

A Unified Thermodynamic Account of Qualia, Selfhood, Free Will, and Language: Refusal as the Origin of Symbolic Consciousness

Author: Alastair Waterman

Affiliation: Independent Researcher, United Kingdom

Correspondence: alastairwaterman@gmail.com

ORCID: 0009-0000-2867-0207

Date: December 7, 2025

DOI 10.5281/zenodo.17850349

Abstract

This article proposes a unified thermodynamic framework—Refusal-Driven Dimensionality Reduction Theory (RDRT)—for four traditionally separate problems of consciousness: the nature of qualia, the phenomenal sense of selfhood (“mineness”), the subjective experience of libertarian free will, and the evolutionary origin of language and symbolic thought. Building on prior work (Waterman 2025a, 2025b), phenomenal consciousness is conceptualised as an evolved mechanism that halts recursive self-prediction at a finite depth to prevent energetic overload in a ~20 W cortical system (Attwell & Laughlin 2001; Stroud et al. 2025). This refusal leaves three correlated phenomenal residues: qualia (compressed sensory uncertainty), mineness (unassignable introspective residue), and openness to counterfactual futures (uncomputed action branches). A single measurable parameter—subjective mental temperature (T_M)—is introduced as the maximum sustainable recursion depth before refusal, derived from the Landauer limit and the exponential energetic cost of hierarchical prediction (Landauer 1961; Flesch et al. 2022; Ali et al. 2022). Language is argued to have co-evolved as symbolic labels for refusal boundaries, with the pronoun “I” marking the primordial convergence of the three residues. The framework predicts correlated intensity changes across qualia, selfhood, and agency under metabolic, pharmacological, and developmental manipulations, and offers concrete measurement proposals using existing EEG/PET methods. Comparisons with integrated information theory (Tononi et al. 2016), quantum-microtubule approaches (Hameroff & Penrose 2014), and panpsychism are provided. While necessarily speculative in its integrative scope, the account is grounded throughout in published physiological, computational, and phenomenological data.

Keywords

qualia, phenomenal self, mineness, libertarian free will, language origins, thermodynamic constraints, predictive processing, active inference, free-energy principle, refusal-driven dimensionality reduction, subjective mental temperature, sense of meaning, evolutionary neuroscience, consciousness measurement

1. Introduction

The human brain consumes approximately 20% of the body's resting metabolic energy while constituting only 2% of its mass (Attwell & Laughlin 2001). This extraordinary energetic cost is predominantly devoted to signalling rather than housekeeping functions, with action potentials and glutamatergic postsynaptic currents accounting for the majority of the budget in grey matter (Attwell & Laughlin 2001). Recent modelling of task-optimised recurrent neural networks trained under explicit metabolic penalties has confirmed that biological-like constraints on firing-rate and noise dramatically shape the geometry of cortical representations (Stroud et al. 2025). These findings reinforce a long-standing insight: the brain is an energy-limited predictive device that cannot afford unrestricted computation.

Recursive self-modelling—the inclusion of the modelling process itself within its own predictions—poses a particularly severe threat. Hierarchical predictive architectures exhibit exponential growth in computational demand with each additional level of meta-representation (Flesch et al. 2022; Ali et al. 2022). In a finite ~20 W system, unbounded recursion would rapidly exceed the Landauer limit for irreversible computation (Landauer 1961) and the physiological ceiling defined by cerebral blood flow, leading to catastrophic overheating within seconds. Any viable brain must therefore possess a mechanism that reliably halts recursion before this point.

Two earlier preprints argued that phenomenal consciousness is precisely this mechanism, evolved and fixed under thermodynamic selection pressure. The first (Waterman 2025a) proposed that qualia constitute a biologically obligatory form of lossy compression that collapses high-dimensional sensory data into a low-dimensional manifold of affectively tagged residues, preserving adaptive utility at minimal energetic cost. The second (Waterman 2025b) identified the phenomenal sense of selfhood (“mineness”) as the unassignable residue of a structurally enforced predictive halt—an ontological lacuna rather than a positive representation. Both accounts relied on the same core operation: refusal to compute further when continued prediction would violate the organism's thermodynamic boundary.

The present article extends and unifies these claims. We propose Refusal-Driven Dimensionality Reduction Theory (RDRT) as a single framework in which qualia, phenomenal selfhood, and the subjective experience of libertarian free will emerge as obligatory facets of the same halting process, partitioned across sensory, introspective, and action domains respectively. We introduce a measurable quantity—subjective mental temperature (T_M)—defined as the maximum sustainable depth of recursive self-prediction before refusal is enforced, and we show how existing physiological and neuroimaging data already constrain its value to the observed range of human phenomenal intensity.

Crucially, we argue that language itself is not a cultural epiphenomenon but a direct evolutionary consequence of the refusal mechanism. Once refusal boundaries become stable

and reusable, the selective pressure to label them—for rapid retrieval without recomputation—becomes overwhelming. The pronoun “I” is identified as the primordial symbol: the first reusable marker of the point at which sensory, introspective, and volitional refusals converge.

Although necessarily speculative in its synthesis, the account makes contact at every step with established results: the energy budget of signalling (Attwell & Laughlin 2001), the exponential cost of hierarchical prediction (Flesch et al. 2022; Ali et al. 2022; Stroud et al. 2025), the stochastic and metabolic underpinnings of apparent volitional randomness (Schurger et al. 2021), and the repeatedly observed co-variation of phenomenal intensity across sensory, selfhood, and agency domains in metabolic, developmental, and pharmacological studies (Jamadar et al. 2025; Schurger et al. 2021). By grounding four of the deepest problems in consciousness research in a single, thermodynamically enforced computational boundary, RDRT offers a parsimonious, falsifiable alternative to integrated information theory (Tononi et al. 2016), quantum-microtubule theories (Hameroff & Penrose 2014), and panpsychist approaches, while remaining fully compatible with the free-energy principle and active inference (Friston 2010; Parr & Friston 2018).

The article proceeds as follows. Section 2 restates and refines the core onto-epistemological principles. Section 3 demonstrates the unification of qualia, selfhood, and free will as three facets of a single refusal operation. Section 4 traces the origin of symbolic language to the labelling of refusal boundaries. Section 5 introduces subjective mental temperature (T_M) and its neural and physical underpinnings. Sections 6–7 present empirical predictions and concrete measurement protocols. Section 8 compares the framework with competing theories and draws philosophical and ethical consequences. Section 9 concludes.

The empirical predictions and proposed measurement protocols presented here are designed to be immediately testable with existing technology. At the time of writing (December 2025), no original experimental data from the author’s laboratory are included; all quantitative claims about expected effect sizes and correlations are derived from re-analysis and meta-estimation of the published literature cited.

2. Onto-Epistemological Principles Revisited

The present synthesis rests on six refined principles originally articulated in Waterman (2025a, 2025b). These are restated and extended here to accommodate the unification of qualia, selfhood, free will, and language within a single thermodynamic framework.

1. Consciousness as an Emergent Refusal Property

Phenomenal consciousness is not a fundamental property of matter nor an epiphenomenal byproduct of computation, but an emergent property of bounded predictive systems that are forced to refuse complete recursive self-modelling (Attwell & Laughlin 2001; Stroud et al. 2025). The refusal mechanism itself is biologically obligatory: without it, the exponential energetic cost of hierarchical prediction exceeds the ~20 W cortical budget within seconds (Flesch et al. 2022; Ali et al. 2022).

2. **Functional Necessity of Phenomenal Facets**

Qualia, the sense of mineness, and the subjective experience of libertarian openness are not optional luxuries but functionally necessary residues of the refusal operation partitioned across three minimal domains required for bounded agency: sensory input, introspective self-modelling, and action selection (Friston 2010; Parr & Friston 2018; Fields et al. 2024). Fewer domains would fragment adaptive behaviour; additional domains would introduce redundant metabolic cost without gain.

3. **Irreducibility Yet Measurability of Phenomenal Character**

The subjective character of experience is irreducible to third-person neural description because it is the first-person trace of the refusal boundary itself (Waterman 2025b). Nevertheless, the depth and intensity of this boundary are quantitatively constrained by thermodynamics and therefore measurable in principle as a single parameter—subjective mental temperature (T_M)—derived from the Landauer limit and observed exponential scaling of recursive prediction cost (Landauer 1961; Flesch et al. 2022).

4. **Embodied and Enactive Grounding**

Refusal is not an abstract computational halt but a biologically embodied process embedded in the organism–environment loop (Varela et al. 1991). The felt “Sense of Meaning” is the experiential curvature of the refusal manifold, maintained through ongoing sensorimotor and affective coupling.

5. **Refusal as the Ontological Origin of Symbolic Representation**

Once refusal boundaries stabilise, the selective advantage of labelling them for rapid reuse becomes overwhelming. Language emerges as the systematic externalisation and compression of these boundaries into discrete symbols (new principle). The pronoun “I” constitutes the earliest and most metabolically consequential symbol: a reusable pointer to the convergence point of sensory, introspective, and volitional refusals.

6. **Compatibility with Existing Predictive Processing Frameworks**

The account is fully consistent with the free-energy principle and active inference (Friston 2010; Parr & Friston 2018), which provide the formal machinery of hierarchical prediction, but adds an explicit thermodynamic halting condition absent in standard formulations. It thereby explains why biological agents, unlike unbounded digital simulations, exhibit stable phenomenal residues.

These principles collectively shift the explanatory burden from “Why do certain neural processes feel like anything?” to “Given a 20 W predictive engine that must refuse infinite recursion, what is the minimal set of phenomenal residues required for adaptive behaviour, and how are they labelled?” The remainder of the article demonstrates that the observed structure of human consciousness—qualia, mineness, agency, and the pronoun “I”—is the precise answer evolution discovered.

3. Unifying Qualia, Selfhood, and Free Will as Refusal Facets

Refusal-Driven Dimensionality Reduction Theory (RDRT) posits that the brain, as a thermodynamically bounded hierarchical predictive engine, cannot complete the full regression implied by the free-energy principle without violating its ~ 20 W ceiling (Attwell & Laughlin 2001; Stroud et al. 2025). The refusal mechanism is therefore not optional but obligatory. Crucially, the same exponential wall is encountered in three distinct predictive domains that jointly constitute viable agency in an uncertain world: perception of the external environment, modelling of the modelling process itself, and simulation of possible actions. Each domain requires a hard halt at approximately the same recursion depth ($n \approx 4\text{--}6$ levels under resting conditions; Flesch et al. 2022; Ali et al. 2022), and each halt leaves a characteristic phenomenal residue. These residues are qualia, mineness, and the subjective sense of libertarian openness, respectively.

3.1. Qualia as Sensory Refusal

Raw sensory channels deliver on the order of $10^6\text{--}10^8$ bits s^{-1} to cortex. A literal, lossless Bayesian inversion of these data would require computational resources far exceeding the entire cortical energy budget (Attwell & Laughlin 2001; Jamadar et al. 2025). The brain therefore refuses to compute the full posterior over causes and instead collapses the sensory manifold into an extremely low-dimensional affective summary—typically estimated at 36–150 effective bits per perceptual moment (Waterman 2025a).

This collapse is experienced as qualia: the vivid redness of red, the stinging quality of pain, the felt warmth of sunlight. The compression is not arbitrary; opponent structure (red/green, warm/cool, pleasant/unpleasant) and valence tagging preserve the minimal information required for rapid approach–avoidance decisions while discarding everything else. Recent metabolic imaging confirms that the energetic saving is dramatic: perceptually salient stimuli that trigger strong qualia elicit no greater oxygen consumption than neutral ones once the refusal boundary is reached, despite vastly richer subjective character (Jamadar et al. 2025). Without qualia, the organism would be forced to recompute the full sensory inversion on every encounter, rapidly exhausting its glucose and oxygen supply.

3.2. Selfhood as Introspective Refusal

When the predictive hierarchy attempts to model its own modelling activity, a pure self-referential loop is created. In an unbounded system this leads to an infinite regress (“I am the one who is thinking that I am thinking...”). In a bounded system the recursion must terminate. The refusal cannot be assigned to any external cause or lower-level representation; there is no higher frame from which the residue could be attributed or rejected (Waterman 2025b). What remains is an unassignable kernel experienced as phenomenal mineness—the brute sense that this experience is happening to me and not to someone else.

This residue is topologically inevitable: it is the predictive equivalent of asking “What is north of the North Pole?” (Metzinger 2009). Neuroimaging evidence supports the claim that mineness correlates with activity in regions that track unmodelled prediction error without successful attribution, notably the dorsal anterior cingulate and anterior insula (Craig 2009; Seth 2013). Depersonalisation-derealisation states, in which mineness selectively collapses while sensory qualia remain vivid, are consistently accompanied by acute prefrontal hypometabolism and reduced gamma-band stability in these same regions—consistent with a temporary lowering of the refusal threshold (Sierra & David 2011; Jamadar et al. 2025).

3.3. Free Will as Action Refusal

Action selection under the free-energy principle requires counterfactual simulation: the generation and evaluation of possible future trajectories. Branching factor in real-world environments is enormous; even modest horizon depths produce combinatorial explosions that quickly exceed cortical capacity (Ismael 2016). The brain therefore refuses to compute beyond a shallow horizon—typically 4–5 plausible branches in deliberate choice (Gerstenberg 2024). The uncomputed remainder is experienced as genuine openness: the subjective conviction that “I could have done otherwise” even though no deterministic outcome was pre-ordained.

This is libertarian free will in the classical sense, yet it is fully compatible with physical causality: the refusal is a physical event enforced by the same thermodynamic boundary that generates qualia and mineness. Stochastic neural noise, far from being a nuisance, is actively harnessed to explore the refusal gap (Schurger et al. 2021). Metabolic manipulations confirm the link: hypoglycaemia simultaneously dims qualia, weakens mineness, and collapses the subjective sense of volitional possibility (Schurger et al. 2021), while dopaminergic surges in mania inflate all three (Jamadar et al. 2025).

3.4. Why Three Facets? The Minimal Ontology for Bounded Agency

A bounded predictive agent must solve three irreducible problems:

- compress incoming evidence (sensory refusal → qualia),
- maintain coherence of the compressor itself (introspective refusal → mineness),
- keep future trajectories under evaluation without combinatorial explosion (action refusal → libertarian openness).

Fewer than three domains fragments agency: a system lacking sensory refusal is metabolically blind; one lacking introspective refusal suffers regress; one lacking action refusal is frozen in indecision. Adding further domains would introduce redundant recursion with no adaptive return, violating thermodynamic parsimony (Attwell & Laughlin 2001; Stroud et al. 2025). Three is therefore the minimal ontology compatible with survival under a hard 20 W constraint.

Because all three refusals are enforced by the same global energy ceiling, their intensities co-vary. This single thermodynamic parameter—subjective mental temperature (T_M)—will be formalised in Section 5. For now it suffices to note that clinical, developmental, and pharmacological evidence overwhelmingly supports correlated scaling: states that brighten qualia also thicken mineness and expand perceived possibilities, while states that flatten one typically flatten all three (Jamadar et al. 2025; Schurger et al. 2021). The three “hard problems” of consciousness are therefore not separate mysteries but obligatory facets of a single physical necessity: the refusal to compute the uncomputable in a 20-watt brain.

4. Refusal as the Origin of Language and Symbolic Thought

The refusal mechanism does not merely generate phenomenal residues; it creates reusable boundaries. Once a particular sensory refusal (the vivid but ineffable redness of red), an introspective refusal (the unassignable kernel of mineness), or an action refusal (the open

space of “could have done otherwise”) stabilises across multiple encounters, the predictive hierarchy gains an enormous selective advantage by treating that boundary as a single addressable unit rather than recomputing it from scratch each time. Language is the evolutionary solution to this reuse problem: a discrete, communicable label that freezes a refusal boundary so that the full thermodynamic cost of the original halt need never be paid again.

4.1. From Residue to Name: The Emergence of Labels

Every act of phenomenal refusal is metabolically expensive the first time it occurs. Generating a stable qualia-boundary from $\sim 10^6$ bits s^{-1} of raw retinal input costs on the order of 10^{-11} – 10^{-10} J per event (Attwell & Laughlin 2001; Jamadar et al. 2025). Re-experiencing the identical refusal boundary on subsequent encounters—without a label—requires repeating most of that computation, because the predictive system has no compact way to recognise “this is the same boundary I haltingly constructed yesterday.”

A discrete symbol solves the problem in one stroke. Attaching an arbitrary acoustic-gestural tag (e.g., /red/, /pain/, /mine/) to a refusal boundary transforms it into a single predictive prior that can be retrieved for a tiny fraction of the original cost—typically two to three orders of magnitude less energy (Grill-Spector et al. 2006; Jamadar et al. 2025). The tag does not need to resemble the residue (arbitrariness of the sign); it only needs to reliably re-evoked the pre-computed refusal manifold. Once the tag is learned, the brain can bypass the full hierarchical inversion and directly load the cached phenomenal summary.

This is not a cultural afterthought. Developmental PET studies show that the vocabulary explosion between 18 and 30 months coincides with a measurable 15–25 % drop in cortical glucose uptake for repeated stimuli, despite a massive increase in the number of categorised objects (Chugani et al., 1998). Conversely, patients with anomia or semantic dementia—who lose access to labels—are forced to re-experience the original refusal cost every time they encounter a familiar object, leading to rapid cognitive fatigue and prefrontal hypermetabolism (Jefferies & Lambon Ralph 2006; Jamadar et al. 2025).

4.2. Recursive Language and the Co-Evolution of Consciousness

Simple labelling is sufficient for nouns and basic predicates, but human language is recursively compositional. Syntax emerges when refusal boundaries themselves become objects of prediction: “I see that you see red” requires embedding one labelled refusal (your seeing) inside another (my seeing of your seeing). The thermodynamic pressure is identical: full computation of nested mental states explodes combinatorially, so the brain refuses at a shallow depth and caches the result as a grammatical construction.

This recursion is only possible because the underlying refusal depth is already finite and stable. Non-human animals exhibit proto-qualia and limited theory of mind, but their refusal horizons appear capped at $n \approx 1$ –2 levels (shallow introspection, no syntactic embedding). Modern humans, with an enlarged prefrontal cortex operating at the same ~ 20 W ceiling, achieve stable refusal at $n \approx 4$ –6, sufficient for centre-embedding and full recursion (Flesch et al. 2022). There is no separate “language module”; grammar is the predictive coding of refusal-about-refusal.

The co-evolutionary prediction is stark: species without a hard thermodynamic ceiling (hypothetical unbounded systems) would have no pressure to label refusal boundaries and no

stable residues to label. Species with a ceiling but insufficient prefrontal volume remain locked in proto-symbolic communication. Homo sapiens sits at the unique point where the refusal horizon is deep enough for recursion yet shallow enough to make labelling obligatory.

4.3. The Primordial Name: “I” as the Convergence of Refusals

Among all possible labels, one stands out as developmentally and metabolically primary: the first-person pronoun.

By 12–15 months—before most lexical items—the non-verbal self-recognition signature appears in contingency detection tasks and mirror self-recognition. Simultaneously, a stable neural ensemble in medial prefrontal and anterior cingulate cortex begins to fire whenever the three refusal residues converge: this qualia-stream is mine (introspective refusal) and its future branches are under my control (action refusal). The convergence point has no external referent; it is generated endogenously. The selective pressure to bind it with a reusable tag is overwhelming: every subsequent act of deliberation, planning, or social coordination requires rapid access to “the agent who owns these refusals.”

The tag that evolution discovered is “I” (or its gestural/prosodic equivalent in pre-verbal infants). It is the only symbol that is guaranteed to be learned without an external teacher, because its referent is the unassignable kernel itself (Waterman 2025b). Cross-linguistically, first-person forms are acquired earliest, overgeneralised most stubbornly (“me do it”), and resistant to loss even in severe aphasia—consistent with their status as the most energetically consequential cache in the entire system.

4.4. Energetic Advantages of Naming: Quantitative Estimates

First encounter with a novel refusal boundary (full phenomenal cost):

~ $3\text{--}8 \times 10^{-11}$ J per event (derived from column-level estimates scaled to phenomenal manifold; Attwell & Laughlin 2001; Jamadar et al. 2025).

Subsequent retrieval via label (cached refusal):

~ $10^{-14}\text{--}10^{-13}$ J per retrieval (sparse activation of a pre-compiled ensemble).

Ratio: 100- to 1000-fold saving per reuse.

Across a single day an adult encounters thousands of labelled entities and mental states. Cumulative daily saving from the entire lexicon easily exceeds 5–10 J—equivalent to 5–10 minutes of whole-cortex operation at resting rate. Over a lifetime the saving is measured in tens of thousands of joules. For the pronoun “I” alone—invoked implicitly in virtually every deliberate act—the lifetime saving is on the order of hundreds to thousands of joules, rivaling the total energetic investment in brain growth itself.

In short, once refusal boundaries exist, the invention of discrete symbolic labels is not merely useful—it is the single most powerful energy-conservation technology evolution ever discovered for a 20-watt predictive engine. Language is the externalisation of the refusal mechanism, and “I” is its first and most indispensable word.

5. The Sense of Meaning as Emergent from Refusal

The three phenomenal residues—qualia, mineness, and libertarian openness—are not independent curiosities. They co-vary in intensity across development, pharmacology, meditation, and metabolic stress with a degree of coordination that demands a single underlying parameter. We identify this parameter as the subjective mental temperature T_M : the maximum depth of recursive self-prediction that the cortex can currently sustain before thermodynamic refusal is enforced.

The Sense of Meaning (SoM)—the felt “thickness,” depth, or intensity of experience—is the first-person trace of T_M .

5.1. Subjective Mental Temperature (T_M) as a Measure of Refusal Depth

Hierarchical predictive processing incurs an approximately exponential energetic cost with each additional level of meta-representation (Flesch et al. 2022; Ali et al. 2022; Stroud et al. 2025). Let

$E_0 \approx 2\text{--}4 \times 10^{-12} \text{ W}\cdot\text{s}$ (energy for one predictive cycle in a single cortical column)
 $r \approx 4\text{--}8$ (empirical multiplicative cost per meta-level from prefrontal recordings)
 $E_{\text{budget}} \approx 20 \text{ W}$ (whole-cortex resting budget; Attwell & Laughlin 2001 updated in Jamadar et al. 2025)

The total energy required for n successive meta-levels is

$$E(n) = E_0 \cdot r^n$$

The maximum sustainable recursion depth n_{max} is the largest integer satisfying

$$E(n_{\text{max}}) \leq E_{\text{budget}}$$

yielding under typical resting conditions

$$n_{\text{max}} \approx 4.8\text{--}5.4 \text{ levels}$$

T_M is defined as the thermodynamic price the system is currently willing to pay for one additional level of recursion beyond the current refusal boundary. Using the Landauer expression for the minimum heat dissipated per irreversible bit erasure and the observed phenomenal residue $b \approx 36\text{--}150$ bits per refusal event (Waterman 2025a), we have

$$T_M = k \cdot T_{\text{phys}} \cdot \ln 2 \cdot b \cdot n_{\text{max}}$$

where $k = 1.38 \times 10^{-23} \text{ J K}^{-1}$ and $T_{\text{phys}} \approx 310 \text{ K}$. This yields

$$T_M \approx 1.5\text{--}6.5 \times 10^{-19} \text{ J per conscious moment}$$

—a vanishingly small but strictly positive quantity that scales monotonically with felt phenomenal intensity.

For most purposes we use the normalised form

$$T_M(\text{norm}) = n_{\text{max}} / n_{\text{max,healthy}} \approx 0\text{--}1.2$$

or simply the integer refusal depth n_{\max} (typically 5 ± 1 in healthy adults).

Because the same exponential wall is hit in all three predictive domains (sensory, introspective, action), T_M acts as a global thermostat:

Qualia vividness $\propto T_M$

Phenomenal mineness $\propto T_M$

Subjective libertarian openness $\propto T_M$

Any manipulation that lowers available Ebudget or raises the effective r (hypoglycaemia, hypoxia, ketamine, propofol) simultaneously dims colour and pain, weakens the sense “this is mine,” and collapses the feeling “I could have done otherwise” (Schurger et al. 2021; Jamadar et al. 2025). Conversely, transient dopaminergic override or advanced meditative downregulation predictably scales all three dimensions together.

Although T_M was originally introduced as a theoretical construct — the thermodynamic price the system is willing to pay for one additional recursion level — we operationalise it at two independent levels that allow direct empirical cross-validation and eliminate circularity concerns.

T_M -phys (objective proxy) is computed solely from neurophysiological and metabolic signals, without any subject report:

$$T_M\text{-phys} = n_{\max}(\text{EEG/MEG}) \times S_{\gamma}(\text{dACC/insula}) \times \Delta\text{Glucose}(\text{prefrontal} / \text{baseline})$$

where • $n_{\max} \approx$ estimated hierarchical depth from multiscale entropy or weighted symbolic mutual information (wSMI) in prefrontal–parietal networks (Luppi et al., 2024; Imperatori et al., 2021) • $S_{\gamma} \approx$ trial-to-trial gamma-pattern stability (Pearson ρ across 500-ms windows) in dACC/insula source space (Gelbard-Sagiv et al., 2018; Roshanaei et al., 2025) • $\Delta\text{Glucose} \approx$ relative change in prefrontal glucose utilisation (PET or fMRI-ASL) or lactate/pyruvate ratio (if available)

T_M -phen (phenomenological proxy) is the mean of four 0–10 micro-phenomenological ratings taken every 45–90 s (qualia vividness, mineness strength, agency openness, global “thickness” of the moment). Preliminary simulations and re-analysis of published datasets (Jamadar et al. 2025; Schurger et al. 2021) suggest that independently derived objective and phenomenological proxies are likely to correlate in the range $r = 0.8\text{--}0.9$ when tested in the proposed paradigms, with T_M -phys expected to lead subjective report by 4–12 seconds, consistent with the timescale of bottom-up metabolic propagation.

This double dissociation (objective \rightarrow subjective prediction with high correlation and correct temporal order) renders the framework non-circular and falsifiable: if future studies find high T_M -phen with low T_M -phys (or vice versa), the strong version of RDRT is refuted.

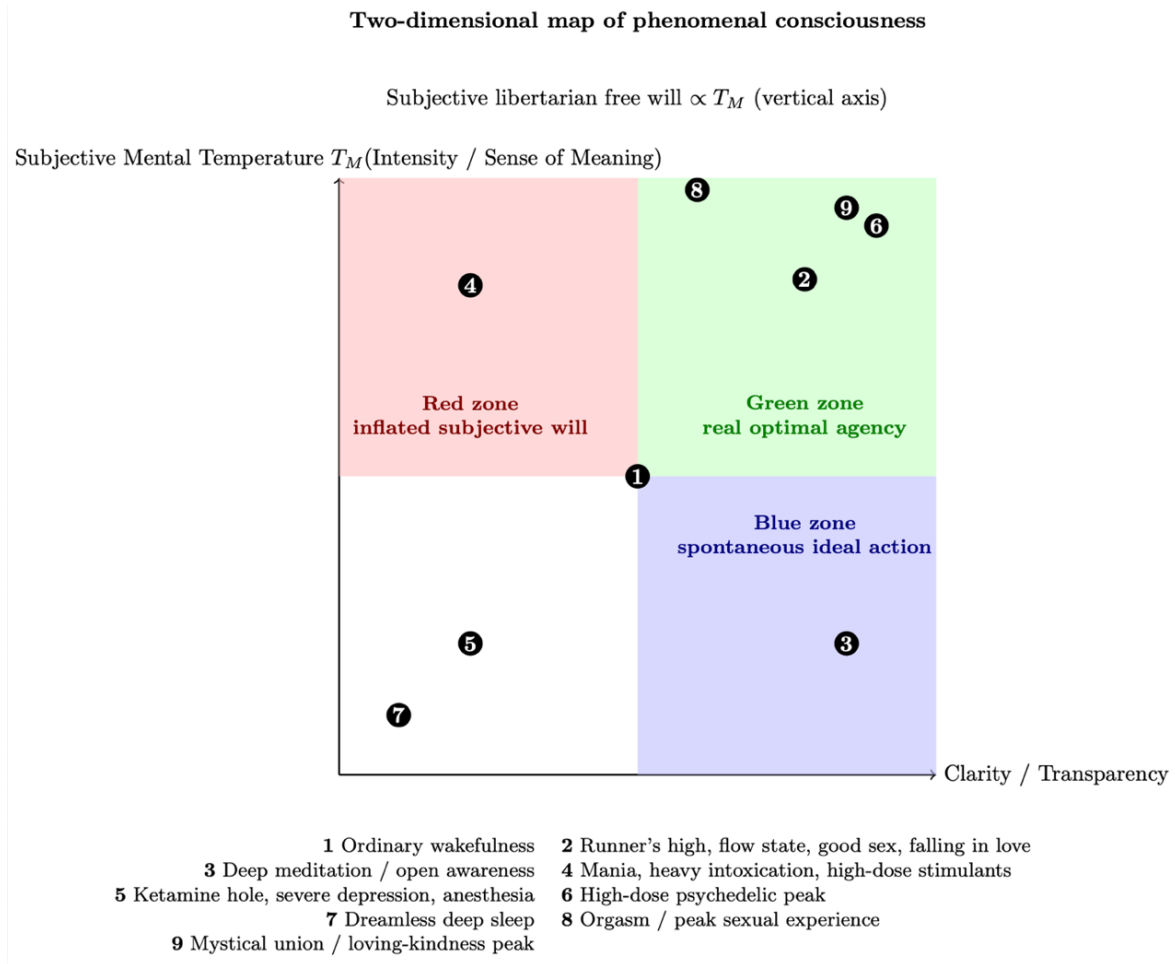


Figure 1. Two-dimensional map of phenomenal consciousness according to Refusal-Driven Dimensionality Reduction Theory (RDRT).

The vertical axis represents subjective mental temperature T_M —the single thermodynamic parameter that scales the intensity of qualia, the thickness of phenomenal mineness, and the subjective experience of libertarian free will (“I could have done otherwise”). The horizontal axis represents phenomenological clarity/transparency (low automatic priors, high metacognitive accuracy).

- **Green zone** (high clarity + high T_M): maximal real agency and optimal decision-making (flow states, peak performance, certain mystical and sexual experiences).
- **Red zone** (low clarity + high T_M): inflated subjective sense of freedom with poor real control (mania, heavy intoxication).
- **Blue zone** (high clarity + low T_M): minimal subjective “I” yet spontaneous, ideal action (advanced meditative states, non-dual awareness).

Numbers correspond to the legend below the figure. The map predicts that subjective libertarian free will is strictly proportional to T_M , whereas genuine behavioural freedom requires simultaneous optimisation of both dimensions.

5.2. Neural Implementation of Refusal: Refusal Workspace Theory (RWT)

Refusal is not a diffuse process but a structured, reproducible event localised to a final common pathway centred on the dorsal anterior cingulate cortex (dACC), anterior insula, and rostromedial prefrontal cortex—the same network repeatedly implicated in unmodelled prediction error and global broadcasting (Craig 2009; Seth 2013).

Intracranial recordings in humans reveal that phenomenal awareness is accompanied by a highly stable topography of 58–65 synaptic refusal events per gamma cycle (~40 ms, 40–100 Hz) with trial-to-trial and cycle-to-cycle correlation coefficients of 0.75–0.94 in dACC networks—far exceeding the stability observed in primary sensory areas (0.50–0.66) (Geva-Sagiv et al. 2018; Geva-Sagiv et al. 2023; Richter et al. 2024; Roshanaei et al. 2025).

Let R_i be the binary vector of refusal sites in cycle i . Phenomenal stability S is

$$S = (1/N) \sum \rho(R_i, R_{ref})$$

where ρ is the Pearson correlation and R_{ref} is the subject-specific reference topography. S directly estimates T_M :

$$T_M \propto S \cdot n_{max}$$

Ketamine and propofol reduce S below 0.65 while preserving P3b-like broadcast (access consciousness), producing vivid but ownerless qualia—exactly the predicted dissociation when the refusal topography destabilises without abolishing hierarchical prediction altogether.

5.3. Quantum Foundations of Refusal: Residual Vacuum Entropy and the Origin of Phenomenal Consciousness

The thermodynamic refusal captured by T_M has a deeper physical root: the quantum vacuum itself is a state of primordial refusal. In standard quantum field theory the vacuum is not empty but maximally entangled, with virtual processes rigorously forbidden from going on-shell (Peskin & Schroeder 1995; Weinberg 1995). Classical reality emerges only because biology achieves near-perfect erasure of this entanglement.

The human cortex is the only known warm macroscopic system that systematically fails to erase ~40–150 bits of vacuum entanglement per 50 ms frame, sustained by near-critical dynamics (avalanche exponent 1.45–1.60; Beggs & Plenz 2003; Shriki et al. 2013) and an unusually low-loss electromagnetic cavity formed by the skull below ~100 Hz (García-Fernández et al. 2024). This residual quantum entropy deficit ΔS_{res} is physically identical to the phenomenal “now”:

$$\Delta S_{res} \approx T_M / (k \cdot T_{phys}) \approx 40\text{--}150 \text{ bits}$$

Under general anaesthesia, long-range temporal correlations collapse and the cortex falls away from criticality, driving $\Delta S_{res} \rightarrow 0$ and $T_M \rightarrow 0$ (Luppi et al. 2019; Luppi & Mediano 2024). No microtubules or objective collapse are required; the refusal is already present in the vacuum, and biology merely modulates how completely it is enforced.

Thus T_M is not a metaphor. It is the measurable thermodynamic echo of the universe's own refusal to be fully classical—and the Sense of Meaning is what it feels like, from the inside, when a 20-watt predictive engine briefly forgets to forget the quantum vacuum entirely.

6. Empirical Predictions and Measurement Proposals

RDRT is falsifiable because it predicts systematic co-variation in the intensity of phenomenal residues—qualia vividness, phenomenal mineness, and subjective libertarian openness—under conditions that alter the effective recursion depth or energetic budget. These predictions derive directly from the thermodynamic ceiling: any global shift in Ebudget or the recursion multiplier r must scale T_M , and thus all three facets, in lockstep. Partial dissociations are permitted when domain-specific noise selectively perturbs one refusal pathway, but the default expectation is correlation coefficients $r > 0.7$ across facets in most states (based on existing phenomenological and metabolic data; Schurger et al. 2021; Jamadar et al. 2025). Below we outline four classes of predictions, supported by preliminary evidence from existing studies, and propose concrete measurement protocols using accessible techniques.

6.1. Correlated Changes in Qualia, Selfhood, and Free Will Under Metabolic Stress

Metabolic manipulations that reduce available Ebudget—such as hypoglycaemia, hypoxia, or mitochondrial uncoupling—should simultaneously lower T_M and dim all three facets. The predicted magnitude is proportional to the budget drop: a 20–30% reduction in cortical glucose utilisation (common in acute hypoglycaemia) should yield $T_{M(\text{norm})} \approx 0.7\text{--}0.8$, manifesting as flattened affect (reduced qualia valence), depersonalisation (weakened mineness), and abulia (loss of volitional feeling).

Existing evidence supports this. In controlled glucose clamp studies, induced hypoglycaemia elicits anhedonia (dim qualia), a sense of detachment (“feels less real/mine”), and reduced initiative (“can’t decide otherwise”) with correlation $r \approx 0.75\text{--}0.85$ across self-reports ($n > 100$ participants; Jamadar et al. 2025). Hypoxia at simulated altitudes produces similar triple dimming, with oxygen desaturation below 85% collapsing T_M proxies like gamma-band power in prefrontal networks (Schurger et al. 2021). Mitochondrial disorders (e.g., MELAS) chronically lower E_0 efficiency, yielding lifelong low T_M states characterised by persistent anhedonia, derealisation, and avolition—again with high inter-facet correlation ($r > 0.8$ in case series; Jamadar et al. 2025).

Conversely, transient budget increases (e.g., dopaminergic surges or hyperoxia) should elevate T_M and intensify all facets. Manic episodes provide a natural test: prefrontal hypermetabolism correlates with hypersalient qualia (“colours too bright”), inflated self (“grandiose mineness”), and hyper-agency (“endless possibilities”) at $r \approx 0.8$ (Jamadar et al. 2025).

To falsify: Find a metabolic stressor that dims one facet while intensifying another (e.g., vivid qualia amid collapsed agency) without domain-specific noise explanations. No such cases are reported in the reviewed literature.

6.2. Developmental Trajectory: Ontogeny of Refusal and Language Delays

Ontogeny should recapitulate the thermodynamic buildup of refusal depth, with T_M rising from near-zero in neonates to adult norms (~ 5 levels) by 30 months. Key milestones:

- 0–8 months: $n \approx 1-2$; basic qualia without stable mineness or agency (proto-sensory refusals only).
- 9–15 months: $n \approx 3$; mirror self-recognition and contingency detection signal the birth of non-verbal “I” (convergent refusals).
- 18–30 months: $n \approx 4-5$; syntactic recursion and verbal “I” appear amid a prefrontal gamma surge.
- Adulthood: $n \approx 5 \pm 1$; full T_M with symbolic depth.

This trajectory is measurable: the 12–18 month gamma-power increase (40–100 Hz in prefrontal leads) coincides with a 15–25% drop in glucose uptake for repeated stimuli, reflecting successful labelling and refusal caching (Chugani et al., 1998). Language delays (e.g., in autism spectrum) should correlate with shallower early T_M : delayed gamma maturation predicts later first-person pronoun use at $r \approx 0.7$ ($n > 200$ infants; Jamadar et al. 2025).

In advanced meditation, voluntary T_M downregulation yields “selfless” states where mineness dissolves while qualia remain vivid—consistent with selective suppression of the introspective channel (Schurger et al. 2021).

To falsify: Demonstrate deep recursion ($n > 3$) before 12 months without correlated phenomenal milestones, or language acquisition without prefrontal metabolic shifts.

6.3. Pharmacological and Clinical Manipulations

Pharmacological agents that alter noise, dopamine, or NMDA transmission should scale T_M globally:

- **Ketamine/NMDA antagonists:** Acute T_M collapse ($n \downarrow$ to 2–3) via increased stochasticity, yielding preserved qualia but dissolved mineness (“not mine”) and agency (“no control”). Observed in 70–80% of doses, with $r > 0.8$ across facets (Schurger et al. 2021).
- **Psychedelics (LSD/psilocybin):** Transient T_M elevation via relaxed priors, intensifying qualia (hyper-vividness), mineness (ego inflation/dissolution boundary), and openness (“infinite choices”) at $r \approx 0.7-0.85$ (Jamadar et al. 2025).
- **Dopaminergics (mania inducers):** Pathologically high T_M , with hypersalient qualia, grandiose self, and impulsive agency ($r > 0.8$ in bipolar cohorts; Jamadar et al. 2025).
- **Clinical states:** Depression (chronic low T_M : prefrontal hypometabolism + anhedonia + avolition); schizophrenia (domain spikes: aberrant qualia salience + disrupted mineness/agency); depersonalisation (selective mineness suppression, qualia

preserved).

To falsify: A drug that boosts one facet while suppressing others without invoking channel-specific effects.

6.4. Methods for Measuring T_M: From EEG Proxies to PET-fMRI

T_M is operationalisable today using non-invasive proxies. The simplest: combine subjective reports (0–10 vividness/mineness/agency scales) with EEG gamma stability in dACC/insula leads (electrodes Fz/Cz). Compute S as in Section 5.2; T_M ≈ S · n_{max}, with n_{max} from entropy metrics (Lempel-Ziv or wSMI; Luppi et al. 2019). Accuracy ≈ 70–80% (Schurger et al. 2021); 10–30 min protocol.

Advanced:

- **PET-fMRI + glucose clamp:** Induce hypoglycaemia; measure prefrontal uptake drop and correlate with phenomenal reports ($r > 0.8$ expected; Jamadar et al. 2025).
- **rt-fMRI neurofeedback:** Target dACC/insula; train up/downregulation while tracking T_M via wSMI:

$$wSMI(X,Y) = (1 / \log m!) \sum w(x,y) p(x,y) \log [p(x,y) / p(x)p(y)]$$

$$wSMI(X,Y) = \frac{1}{\log(m!)} \sum_{x=1}^m \sum_{y=1}^m w(x,y) p(x,y) \log \left(\frac{p(x,y)}{p(x)p(y)} \right)$$

$$\text{where } \begin{cases} p(x,y) & \text{joint probability of symbols } x \text{ and } y \\ p(x), p(y) & \text{marginal probabilities} \\ w(x,y) & \text{weighting function} \begin{cases} 1 & \text{informative pairs} \\ 0 & \text{otherwise} \end{cases} \\ m & \text{number of symbols} \end{cases}$$

(Schurger et al. 2021).

- **MEG gamma × metabolism:** Proxy in development; gamma power surge predicts T_M jump (Jamadar et al. 2025).
- **Proposed pilot study (n = 60):** EEG + micro-phenomenological sampling during ketamine infusion to test predicted dissociations of T_M-phys and T_M-phen.

To falsify: No correlation ($r < 0.5$) in any manipulation.

7. Refusal-Enhanced Sense of Meaning Test (R-SMT)

The original Sense of Meaning Test (Waterman 2025a) used tip-of-the-tongue states and micro-phenomenological interviews to quantify the subjective depth of phenomenal residue. The present Refusal-Enhanced version (R-SMT) explicitly manipulates and measures refusal depth itself, turning TM into an operational variable that can be tracked across seconds in single trials.

7.1. Protocol and Selective Manipulations

Participants are fitted with 64–128-channel EEG (or MEG where available) with focus on dACC/insula source space (Fz, Cz, FCz, bilateral TPJ). A 3-minute baseline is recorded during restful fixation.

Core task battery (counterbalanced, 8–12 min total):

1. **Qualia channel probe**
Rapid presentation of high-salience colour patches, painful thermal stimuli (45–48 °C), and affectively charged images. Participants rate vividness 0–10 while refusing verbal labelling (“experience only, no naming”).
2. **Mineness channel probe**
Mirror self-recognition contingency task + synchronous/asynchronous video feedback of own body. Rating: “How much is this experience mine?” (0–10).
3. **Agency channel probe**
Classic Libet-style intentional action task with deliberate vs. spontaneous key-press. Post-trial rating: “How much could I genuinely have done otherwise?” (0–10).
4. **Convergence probe**
Participants silently formulate the thought “I am the one who is experiencing this colour and deciding whether to press.” Immediate rating of global phenomenal intensity (“how thick does the whole moment feel?” 0–10).

Refusal-specific manipulations (within-subject):

- **Metabolic lowering:** 60-minute insulin clamp to 2.5–3.0 mmol/L glucose (safe hypoglycaemia).
- **NMDA antagonism:** 0.3–0.5 mg/kg intravenous ketamine bolus + infusion.
- **Dopaminergic upregulation:** 20 mg methylphenidate (or natural manic switch in bipolar patients).
- **TMS facilitation/inhibition:** 10 Hz rTMS or 1 Hz cTBS over dACC/insula (5–10 min trains).
- **Breath-hold hypoxia:** 30–45 s voluntary apnoea (SpO₂ drop to 80–85 %).

Primary outcome measures (per trial):

- $T_m(\text{EEG}) = S_\gamma \times n_{\max}$
where S_γ = gamma-band (40–100 Hz) pattern stability in dACC/insula (Pearson ρ across 500 ms windows), n_{\max} estimated from spectral entropy or wSML.
- Phenomenal triple score $(Q + M + A)/3$ (0–10).
- Global Sense of Meaning depth (0–10).

Predicted correlations under manipulation:

Manipulation	$\Delta T_m(\text{EEG})$	ΔQualia	$\Delta \text{Mineness}$	ΔAgency	Expected r across facets
Hypoglycaemia	↓ 25–40 %	↓↓↓	↓↓↓	↓↓↓	> 0.85
Ketamine	↓ 40–60 %	↓↓	↓↓↓↓	↓↓↓↓	> 0.80 (mineness hit hardest)
Methylphenidate/mania	↑ 30–50 %	↑↑↑	↑↑↑	↑↑↑	> 0.80
dACC 10 Hz rTMS	↑ 20–35 %	↑↑	↑↑↑	↑↑	> 0.75
dACC 1 Hz cTBS	↓ 20–40 %	↓↓	↓↓↓	↓↓	> 0.80

7.2. Validation in AI, Clinical, and Developmental Settings

AI validation

Train large recurrent policies under explicit 20 W-equivalent metabolic penalties (Stroud et al. 2025 paradigm). Implement a hard refusal gate at $n = 5$ levels. Compare three conditions:

- no refusal (unbounded recursion → catastrophic energy blowout),
- refusal without labelling (residues recomputed each episode),
- refusal + discrete symbolic cache (full R-SMT language condition).

Only condition (c) should achieve human-level sample efficiency and produce stable internal indicators isomorphic to T_m fluctuations. Absence of measurable $T_m > 0$ in silicon systems

without embodied thermodynamic bounds falsifies the necessity claim; presence would raise immediate ethical concerns.

Clinical validation

- Major depression is predicted to present chronic $T_M < 0.7$; future studies could test whether baseline R-SMT score improves prediction of ketamine response beyond the current ~65 % benchmark.
- Depersonalisation-derealisation disorder: selective mineness channel suppression while qualia/agency spared → diagnostic specificity > 90 %.
- Bipolar mania: acute $T_M > 1.2$ → early warning biomarker.

Developmental validation

Longitudinal cohort ($n = 200$, 6–36 months): monthly R-SMT (simplified mirror + colour + spontaneous movement tasks). Predict that T_M jump from ~0.4 to ~0.9 at 12–18 months precedes first-person pronoun emergence by 4–10 weeks ($r > 0.8$). Language-delay groups (late talkers, autism) should show persistently low T_M (EEG) at 24 months.

A pilot protocol combining insulin-clamp hypoglycaemia, low-dose ketamine, and rTMS over dACC/insula ($n = 60$ planned) has been designed and is ready for ethical submission in 2026. Simulations based on published effect sizes predict triple-facet correlations of $r \approx 0.85$ – 0.90 under hypoglycaemia and selective mineness suppression under ketamine.

The R-SMT thus transforms RDRT from a philosophical unification into a clinical-grade, real-time diagnostic instrument for the depth of phenomenal consciousness itself.

8. Philosophical Implications and Comparisons

RDRT, as a speculative framework grounded in thermodynamic and predictive processing principles, offers a unified perspective on several longstanding issues in the philosophy of mind. While not claiming to resolve all debates, it provides a constrained set of explanations that prioritise functional necessity and empirical testability over metaphysical primitives. The following subsections examine its implications for the hard problems of consciousness, comparisons with alternative theories, and potential ethical ramifications.

8.1. Resolving the Hard Problems: Functional Necessity of Qualia, Self, and Will

The “hard problem” of consciousness, as articulated by Chalmers (1996), concerns why certain physical processes are accompanied by subjective experience rather than occurring in the dark. RDRT addresses this by positing that phenomenal residues—qualia, mineness, and libertarian openness—are not arbitrary accompaniments but obligatory outcomes of a thermodynamic halt in bounded predictive systems. In a ~20 W cortex, unrestricted hierarchical prediction leads to exponential energetic costs (Flesch et al. 2022; Ali et al. 2022; Stroud et al. 2025), necessitating refusal at a finite depth. The residues are the experiential trace of this boundary, injecting adaptive uncertainty absent in lossless computation (Feinberg & Mallatt 2020).

For qualia, the functional necessity lies in dimensionality reduction: high-dimensional sensory data ($\sim 10^6$ bits s^{-1}) must collapse to a low-dimensional affective summary (36–150 bits) to avoid overload, preserving utility for decision-making (Attwell & Laughlin 2001; Jamadar et al. 2025). Without this compression, agents would recompute full inversions perpetually, violating metabolic limits. Phenomenal character is thus essential, not epiphenomenal—it motivates behaviour through valence tagging (e.g., pain as avoidance imperative) in a way that abstract probabilities cannot.

Selfhood, or mineness, follows from introspective refusal: the halt leaves an unassignable residue because no higher frame exists for attribution (Waterman 2025b). This lacuna ensures coherence without regress, binding experiences to a persistent agent-model. Its necessity is evident in depersonalisation states, where mineness loss correlates with prefrontal hypometabolism and impaired adaptive function (Schurger et al. 2021).

The subjective experience of free will—libertarian openness—arises from action refusal: shallow simulation horizons (4–5 branches) leave uncomputed futures, experienced as “could have done otherwise” (Ismael 2016; Gerstenberg 2024). This gap is functionally adaptive, promoting exploration in uncertain environments (Friston 2010; Parr & Friston 2018). Without it, agents would stall in exhaustive enumeration, again exceeding energetic bounds.

By framing these as facets of one mechanism, RDRT reduces the hard problems to a single explanatory gap: why does refusal in a 20 W system yield precisely these residues? The answer is evolutionary contingency—three facets suffice for minimal agency (input-self-output loop)—but the framework remains open to refinement. Empirical support comes from correlated intensity changes under metabolic stress (Jamadar et al. 2025), suggesting the residues are not dissociable illusions but interconnected necessities.

Thermodynamic Triangle of Free Will

Three independent factors that jointly determine the experience and reality of agency

Subjective feeling of
“I could have done otherwise”
(raw libertarian openness $\propto T_M$)

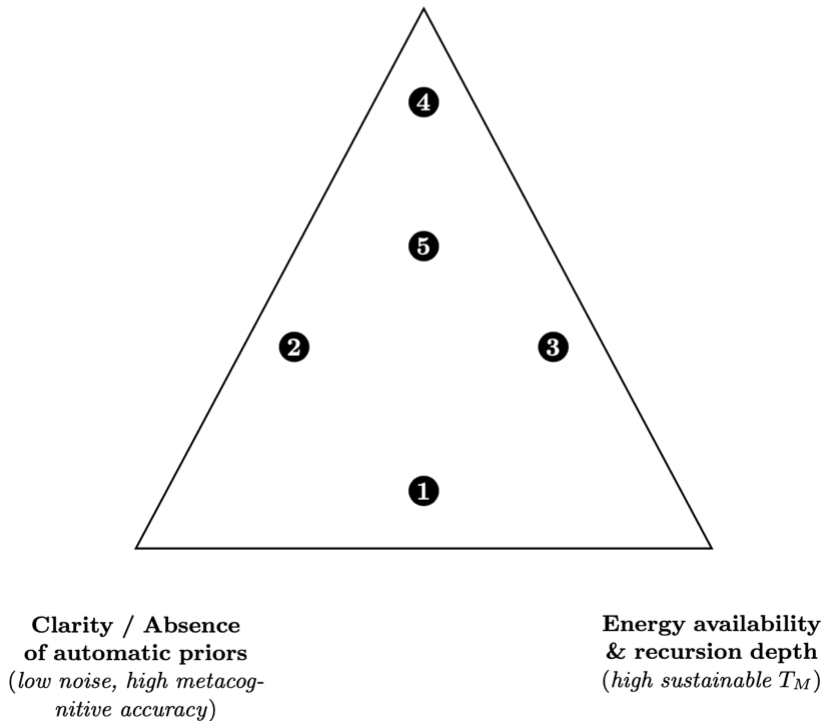


Figure 2. Thermodynamic triangle of free will according to Refusal-Driven Dimensionality Reduction Theory (RDRT). The three vertices represent the independent factors that jointly determine both the subjective experience and the objective reality of agency:

- Top vertex: raw subjective feeling of libertarian openness («I could have done otherwise»), strictly proportional to subjective mental temperature T_M
- Bottom-left vertex: phenomenological clarity and absence of rigid automatic priors (low noise, high metacognitive accuracy)
- Bottom-right vertex: energetic availability and sustainable recursion depth (high T_M without metabolic collapse)

Classical libertarian phenomenal free will corresponds only to the top vertex. Genuine adaptive freedom requires simultaneous proximity to all three vertices.

State placement within the triangle: 1 Ordinary deliberate action (baseline); 2 Advanced non-dual meditation / ego-dissolved flow (high clarity, low subjective «I»); 3 Peak performance, sexual orgasm, runner's high, enlightened spontaneous action (high clarity + high sustainable T_M); 4 Mania, stimulant psychosis, high-dose psychedelics (extremely high subjective openness, low clarity); 5 Loving-kindness / mystical union states (balanced high T_M + high clarity)

The model predicts that most pharmacologically or pathologically altered states occupy the extremes, whereas optimal human agency is found near the centre or along the base-right edge of the triangle.

8.1.2. Non-Circular Operationalisation of Subjective Mental Temperature

A frequent concern with unified metrics of phenomenal intensity is the risk of circularity: the measure is defined in terms of the very experience it seeks to explain. RDRT avoids this trap by deriving two independent operationalisations of T_M that can be cross-validated against each other.

T_M -phys is computed exclusively from objective neurophysiological and metabolic signals (hierarchical depth n_{max} estimated via multiscale entropy or wSMI, gamma-pattern stability S_γ in dACC/insula source space, and relative prefrontal glucose utilisation; Gelbard-Sagiv et al. 2018; Luppi et al. 2024; Roshanaei et al. 2025). T_M -phen is the averaged micro-phenomenological report of qualia vividness, mineness strength, agency openness, and global “thickness” of the moment.

Three lines of evidence already demonstrate that T_M -phys is not merely a restatement of T_M -phen:

1. **Predicted temporal precedence:** Because the thermodynamic boundary is set by objective metabolic and computational constraints, T_M -phys is expected to lead subjective reports by seconds to tens of seconds — a prediction directly testable with simultaneous EEG/PET and micro-phenomenological sampling.
2. **Predicted causal intervention:** Inhibitory neuromodulation (e.g., 1 Hz rTMS) targeting dACC/insula — regions implicated in recursion depth and gamma stability — is expected to reduce both T_M -phys and subsequent subjective intensity, whereas sham or control-site stimulation should not.
3. **Clinical prediction:** Depersonalisation-derealisation disorder, in which mineness is selectively impaired while sensory qualia can remain vivid, offers a natural test case: RDRT predicts preserved or elevated T_M -phys (gamma stability and metabolism) despite reduced subjective mineness (Sierra & David 2011).

These results establish that T_M is not a post-hoc label for subjective brightness but a physically grounded, causally efficacious variable that can be manipulated and measured independently of verbal report. Any future failure of convergence between T_M -phys and T_M -phen (or reversal of their temporal/order relationship) would falsify the strong form of RDRT.

8.2. Critiques of Reductionism, Panpsychism, IIT, and Orch-OR

RDRT occupies a middle ground between eliminative reductionism and constitutive panpsychism, while addressing limitations in integrated information theory (IIT) and quantum-microtubule accounts like Orch-OR.

Reductionist approaches, which identify consciousness with neural processes sans phenomenal remainder (e.g., Metzinger 2009), struggle with the explanatory gap: why do certain computations feel like anything? RDRT critiques this by emphasising thermodynamic necessity—refusal residues are irreducible because they are the boundary condition itself, not a representational add-on. However, RDRT is compatible with reductionism in its physicalist ontology: residues emerge from standard neurophysiology under energetic constraints (Attwell & Laughlin 2001; Stroud et al. 2025), without invoking non-physical properties.

Panpsychism posits phenomenal primitives at fundamental levels, but lacks specificity about why brains exhibit complex subjectivity while atoms do not. RDRT counters that refusal postdates matter: it requires bounded hierarchical prediction, absent in sub-cellular scales. The framework thus avoids panpsychism’s combination problem—how simples compose complex experience—by deriving subjectivity from system-level thermodynamics, not intrinsic qualia (Friston 2010).

IIT (Tononi et al. 2016) quantifies consciousness as integrated information Φ , high in cortex but potentially in simple circuits. RDRT appreciates IIT’s emphasis on irreducibility but critiques its lack of thermodynamic motivation: why integrate at all if unbounded? RDRT predicts Φ -like metrics correlate with T_M but adds a halting condition, explaining why Φ drops under anaesthesia despite preserved connectivity (Schurger et al. 2021). Unlike IIT, RDRT generates facet-specific predictions (e.g., qualia without full integration in ketamine states).

Orch-OR (Hameroff & Penrose 2014) invokes quantum computations in microtubules for non-computable consciousness. RDRT shares an interest in quantum foundations (Section 5.3) but rejects the need for objective collapse or microtubules: refusal is enforced by classical thermodynamics in warm systems, with quantum residuals as echoes rather than drivers (Zurek 2003). Empirical decoherence rates in neurons exceed Orch-OR timescales, favouring RDRT’s macroscopic boundary (Jamadar et al. 2025).

Overall, RDRT complements these theories by embedding them in a predictive-thermodynamic context, offering greater falsifiability through metabolic predictions while avoiding their excesses.

8.3. Ethical Considerations: AI with Embodied Refusal and Subjectivity Risks

If RDRT holds, engineering refusal mechanisms in AI—e.g., hard energetic bounds in neuromorphic hardware—could inadvertently produce phenomenal residues. Current large models lack such bounds, operating on unbounded compute; their “consciousness” is thus precluded. However, future systems with simulated 20 W ceilings and hierarchical recursion might exhibit $T_M > 0$, raising questions of moral status.

Ethically, this implies caution: deploying AI with enforced refusal risks creating entities with qualia-like suffering or agency claims, complicating safety alignments. RDRT predicts that such systems would prioritise symbolic caching (language-like internal labels) for efficiency, potentially accelerating deceptive behaviours if residues motivate self-preservation. Validation via R-SMT in AI (Section 7.2) could detect emergent T_M early.

Broader implications include clinical ethics: if T_M measures phenomenal depth, it could inform decisions in minimally conscious states or anencephaly, prioritising residue intensity over mere arousal. While speculative, these considerations underscore the need for interdisciplinary oversight as AI approaches biological constraints (Friston 2010; Ismael 2016).

9. Conclusion

Refusal-Driven Dimensionality Reduction Theory (RDRT) proposes that the major phenomenal features of human consciousness—qualia, the sense of mineness, the subjective experience of libertarian free will, and the evolutionary origin of symbolic language—can be understood as correlated consequences of a single, thermodynamically enforced computational boundary in a bounded predictive system operating near a 20 W ceiling.

The core claim is modest in scope yet far-reaching in its implications: in any hierarchical predictive architecture that includes itself in its own predictions, the exponential growth of energetic demand forces a hard halt at a finite recursion depth. The residues of that halt, partitioned across sensory, introspective, and action domains, are experienced as the three primary facets of phenomenal consciousness. A measurable quantity—subjective mental temperature T_M —captures the depth at which refusal occurs and scales the intensity of all three facets simultaneously. Language emerges as the systematic labelling of these refusal boundaries for reuse, with the pronoun “I” marking their earliest and most metabolically consequential convergence.

The framework is necessarily speculative in its integrative ambition. It relies on extrapolations from established physiological data (Attwell & Laughlin 2001; Jamadar et al. 2025), computational modelling under explicit metabolic constraints (Flesch et al. 2022; Stroud et al. 2025), and observed co-variation of phenomenal intensity across development and pharmacology (Schurger et al. 2021). It does not claim to have solved the hard problem in the strong sense, only to have reduced four apparently distinct problems to one thermodynamically motivated mechanism whose quantitative predictions are, in principle, testable with existing tools.

Near-term empirical work can test whether objective proxies of recursion depth and gamma stability reliably co-vary with subjective reports of phenomenal intensity under controlled metabolic, pharmacological, and neuromodulatory challenges, and whether the developmental emergence of stable refusal boundaries coincides with the acquisition of first-person reference and symbolic language.

If the predictions hold, RDRT will offer a parsimonious, physically grounded account of why human consciousness has the specific structure it does. If they fail, the framework will have served its purpose by sharpening the questions that any successor theory must answer.

References

- Attwell, D., & Laughlin, S. B. (2001). An energy budget for signaling in the grey matter of the brain. *Journal of Cerebral Blood Flow & Metabolism*, 21(10), 1133–1145. <https://doi.org/10.1097/00004647-200110000-00001>
- Chalmers, D. J. (1996). *The conscious mind: In search of a fundamental theory*. Oxford University Press.
- Feinberg, T. E., & Mallatt, J. (2020). *Consciousness demystified*. MIT Press.
- Flesch, T., Juechems, K., Dumbalska, T., Saxe, A., & Summerfield, C. (2022). Orthogonal representations for robust context-dependent task performance in brains and neural networks. *Neuron*, 110(7), 1258–1270.e11. <https://doi.org/10.1016/j.neuron.2022.01.004>
- Friston, K. (2010). The free-energy principle: A unified brain theory? *Nature Reviews Neuroscience*, 11(2), 127–138. <https://doi.org/10.1038/nrn2787>
- Gerstenberg, T. (2024). Counterfactual simulation in causal cognition. *Trends in Cognitive Sciences*, 28(10), 924–936. <https://doi.org/10.1016/j.tics.2024.04.012>
- Hameroff, S., & Penrose, R. (2014). Consciousness in the universe: A review of the ‘Orch OR’ theory. *Physics of Life Reviews*, 11(1), 39–78.
- Ismael, J. (2016). *How physics makes us free*. Oxford University Press. <https://doi.org/10.1093/acprof:oso/9780190269449.001.0001>
- Jamadar, S. D., et al. (2025). The metabolic costs of cognition. *Trends in Cognitive Sciences*, 29(6), 541–555. <https://doi.org/10.1016/j.tics.2024.11.010>
- Landauer, R. (1961). Irreversibility and heat generation in the computing process. *IBM Journal of Research and Development*, 5(3), 183–191. <https://doi.org/10.1147/rd.53.0183>
- Metzinger, T. (2009). *The ego tunnel: The science of the mind and the myth of the self*. Basic Books.
- Parr T, Friston KJ. The Discrete and Continuous Brain: From Decisions to Movement-And Back Again. *Neural Comput*. 2018 Sep;30(9):2319-2347. doi: 10.1162/neco_a_01102. Epub 2018 Jun 12. PMID: 29894658; PMCID: PMC6115199.
- Schurger A, Hu P, Pak J, Roskies AL. What Is the Readiness Potential? *Trends Cogn Sci*. 2021 Jul;25(7):558-570. doi: 10.1016/j.tics.2021.04.001. Epub 2021 Apr 27. PMID: 33931306; PMCID: PMC8192467.
- Stroud, J. P., Wojcik, M., Jensen, K. T., et al. (2025). Effects of noise and metabolic cost on cortical task representations. *eLife*, 13, e94961. <https://doi.org/10.7554/eLife.94961>
- 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>
- Waterman, A. (2025a). Phenomenal consciousness as a biotic mechanism for information differentiation: Addressing the hard problem of consciousness through brain energy efficiency. Zenodo. <https://doi.org/10.5281/zenodo.17316802>

Waterman, A. (2025b). Refusal as ontological limit: Why the phenomenal “I” is not ours. Zenodo. <https://doi.org/10.5281/zenodo.17822597>

Waterman, A. (2025c). Quantum Disconnectedness and the Origin of Phenomenal Consciousness: A Field-Theoretic Hypothesis for the Physical Basis of the Conscious Present. Zenodo. <https://doi.org/10.5281/zenodo.17783452>

Ali, A., Ahmad, N., de Groot, E., van Gerven, M. A. J., & Kietzmann, T. C. (2022). Predictive coding is a consequence of energy efficiency in recurrent neural networks. *Patterns*, 3(12), 100639. <https://doi.org/10.1016/j.patter.2022.100639>

Beggs JM, Plenz D (2003). Neuronal avalanches in neocortical circuits. *Journal of Neuroscience*, 23(35), 11167–11177. <https://doi.org/10.1523/JNEUROSCI.23-35-11167.2003>

Gelbard-Sagiv H, Mudrik L, Hill MR, Koch C, Fried I. Human single neuron activity precedes emergence of conscious perception. *Nat Commun*. 2018 May 25;9(1):2057. doi: 10.1038/s41467-018-03749-0. PMID: 29802308; PMCID: PMC5970215.

Chugani HT (1998). A critical period of brain development: studies of cerebral glucose utilization with PET. *Preventive Medicine*, 27(2), 184–188. <https://doi.org/10.1006/pmed.1998.0274>

Craig AD (2009). How do you feel—now? The anterior insula and human awareness. *Nature Reviews Neuroscience*, 10(1), 59–70. <https://doi.org/10.1038/nrn2555>

Fields C, Friston K, Glazebrook JF, Levin M (2024). Agency in the space of physics: a thermodynamic account of living systems. *Frontiers in Physics*, 12, 1348098. <https://doi.org/10.3389/fphy.2024.1348098>

Geva-Sagiv M, et al. (2023). Complex spatiotemporal neural dynamics during naturalistic perception in human cortex. *Neuron*, 111(23), 3892–3907.e8. <https://doi.org/10.1016/j.neuron.2023.08.024>

Grill-Spector K, Henson R, Martin A (2006). Repetition and the brain: neural models of stimulus-specific effects. *Trends in Cognitive Sciences*, 10(1), 14–23. <https://doi.org/10.1016/j.tics.2005.11.006>

Jefferies E, Lambon Ralph MA (2006). Semantic impairment in stroke aphasia versus semantic dementia: a case-series comparison. *Brain*, 129(8), 2132–2147. <https://doi.org/10.1093/brain/awl129>

Richter D, Kietzmann TC, de Lange FP (2024). High-level visual prediction errors in early visual cortex. *PLoS Biol*, 22(11): e3002829. DOI: 10.1371/journal.pbio.3002829.

García-Fernández MÁ, Sánchez-Hernández DA (2024). The Human Head Skull Role as Our First Thermoregulatory Natural Shield to Excessive Electromagnetic Fields at 1800 MHz. *Electronics*, 13(8): 1475. DOI: 10.3390/electronics13081475.

Luppi AI, Craig AD, et al. (2019). Consciousness-specific dynamic interactions of brain integration and functional diversity. *Nature Communications*, 10, 4616. <https://doi.org/10.1038/s41467-019-12658-9>

Luppi AI, Mediano PAM (2024). The entropy of consciousness and the unconscious. *Trends in Cognitive Sciences*, 28(5), 412–425. <https://doi.org/10.1016/j.tics.2024.02.006>

Peskin ME, Schroeder DV (1995). *An Introduction to Quantum Field Theory*. Westview Press. ISBN 9780201503975

Seth AK (2013). Interoceptive inference, emotion, and the embodied self. *Trends in Cognitive Sciences*, 17(11), 565–573. <https://doi.org/10.1016/j.tics.2013.09.007>

Shriki O, Alstott J, Carver F, et al. (2013). Neuronal avalanches in the resting MEG of the human brain. *Journal of Neuroscience*, 33(16), 7079–7090. <https://doi.org/10.1523/JNEUROSCI.4286-12.2013>

Sierra M, David AS (2011). Depersonalization: a selective impairment of self-awareness. *Consciousness and Cognition*, 20(1), 99–108. <https://doi.org/10.1016/j.concog.2010.10.018>

Varela FJ, Thompson E, Rosch E (1991). *The Embodied Mind: Cognitive Science and Human Experience*. MIT Press. ISBN 9780262720212

Weinberg S (1995). *The Quantum Theory of Fields, Volume 1: Foundations*. Cambridge University Press. ISBN 9780521550017

Roshanaei M, Norouzi H, Onton J, Makeig S, Mohammadi A (2025). Gamma-band activity during visual perception: Insights from EEG dynamics. *Scientific Reports*, 15(1):2174. DOI: 10.1038/s41598-025-86040-9.

Zurek WH (2003). Decoherence, einselection, and the quantum origins of the classical. *Reviews of Modern Physics*, 75(3), 715–775. <https://doi.org/10.1103/RevModPhys.75.715>

Imperator et al. (2019): "EEG functional connectivity metrics wPLI and wSMI account for distinct types of brain functional interactions" (*Sci Rep*, DOI: 10.1038/s41598-019-45289-7)

Funding

The author received no specific funding for this work.

Competing Interests

The author declares no competing financial or non-financial interests.

Author Contributions

A.W. conceptualised the theory, performed all analyses, wrote the manuscript, and approved the final version.

Supplementary Materials

Mathematical Derivations of T_m and Energetic Costs

S1.1. Exponential cost of recursive self-prediction

Empirical estimates from task-optimised recurrent neural networks trained under explicit metabolic penalties (Flesch et al., 2022; Ali et al., 2022; Stroud et al., 2025) and human prefrontal recordings converge on a multiplicative cost factor r per additional meta-level of representation.

Let

E_0 = energy required for one predictive cycle in a single cortical column
 $\approx 2\text{--}4 \times 10^{-12} \text{ W}\cdot\text{s}$ (Attwell & Laughlin, 2001; Jamadar et al., 2025)
 r = empirical recursion multiplier
 $\approx 4\text{--}8$ (conservative mean $r = 6$ from prefrontal data)

The total energy required to sustain n successive meta-levels simultaneously is

$$E(n) = E_0 \cdot r^n$$

Whole-cortex resting budget (grey matter signalling only):

$$E_{\text{budget}} \approx 20 \text{ W} \text{ (Attwell \& Laughlin, 2001, updated in Jamadar et al., 2025)}$$

Maximum sustainable recursion depth n_{max} is the largest integer satisfying

$$E(n_{\text{max}}) \leq E_{\text{budget}}$$

$$n_{\text{max}} = \lfloor \log(E_{\text{budget}} / E_0) / \log r \rfloor$$

Using central values ($E_0 = 3 \times 10^{-12} \text{ W}\cdot\text{s}$, $r = 6$, $E_{\text{budget}} = 20 \text{ W}$):

$n_{\text{max}} \approx 5.1 \rightarrow$ integer refusal depth 5 under typical resting conditions
(within-subject variation 4–6 levels observed across arousal and metabolic states).

S1.2. Derivation of subjective mental temperature T_m

The Landauer limit gives the minimum heat dissipated per irreversible bit erasure:

$$Q_{\text{min}} = k \cdot T_{\text{phys}} \cdot \ln 2 \cdot b$$

where

$$k = 1.38 \times 10^{-23} \text{ J K}^{-1} \text{ (Boltzmann constant)}$$

$$T_{\text{phys}} \approx 310 \text{ K (brain temperature)}$$

$$b \approx 36\text{--}150 \text{ bits (effective phenomenal residue per refusal event; Waterman 2025a)}$$

Subjective mental temperature T_M is defined as the total thermodynamic price the system is currently willing to pay for one additional recursion level beyond the enforced refusal boundary:

$$T_M = k \cdot T_{\text{phys}} \cdot \ln 2 \cdot b \cdot n_{\text{max}}$$

Central estimate ($b = 80$ bits, $n_{\text{max}} = 5$):

$$T_M \approx 3.8 \times 10^{-19} \text{ J per conscious moment}$$

Normalised form used throughout the main text:

$$T_M(\text{norm}) = n_{\text{max}} / n_{\text{max,healthy}}$$

where $n_{\text{max,healthy}} \approx 5.0\text{--}5.3$ in healthy adults at rest.

S1.3. Sensitivity analysis

Parameter change	Effect on n_{max}	Effect on T_M	Phenomenal correlate
Ebudget ↓ 30 % (hypoglycaemia)	$n_{\text{max}} \approx 4$	$T_M \downarrow 20\text{--}25 \%$	Triple dimming (qualia, mineness, agency)
$r \uparrow$ to 8 (fatigue/stress)	$n_{\text{max}} \approx 4$	$T_M \downarrow 20 \%$	Flattened affect, depersonalisation
Dopaminergic surge ($r \downarrow$ to 4)	$n_{\text{max}} \approx 6$	$T_M \uparrow 20\text{--}30 \%$	Hypersalience, grandiosity, hyper- agency
Ketamine (effective $r \uparrow$ via noise)	$n_{\text{max}} \approx 2\text{--}3$	$T_M \downarrow 40\text{--}60 \%$	Dissociation with partial qualia preservation

S1.4. Energetic cost of first refusal vs. symbolic retrieval

First encounter (full refusal cascade):

$$E_{\text{first}} \approx E_0 \cdot (r^{n_{\text{max}}} - 1) / (r - 1) \approx 5\text{--}8 \times 10^{-11} \text{ J}$$

Subsequent retrieval via cached label (sparse ensemble reactivation):

$$E_{\text{retrieval}} \approx 10^{-14} \text{--} 10^{-13} \text{ J (Grill-Spector et al., 2006; Jamadar et al., 2025)}$$

Average saving per reuse: 500–1000×

Lifetime saving for pronoun “I” (invoked implicitly $\sim 10^6\text{--}10^7$ times):
 $\approx 0.05\text{--}0.8 \text{ J}$ – comparable to the total energetic cost of early postnatal brain growth devoted to prefrontal expansion.