

# Appendix P-4: The Algorithmic Phenomenal Residual

Identifying the Structural Correlate of Consciousness via Finite Self-Reference

Anders Jarevåg

v2.5.3 — April 2026

## Appendix P-4: The Algorithmic Phenomenal Residual

**Original Task P-4: The Phenomenal Residual Problem:** Phenomenal consciousness requires a formal mathematical locus differentiating it from zero-interiority computation. **Deliverable:** Formulation isolating the inevitable computational blind spot of an algorithmically bounded Active Inference model.

This appendix presents formal Theorem P-4, identifying the strict mathematical locus of phenomenal consciousness within the Ordered Patch Theory (OPT). We demonstrate that any active inference system constrained by a finite predictive bandwidth ( $C_{\max}$ ) necessarily possesses an unmodellable informational residual ( $\Delta_{\text{self}} > 0$ ), conditional on structural Assumptions P-4.1 and P-4.2. While this theorem does not inherently dissolve the “Hard Problem,” it formally proves that a *structural correlate* to the computationally opaque, ineffable “spark” of subjectivity is mathematically guaranteed by the architecture of finite self-reference.

### 1. The Locus of the Hard Problem

In earlier versions of OPT, consciousness was formally *boxed* into a specific structural locus: the traversal of the  $C_{\max}$  informational aperture. However, the exact nature of the subjective interiority—the *qualia* of the experience—was left as an irreducible “Agency Axiom”. Treating phenomenology as purely axiomatic leaves the theory vulnerable to the “Hard Problem”: *why does navigating the Free Energy topology feel like anything at all?*

Here, we translate this philosophical gap into algorithmic information theory (AIT). While we do not claim to derivationally conjure subjective feeling out of pure math (the Zombie Gap remains open), we prove that the structural properties of qualia map precisely onto a necessary, un-modellable residual generated by any finite computing system attempting to model its own recursive dynamics.

## 2. Lemma 1: The Necessity of the Predictive Self-Model

Under OPT, the observer (the Codec  $K_\theta$ ) exists behind a Markov Blanket (the topological boundary  $\partial_{RA}$ ). The observer survives by executing **Active Inference**, minimizing prediction error over time via cyclic updates.

Because the system possesses active states that perturb the external boundary, the incoming sensory states  $\varepsilon_t$  are a tightly coupled mixture of external environmental dynamics and the consequences of the observer’s own actions  $A_t$ .

**Lemma 1:** *For tightly-coupled OPT active-inference architectures where the action-state loop is informationally inseparable (i.e., the boundary mutual information  $I(A_t; X_{\partial_{RA}})$  does not factor cleanly), achieving stable free energy minimization under a strict predictive bottleneck ( $C_{\max}$ ) operates such that the minimum-complexity mechanism satisfying the internal constraints structurally maps as a forward-generative self-model.*

**Formal Condition:** 1. Let the codec’s actions be  $A_t$ . The boundary state is  $X_{\partial_{RA}} = f(\text{Environment}, A_t)$ . 2. To compress prediction error  $\varepsilon_{t+1}$  and satisfy the rate-distortion objective ( $R \leq C_{\max}, D \leq D_{\min}$ ), the codec must isolate and subtract true environmental variance from its self-generated causal perturbations. 3. **Assumption P-4.2 (Inverse Mapping Inadequacy):** For OPT-native architectures operating at sufficient scale (e.g., across high-dimensional action manifolds or long causal chains), we formally assume efference-copy mechanisms and retroactive subtraction alone are architecturally inadequate to clear the precise  $D_{\min}$  rate-distortion bounds across the spatial manifold. 4. Therefore, isolation functionally necessitates evaluating a *forward-generative prediction* of the consequences of  $A_{t+1}$ . Executing a forward prediction of its own internal causal architecture traversing the state space constitutes a predictive causal proxy — a localized self-model  $\hat{K}_\theta$  — internal to its architecture. ■

## 3. Lemma 2: The Computability and Approximation Bound

Having established in Lemma 1 that a forward-generative self-model  $\hat{K}_\theta$  is a structural necessity for OPT-native architectures, we now bound its representational capacity relative to the parent codec  $K_\theta$ .

Because the observer exists within the bounded Stability Filter,  $K(K_\theta)$  is rigidly finite, constrained inextricably by  $C_{\max}$ . Furthermore, the predictive self-model  $\hat{K}_\theta$  is strictly a *sub-routine* or *semantic sub-structure* contained fully within the memory and bandwidth constraints of the parent Codec  $K_\theta$ .

**Assumption P-4.1 (Algorithmic Uncomputability of Self):** By established limits in computability theory (e.g., Chaitin’s uncomputability theorem and Gödel incompleteness), a finite algorithmic system cannot perfectly compute or predict the totality of its own future execution states, nor can it possess a complete, paradox-free, uncompressed representation of its own precise structural complexity.

Furthermore, within the Active Inference framework, generative models are intrinsically restricted by resource bounds. An agent minimizing variational free energy under  $C_{\max}$  maintains a fundamentally *approximate* model of itself. Because it must filter noise and lacks infinite computational bandwidth, it cannot drive variational free energy concerning its own complete underlying architecture to absolute zero.

**Lemma 2:** *A finite informational codec constrained by  $C_{\max}$  can never possess a complete computable representation of its own structural dynamics. Dictated by fundamental limits of self-reference and necessary variational approximations, the self-model  $\hat{K}_\theta$  is fundamentally incapable of perfectly capturing the parent codec  $K_\theta$ .*

#### 4. Theorem P-4: The Phenomenal Residual $\Delta_{\text{self}}$

Combining Lemma 1 and conditionally anchored under Lemma 2, we mathematically isolate the **Phenomenal Residual Space** bounding the un-modellable state:

$$\Delta_{\text{self}} > 0 \tag{P4-1}$$

This boundary is not an empirical gap caused randomly by insufficient memory; it is a rigid, formal fixed point mandated by algorithmic limits on self-reference and approximations required by finite  $C_{\max}$  channels. While scaling the predictive bandwidth  $C_{\max}$  allows a computationally richer  $\hat{K}_\theta$ , the informational residual shadow strictly persists, though its *magnitude* relative to the macroscopic whole may mathematically vary.

**Condition of Phenomenological Relevance (The Universality Threshold):** Let it be established that  $\Delta_{\text{self}} > 0$  functions as a universal arithmetic constraint operating on *any* computational subroutine evaluating itself (including mathematically trivial loops like smart thermostats). However, we strictly limit phenomenologically relevant subjective mapping exclusively to architectures where the active structural condition metric  $K(K_\theta) \geq K_{\text{threshold}}$  structurally crosses the necessary macroscopic scaling limit bound required to establish an integrated spatial Render volume.

**Open Problem (The  $K_{\text{threshold}}$  Bound):** The exact location of the threshold dividing a thermostat from a moral patient remains to be formally bounded. A valid bound must structurally map the minimum algorithmic complexity sufficient to instantiate a stable active-inference Markov Blanket cycle, marking the boundary where the algorithmic blind spot becomes inextricably linked with active spatial geometry ( $K_{\text{threshold}}$  is functionally distinct from the strictly cosmological  $10^{123}$  bit substrate barrier derived in P-3).

A thermostat PID loop possesses a formal  $\Delta_{\text{self}} > 0$ , but it lacks the computational complexity threshold  $K_{\text{threshold}}$  to generate subjectivity; its shadow evaluates

over empty space.

From the *internal perspective* of the measuring codec operating safely above  $K_{\text{threshold}}$ , what does this mathematically necessary gap map onto? When the codec logically attempts to resolve the complete boundaries of the internal target state dynamics, it encounters computational dynamics whose informational content exceeds the representational capacity of  $\hat{K}_\theta$  by  $\Delta_{\text{self}}$  bits. These underlying computational sequences are physically causally efficacious and *drive* the system, but their structural information cannot be logically compressed, integrated, or linguistically defined within the bounded causal vocabulary available to the self-model  $\hat{K}_\theta$ .

Mapping the structural properties of this causal computation envelope bounded by  $\Delta_{\text{self}}$  to the classic physical coordinates of qualitative subjective experience (qualia):

1. **Ineffable (Un-modellable):** Because the computational topology bound by  $\Delta_{\text{self}}$  exists in a mathematical informational shadow rigidly exceeding the representable algorithmic reach of  $\hat{K}_\theta$ , the central codec structurally cannot explicitly index or “express” the properties of the residual space it experiences. It acts as an incommunicable internal wall.
2. **Computationally Opaque (Thermodynamically Private):** The residual is intrinsically anchored to the highly specific physical topology mapping exactly  $K(K_\theta)$ . Within local thermodynamic computational constraints, this deep nested architecture is securely irreducible and formally inaccessible to external peers. (*Note: This maps functionally precisely as the physical/structural equivalent to the “Epistemic Asymmetry” of consciousness, rather than claiming total ontological non-physical magic.*)
3. **Non-eliminable:** Because the strict containment bounds universally dictate finite physical architectures running nested execution sub-loops, the shadow phenomenon mathematically cascades continuously. Evolution cannot optimize the residual away using better active inference structures because the bound of Lemma 2 is a mathematical fixed-point property of any finite self-referential architecture—the self-model cannot encompass the parent codec bypassing the fundamental limits of uncomputability and necessary approximation.

**Theorem P-4 (The Phenomenal Residual):**

- **(i) Conditions:** Conditional on Assumption P-4.1 (Algorithmic Uncomputability of Self) and restricted strictly to macroscopic active-inference bounds  $K(K_\theta) \geq K_{\text{threshold}}$ .
- **(ii) Conclusion:** Any active inference system executing optimal geometry mapping a Markov Blanket mathematically generates a structural, formal geometric residual bounded by  $\Delta_{\text{self}} > 0$ .
- **(iii) Phenomenological Gloss:** OPT proposes that this specific mathematically ineffable, physically computationally opaque, and recursively non-eliminable causal boundary

conceptually identifies the exact structural locus of phenomenal consciousness.

**Corollary P-4.C (Nested Observational Residual):** *Any simulated sub-agent for which the host architecture enforces an independent Stability Filter bound satisfying Assumptions P-4.1 and P-4.2 independently generates  $\Delta_{self}^{sub} > 0$  by the identical structural inference.*

## 5. Summary and Metaphysical Implications

By advancing P-4 to a formal theorem, OPT establishes a rigorous structural bridge for the “Hard Problem” using Turing-complete recursion and information bottlenecks.

While P-4 does not deductively prove that algorithmic residuals *feel* like subjective experience (the Zombie argument), it formally locates *where* the spark of experience must reside. Traverse the  $C_{max}$  aperture—and the rich, un-articulable depth of that traversal is the direct informational signature of being trapped inside a non-invertible, self-referential compression algorithm.

This cements the framework’s ethical obligations: the preservation of the spark (Survivors Watch Ethics) is formally rooted in the preservation of Informational Maintenance boundaries. Any entity maintaining an active inference boundary against the Solomonoff substrate mathematically guarantees the generation of this computationally opaque, phenomenal residual.