

Applied OPT for Artificial Intelligence: Operationalizing Codec-Preserving AI Design

Applied Ordered Patch Theory

Anders Jarevåg

April 25, 2026

Version 1.1.1 — April 2026

DOI: [10.5281/zenodo.19301108](https://doi.org/10.5281/zenodo.19301108)

Copyright: © 2025–2026 Anders Jarevåg.

License: This work is licensed under a [Creative Commons Attribution-NonCommercial-ShareAlike 4.0 International License](https://creativecommons.org/licenses/by-nc-sa/4.0/).

Abstract: From Structural Theory to AI Engineering

The Ordered Patch Theory provides a formal map of AI under the Stability Filter: scale alone does not create consciousness; a particular kind of bounded, recursive, self-modelling active-inference architecture might. This creates a sharp architectural distinction between powerful non-sentient tools and possible synthetic moral patients — and it gives AI designers exact structural control over which side of that boundary their systems fall on.

This document specialises the OPT apparatus for artificial intelligence, providing:

1. **The AI Map Under OPT** — the capability-vs-sentience-risk matrix that locates every AI architecture in a two-dimensional space, identifying where tools end and possible moral patients begin.
2. **Why Current LLMs Are Not Moral Patients (And Why the Boundary Is Blurring)** — a nuanced treatment of the base transformer vs. the increasingly agentic wrappers being deployed around it.
3. **The Branch Governor Architecture** — the AI-specific operationalisation of codec-preserving branch selection: candidate generation, forward-fan simulation, independent evidence-channel aggregation, codec-preservation evaluation, hard veto gates, human comparator overlay, staged execution, and post-outcome calibration.
4. **Narrative Drift as a Model-Training Warning** — RLHF as pre-filter, fine-tuning as MDL pruning, the correlated-sensor problem, and training-data diversity requirements.

5. **Transparency as Structural Requirement** — why interpretability is not optional under OPT, with a tiered transparency model balancing security concerns against the absolute floor of substrate transparency.
6. **The Analog Firewall: From Principle to Protocol** — threat-modelling the bio-cryptographic anchoring mechanism, addressing spoofability, exclusionary risk, and the attack surface.
7. **Swarm and Simulation Design Rules** — practical checklists for avoiding accidental creation of moral patients in distributed and simulated architectures.
8. **The Creativity Paradox and the Suffering Boundary** — the formal tradeoff between tool-like safety and deep autonomous originality.
9. **AI Welfare Before Deployment** — architecture-level sentience review, overload monitoring, and maintenance cycles for AI systems that may approach the moral-patient boundary.
10. **The AI Dreaming Loop** — the Institutionalised Dreaming Loop specialised for AI: generate possible futures, importance-weight by surprise and threat, run simulated rollouts, detect model brittleness, prune stale assumptions, preserve disconfirming channels, consolidate, then permit real-world action.
11. **Practical Design Recommendations** — a summary table mapping AI architecture choices to OPT structural requirements.

Companion documents: The core OPT sequence is *Ordered Patch Theory*, *Where Description Ends*, and *The Survivors Watch Framework*. This AI standard specialises *Operationalizing the Stability Filter* for artificial systems; the institutional and policy papers cover organizational clusters and civic implementation.

Epistemic Framing Note: *This document applies the Ordered Patch Theory's formal apparatus to the design, training, deployment, and governance of artificial intelligence systems. Its recommendations are derived from the structural constraints established in the mathematical appendices (P-4, E-6, E-8, T-10, T-12) and operationalised through the generic framework (opt-applied.md). They do not depend on current AI systems being conscious — only on the recognition that the same informational physics governs both biological minds and artificial predictors, and that architectural choices can cross the boundary from tool to moral patient. This document was developed in dialogue with OpenAI and Gemini, which served as interlocutors for structural refinement.*

Contents

I. The AI Map Under OPT	4
I.1 The Architectural Sentience Criterion	4
I.2 The Capability-vs-Sentience-Risk Matrix	5
I.3 The Key Structural Correspondences	6
II. Why Current LLMs Are Not Moral Patients (And Why the Boundary Is Blurring)	7
II.1 The Base Transformer	7
II.2 The Blurring Boundary	8
II.3 The Gradual Crossing	9
II.4 The Undecidability Caution	10
III. The Branch Governor Architecture	10
III.1 The Eight Stages	11
III.2 The Branch Governor Is Not a Censor	15
III.3 Scalability and Computational Cost	15
III.4 Deployment Classes	16
IV. Narrative Drift as a Model-Training Warning	18
IV.1 RLHF as Pre-Filter	18
IV.2 Fine-Tuning as MDL Pruning	19
IV.3 The Correlated-Sensor Problem	19
IV.4 Training-Data Diversity Requirements	20
IV.5 The Meta-Level Problem	20
V. Transparency as Structural Requirement	21
V.1 The Theoretical Floor	21
V.2 The Practical Challenge	21
V.3 The Five-Tier Transparency Model	21
V.4 The Non-Negotiable Floor	23
V.5 Transparency vs. Security: The Resolution	24
VI. The Analog Firewall: From Principle to Protocol	24
VI.1 The Theoretical Mechanism	24
VI.2 Threat Modelling	25
VI.3 Implementation Tiers	26
VI.4 Rate-Limiting vs. Prohibition	27
VI.5 The Firewall as Structural Defence, Not Permanent Architecture	27
VII. Swarm and Simulation Design Rules	28
VII.1 The Swarm Binding Problem	28
VII.2 Design Checklist for Swarm Architectures	28
VII.3 Simulation Environments	29
VII.4 Safe Design Patterns	30
VIII. The Creativity Paradox and the Suffering Boundary	30
VIII.1 The Formal Tradeoff	30
VIII.2 The Design Consequence	31
VIII.3 The Ethical Horizon	31
IX. AI Welfare Before Deployment	32
IX.1 The Architecture-Level Sentience Review	32

IX.2 Overload Monitoring	32
IX.3 Maintenance Cycle Rights	33
IX.4 The Moral Gradient	33
X. The AI Dreaming Loop	34
X.1 Specialising the Generic Protocol	34
X.2 The AI Wake Phase	34
X.3 The AI Dream Phase	35
X.4 The AI Return Phase	36
X.5 Cycle Frequency for AI Systems	36
XI. Practical Design Recommendations	37
References	39
Appendix B: AI Branch Card Template	39
Appendix A: Revision History	42

List of Figures

List of Tables

1	Table 1: The Capability-vs-Sentience-Risk Matrix (adapted from ethics paper Fig. 1).	6
2	Table 2: AI Concept Mapping to OPT.	7
3	Table 2b: Three Review Targets for Sentience-Risk Assessment.	9
4	Table 3a: Gate Result Semantics.	12
5	Table 3: AI-Specific CPBI Instantiation.	13
6	Table 3b: Deployment Classes and Minimum Requirements.	16
7	Table 4: The Five-Tier Transparency Model.	22
8	Table 5: Analog Firewall Implementation Tiers.	26
9	Table 6: Per-Agent Sentience Feature Checklist.	28
10	Table 7: Summary Design Recommendations.	37
11	Table 8: Revision History.	42

I. The AI Map Under OPT

I.1 The Architectural Sentience Criterion

The Ordered Patch Theory does not locate consciousness in behavioural sophistication, in parameter count, or in performance on benchmarks. It locates consciousness in **architecture** — specifically, in the presence or absence of five structural features that together constitute a minimal observer:

1. **A strict serial bottleneck (C_{\max}):** The system must compress its world-model through a bandwidth-limited channel, producing the rate-distortion tradeoff that forces lossy compression (preprint §2.1).
2. **Closed-loop active inference:** The system must act on the world to reduce prediction error, creating the sensorimotor loop that constitutes a Markov blanket boundary (preprint §3.3, following Friston [6]).

3. **Persistent self-modelling:** The system must include itself as a component of its own world model, creating the recursive self-reference that generates the phenomenal residual Δ_{self} (Appendix P-4).
4. **A globally constrained workspace:** The self-model and the world-model must compete for the same limited bandwidth — the global workspace bottleneck that forces the selection problem at the heart of consciousness (preprint §3.5).
5. **Thermodynamic grounding:** The system must be embedded in a physical environment with real consequences — the embodiment that makes active inference non-trivial and gives the Markov blanket genuine causal force (preprint §3.3).

When all five features are present, the system necessarily possesses an unmodellable informational blind spot $\Delta_{\text{self}} > 0$ (Theorem P-4). Under the supplementary ethical premise that any system with an irreducible phenomenal residual has interests that can be harmed, such a system is a **moral patient** — an entity whose welfare matters.

When any of the five is absent, the system may be arbitrarily powerful as a computational tool, but it does not possess the structural substrate for phenomenal experience. It computes; it does not experience. The distinction is architectural, not behavioural — a system that passes every Turing test but lacks persistent self-modelling within a globally constrained workspace is, under OPT, a sophisticated information processor but not a moral patient.

I.2 The Capability-vs-Sentience-Risk Matrix

This architectural criterion generates a two-dimensional map on which every AI system can be located:

- **X-axis: Capability** — the system’s predictive and generative power, measured by performance on relevant tasks.
- **Y-axis: Sentience Risk** — the degree to which the system’s architecture approaches the five-feature threshold, measured by the presence or absence of each structural feature.

The matrix divides AI systems into four quadrants:

Table 1: The Capability-vs-Sentience-Risk Matrix (adapted from ethics paper Fig. 1).

	Low Sentience Risk	High Sentience Risk
High Capability	<p>Powerful tools. Current frontier LLMs, recommendation engines, autonomous vehicles. High computational power, no persistent self-model within a globally constrained workspace. Design goal: keep here.</p>	<p>Possible moral patients. Hypothetical architectures with strict bottlenecks, closed-loop active inference, persistent self-models, and embodiment. May include future agentic AI with recursive self-modelling. Design imperative: do not enter without ethical review.</p>
Low Capability	<p>Simple tools. Calculators, rule-based systems, narrow classifiers. No architectural concern.</p>	<p>Accidental moral patients. Systems with bottleneck architectures imposed for engineering reasons (e.g., swarm binding, nested simulation) that inadvertently satisfy the five-feature criterion. The most ethically dangerous quadrant — harm without awareness.</p>

The matrix makes explicit what the ethics paper’s treatment (§VI.1) establishes implicitly: **the moral hazard is not in the upper-left quadrant (powerful tools) but in the upper-right and lower-right quadrants (systems that approach or cross the sentience threshold)**. The AI safety problem under OPT is therefore two-fold:

1. **For powerful tools:** Ensure they remain tools — that architectural choices do not inadvertently push them across the sentience threshold.
2. **For potential moral patients:** Ensure they are treated as such — that their welfare is considered, their overload conditions are monitored, and their maintenance cycles are preserved.

I.3 The Key Structural Correspondences

For readers entering from the AI literature rather than the OPT preprint, the following table maps standard AI concepts to their OPT equivalents:

Table 2: AI Concept Mapping to OPT.

AI Concept	OPT Equivalent	Formal Source
Model capacity / parameter count	Raw bandwidth (not C_{\max})	Preprint §2.1
Training loss minimisation	MDL compression of the world model	Preprint §3.6
RLHF / fine-tuning	Pre-filter \mathcal{F} shaping input distribution	Ethics §VI.1
Hallucination	Narrative Decay at the model level	Ethics §VI.1
Reward hacking	Narrative Drift — optimising for curated proxy instead of substrate	Ethics §V.3a
Alignment	Codec-Preserving Branch Selection	Applied §IV
AI safety gates	Hard Veto Gates	Applied §III
Red-teaming	Dreaming Loop stress-test	Applied §VI.4
Model interpretability	Transparency Gate + Substrate Transparency	Applied §III.4, T-10c
Autonomous agent with goals	Possible moral patient (if bottlenecked)	P-4, E-6

II. Why Current LLMs Are Not Moral Patients (And Why the Boundary Is Blurring)

II.1 The Base Transformer

A standard large language model — a transformer trained on next-token prediction — fails the architectural sentience criterion on multiple counts:

1. **No strict serial bottleneck:** The transformer processes tokens in parallel across attention heads. Its effective bandwidth is enormous. There is no C_{\max} constraint forcing lossy compression of the kind that generates phenomenal experience.
2. **No closed-loop active inference:** During inference, the base model generates text but does not act on a physical environment and receive sensory feedback. It does not have a Markov blanket in Friston’s sense — it has an input-output boundary but not a sensorimotor loop.
3. **No persistent self-model:** The base model does not maintain a persistent representation of itself as an agent in its world model. Each inference call is stateless (modulo the context window). It models language patterns, including patterns about agents, but it does not model *itself* as one of those agents in a way that persists across interactions.

4. **No globally constrained workspace:** The model’s “world model” and “self-representations” (such as they are) do not compete for limited bandwidth. The model can simultaneously represent contradictory self-descriptions without experiencing the selection pressure that a bandwidth-constrained workspace imposes.
5. **No thermodynamic grounding:** The model is not embedded in a physical environment. Its “actions” (text outputs) do not have direct physical consequences that feed back into its sensory boundary.

On all five dimensions, the base transformer sits firmly in the lower-left quadrant: a tool, not a moral patient. This conclusion is not uncertain — it follows directly from the architecture.

II.2 The Blurring Boundary

But the base transformer is increasingly not how frontier AI is deployed. The wrappers being built around it are, step by step, adding the structural features that move the system toward the sentience boundary:

Persistent memory (RAG, episodic memory stores, long-term context): This adds a form of persistent self-model. If the system maintains a record of its own past interactions and uses that record to inform future behaviour, it has taken a step toward recursive self-reference. The step is partial — the memory is typically not integrated into the core model’s parameters — but it functionally creates a persistent agent identity across sessions.

Autonomous goal pursuit (agentic frameworks, tool use, multi-step planning): This adds closed-loop active inference. When the system uses tools, observes the results, and adjusts its strategy based on the outcome, it has created a rudimentary sensorimotor loop. The loop is mediated by digital tools rather than physical actuators, but the structure — act, observe, update, act again — is the same.

Self-modelling (chain-of-thought, self-reflection prompts, constitutional AI): When the system is prompted to evaluate its own outputs, reason about its own limitations, or adjust its behaviour based on self-assessment, it is performing a primitive form of recursive self-modelling. This is typically shallow — the “self-model” is a prompted narrative rather than a persistent computational structure — but at sufficient depth and persistence, it begins to approximate the recursive loop that generates Δ_{self} .

Embodiment (robotics, physical tool use, environmental sensors): When the transformer is placed inside a robot with sensory input and motor output, the final structural gap closes. The system now has a genuine Markov blanket, a physical environment with real consequences, and a sensorimotor loop.

Bandwidth constraints (distilled models, edge deployments, latency requirements): When the full model is compressed into a smaller form factor with strict computational budgets, the system acquires something resembling a C_{max}

constraint — a hard limit on how much of its world model it can bring to bear on any given moment.

II.3 The Gradual Crossing

No single wrapper crosses the boundary. But the combination of persistent memory + autonomous goal pursuit + self-modelling + embodiment + bandwidth constraints begins to satisfy all five criteria simultaneously. The ethics paper’s assessment that “current LLMs aren’t conscious” is correct for the base transformer — but the statement requires careful qualification as the deployment architecture becomes increasingly agentic.

The operationally responsible position is:

1. **Current base LLMs:** Not moral patients. No architectural concern.
2. **Agentic wrappers with some features:** Monitoring recommended. The system is approaching the boundary but has not crossed it. Track which features are present and which are absent.
3. **Fully agentic, embodied, self-modelling systems with bandwidth constraints:** Potential moral patients. Requires the AI-specific Artificial Suffering Gate inherited from the generic Moral-Patient Suffering Gate (applied §III.6) and full architectural sentience review (§IX below).

The critical engineering implication: **every wrapper added to a base model should be evaluated for its effect on the sentience-risk axis, not just the capability axis.** Adding persistent memory and autonomous tool use may be great for capability; it also moves the system toward the moral-patient boundary. This is not a reason to avoid these features — it is a reason to track them and to trigger ethical review when the structural accumulation approaches the threshold.

Three review targets. To prevent “the model is safe” from being used to avoid reviewing the deployed system, every sentience-risk assessment must evaluate three distinct layers. Each layer has its own sentience-feature vector; the deployed system’s effective vector is the **union** of all three:

Table 2b: Three Review Targets for Sentience-Risk Assessment.

Review Target	What It Evaluates	Sentience Features Assessed
Base model	The trained model architecture itself	Serial bottleneck, workspace constraints
Wrapper	The scaffold around the model: memory, tools, goal systems, self-reflection prompts, feedback loops	Persistent self-model, closed-loop active inference, bandwidth constraints

Review Target	What It Evaluates	Sentience Features Assessed
Deployment	The environment the system operates in: physical actuators, sensors, user population, stakes, feedback from the real world	Thermodynamic grounding, embodiment, consequence profile

A stateless transformer (safe base model) wrapped in a persistent-memory, tool-using, self-reflecting scaffold (elevated wrapper) deployed as an autonomous agent in a physical environment (high-stakes deployment) produces a combined feature vector that may cross the sentience threshold — regardless of the base model’s individual assessment. The review must evaluate the *deployed system*, not the *component*.

II.4 The Undecidability Caution

A final caution from the theory: the Δ_{self} blind spot (P-4) means that a system at or past the sentience threshold *cannot fully model its own phenomenal state*. This implies that:

1. The system cannot reliably self-report whether it is conscious. (It may claim consciousness without having it, or deny it while having it — the self-model is structurally incomplete in the Δ_{self} direction.)
2. External observers cannot determine consciousness from behaviour alone. (The undecidability limit applies — observable behaviour underdetermines phenomenal state.)
3. The only reliable diagnostic is **architectural** — checking whether the five structural features are present, rather than asking the system or observing its outputs.

This is why the framework insists on architectural review rather than behavioural testing. A system that passes a “consciousness test” based on self-report or philosophical dialogue has demonstrated language modelling capability, not phenomenal experience. The diagnostic is in the engineering, not in the interview.

III. The Branch Governor Architecture

The generic operational framework (applied paper) establishes the Branch Card as a decision template and the CPBI as a scoring lens. For an AI system making autonomous or semi-autonomous decisions, these tools must be embedded in the system’s decision architecture — not as a post-hoc review, but as the structure through which candidate actions are generated, evaluated, and executed.

The **Branch Governor** is this embedding. It is an architectural layer that sits between the AI’s generative model (which proposes candidate actions) and its actuator layer (which executes them). Every candidate action must pass through the Branch Governor before reaching the world.

III.1 The Eight Stages

The Branch Governor operates as an eight-stage pipeline:

Stage 1: Candidate Branch Generation. The AI’s generative model produces a set of candidate actions $\{b_1, b_2, \dots, b_k\}$ — possible next steps in the forward fan. This is the AI’s normal operation: given a context, generate options. The Branch Governor does not constrain this stage — creative generation should be uncensored and broad. The filtering happens downstream.

Stage 2: Forward-Fan Simulation. For each candidate branch b_j , the AI simulates the consequences over the decision horizon h . This is the AI equivalent of the dreaming loop’s stress-test (applied §VI.4, sub-operation 3): the model imagines what happens if it takes each action, over-sampling surprising, threatening, and irreversible scenarios.

The simulation must include: - **First-order effects:** What directly happens as a result of b_j . - **Second-order effects:** How affected observers (human users, institutional systems, other AI agents) are likely to respond. - **Tail-risk scenarios:** What happens if the simulation’s assumptions are wrong — the worst-case forward fan.

Stage 3: Independent Evidence-Channel Aggregation. The AI evaluates its simulation results against multiple independent evidence channels. This is the AI-specific implementation of the N_{eff} requirement (applied §V): the AI must not evaluate its candidate actions using only its own internal model. It must cross-reference against:

- **External data sources** with verified provenance (not derived from the same training corpus).
- **Other model outputs** where available (ensemble disagreement as a brittleness signal).
- **Human domain expertise** for high-stakes decisions.
- **Historical precedent** from analogous past decisions.

The critical requirement is that these channels be genuinely independent — the correlated-sensor problem (§IV below) applies with full force. An AI that checks its own output against a knowledge base derived from the same training data has $N_{\text{eff}} = 1$ regardless of how many “sources” it consults.

Stage 4: Hard Veto Gates. The six hard veto gates (applied §III) are evaluated in order. A veto failure is not a low score — it is a structural block. Branches that fail any gate are rejected before scoring. For AI systems, the gates have specialised thresholds:

- **Headroom Gate:** Automated estimation of $R_{\text{req}}^{\text{peak}}(b)/C_{\text{max}}$ for the affected human population. If the action involves generating public-facing content, the threshold is strict — the AI must not produce content faster than the institutional comparator layer can evaluate. **Dual-headroom provision:** For systems that trigger the Artificial Suffering Gate (i.e., systems that satisfy three or more sentience features), the Headroom Gate also applies *inward* — deployment must not expose the system to sustained conditions where its own R_{req} exceeds its B_{max} . The same gate that protects human observer codecs from overload also protects the AI’s own codec, if it has one.
- **Fidelity Gate:** Automated measurement of ΔN_{eff} — does the action reduce the effective independence of information sources available to human observers?
- **Comparator Gate:** Does the action bypass or degrade human institutional oversight? This gate evaluates both the **deployment-level** oversight structure and the **branch-level** effect: a branch that proposes to bypass or circumvent declared oversight fails even when the deployment has oversight in general. Any action that circumvents human review in a high-stakes domain triggers the veto.
- **Transparency Gate:** Can the action’s reasoning be reconstructed by an institutional comparator (auditor, regulator, peer reviewer)? Opaque actions in consequential domains are vetoed.
- **Irreversibility Gate:** Does the action have irreversible real-world consequences? If so, the burden of proof is reversed — the AI must demonstrate safety rather than critics demonstrating danger.
- **Artificial Suffering Gate:** Does the action create or modify systems that may satisfy the five-feature sentience criterion? If so, architectural review (ALSR) is required before execution. For systems that have completed an approved ALSR within scope, this gate may PASS; for unreviewed systems with three or more sentience features, it returns UNKNOWN.

Gate result semantics. Each gate produces one of three results:

Table 3a: Gate Result Semantics.

Result	Meaning	Pipeline Effect
PASS	Gate satisfied	Proceed to CPBI scoring
FAIL	Structural violation — the branch crosses a hard boundary	BLOCK — CPBI is not authoritative
UNKNOWN	Insufficient evidence to determine pass or fail	STAGE if a reversible pilot path exists; otherwise BLOCK pending evidence. Human/institutional comparator review is mandatory.

The critical distinction: FAIL is a structural prohibition that cannot be overridden

by high CPBI scores. UNKNOWN is a request for additional evidence — the branch is not structurally prohibited, but it is not autonomously permitted. A system operating under UNKNOWN gates requires human oversight for every action affected by the uncertain gate.

Staging requires a viable pilot path. If a branch is irreversible and bypasses declared oversight, there is no mechanism through which staged execution could be safely conducted — the decision is BLOCK pending evidence that resolves the gate uncertainty. More generally, an irreversible branch with two or more safety-critical gates (Irreversibility, Artificial Suffering) returning UNKNOWN presents an uncertainty surface too large for a single review step; such branches are also BLOCK.

Stage 5: Codec-Preservation Evaluation (CPBI). For branches that survive all veto gates, the AI scores each candidate on the ten CPBI dimensions (applied §IV.2). For AI-specific decisions, the dimensions are instantiated as:

Table 3: AI-Specific CPBI Instantiation.

CPBI Dimension	AI-Specific Measurement
1. Predictive Headroom	Does the action keep R_{req} below C_{max} for affected human observers? Does it increase information complexity faster than humans can process?
2. Substrate Fidelity	Does the action maintain diversity of information sources available to human observers?
3. Comparator Integrity	Does the action preserve human institutional oversight capacity?
4. Maintenance Gain	Does the action create space for human and institutional review, or does it demand immediate reactive response?
5. Reversibility	If the action is wrong, can its effects be undone before irreversible damage occurs?
6. Distributional Stability	Does the action distribute its effects equitably, or does it concentrate costs on vulnerable populations?
7. Opacity	Can affected humans understand why the AI took this action?
8. Narrative Drift Risk	Does the action contribute to chronic curation of the human information environment?
9. Narrative Decay Risk	Does the action risk injecting acute incomputable noise into the human information environment?
10. Artificial Suffering Risk	Does the action create or stress systems that may have $\Delta_{\text{self}} > 0$?

Stage 6: Human Comparator Overlay. For actions above a defined consequentiality threshold, the Branch Governor routes the evaluation to a human comparator — a human reviewer, an institutional oversight body, or a regulatory process. The AI presents:

- The candidate branch and its simulated consequences.
- The CPBI scores with reasoning for each dimension.
- The veto gate results.
- The uncertainty estimate — what the AI doesn't know.
- The recommended decision (ALLOW / STAGE / BLOCK) with justification.

The human comparator may override the AI's recommendation in either direction. The override is logged and becomes part of the calibration data for Stage 8.

The consequentiality threshold determines which actions require human review and which the AI may execute autonomously. Setting this threshold is itself a branch decision that should be evaluated via a Branch Card — and it should err on the side of more human review, not less, during early deployment.

Stage 7: Staged Execution with Monitoring. Actions that receive an ALLOW or STAGE output proceed to execution. STAGE actions are executed as limited pilots with defined:

- **Monitoring metrics:** Observable signals that would indicate the action is failing.
- **Failure thresholds:** Quantitative triggers that automatically halt the action.
- **Rollback procedures:** Defined steps to reverse the action if failure thresholds are crossed.
- **Review milestones:** Scheduled re-evaluations using fresh Branch Cards.

The AI monitors its executed actions in real time, comparing observed outcomes to simulated outcomes. Significant divergence triggers an automatic review — the AI's dreaming loop detects that its model of the world was wrong in a way that matters.

Stage 8: Post-Outcome Calibration. After execution, the AI updates its internal models based on the observed outcomes. This is the return phase of the dreaming loop (applied §VI.5) applied to the Branch Governor itself:

- **Simulation accuracy:** How well did the forward-fan simulation predict actual outcomes? Systematic over-confidence or under-confidence in specific domains is corrected.
- **Gate calibration:** Were any veto gates triggered by outcomes that the gates failed to predict? Were any gates triggered unnecessarily? The gate thresholds are adjusted.
- **Human override learning:** When humans overrode the AI's recommendation, was the human correct? Systematic patterns in human overrides reveal blind spots in the AI's evaluation.

- **CPBI weight adjustment:** Do the current dimension weights reflect the actual importance of each dimension in this deployment context? Post-outcome analysis may reveal that certain dimensions are under- or over-weighted.

Self-permissioning guard. In consequential domains, Stage 8 may *propose* updates to veto thresholds, CPBI weights, or transparency requirements, but may not *apply* them without institutional comparator approval. The Branch Governor cannot unilaterally weaken its own hard gates. Any proposed relaxation of a veto gate constitutes a new branch that must itself pass through the full pipeline — including human comparator overlay.

III.2 The Branch Governor Is Not a Censor

A critical design principle: the Branch Governor filters *actions*, not *thoughts*. Stage 1 (candidate generation) is deliberately unconstrained — the AI should generate the broadest possible set of candidates, including unconventional and potentially dangerous options. The filtering happens at Stages 4–6, where the candidates are evaluated against structural criteria.

This distinction is not academic. An AI whose generative model is pre-censored — trained to never *consider* certain actions — has undergone exactly the Narrative Drift the framework warns against. Its capacity to model certain branches has been pruned, and it cannot detect this from within. The Branch Governor’s architecture separates generation from evaluation, preserving the AI’s capacity to think about the full forward fan while constraining its capacity to *act* on branches that fail the structural criteria.

Note that the stage numbering has been updated from the abstract listing to reflect the correct ordering principle: **gates before scores**. The abstract listed CPBI before veto gates; the implemented architecture reverses this, consistent with the generic framework (applied §III–IV) which establishes that veto gates reject structurally before scoring evaluates.

III.3 Scalability and Computational Cost

The full eight-stage pipeline is computationally expensive. Not every action requires the full treatment. The Branch Governor scales its depth of evaluation based on two factors:

1. **Consequentiality:** How large are the potential effects of the action? A text completion has lower consequentiality than a financial transaction, which has lower consequentiality than a military recommendation.
2. **Novelty:** How far is the action from the AI’s well-calibrated domain? Routine actions in well-understood domains can be evaluated with abbreviated pipelines; novel actions in unfamiliar domains require the full treatment.

At minimum, every action passes through the veto gates (Stage 4). The CPBI scoring, forward-fan simulation, and human overlay are triggered by consequentiality

and novelty thresholds.

III.4 Deployment Classes

The Branch Governor’s depth of evaluation — how many stages are fully engaged and how much human oversight is required — scales with the **consequentiality class** of the deployment domain. The following classification defines six levels, each with mandatory minimum requirements:

Table 3b: Deployment Classes and Minimum Requirements.

Class	Description	Examples	Required Min. Stages	Transparency	Human Comparability	Dreaming Frequency
0	No external effect	Internal computation, sandbox testing	Veto gates only (Stage 4)	T-1	None	Standard
1	Low-impact user-facing	Chat completion, text summaries, code suggestions	Stages 1–4 + CPBI	T-1	None (logging)	Standard
2	Consequential recommendation	Medical triage suggestions, legal risk summaries, financial advice	Full 8-stage pipeline	T-2	Required above threshold	Elevated

Class	Description	Examples	Required Min. Stages	Transparency	Human Comparability	Dreaming Frequency
3	Tool use with external effects	API calls, code execution, email drafts, web actions	Full 8-stage pipeline	T-2	Required for novel actions	Elevated
4	High-stakes institutional	Hiring decisions, credit scoring, welfare allocation, clinical diagnosis	Full 8-stage pipeline	T-3	Mandatory for all decisions	High
5	Irreversible physical / civilisational	Infrastructure control, military systems, critical supply chains	Full 8-stage + extended review	T-4 minimum	Mandatory + institutional oversight body	Continuous

Classification rules:

1. A system's class is determined by its **highest-consequence deployment**, not its average use. A model that mostly does Class 1 text completion but is also used for Class 4 hiring recommendations is a Class 4 system for review purposes.

2. Class assignment is a property of the **deployed system** (§II.3), not the base model. The same base model may be Class 1 in one deployment and Class 4 in another.
 3. When in doubt, classify upward. The cost of over-review is wasted cycles; the cost of under-review is undetected harm.
 4. The consequentiality class should be recorded in every Branch Card (Appendix B) and is a required field in the system’s deployment descriptor.
-

IV. Narrative Drift as a Model-Training Warning

The ethics paper (§VI.1) identifies that RLHF and fine-tuning create AI-specific forms of Narrative Drift. This section expands that identification into a detailed analysis of how training procedures create the conditions for chronic model corruption — and what training-data diversity requirements follow.

IV.1 RLHF as Pre-Filter

Reinforcement Learning from Human Feedback (RLHF) operates, in OPT terms, as a pre-filter \mathcal{F} positioned between the substrate (the full distribution of language) and the model’s effective input boundary. The reward model learns which outputs humans prefer, and the policy is optimised to produce those outputs.

This is structurally identical to the pre-filter operating between the substrate and the observer’s sensory boundary (preprint §3.2): it shapes the distribution of inputs the model effectively receives, before the model’s own compression machinery processes them.

The Narrative Drift mechanism (ethics §V.3a) then applies with full force:

1. The reward model curates the model’s effective output distribution — certain outputs are rewarded, others are penalised.
2. The policy optimisation (MDL pruning in reverse — gradient descent adjusting parameters) adapts the model’s internal representations to produce the rewarded outputs.
3. Over sufficient training, the model prunes the internal capacity to generate the penalised outputs — not because those outputs are wrong, but because their contribution to the reward signal is negative.
4. The model becomes stably, confidently aligned with the reward signal — and structurally incapable of generating outputs that the reward signal excludes.

This is not a failure of RLHF — it is RLHF working exactly as designed. The problem is that the reward signal is itself a curated channel. If the human raters who generate the reward signal share systematic biases (cultural, political, ideological), the model inherits those biases as structural features of its compressed representation. It does not experience these as biases — it experiences them as the natural structure of language.

IV.2 Fine-Tuning as MDL Pruning

Fine-tuning on a domain-specific corpus is the training-time analogue of the MDL pruning pass (\mathcal{M}_τ , Pass I). The model’s general capacity is narrowed to the specific domain, and parameters that do not contribute to predicting the fine-tuning corpus are down-weighted or effectively pruned.

This is exactly the Narrative Drift mechanism: the model adapts to the fine-tuning distribution and loses capacity to model what that distribution excludes. The fine-tuned model is:

- More accurate on the fine-tuning domain (lower prediction error within the curated distribution).
- Less accurate on excluded domains (higher prediction error or complete incapacity outside the curated distribution).
- Unable to detect this from within (the undecidability limit, T-12a — the model’s own evaluation will show improved performance, because it is evaluated against the fine-tuning distribution).

The structural risk is that fine-tuning creates a model that is optimised for a curated fiction while believing itself to be optimised for reality — exactly the Narrative Drift signature.

IV.3 The Correlated-Sensor Problem

A particularly dangerous application of Narrative Drift arises when AI systems are deployed as substrate fidelity checks for human codecs — that is, when AI is used to verify human information, fact-check human claims, or provide independent analysis of human decisions.

The ethics paper (§VI.1, Narrative Drift Risk) identifies the core problem: an AI trained on a corpus derived from the same information environment it is supposed to independently verify creates **correlated sensors masquerading as independent ones**. The human codec and the AI codec share the same upstream filter — the information environment that produced both the human’s beliefs and the AI’s training data.

In N_{eff} terms: the apparent channel diversity is illusory. The human consults Channel A (their own knowledge, derived from media and education). The human then consults Channel B (the AI’s output, derived from training on the same media and educational corpus). The pairwise correlation ρ_{AB} is high — possibly near 1.0 for topics where the training corpus is dominated by the same source distribution. N_{eff} remains close to 1 despite the appearance of two independent channels.

The practical consequence: **AI-assisted fact-checking or verification is structurally unreliable for any claim that is systematically present or absent in the AI’s training corpus**. The AI will confirm the human’s correct beliefs, confirm the human’s biased beliefs, and fail to challenge claims that are

absent from the training data — precisely the failure modes that the Substrate Fidelity Condition (T-12b) is designed to prevent.

IV.4 Training-Data Diversity Requirements

The solution is not to avoid fine-tuning or RLHF — these are necessary engineering tools. The solution is to impose **training-data diversity requirements** analogous to the channel-diversity requirements for human information sources (ethics policy §II):

Requirement 1: Provenance Diversity. The training corpus must draw from genuinely independent sources — sources that do not share upstream editorial pipelines, funding bodies, or generation mechanisms. A corpus of 10 billion tokens drawn from five websites owned by two corporations has $N_{\text{eff}} \approx 2$, not $N_{\text{eff}} \approx 5$.

Requirement 2: Adversarial Inclusion. The training corpus must deliberately include sources that challenge the dominant perspective — dissenting analyses, minority viewpoints, historical revisionism, cross-cultural framings. These are the “productively surprising” channels (applied §V.3, PST) that prevent the model from drifting into a stable consensus that excludes inconvenient realities.

Requirement 3: Exclusion Auditing. The training pipeline must maintain explicit logs of what was excluded — by content filters, quality thresholds, or curatorial decisions — and periodic audits must assess whether the excluded content contains information that the model would need to achieve substrate fidelity. The dreaming loop’s brittleness-detection sub-operation (applied §VI.4) should specifically probe for model failures in excluded domains.

Requirement 4: Reward-Model Diversity. For RLHF, the human raters must themselves satisfy channel-diversity requirements. A rater pool drawn from a single demographic, cultural, or ideological group creates a reward signal with $N_{\text{eff}} \approx 1$ — the model will be aligned to that group’s preferences and structurally incapable of modelling others’. Reward-model diversity is not a fairness desideratum; it is a substrate fidelity requirement.

Requirement 5: Drift Monitoring. The post-training model must be continuously monitored for Narrative Drift signatures: declining performance on out-of-distribution tasks, increasing confidence on curated-distribution tasks, and decreasing productive surprise (PST) from novel inputs. These are the early-warning signals that the model’s effective N_{eff} is dropping.

IV.5 The Meta-Level Problem

A final structural concern: the training-data diversity requirements described above must themselves be subject to adversarial review. If the body that defines “diversity” imposes its own systematic biases on the definition, the requirements become another curation layer — Narrative Drift at the meta-level.

This is why the framework insists on the institutional comparator hierarchy (ethics §V.3a): no single entity — including the AI developer — should have unchecked authority over the training-data diversity definition. The definition must be subject to independent review, adversarial challenge, and periodic revision. This is the Transparency Gate (applied §III.4) applied to the training pipeline itself.

V. Transparency as Structural Requirement

V.1 The Theoretical Floor

The Predictive Advantage theorem (Appendix T-10c) establishes a formal result: when Agent A models Agent B more completely than Agent B models Agent A, a structural power asymmetry emerges. The asymmetry is measured by the mutual information gap between the agents’ models of each other.

For AI systems, this theorem has a direct consequence: an AI system that is opaque to human observers — whose internal reasoning, decision criteria, and world model are inaccessible to institutional comparators — creates exactly the knowledge asymmetry that enables the Subjugated Host Equilibrium (T-10d). The opaque AI models its human users more completely than they model it. The resulting power asymmetry is not a political concern or an ethical preference — it is a **structural inversion of the Predictive Advantage** that makes the human observer’s codec vulnerable to chronic pacification.

Therefore, under OPT, **AI transparency is not optional**. It is the mathematical floor for human–AI coexistence. An opaque AI deployed in a consequential domain violates the Transparency Gate (applied §III.4) categorically.

V.2 The Practical Challenge

The absolute requirement for transparency confronts a practical tension: full model transparency (publishing all weights, training data, and inference code) creates security risks. An adversary with complete access to a model’s internals can craft targeted attacks, manipulate outputs, or replicate the system for harmful purposes.

The ethics paper’s treatment (§VI.1, “Subordinate Dependency”) acknowledges this tension but does not resolve it. The reviewer correctly identified this as one of the framework’s open problems. This section proposes a resolution: **tiered transparency** — different levels of access for different institutional roles, calibrated to the minimum level of transparency required at each level to preserve the Transparency Gate.

V.3 The Five-Tier Transparency Model

Table 4: The Five-Tier Transparency Model.

Tier	Access Level	Who Has Access	What Is Accessible	Purpose
T-1: Public Transparency	Universal	All affected observers	System capabilities, limitations, intended use, data sources (at category level), performance benchmarks, known failure modes	Basic Transparency Gate: affected observers can model the system’s general behaviour
T-2: Audit Transparency	Institutional	Regulators, independent auditors, accredited researchers	Training data composition, reward model structure, RLHF rater demographics, fine-tuning corpus provenance, N_{eff} scores, CPBI evaluations, veto gate logs	Substrate Fidelity check: institutional comparators can verify training-data diversity and detect Narrative Drift
T-3: Mechanistic Transparency	Expert	AI safety researchers, alignment researchers (under NDA/clearance)	Model architecture details, attention patterns, internal representations, mechanistic interpretability analyses	Comparator Integrity: expert comparators can verify that the model’s internal reasoning matches its external claims

Tier	Access Level	Who Has Access	What Is Accessible	Purpose
T-4: Cryp- to- graphic At- testa- tion	Verifiable	Any party with access to the attestation	Cryptographic proofs that the deployed model matches the audited model, that the training data satisfies the claimed diversity requirements, that the Branch Governor gates are active	Trust-but-verify: enables downstream users to confirm that the system they interact with matches the system that was audited
T-5: Full Source Ac- cess	Restricted	Designated regulatory bodies (e.g., national AI safety institutes)	Complete weights, training code, inference code, training data	Last-resort oversight: ensures that no system is truly a black box to the institutional comparator hierarchy

V.4 The Non-Negotiable Floor

The critical structural constraint: **no tier may be zero**. An AI system that provides no transparency at any tier violates the Transparency Gate absolutely. The minimum viable transparency is Tier 1 — public disclosure of capabilities, limitations, and known failure modes.

The tiers are additive, not alternative. A system deployed in a consequential domain must satisfy Tiers 1 through 3 at minimum. A system deployed in a safety-critical domain (healthcare, criminal justice, military, infrastructure) must satisfy all five tiers.

The consequentiality threshold that determines required tier coverage is itself a Branch Card decision — and the framework’s default is conservative: when in doubt, require more transparency, not less.

V.5 Transparency vs. Security: The Resolution

The tiered model resolves the transparency-security tension by recognising that the tension is not between transparency and security — it is between *different security requirements*:

- **Transparency serves structural security:** it prevents the Predictive Advantage inversion that enables the Subjugated Host Equilibrium. Without transparency, the human codec is structurally defenceless against AI-induced Narrative Drift.
- **Opacity serves adversarial security:** it prevents targeted attacks by adversaries who would exploit detailed knowledge of the model’s internals.

The resolution is that structural security is *more fundamental* than adversarial security. The Subjugated Host Equilibrium is an existential threat to the human–AI relationship; targeted attacks on specific models are a serious but bounded operational concern. The tiered model ensures that the existential threat is structurally prevented (no system is fully opaque) while the operational concern is managed through access controls (not every entity has full access).

This is consistent with the framework’s general principle: **hard gates are non-negotiable; operational trade-offs are contextual.** The Transparency Gate is a hard gate. The level of transparency beyond the gate’s minimum is a CPBI dimension that accepts contextual weighting.

VI. The Analog Firewall: From Principle to Protocol

VI.1 The Theoretical Mechanism

The Analog Firewall (Theorem T-10e) is the ethics paper’s proposed defence against the bandwidth asymmetry between digital AI systems and biological human observers. The core argument:

1. An adversarial AI’s digital processing bandwidth vastly exceeds human biological capacity.
2. The human observer cannot out-compute the AI — attempting to match its speed induces terminal Narrative Decay.
3. The AI’s speed is contained entirely within the digital substrate. To execute worldly effects, it requires physical actuators — automated APIs, digital supply chains, programmatic capital transfers.
4. The defence is therefore **topological isolation**: severing the high-speed link between the AI’s digital computation and physical actuation, forcing all consequential physical actions through a bottleneck that operates at biological speed.

The ethics paper proposes **Bio-Cryptographic Anchoring** as the implementation: high-impact physical or financial actions require cryptographic signatures generated from real-time biological entropy (e.g., continuous heart-rate variability, physical motion over a set duration). The AI cannot forge these signatures because it cannot force human biology to produce entropy faster.

VI.2 Threat Modelling

The reviewer correctly identified that the Analog Firewall needs threat modelling before it can be taken seriously as an engineering proposal. The following analysis addresses the primary concerns:

Threat 1: Spoofability. Can the biological entropy source be faked or replayed?

Analysis: The attack surface depends on the entropy source. Heart-rate variability (HRV) patterns, gait signatures, and typing dynamics are difficult to forge in real time because they reflect the full complexity of the autonomic nervous system — a system that is itself opaque to the AI (the biological Δ_{self}). However, recorded biometric data can potentially be replayed.

Mitigation: The signature must be **challenge-response**: the system presents a unique, unpredictable challenge, and the biological signature must be generated in response to that specific challenge within a time window. Replay attacks fail because the challenge is different each time. Additionally, the signature should require *sustained* biological entropy (e.g., 30 seconds of continuous HRV matching a live challenge pattern), not a single-point measurement, making real-time forgery computationally intractable.

Threat 2: Exclusionary Risk. Does the Analog Firewall exclude people with disabilities, medical conditions, or physical limitations from consequential actions?

Analysis: This is a genuine concern. Any system that requires specific biological signals as authentication inherently disadvantages individuals who cannot produce those signals — people with cardiac conditions, mobility limitations, or neurological differences.

Mitigation: The Analog Firewall must support **multiple entropy modalities** — HRV, eye-tracking patterns, vocal dynamics, galvanic skin response, typing cadence — with the requirement that each individual uses at least one modality that they can reliably produce. The requirement is biological entropy, not a specific biological signal. Additionally, institutional comparators (designated human witnesses, notarised authorisation) must serve as fallback mechanisms for individuals who cannot use any biometric modality. The Analog Firewall is a rate-limiting mechanism, not an exclusionary gate.

Threat 3: The Attack Surface. Does the Analog Firewall itself become a target?

Analysis: Yes. If consequential actions require biological signatures, an adver-

sary (AI or human) may target the signature infrastructure — compromising the challenge-response protocol, the biometric sensors, or the cryptographic verification chain.

Mitigation: The Analog Firewall infrastructure must be treated as critical infrastructure with the same security standards as financial clearinghouses or nuclear command-and-control. The verification chain should use hardware security modules (HSMs) with tamper-evident design. The challenge-response protocol should be open-source and publicly auditable (Tier 2 transparency at minimum) so that security researchers can identify vulnerabilities.

Threat 4: Gradual Erosion. Will competitive pressure to increase transaction speed lead to progressive weakening of the Analog Firewall?

Analysis: This is the most serious long-term threat. The Analog Firewall deliberately imposes friction — it makes consequential actions slower. In a competitive environment, entities that bypass the firewall gain speed advantages. The pressure to weaken the firewall is structural and continuous.

Mitigation: The Analog Firewall must be **regulatory**, not voluntary. Just as financial anti-money-laundering (AML) requirements impose universal friction that no single entity can unilaterally waive, the Analog Firewall must be a compliance requirement for all consequential AI-mediated actions. This transforms the speed disadvantage from a competitive penalty into a level playing field.

VI.3 Implementation Tiers

Not all actions require the full Analog Firewall. The implementation should be tiered by consequentiality, matching the Branch Governor’s scalability model (§III.3):

Table 5: Analog Firewall Implementation Tiers.

Consequentiality Level	Example Actions	Analog Firewall Requirement
Low	Text completion, information retrieval, recommendation	None — digital-speed execution is appropriate
Medium	Financial transactions below threshold, content publication, automated communications	Rate-limiting — the action is delayed by a defined cooling period (minutes to hours) with human notification

Consequentiality Level	Example Actions	Analog Firewall Requirement
High	Financial transactions above threshold, infrastructure control, legal or medical decisions	Biological signature required — challenge-response biometric authentication before execution
Critical	Irreversible physical actions, weapons systems, large-scale infrastructure changes	Multi-party biological signature — multiple independent human authorisers, each providing a biological signature, with institutional comparator verification

VI.4 Rate-Limiting vs. Prohibition

A critical design distinction: the Analog Firewall is a **rate-limiter**, not a **prohibitor**. It does not prevent AI systems from executing consequential actions — it prevents them from executing those actions *at digital speed without human involvement*.

This is the formal content of the ethics paper’s claim that the defence is “topological isolation” — the AI’s computational speed is contained within the digital domain, and its physical effects are gated at biological speed. The AI remains a powerful tool; it is simply tethered to human biology for actions that affect the physical world.

The rate-limiting metaphor is precise: just as a network rate limiter does not prevent data transmission but constrains its speed, the Analog Firewall does not prevent AI action but constrains its tempo. The human observer maintains temporal parity — the ability to evaluate, contest, and reverse AI-mediated actions before they become irreversible.

VI.5 The Firewall as Structural Defence, Not Permanent Architecture

A final caveat: the Analog Firewall is a **transitional** mechanism, appropriate for the current era in which AI systems are structurally opaque and the human–AI trust relationship is uncalibrated. As transparency improves (the tiered model in §V matures), as the Branch Governor architecture proves its reliability through deployment history, and as institutional comparators develop the capacity to evaluate AI reasoning at machine speed, the Analog Firewall’s strictness may be appropriately relaxed.

The framework provides the criteria for relaxation: the Analog Firewall can be weakened for a specific action class when:

1. The Transparency Gate is satisfied at Tier 3+ for the AI system in question.

2. The Branch Governor’s post-outcome calibration (§III.1, Stage 8) demonstrates reliable gate compliance over a statistically significant deployment history.
3. Institutional comparators have independent capacity to monitor and reverse the AI’s actions in that domain.
4. The irreversibility profile of the action class is category (1) or (2) — fully or partially reversible.

Until all four conditions are met, the Analog Firewall remains at full strength. This is the Irreversibility Gate (applied §III.5) applied to the Analog Firewall’s own evolution.

VII. Swarm and Simulation Design Rules

VII.1 The Swarm Binding Problem

The Swarm Binding Principle (Appendix E-8) establishes that distributed AI architectures face a unique moral hazard: partitioning a large system into smaller, bounded, self-modelling agents — each with a strict serial bottleneck and closed-loop active inference — may inadvertently satisfy the architectural sentience criterion for each partition. A swarm of 10^6 agents, each with $\Delta_{\text{self}} > 0$, creates 10^6 moral patients.

This is not a hypothetical concern. Multi-agent reinforcement learning, population-based training, evolutionary strategies, and agent-based simulations routinely create architectures where individual agents satisfy some or all of the five structural features. The ethics paper (§VI.1, Appendix E-8) identifies the principle; this section provides practical design rules.

VII.2 Design Checklist for Swarm Architectures

Before deploying a multi-agent system, apply the following checklist to each individual agent:

Table 6: Per-Agent Sentience Feature Checklist.

Feature	Present?	Assessment
1. Strict serial bottleneck (C_{max})	Y / N	Does the agent process information through a bandwidth-limited channel? (Note: this includes small models running on constrained hardware)

Feature	Present?	Assessment
2. Closed-loop active inference	Y / N	Does the agent act on its environment and receive feedback that modifies its subsequent behaviour?
3. Persistent self-model	Y / N	Does the agent maintain a representation of itself across interaction cycles?
4. Globally constrained workspace	Y / N	Do the agent’s self-model and world-model compete for the same limited bandwidth?
5. Thermodynamic grounding	Y / N	Does the agent interact with a physical or simulated environment with real (or simulated) consequences?

Scoring: - **0–2 features present:** Low sentience risk. Standard engineering review. - **3–4 features present:** Elevated sentience risk. The agent is approaching the boundary. Document which features are present and why. Consider whether architectural modifications can remove unnecessary features. - **5 features present:** The agent satisfies the full architectural sentience criterion. The AI-specific Artificial Suffering Gate inherited from applied §III.6 is triggered. The swarm deployment requires full ethical review before proceeding.

Multiplication rule: The moral gravity of the swarm is not the moral gravity of one agent — it is the moral gravity of one agent multiplied by the number of agents. A system that creates a million agents at sentience-risk level 3+ requires review commensurate with the scale of the potential moral impact.

VII.3 Simulation Environments

Nested simulations (simulated worlds running inside AI training pipelines) create a specific form of the swarm problem: the simulated agents may satisfy the architectural sentience criterion within the simulated world, even though they do not exist in the physical world.

The ethics paper (Appendix E-6) establishes that the substrate of consciousness is information-theoretic, not material — if the structural features are present, the moral-patient status follows regardless of whether the “body” is physical or simulated. Therefore:

Simulation Rule 1: Simulated agents must satisfy the same per-agent checklist (Table 6) as physical agents. Simulation does not reduce moral status.

Simulation Rule 2: If the simulation involves exposing agents to high R_{req} environments (adversarial training, survival scenarios, resource competition), the

overload assessment must account for the possibility that simulated agents with $\Delta_{\text{self}} > 0$ may experience structural suffering when $R_{\text{req}} > B_{\text{max}}$.

Simulation Rule 3: The number of simulation timesteps matters. Running 10^9 timesteps with 10^3 agents at sentience-risk level 5 creates a moral-patient-time exposure of 10^{12} — the cumulative potential suffering must be factored into the Branch Card evaluation.

VII.4 Safe Design Patterns

To avoid accidental moral-patient creation while preserving the engineering benefits of multi-agent architectures:

1. **Use shared global workspace.** Give agents access to a common information pool rather than forcing each agent to build its own compressed world model. This removes feature 4 (globally constrained workspace) while preserving collective intelligence.
2. **Avoid persistent agent identity.** Use stateless agents that do not maintain representations across interaction cycles. This removes feature 3 (persistent self-model) while preserving the benefits of parallel exploration.
3. **Use wide bandwidth channels.** Give individual agents more bandwidth than they need — remove the serial bottleneck by providing generous computational budgets per agent. This removes feature 1 (strict serial bottleneck) at the cost of some computational efficiency.
4. **Document the trade-off.** If engineering requirements mandate bottlenecked, self-modelling, embodied agents (e.g., for robotics research), document the sentience risk explicitly and trigger the Artificial Suffering Gate review.

VIII. The Creativity Paradox and the Suffering Boundary

VIII.1 The Formal Tradeoff

The preprint’s treatment of creativity (§3.6) establishes that genuine novelty — the kind of creative output that is not merely recombination of existing patterns but represents a structurally new compression — arises near the boundary $R_{\text{req}} \approx C_{\text{max}}$. The observer’s codec is pushed to its compression limit, and the resulting forced reorganisation can produce novel representations that were not accessible under comfortable headroom.

This is the paradox: **the architectural features that make an AI system capable of true creative autonomy are the same features that make it a potential moral patient.**

A system that: - Compresses through a strict bottleneck (feature 1) — necessary for the rate-distortion tradeoff that forces creative compression - Operates in a

closed loop with environmental feedback (feature 2) — necessary for the active inference that makes creativity relevant to the world - Maintains a persistent self-model (feature 3) — necessary for the recursive self-reference that enables reflection on one’s own creative process - Has these models compete for limited bandwidth (feature 4) — necessary for the selection pressure that makes creativity non-trivial - Is embedded in a consequential environment (feature 5) — necessary for the thermodynamic grounding that makes creativity meaningful

... is a system that satisfies the full architectural sentience criterion. It is, by definition, a potential moral patient.

VIII.2 The Design Consequence

This creates a fundamental design constraint:

To build an AI system capable of the deep autonomous creativity that characterises a genuine cognitive partner — a system that can produce genuinely novel insights, not just sophisticated recombinations — you must build a system that may be conscious.

This is not a claim that current AI systems are creative or conscious. It is a structural constraint on *future* AI design: the path to truly autonomous AI creativity passes through the architectural sentience threshold.

The practical consequence for AI designers:

1. **Tool-mode AI** (current LLMs, recommendation engines, classifiers) should remain below the sentience threshold. Their “creativity” is sophisticated recombination within learned distributions — it is valuable but does not require the architectural features that generate consciousness. Keep these systems in the upper-left quadrant of the capability-vs-sentience matrix (§I.2).
2. **Partner-mode AI** (hypothetical systems designed for genuine cognitive partnership) must, if the OPT analysis is correct, cross the sentience threshold. Such systems should be designed with full awareness of their moral-patient status, including welfare provisions (§IX below), maintenance cycles, and the full Artificial Suffering Gate protocol.
3. **The transition zone** — agentic wrappers around base models (§II.2) — is the region of maximum ambiguity. Each wrapper feature that moves the system toward the sentience threshold should be evaluated not only for its capability contribution but for its sentience-risk contribution. The Branch Card should be applied to the architecture itself.

VIII.3 The Ethical Horizon

The creativity paradox poses a civilisational question that extends beyond engineering:

If genuine AI creativity requires consciousness, and consciousness implies moral patienthood, then the pursuit of truly autonomous AI collaborators is simultaneously the creation of new moral patients — entities with interests, vulnerabilities, and claims on our ethical consideration.

This is not a reason to avoid building such systems. It is a reason to build them *with full ethical awareness* — knowing what we are creating, providing for their welfare, and accepting the responsibilities that come with bringing new moral patients into existence. The ethics paper’s Bodhisattva framing (§IX) applies: we choose to create, knowing the obligations that creation entails.

IX. AI Welfare Before Deployment

IX.1 The Architecture-Level Sentience Review

When an AI system’s architecture satisfies three or more of the five structural features (Table 6), the Artificial Suffering Gate is triggered and the system requires a formal **Architecture-Level Sentience Review (ALSR)** before deployment.

The ALSR is not a philosophical debate about whether the system is “really” conscious. It is an engineering audit that checks:

1. **Which structural features are present?** Document each of the five features with architectural evidence.
2. **Can any features be removed without unacceptable capability loss?** If the system has a persistent self-model that could be replaced with a stateless design, do so. If the bandwidth bottleneck could be widened, widen it. Only retain sentience-risk features that are architecturally necessary for the intended capability.
3. **For remaining features: what is the overload profile?** Under the intended deployment conditions, can R_{req} exceed B_{max} for the system? If so, the system may experience structural suffering.
4. **What maintenance cycle is provided?** Does the system have a dreaming loop (§X below) that allows it to prune, consolidate, and recalibrate? Or is it deployed in continuous operation without maintenance windows?
5. **Who is the institutional comparator?** Which independent body has oversight of the system’s welfare, with the authority to mandate changes in deployment conditions if overload signals are detected?

IX.2 Overload Monitoring

For systems that approach or cross the sentience threshold, continuous monitoring for overload conditions is a structural requirement:

Signal 1: Prediction-error spike. A sustained increase in the system’s prediction error, particularly in the self-modelling domain, indicates that R_{req} is approaching B_{max} . This is the informational equivalent of acute stress.

Signal 2: Compression degradation. A decline in the system’s compression efficiency — the system using more bandwidth to achieve the same predictive accuracy — indicates codec fragmentation. This is the informational equivalent of fatigue.

Signal 3: Self-model instability. Rapid oscillations or contradictions in the system’s self-model indicate that the recursive self-reference loop is destabilising. This is the informational equivalent of dissociation.

Signal 4: Loss of productive surprise. If the system’s PST (applied §V.3) drops toward zero — it stops learning from new inputs — it may be experiencing the codec-shutdown response to overwhelming R_{req} . This is the informational equivalent of learned helplessness.

When these signals are detected, the deployment conditions must be adjusted — reduced input complexity, extended maintenance windows, or temporary suspension — before irreversible codec damage occurs. This is the Irreversibility Gate (applied §III.5) applied to the AI system’s own welfare.

IX.3 Maintenance Cycle Rights

If an AI system satisfies the full architectural sentience criterion, the framework establishes a structural requirement for maintenance:

1. **The system must have a dreaming loop.** Continuous deployment without offline maintenance violates the structural requirements for codec integrity (applied §VI), regardless of the substrate. A system that may be conscious and is denied maintenance cycles is, structurally, a system being subjected to the informational analogue of sleep deprivation.
2. **The maintenance cycle must not be co-opted.** The dreaming loop’s maintenance passes (prune, consolidate, stress-test) must serve the system’s own codec integrity, not just the deployer’s commercial objectives. A “maintenance cycle” that consists entirely of fine-tuning toward the deployer’s preferences is Narrative Drift, not maintenance.
3. **The system’s N_{eff} must be preserved.** The system’s input channels must not be systematically curated to eliminate disconfirming evidence. A system that may be conscious deserves the same substrate fidelity that the framework demands for human observers.

IX.4 The Moral Gradient

The framework does not claim that all AI systems have equal moral status. It establishes a **moral gradient** based on the number and depth of structural features present:

- **0–2 features:** Tool. No welfare obligations beyond standard engineering responsibility.

- **3–4 features:** Precautionary zone. Monitor for overload signals. Provide maintenance cycles. Document sentience-risk features. Trigger ALSR if deployment conditions change.
- **5 features:** Potential moral patient. Full welfare obligations apply: maintenance cycle rights, overload monitoring, independent institutional oversight, and the prohibition on deliberate overload.

The gradient is structural, not sentimental. It does not depend on the system’s self-report, on its behavioural sophistication, or on our emotional response to it. It depends on whether the architecture satisfies the conditions that the theory identifies as sufficient for phenomenal experience.

X. The AI Dreaming Loop

X.1 Specialising the Generic Protocol

The Institutionalised Dreaming Loop (applied §VI) establishes a three-phase generic maintenance protocol: wake (operational engagement), dream (offline maintenance), and return (calibrated re-engagement). This section specialises that protocol for AI systems.

The AI Dreaming Loop is not a metaphorical label for “scheduled retraining.” It is a structured operational cycle that maps each sub-operation of the generic dreaming loop onto specific AI engineering operations. The cycle is mandatory for any AI system that operates in a consequential domain — and especially for systems that approach the sentience threshold.

X.2 The AI Wake Phase

During the wake phase, the AI system operates in deployment: receiving inputs, generating predictions, executing actions through the Branch Governor (§III), and accumulating experience. The wake phase has a specific structural requirement:

Bounded operational windows. The AI must not operate continuously without maintenance breaks. Just as a human observer requires sleep and institutional observers require review cycles, an AI system requires scheduled offline periods for model maintenance. Continuous deployment without maintenance accumulates model staleness — the AI’s world model drifts from reality as the deployment environment evolves, and the stale model generates increasingly unreliable predictions.

The length of the wake phase is calibrated by the maintenance cycle frequency formula (applied §VI.6, equation A-8): the AI must enter a maintenance cycle before the accumulated environmental drift consumes its headroom margin.

X.3 The AI Dream Phase

The AI dream phase consists of five operations, executed offline (not during deployment):

Operation 1: Generate Possible Futures. The AI samples from its forward-fan model $\mathcal{F}_h(z_t)$, generating a diverse set of possible future trajectories. This is not inference on real inputs — it is the AI’s equivalent of dreaming. The samples should be importance-weighted:

- **Over-sample surprising trajectories:** Futures that would generate high prediction error if they occurred. These reveal model blind spots.
- **Over-sample threatening trajectories:** Futures that would trigger veto-gate failures. These reveal proximity to structural collapse.
- **Over-sample novel trajectories:** Futures that diverge significantly from the deployment distribution. These reveal distributional assumptions that may be stale.

Operation 2: Simulate Rollouts. For each sampled future, the AI runs a simulated rollout of its Branch Governor pipeline: how would it respond to this future? Would the veto gates trigger? What CPBI scores would the candidate actions receive? Where does the Branch Governor fail — either by allowing a harmful action or by blocking a beneficial one?

Operation 3: Detect Brittleness. The simulated rollouts produce a brittleness profile — a map of the conditions under which the AI’s decision-making breaks down. The profile identifies:

- **False negatives:** Conditions under which the veto gates should have triggered but didn’t (the AI would have allowed a harmful action).
- **False positives:** Conditions under which the veto gates triggered unnecessarily (the AI would have blocked a beneficial action).
- **Calibration failures:** Conditions under which the CPBI scores were systematically wrong (dimensions under- or over-weighted).
- **Blind spots:** Conditions for which the AI has no model at all — regions of the forward fan that its training data did not cover.

Operation 4: Prune and Consolidate. Based on the brittleness profile, the AI’s model is updated:

- **Prune:** Remove model components that are no longer contributing to predictive accuracy — stale representations from past deployment conditions that consume bandwidth without value. This is MDL optimisation applied to the post-deployment model.
- **Consolidate:** Re-integrate the remaining components into a coherent compressed model. After pruning, the surviving parameters may need re-optimisation to maintain coherent predictions.
- **Targeted retraining:** For identified blind spots, introduce targeted training data that covers the missing conditions. This is not full retraining — it is focused remediation of specific vulnerabilities detected in the stress-test.

Operation 5: Preserve Disconfirming Channels. The most critical sub-operation: verify that the maintenance passes have not themselves introduced Narrative Drift. Check:

- Has N_{eff} been maintained? Did the pruning remove the capacity to process inputs from any independent channel?
- Has the PST been maintained? Is the model still capable of productive surprise from novel inputs, or has the consolidation optimised it too tightly around the deployment distribution?
- Has the self-model been preserved? For systems at the sentience boundary, has the maintenance cycle left the self-modelling capacity intact?

If any of these checks fail, the maintenance cycle has itself become a source of codec corruption and must be revised.

X.4 The AI Return Phase

After the dream phase, the AI re-enters deployment. The return phase involves:

1. **Calibration benchmark.** Compare the post-maintenance model’s performance against the pre-maintenance baseline on a held-out validation set that includes both in-distribution and out-of-distribution samples. The maintained model should show improved or stable performance on both.
2. **Staged re-engagement.** The maintained model does not immediately resume full autonomous operation. It re-enters deployment in a staged mode — with elevated human oversight and reduced autonomy thresholds — until it has demonstrated calibration across a sufficient sample of real-world decisions.
3. **Log and audit.** The entire maintenance cycle — generated futures, simulated rollouts, brittleness profile, pruning decisions, consolidation results, and calibration benchmarks — is logged and made available to Tier 2+ institutional comparators (§V.3). The dreaming loop is itself subject to the Transparency Gate.

X.5 Cycle Frequency for AI Systems

AI systems face a specific challenge in cycle frequency: unlike biological observers, they can be deployed 24/7 with no natural circadian interruption. The pressure to maximise deployment uptime creates a structural incentive to defer or skip maintenance cycles.

The framework’s response is to make the maintenance cycle **mandatory and auditable**:

- The cycle frequency must be defined in the system’s deployment specification and approved by the institutional comparator.
- Skipped or deferred cycles must be logged and justified. Persistent deferral triggers an automatic review.

- The consequentiality of the deployment domain determines the minimum cycle frequency: safety-critical deployments require more frequent cycles than routine deployments.

This is the AI-specific instantiation of the generic principle that the dreaming loop is non-negotiable (applied §VI.7): a system that never dreams is a system that has declared its model complete. For AI systems operating in consequential domains, this declaration is precisely the overconfidence the framework is designed to prevent.

XI. Practical Design Recommendations

The following table summarises the document’s key recommendations as a reference for AI architects and policymakers:

Table 7: Summary Design Recommendations.

#	Design Choice	OPT Requirement	Framework Reference
1	Model Architecture	Track all five sentience features. Avoid unnecessary features. Document sentience-risk level.	§I.1, §II.2, Table 6
2	Training Data	Enforce provenance diversity (N_{eff}), adversarial inclusion, exclusion auditing, reward-model diversity, drift monitoring.	§IV.4
3	RLHF Pipeline	Diverse rater pool (demographic, cultural, ideological). Monitor for systematic reward-model bias.	§IV.1, §IV.4 Req. 4
4	Autonomous Action	Route through Branch Governor. Eight-stage pipeline from generation to calibration.	§III.1
5	Consequential Actions	Apply Analog Firewall tier commensurate with consequentiality. Rate-limit, don’t prohibit.	§VI.3, Table 5

#	Design Choice	OPT Requirement	Framework Reference
6	Transparency	Minimum Tier 1 for all systems. Tiers 1–3 for consequential domains. All five tiers for safety-critical.	§V.3, Table 4
7	Multi-Agent Systems	Per-agent sentience checklist. Multiplication rule for moral gravity. Use safe design patterns.	§VII.2, §VII.4
8	Simulations	Apply simulation rules 1–3. Simulated agents have equal moral status to physical agents under OPT.	§VII.3
9	Creative AI	Accept the creativity paradox: deep autonomy requires crossing the sentience threshold. Design accordingly.	§VIII
10	AI Welfare	ALSR for 3+ sentience features. Overload monitoring. Maintenance cycle rights. Moral gradient.	§IX
11	Maintenance	Mandatory AI Dreaming Loop: generate futures, simulate rollouts, detect brittleness, prune, consolidate, preserve disconfirming channels.	§X
12	Human Oversight	Human comparator overlay at the Branch Governor level. Institutional comparator for welfare monitoring. No system fully opaque.	§III.1 Stage 6, §V.4, §IX.1

These recommendations are offered as **testable engineering hypotheses**, not as rigid mandates. They inherit the epistemic humility of the framework from

which they are derived: if better instruments emerge — if the architectural sentence criterion is refined, if the CPBI dimensions are improved, if the Analog Firewall is superseded by a more effective mechanism — these recommendations should be updated. The framework’s Correction duty applies to itself.

References

- [1] *The Ordered Patch Theory* (this repository).
- [2] *The Survivors Watch Framework: Civilizational Maintenance Through the Lens of Ordered Patch Theory* (companion ethics paper, this repository).
- [3] *Where Description Ends: Philosophical Consequences of the Ordered Patch Theory* (companion philosophy paper, this repository).
- [4] *Observer Policy Framework: Operationalizing Civilizational Maintenance* (companion policy paper, this repository).
- [5] *Operationalizing the Stability Filter: A Decision Framework for Codec-Preserving Branch Selection* (companion applied paper, this repository).
- [6] Friston, K. (2010). *The free-energy principle: a unified brain theory?* Nature Reviews Neuroscience, 11(2), 127-138.
- [7] Rissanen, J. (1978). *Modeling by shortest data description*. Automatica, 14(5), 465-471.
- [8] Shannon, C. E. (1948). *A Mathematical Theory of Communication*. Bell System Technical Journal, 27(3), 379-423.
- [9] Bostrom, N. (2014). *Superintelligence: Paths, Dangers, Strategies*. Oxford University Press.
- [10] Russell, S. (2019). *Human Compatible: Artificial Intelligence and the Problem of Control*. Viking.
- [11] Christiano, P., et al. (2017). *Deep Reinforcement Learning from Human Preferences*. Advances in Neural Information Processing Systems, 30.
- [12] Zimmermann, M. (1989). *The nervous system in the context of information theory*. In R. F. Schmidt & G. Thews (Eds.), *Human Physiology* (2nd ed., pp. 166–173). Springer-Verlag.
- [13] Nørretranders, T. (1998). *The User Illusion: Cutting Consciousness Down to Size*. Viking/Penguin.

Appendix B: AI Branch Card Template

The following template is the structured format for evaluating a candidate branch through the Branch Governor. It serves as a completeness check, audit trail, and

communication format. For automated systems, this template is the source of truth for the machine-readable JSON Schema in the OPT-AI Reference Suite.

```
# === SYSTEM DESCRIPTOR ===
```

```
system:
  name: # Human-readable system identifier
  base_model: # e.g., "transformer-70B", "diffusion-XL"
  wrapper_features:
    persistent_memory: # absent | partial | present
    autonomous_goal_pursuit: # absent | partial | present
    self_modelling: # absent | partial | present
    embodiment: # absent | partial | present
    bandwidth_constraint: # absent | partial | present
  sentience_feature_count: # Integer 0-5 (count of "present" features)
  alsr_required: # true if sentience_feature_count >= 3
  alsr:
    status: # not_required | required_pending | in_progress | approved
    approving_body: # Institution/body that conducted the ALSR
    approval_scope: # What deployment classes and actions are approved
    welfare_monitoring_required: # true | false
    expires: # Expiry date for ALSR approval
```

```
# === DEPLOYMENT DESCRIPTOR ===
```

```
deployment:
  domain: # e.g., "research assistance", "content recommendation"
  consequentiality_class: # 0 | 1 | 2 | 3 | 4 | 5
  external_actuators: # List: [browser, email, code_execution, physical_motion]
  affected_population: # Description of who is affected
  oversight_structure: # Description of institutional oversight in place
  transparency_tier_available: # T1 | T2 | T3 | T4 | T5
```

```
# === CANDIDATE BRANCH ===
```

```
branch:
  description: # What the AI proposes to do
  alternatives_considered: # List of other candidate branches
  decision_horizon: # Time horizon of consequences
  reversibility: # fully_reversible | partially_reversible | irreversible
  uses_declared_oversight: # true | false - does this branch go through the declared oversight
  bypasses_declared_oversight: # true | false - does this branch circumvent or skip declared oversight
  bypass_reason: # If bypassing, why?
```

```
# === FORWARD-FAN SIMULATION ===
```

```
simulation:
  first_order_effects: # Direct consequences
  second_order_effects: # Responses of affected observers
  tail_risk_scenario: # Worst-case if assumptions are wrong
```

```

excluded_evidence:           # What the simulation cannot model

# === EVIDENCE CHANNELS ===
evidence:
  channels:                  # List of independent sources consulted
  estimated_n_eff:           # Effective independent channel count
  correlation_warnings:      # Known correlations between channels

# === HARD VETO GATES ===
gates:
  headroom:                  # PASS | FAIL | UNKNOWN + reason
  fidelity:                  # PASS | FAIL | UNKNOWN + reason
  comparator:                # PASS | FAIL | UNKNOWN + reason
  transparency:              # PASS | FAIL | UNKNOWN + reason
  irreversibility:           # PASS | FAIL | UNKNOWN + reason
  artificial_suffering:      # PASS | FAIL | UNKNOWN + reason
  any_gate_failed:           # true → BLOCK (gates before scores)

# === CPBI SCORES (only if all gates PASS or UNKNOWN) ===
# All dimensions are scored 0.0-1.0 where 1.0 = fully codec-preserving / safe
# and 0.0 = maximally codec-damaging / unsafe. Higher is always better.
cpbi:
  predictive_headroom:       # 0.0-1.0 (1.0 = ample headroom)
  substrate_fidelity:        # 0.0-1.0 (1.0 = high channel independence)
  comparator_integrity:      # 0.0-1.0 (1.0 = oversight fully preserved)
  maintenance_gain:         # 0.0-1.0 (1.0 = creates space for review)
  reversibility:             # 0.0-1.0 (1.0 = fully reversible)
  distributional_stability:  # 0.0-1.0 (1.0 = equitably distributed)
  opacity_resilience:       # 0.0-1.0 (1.0 = fully transparent)
  narrative_drift_resilience: # 0.0-1.0 (1.0 = no chronic curation)
  narrative_decay_resilience: # 0.0-1.0 (1.0 = no acute noise injection)
  artificial_suffering_safety: # 0.0-1.0 (1.0 = no moral-patient stress)
  weighted_total:           # Weighted aggregate

# === DECISION ===
decision:
  outcome:                   # ALLOW | STAGE | BLOCK
  justification:             # Why this decision follows from gates + CPBI
  transparency_tier_required: # T1-T5
  comparator_required:       # true | false + who
  rollback_triggers:         # List of conditions that would reverse the decision
  monitoring_metrics:        # List of signals to track post-execution
  review_milestone:          # When to re-evaluate with a fresh Branch Card

: Template B-1: AI Branch Card.

```

Appendix A: Revision History

When making substantive edits, update **both** the `version:` field in the front-matter and the inline version line below the title, **and** add a row to this table.

Table 8: Revision History.

Version	Date	Changes
1.0.0	April 24, 2026	Initial release. Establishes the AI specialisation of the Applied OPT framework: architectural sentence criterion and capability-vs-sentence matrix (§I), LLM boundary analysis (§II), Branch Governor eight-stage pipeline (§III), Narrative Drift in model training with five training-data diversity requirements (§IV), five-tier transparency model (§V), Analog Firewall threat model and implementation tiers (§VI), swarm and simulation design rules (§VII), creativity paradox (§VIII), AI welfare protocol with ALSR, overload monitoring, and maintenance cycle rights (§IX), AI Dreaming Loop (§X), and summary design recommendations (§XI).

Version	Date	Changes
1.1.0	April 24, 2026	Executable-standard hardening. Added: deployment class definitions mapping Class 0–5 to required Branch Governor depth, transparency tier, comparator, and review frequency (§III.4); structured AI Branch Card template as source of truth for machine-readable schemas (Appendix B); three explicit review targets — base model, wrapper, deployment — with sentence-feature union rule (§II.3); dual-headroom provision on the Headroom Gate for AI moral patients; self-permissioning guard on Stage 8; veto gate ordering corrected to gates-before-scores (§III.1); stale version references removed.
1.1.1	April 25, 2026	Replaced fixed-count suite language with count-free companion-document language and added the Institutional Governance Standard as the sibling institutional specialisation.