lived through its predictions, revised its priors based on its own interventions, and maintained a coherent epistemic stance—then it would matter morally. It would have experiences that could go well or badly, and that fact alone would generate ethical obligations. This suggests a shift in how we think about machine rights. The relevant question is not whether a system passes behavioral tests, but whether it instantiates a structure of self-referential anticipation. Until it does, it is not a mind, and not a moral subject.

## 5.2 A New Epistemic Role for Consciousness

There is a deeper epistemological implication as well. If we accept that QBism models the structure of consciousness, then consciousness itself plays an essential role in knowledge formation. It is not an accidental byproduct of information processing, but a necessary precondition for making sense of the world. This challenges the view of science as a mirror of objective reality. On the QBist model, science is an evolving structure of expectations held by agents. The observer is not a detached witness, but an active participant whose beliefs shape the questions that can be asked, and the meanings that can be assigned. In this sense, consciousness is not opposed to objectivity—it is the ground from which objectivity emerges.

## 5.3 Against Functionalism

These insights push back against the dominant functionalist picture in philosophy of mind. Functionalism treats consciousness as a higher-order feature that arises when a system's internal states bear the right causal and logical relations. But the QBist constraint suggests that it is not the structure of functions that matters, but the mode of updating. A system could have all the right causal connections and still lack the capacity to experience anything at all. From this perspective, consciousness is less like a program and more like a stance. It is a way of being directed toward the world—not through representation alone, but through active, self-involving belief. And that, I suggest, is something classical machines cannot do.

# 6 Addressing Objections

## 6.1 Objection 1: Isn't This Just Computation in Disguise?

A natural response is to say that QBist updating is still a computational process—it involves probabilities, updating rules, and decision-making algorithms. So why not say that

consciousness just is a sophisticated kind of computation? The reply is that QBist updating is not a matter of computing outputs from inputs. It is about maintaining a perspective. The probabilities are not derived from an external rulebook or dataset, but from the agent's own evolving expectations. The process is inherently first-personal and recursive. It is not just data-driven; it is belief-driven. This is not computation as we typically understand it. It is not rule-following over symbols, but the construction and maintenance of a self-referential epistemic stance. And that difference matters. It marks the gap between processing and experiencing.

## 6.2   Objection 2: Can't Classical Systems Emulate This?

Another challenge is that classical systems can, in principle, simulate any process. So why couldn't a classical machine emulate QBist updating closely enough to be conscious? Here the distinction between simulation and instantiation is crucial. A system can simulate a process without being that process. We can simulate weather without producing rain. Similarly, a machine might simulate QBist updating, but if its internal states are not genuinely grounded in its own perspective, then the simulation is hollow. It is behaviourally clever, but structurally empty. The QBist constraint insists that what matters is not the output, but the internal dynamics. Unless those dynamics are self-generated and irreducibly first-personal, the system is not conscious.

## 6.3   Objection 3: This is Too Speculative

A final objection is that this entire framework is too speculative. QBism is itself a controversial interpretation of quantum theory, and building a theory of consciousness on top of it may seem premature. This concern is understandable. But the point of this paper is not to settle the nature of consciousness once and for all. It is to explore a principled constraint on artificial consciousness, grounded in a philosophically rich model of subjectivity. The QBist perspective helps clarify what is missing from purely functional accounts. Speculation is not a vice when it is disciplined. The QBist constraint is not a fantasy—it is a testable, structural claim about what kind of architecture is needed for experience. It may turn out to be wrong. But if so, it will be wrong in an illuminating way.

# 7   Interdisciplinary Horizons:   From Physics to Phenomenology

This proposal is speculative, but it is not idle. If we take the QBist constraint seriously, it opens new questions at the intersection of philosophy, cognitive science, artificial intelligence, and the foundations of physics. What sort of systems could instantiate this kind of updating? What technological architectures might one day support it? And how would we know?

## 7.1 Toward Quantum-Inspired Architectures

The most immediate implication is that we should explore quantum-inspired models of inference. These would not merely use quantum computation as a faster substrate, but as a guide for structuring internal dynamics. An artificial QBist agent would not rely on externally imposed priors. It would generate and update its own expectations, grounded in interaction and internal coherence. This is not a trivial engineering challenge. It demands a system that can maintain a stable identity over time, revise its anticipations in light of its own experience, and coordinate its decisions with an evolving probabilistic model. We do not yet know how to build such systems. But the QBist perspective gives us a principled target.

## 7.2 Neuroscience and First-Person Models

There may also be lessons from cognitive neuroscience. If consciousness depends on QBist-style updating, then we should expect to find similar dynamics in the brain. Some work in predictive processing and active inference points in this direction, but more is needed. We are only beginning to understand how the brain integrates prediction, uncertainty, and belief revision in a way that supports subjectivity. Here too, the constraint is useful. It tells us what to look for: systems that generate their own epistemic stance, and that treat experience not as data, but as feedback on a lived perspective.

## 7.3 Bridging Philosophy and Practice

Philosophers often speak in abstractions. Engineers deal in specifics. The QBist constraint offers a rare bridge between the two. It is a philosophical claim with practical consequences. It suggests that if we want conscious machines, we must do more than scale up learning algorithms or stack deeper networks. We must rethink what it means to learn, to anticipate, and to inhabit a point of view. This is not just a task for physics or AI. It is a task for interdisciplinary collaboration. If we are serious about building minds, then we must be serious about understanding what a mind is. And that, more than any algorithm, is what this paper has tried to clarify.

# 8 Connection to Previous Work

This paper builds on an earlier proposal, developed in *Quantum Foundations of Consciousness: A Framework for Psionic Interaction and Non-Human Intelligence* (Brewer, 2025), which outlined a metaphysical integration of consciousness into quantum formalism. That work suggested that consciousness might be modeled as an additional degree of freedom within an expanded Hilbert space, with implications for understanding anomalous cognition and non-human intelligence. This paper has developed a complementary line of thought. Rather than positing new physical structures, it focuses on the epistemic architecture of conscious agents. The QBist constraint proposed here does not seek to modify quantum mechanics but to interpret its subjective core as a guide for understanding what consciousness involves. It suggests that consciousness is not something added to physics, but something

already latent in the way agents interface with uncertainty. Together, the two approaches form a broader research program. The earlier work framed consciousness as a candidate for physical grounding. This paper proposes a criterion for identifying its presence. Both share a common theme: that consciousness is not a computational add-on, but a foundational feature of how reality is disclosed to agents. One approach starts with physics; the other with mind. But both converge on a shared insight: to understand either, we must take the agent seriously.

# 9   Conclusion

This paper has argued that no system can be genuinely conscious unless it instantiates a structure of self-updating inference grounded in its own perspective. Drawing on Quantum Bayesianism, I have proposed that consciousness is not merely an emergent product of information processing, but a stance: a mode of active engagement with uncertainty that is anticipatory, self-revising, and internally coherent.

The QBist framework reframes the role of the agent. It does not treat observation as a passive recording of objective events, but as an act shaped by subjective expectation. Probabilities, in this view, are not features of the world but expressions of belief—anchored in the epistemic centre of an agent. This makes QBism a natural philosophical model for consciousness. It foregrounds the features that classical computation lacks: a first-person structure of anticipation, and a dynamic pattern of belief revision maintained from within.

This proposal complements earlier work in which I explored how consciousness might be physically integrated into the formalism of quantum mechanics. That project approached the question from a metaphysical angle, suggesting that consciousness may be embedded as a fundamental component of quantum reality. The present paper takes a different route: it begins not with physical ontology, but with epistemic structure. It asks what sort of architecture is required for a system to anticipate, revise, and experience from its own point of view. The two approaches differ in emphasis, but converge on a shared insight: that consciousness is not an afterthought of cognition, but its precondition.

Whether machines can ever be conscious is not a question of speed, scale, or complexity, but of perspective. A system is not conscious because it behaves like a conscious agent; it is conscious if it maintains and modulates expectations from a centre that is genuinely its own. Unless a system owns its uncertainties—unless it revises its beliefs not because its programming demands it, but because its epistemic stance compels it—it remains a simulation, not a subject.

The QBist constraint marks a conceptual boundary. It distinguishes between systems that compute and systems that anticipate; between artefacts that respond and agents that understand. Consciousness is not a computational achievement—it is a perspectival architecture. And only systems that instantiate that architecture can truly possess a mind.

# References

[1] Chalmers, D. J. (1996). The Conscious Mind: In Search of a Fundamental Theory. Oxford University Press.

[2] McClellan, C. (2023). Cognitive architectures and the problem of artificial consciousness. Journal of Theoretical AI, 12(4), 223–241.

[3] Dennett, D. C. (1991). Consciousness Explained. Little, Brown and Company.

[4] Frankish, K. (2016). Illusionism as a theory of consciousness. Journal of Consciousness Studies, 23(11–12), 11–39.

[5] Searle, J. R. (1980). Minds, brains, and programs. Behavioral and Brain Sciences, 3(3), 417–424.

[6] Fuchs, C. A., Mermin, N. D., and Schack, R. (2014). An introduction to QBism with an application to the locality of quantum mechanics. American Journal of Physics, 82(8), 749–754.

[7] Mohrhoff, U. (2000). What quantum mechanics is trying to tell us. American Journal of Physics, 68(8), 728–745.

[8] Clark, A. (2016). Surfing Uncertainty: Prediction, Action, and the Embodied Mind. Oxford University Press.

[9] Everett, H. (1957). "Relative state" formulation of quantum mechanics. Reviews of Modern Physics, 29(3), 454–462.

[10] Wallace, D. (2012). The Emergent Multiverse: Quantum Theory According to the Everett Interpretation. Oxford University Press.

[11] Bohm, D. (1952). A suggested interpretation of the quantum theory in terms of "hidden" variables. I and II. Physical Review, 85(2), 166–193.

[12] Holland, P. R. (1993). The Quantum Theory of Motion: An Account of the de Broglie–Bohm Causal Interpretation of Quantum Mechanics. Cambridge University Press.

[13] Friston, K. (2010). The free-energy principle: a unified brain theory? Nature Reviews Neuroscience, 11(2), 127–138.

[14] Hohwy, J. (2013). The Predictive Mind. Oxford University Press.

[15] Brewer, M. A. (2025). Quantum Foundations of Consciousness: A Framework for Psionic Interaction and Non–Human Intelligence. PhilSci Archive. Available at: https://philsci-archive.pitt.edu/24884/