# Fisher Information Geometry and the Thermodynamic Cost
# of Self-Referential Calibration

Sidong Liu, PhD

iBioStratix Ltd

sidongliu@hotmail.com

February 2026

## Abstract

Papers I and II of the T-DOME series [1, 2] established that persistent agents must carry non-Markovian memory (Paper I) and must spontaneously break the gauge symmetry of their internal Clifford algebra $Cl(V, q)$ to form a compressed reference frame—the "ego" $\mathfrak{E} = (\mathcal{F}^*, V_{\mathrm{fg}}^*)$ (Paper II). Paper II concluded with the **Delusion Trap**: under environmental drift, a fixed reference frame decouples from the optimal gauge on the logarithmic timescale $t_{\mathrm{del}} = \Lambda^{-1} \ln(\pi/4\theta_0)$, and the agent cannot detect this failure from within its own foreground subspace $V_{\mathrm{fg}}$.

In this final work we derive the theory of **self-referential calibration**. We show that while the agent cannot observe the background subspace $V_{\mathrm{bg}}$ directly, it can measure the **Fisher information** of its own prediction-residual stream with respect to its frame parameters $\sigma$. We prove three main results:

1. **Drift Detectability** (Theorem 14): environmental drift generates a quadratically growing signal in the self-referential Fisher information $\mathcal{I}_F(\sigma)$, detectable before the Delusion Trap closes.

2. **Self-Referential Cramér–Rao Bound** (Theorem 18): the agent's drift-estimation error is bounded below by $1/(n_{\mathrm{eff}}\,\mathcal{I}_F + \mathcal{I}_{\mathrm{ego}})$, where $\mathcal{I}_{\mathrm{ego}}$ quantifies the rigidity of the ego prior.

3. **Thermodynamic Cost of the Loop** (Theorem 28): the minimum dissipation rate for self-referential calibration is $\dot{W}_{\mathrm{loop}} \geq k_B T \ln 2\,[h_\mu\,k^* + \mathcal{C}_{\mathrm{meta}}] + \mathcal{L}^2/\tau_{\mathrm{recalib}}^2$, where $\mathcal{L}$ is the thermodynamic length of the frame update and $\tau_{\mathrm{recalib}}$ is the recalibration time.

The calibration loop satisfies a Lyapunov tracking bound (Theorem 23), keeping the mismatch within a neighbourhood whose size is set by the ratio of environmental drift speed to adaptation rate. We identify this loop as the minimal physical realisation of *reflexivity*—estimating drift from residual statistics and correcting the frame via Lyapunov-monitored natural gradient descent. Combining with Papers I and II, we state a **Four-Part Structure Proposition** (Proposition 27): within the class

of agents satisfying (C1)–(C5), a sufficient architecture for persistence under drift requires (1) an external observable geometry, (2) an internal control algebra, (3) a self-monitoring Lyapunov function, and (4) biased non-Markovian memory.

# 1 Introduction

## 1.1 Context: The Delusion Trap

Paper II of this series [2] established that persistent agents under bounded computation must spontaneously break the gauge symmetry of their internal algebra $Cl(V, q)$, selecting a privileged reference frame $\mathcal{F}^*$ that compresses the memory kernel into a tractable $k^*$-dimensional foreground subspace $V_{\mathrm{fg}}$. This gauge fixing—the "ego" $\mathfrak{E} := (\mathcal{F}^*, V_{\mathrm{fg}}^*)$—is not an additional hypothesis but the survival-optimal strategy under bounded rationality.

However, Paper II's final theorem revealed a fatal consequence. Under environmental drift (spectral-density parameters changing at rate $\varepsilon$), the mismatch angle between the agent's fixed frame and the instantaneous optimal frame grows as $\theta(t) = \theta_0\, e^{\Lambda t}$ (Paper II, Definition 27), where $\Lambda \sim \varepsilon/\tau_{\mathrm{adapt}}$ is the environmental Lyapunov exponent. Beyond the *delusion time*

$$t_{\mathrm{del}} = \frac{1}{\Lambda}\, \ln\left(\frac{\pi/4}{\theta_0}\right), \tag{1}$$

three catastrophic failures occur simultaneously (Paper II, Theorem 29):

1. The hidden survival component dominates: $|\mathcal{S}_{\mathrm{hid}}| > |\mathcal{S}_{\mathrm{vis}}|$.

2. The agent's update direction anti-correlates with the true survival gradient: $\langle \nabla_u \mathcal{S}_{\mathrm{vis}},\, \nabla_u \mathcal{S}_{\mathrm{full}} \rangle < 0$.

3. All four bias terms ($\mathcal{B}_{\mathrm{select}}$, $\mathcal{B}_{\mathrm{frame}}$, $\mathcal{B}_{\mathrm{center}}$, $\mathcal{B}_{\mathrm{inc}}$) operate within $V_{\mathrm{fg}}$ and cannot register changes in the background $V_{\mathrm{bg}}$.

Paper II further showed (Remark 31) that "dithering"—randomly probing the background subspace—fails because the agent has no gradient signal to indicate *when* or *where* to probe. The exponential divergence in $V_{\mathrm{bg}}$ is invisible until it dominates, at which point it is too late.

*The present paper provides the escape.*

## 1.2 Position within Papers I–III

This paper is the third and final of the T-DOME framework, closing the three-paper sequence.

| Framework | Question | Result | Status |
|---|---|---|---|
| HAFF [17, 18] | How does geometry emerge? | Algebra → Geometry | Complete |
| Q-RAIF [20, 21] | What algebra must an observer have? | $Cl(V, q) \hookrightarrow Cl(1, 3)$ | Complete |
| T-DOME I [1] | Why must agents carry memory? | Markovian ceiling; memory as necessity | Complete |
| T-DOME II [2] | Why must agents break symmetry? | Reference-frame selection under bounded computation | Complete |
| **T-DOME III** (this work) | How does self-calibration arise? | Fisher self-referential bound; thermodynamic cost of reflexivity | **This paper** |

The three T-DOME papers form an irreversible logical chain:

1. **Paper I:** Without memory, a system is trapped in the Markovian present. Memory breaks this trap but floods the system with unbounded historical data.

2. **Paper II:** Unbounded memory under finite computational resources causes processing collapse. Spontaneous symmetry breaking resolves the overload but introduces systematic bias.

3. **Paper III (this work):** Uncorrected bias diverges from a changing environment. A self-referential calibration loop—monitoring the Fisher information of one's own prediction stream—resolves the bias but requires a second-order control structure and an irreducible thermodynamic cost.

Each resolution creates the precondition for the next crisis: memory enables overload, compression enables bias, and bias demands calibration. Only the complete closure *Paper I + Paper II + Paper III* allows a system to persist under the Second Law in a drifting environment.

## 1.3   The Information-Geometric Insight

The key observation that resolves the Delusion Trap is subtle: *while the agent cannot observe $V_{\text{bg}}$ directly, it can observe the statistical properties of its own prediction residuals in $V_{\text{fg}}$.*

The prediction residual $e(t) := \mathcal{S}_{\text{vis}}(t) - \mathcal{S}_{\text{vis}}^{(\text{pred})}(t)$ lies in $V_{\text{fg}}$ by construction. Its *value* carries no information about the background. But its *distribution*—the probability law $p(e \,|\, \sigma)$, parametrised by the gauge-fixing parameter $\sigma$—does depend on $\sigma$, because the projection $\Pi_{\mathcal{F}}(\sigma)$ determines which environmental correlations are captured and which are discarded.

When the frame $\sigma$ drifts away from the optimal $\sigma^*$, the residual distribution shifts. The *Fisher information metric*

$$g_{ij}(\sigma) = \mathbb{E}_\sigma \left[ \frac{\partial \log p(e \,|\, \sigma)}{\partial \sigma^i} \frac{\partial \log p(e \,|\, \sigma)}{\partial \sigma^j} \right] \tag{2}$$

measures the sensitivity of this distribution to changes in $\sigma$. A spike in $g_{ij}$—a "stress" in the agent's internal geometry—is the signal that the reference frame is becoming stale.

This is the mathematical realisation of the "second-order operation" demanded by Paper II, Section 7.5: the agent does not need to see the truth (the full $Cl(V, q)$), but only the *rate of change of its own prediction error* as a function of its frame parameters. Fisher information is precisely this quantity.

## 1.4 Relation to Architectural Incompleteness

The architectural incompleteness result [19] established *architectural incompleteness*: the observable-algebra framework cannot self-ground. Paper II provided a partial operational response (the ego as gauge fixing under bounded computation). The present paper provides the final operational response: the self-referential calibration loop cannot *eliminate* architectural incompleteness, but it can *track* the consequences of incompleteness in real time. The Lyapunov function $V(\sigma)$ monitors the distance between the agent's frame and the optimal frame without requiring access to the "complete" description—it operates entirely within the agent's own predictive statistics.

## 1.5 Scope and Disclaimers

1. *Reflexivity* refers throughout to second-order control: the ability of a system to monitor and adjust its own monitoring process. It carries *no* implication of phenomenal consciousness, subjective experience, or qualia.

2. The self-referential calibration loop does not *eliminate* the ego's bias; it tracks and compensates for drift in the bias. The four bias terms of Paper II persist in the calibrated phase.

3. The thermodynamic cost bounds are information-theoretic lower bounds, not claims about specific physical implementations.

4. The framework applies to systems satisfying (C1)–(C5) (Section 2.6). It is not a universal theory of agency.

**Related work.** The Fisher information metric on statistical manifolds was introduced by Rao [3] and shown to be unique by Čencov [7]. The natural gradient and information geometry were developed by Amari [6, 5]. Thermodynamic length and optimal finite-time transformations were established by Crooks [8] and Sivak–Crooks [9]. The connection between Fisher information and entropy production was formalised by Ito [12] and Barato–Seifert [11]. Second-order cybernetics originates with Ashby [13] and von Foerster [14]. Adaptive control and self-tuning regulators are treated in [15]. The Bayesian Cramér–Rao bound (van Trees inequality) is from [10].

**Summary of contributions.** This paper establishes three main results:

1. **Drift Detectability** (Theorem 14): the self-referential Fisher information of the prediction-residual stream grows quadratically with accumulated drift, providing a detectable signal before the Delusion Trap closes.

2. **Self-Referential Cramér–Rao Bound** (Theorem 18): drift-estimation precision is bounded by the sum of data Fisher information and ego rigidity.

3. **Thermodynamic Cost** (Theorem 28): the self-calibration loop requires a minimum dissipation rate with three distinct components (sensing, computing, actuating).

# 2 Mathematical Preliminaries

## 2.1 Inherited Framework from Papers I and II

We briefly recall the key objects; the reader is referred to Papers I and II for full definitions and proofs.

**From Paper I [1].**

- **Survival functional.** $\mathcal{S}[\Lambda, \tau] := \Delta F - W[0, \tau]$ (Paper I, Eq. (9)).

- **Markovian Ceiling.** $\mathcal{S}[\Lambda^{\mathrm{M}}, \tau] \leq 0$ for all $\tau \geq 0$.

- **Memory kernel.** $\mathcal{K}(t, s)$: the non-Markovian superoperator encoding system–environment correlations.

- **Entropy rate.** $h_\mu := \lim_{T \to \infty} T^{-1} H(X_{0:T})$ (bits per unit time per algebraic component).

- **Predictive information.** $I_{\mathrm{pred}} := I(\overleftarrow{X}; \overrightarrow{X})$.

**From Paper II [2].**

- **Internal algebra.** $\mathcal{O}_{\mathrm{int}} = Cl(V, q)$, $D = \dim Cl(V, q) = 2^n$, gauge group $G = \mathrm{Aut}(Cl(V, q))$.

- **Gauge bundle.** $\pi : P \to M$, structure group $G$; a section $\sigma : M \to P$ is a reference frame.

- **Ego.** $\mathfrak{E} := (\mathcal{F}^*, V_{\mathrm{fg}}^*)$ with $k^* = \lfloor \mathcal{C}_{\mathrm{budget}} / h_\mu \rfloor$ foreground components.

- **Projected kernel.** $\mathcal{K}_\mathcal{F}(t, s) = \Pi_\mathcal{F} \mathcal{K}(t, s) \Pi_\mathcal{F}$.

- **Survival decomposition.** $\mathcal{S} = \mathcal{S}_{\mathrm{vis}}(\mathcal{F}) + \mathcal{S}_{\mathrm{hid}}(\mathcal{F})$.

- **Four bias terms.** $\mathcal{B}_{\mathrm{select}}, \mathcal{B}_{\mathrm{frame}}, \mathcal{B}_{\mathrm{center}}, \mathcal{B}_{\mathrm{inc}}$ (Paper II, Proposition 18, Table 2).

- **Delusion Trap.** $t_{\mathrm{del}} = \Lambda^{-1} \ln(\pi / 4\theta_0)$ (Paper II, Theorem 29).

- **Information-objects convention.** $I(\mathcal{K}_\mathcal{F}; \mathcal{K}) \equiv I(\hat{X}; X)$ on induced record processes (Paper II, Remark 15).

## 2.2 Fisher Information Metric

**Definition 1** (Fisher information matrix). *Let $\{p(x\,|\,\theta) : \theta \in \Theta \subset \mathbb{R}^d\}$ be a parametric family of probability densities satisfying standard regularity conditions (interchange of differentiation and integration). The* Fisher information matrix *is*

$$g_{ij}(\theta) := \mathbb{E}_\theta\left[\frac{\partial \log p(x\,|\,\theta)}{\partial \theta^i}\frac{\partial \log p(x\,|\,\theta)}{\partial \theta^j}\right] = -\mathbb{E}_\theta\left[\frac{\partial^2 \log p(x\,|\,\theta)}{\partial \theta^i\,\partial \theta^j}\right]. \tag{3}$$

*The pair $(\Theta, g)$ is a Riemannian manifold called the* statistical manifold.

**Remark 2** (Uniqueness). *By Čencov's theorem [7], the Fisher–Rao metric $g^{\mathrm{FR}}$ is, up to a positive scalar multiple, the unique Riemannian metric on the space of probability distributions that is invariant under all Markov morphisms (sufficient-statistic embeddings). This uniqueness guarantees that the Fisher metric is the* canonical *choice for measuring drift on the statistical manifold of the agent's predictive model—it is not a design choice but a mathematical necessity.*

**Proposition 3** (Cramér–Rao bound). *For any unbiased estimator $\hat{\theta}$ of $\theta$ based on $n$ independent observations:*

$$\mathrm{Cov}(\hat{\theta}) \succeq \frac{1}{n}\left[g(\theta)\right]^{-1} \tag{4}$$

*in the Löwner order. The scalar case reads $\mathrm{Var}(\hat{\theta}) \geq 1/\big(n\,g(\theta)\big)$.*

**Remark 4** (Effective independence). *Throughout this paper, references to "independent observations" in the context of continuous-time residual streams should be read as* effective independence *after thinning by the environmental decorrelation time $\tau_E$, yielding an effective sample size $n_{\mathrm{eff}} \approx T/\tau_E$. In particular, the sample count $n$ in (4) becomes $n_{\mathrm{eff}}$ in the self-referential setting of Section 4.2.*

**Remark 5** (Fisher metric and KL divergence). *The Fisher metric arises as the Hessian of the Kullback–Leibler divergence [16]:*

$$D_{\mathrm{KL}}\big(p_\theta \,\|\, p_{\theta+d\theta}\big) = \tfrac{1}{2}\,g_{ij}(\theta)\,d\theta^i\,d\theta^j + O\big(|d\theta|^3\big). \tag{5}$$

*This identifies the Fisher metric as the infinitesimal measure of statistical distinguishability.*

## 2.3 Information Geometry

Following Amari [4, 5], the statistical manifold $(\Theta, g)$ carries additional geometric structure beyond the Riemannian metric.

**$\alpha$-connections.** For each $\alpha \in [-1, 1]$, Amari defines an affine connection $\nabla^{(\alpha)}$ on $\Theta$. The cases $\alpha = 1$ (exponential connection, $\nabla^{(e)}$) and $\alpha = -1$ (mixture connection, $\nabla^{(m)}$) are *dual* with respect to $g$: $\partial_k\,g(X, Y) = g(\nabla_k^{(e)}X,\, Y) + g(X,\, \nabla_k^{(m)}Y)$. For exponential families, $\nabla^{(e)}$ is flat in natural parameters and $\nabla^{(m)}$ is flat in expectation parameters—the *dually flat structure*. The case $\alpha = 0$ recovers the Levi-Civita connection of the Fisher metric.

**Natural gradient.** Standard gradient descent in parameter space ignores the curvature of the statistical manifold. The *natural gradient* [6]

$$\dot{\theta} = -\eta\, g^{-1}(\theta)\, \nabla_\theta L(\theta), \tag{6}$$

where $\eta > 0$ is the learning rate and $L(\theta)$ is a loss function, provides the steepest descent direction in the Fisher metric. It is reparametrisation-invariant and Fisher-efficient (achieves the Cramér–Rao bound asymptotically).

**Pythagorean theorem.** In a dually flat space, the KL divergence satisfies a generalised Pythagorean relation: $D_{\mathrm{KL}}(p\,\|\,r) = D_{\mathrm{KL}}(p\,\|\,q) + D_{\mathrm{KL}}(q\,\|\,r)$ when $q$ is the $m$-projection of $p$ onto a submanifold containing $r$. This decomposition will be applied to separate the foreground-recoverable and background-irrecoverable components of drift.

## 2.4 Thermodynamic Length

**Definition 6** (Thermodynamic length). *Let $\lambda(t)$ for $t \in [0, \tau]$ be a path through control parameter space, and let $\zeta_{ij}(\lambda)$ be the* friction tensor *(the time-integrated equilibrium force–force correlation function at $\lambda$). The* thermodynamic length *of the path [8] is*

$$\mathcal{L} := \int_0^\tau \sqrt{\zeta_{ij}(\lambda)\, \dot{\lambda}^i\, \dot{\lambda}^j}\; dt. \tag{7}$$

**Proposition 7** (Sivak–Crooks bound). *The excess (dissipated) work during a finite-time transformation of duration $\tau$ satisfies [9]*

$$W_{\mathrm{ex}} \;\geq\; \frac{\mathcal{L}^2}{\tau}. \tag{8}$$

*The minimum is achieved by the geodesic of the friction tensor $\zeta$. In the linear-response regime, the friction tensor is related to the Fisher metric of the equilibrium distribution at $\lambda$ by $\zeta_{ij}(\lambda) \sim \tau_{\mathrm{relax}}\, g_{ij}^{\mathrm{Fisher}}(\lambda)$, where $\tau_{\mathrm{relax}}$ is the relaxation time.*

## 2.5 Second-Order Cybernetics

Von Foerster [14] distinguished two levels of control:

- **First-order cybernetics**: feedback control of observed systems. The controller adjusts its actions based on the output of a sensor. Paper II's ego is a first-order structure: it processes environmental data within a fixed frame.

- **Second-order cybernetics**: feedback control of the *observing* system itself. The controller adjusts the *sensor*—or equivalently, the reference frame within which the sensor operates. This is what Paper III provides.

Ashby's Law of Requisite Variety [13] provides a lower bound on the complexity of the meta-controller:

$$\mathrm{dim}(\text{meta-controller state space}) \;\geq\; \mathrm{dim}(\text{environmental drift subspace}). \tag{9}$$

7

The meta-observer must have at least as many adjustable parameters as there are independent modes of environmental drift.

In adaptive control theory [15], the analogous result is the *persistent excitation* condition: parameter estimates converge if and only if the input signal is "rich enough" to excite all modes of the system. In our framework, persistent excitation corresponds to $h_\mu > 0$—the environment must continue to generate novelty for the self-calibration loop to function.

**Remark 8** (Operational content). *The second-order cybernetic structure in this paper is* not *a philosophical metaphor. It has concrete operational content: the natural gradient update (6) is a specific algorithm that takes as input the Fisher information of the residual stream and produces as output an update to the frame parameter $\sigma$. This algorithm can be implemented by any physical system capable of accumulating second-moment statistics of its own prediction errors over a window of length $T \geq \tau_E$.*

## 2.6 Standing Assumptions

**Definition 9** (Standing Assumptions). *Throughout this paper, the following conditions are assumed:*

(C1) **Inherited framework.** *All assumptions (B1)–(B5) of Paper II [2] remain in force. This transitively includes (A1)–(A5) of Paper I [1] (open quantum system, thermal bath, well-defined free energy, finite Hilbert space, weak coupling) and the realizability embedding $\phi : Cl(V, q) \hookrightarrow Cl(1, 3)$ (Q-RAIF [22]). We invoke this embedding strictly as a structural inheritance from the earlier papers; no new physical claims about $Cl(1, 3)$ spacetime are introduced here. Additionally, the Delusion Trap is active: $\tau_{\mathrm{mem}} > \tau_{\mathrm{par}}$ and $\Lambda > 0$.*

(C2) **Environmental drift.** *The instantaneous optimal frame $\mathcal{F}^*(t)$ rotates continuously in $G$ at a rate characterised by the Lyapunov exponent $\Lambda > 0$ (Paper II, Eq. (37)).*

(C3) **Finite meta-observer budget.** *The self-calibration loop has a computational budget $\mathcal{C}_{\mathrm{meta}} < \infty$ (bits per unit time), distinct from the ego's processing budget $\mathcal{C}_{\mathrm{budget}}$.*

(C4) **Regularity.** *The agent's predictive family $\{p(e \,|\, \sigma) : \sigma \in G/H\}$ satisfies standard Fisher information regularity: full rank, finite Fisher matrix, and interchange of differentiation and integration. This extends (A5) from Paper I.*

(C5) **Persistent excitation.** *The environmental entropy rate satisfies $h_\mu > 0$ for all $t$. The environment generates new information indefinitely; no "frozen" regimes occur.*

# 3 The Drift Detection Problem

## 3.1 Why First-Order Control Fails

**Theorem 10** (First-Order Insufficiency). *Under assumptions (C1)–(C5), decompose the prediction residual as $e(t) = e_{\mathrm{drift}}(t) + \xi(t)$, where $e_{\mathrm{drift}}$ is the deterministic drift-induced*

*component (second-order in $\theta$) and $\xi(t)$ is the innovation noise, whose distribution is symmetric on $V_{\text{fg}}$ under (C5). No first-order controller—one that updates $\dot{\sigma} = f(e(t))$ based on the instantaneous residual without computing statistical properties of the error stream—can uniformly reduce the drift. Specifically: for any deterministic update function $f$, there exists a measurable event $\mathcal{E} \subset V_{\text{fg}}$ with $\mathbb{P}(\mathcal{E}) \geq 1/2 - O(\text{SNR})$ under the symmetric innovation distribution $p(\xi)$, such that for all $\xi \in \mathcal{E}$ the update direction satisfies $\langle \dot{\sigma}, \dot{\sigma}^* \rangle \leq 0$, where $\text{SNR} \sim \theta^4/h_\mu$.*

*Proof. Probability space.* The probability is taken over the innovation sequence $\{\xi(t)\}_{t \geq 0}$ under the symmetric distribution induced by the bath coupling (C5). All expectations below are over $p(\xi)$.

*Signal-to-noise separation.* The prediction error $e(t)$ lies in $V_{\text{fg}}$ by construction. Frame drift manifests as a rotation of the optimal frame $\mathcal{F}^*(t)$ in the gauge group $G$, shifting survival weight from $V_{\text{fg}}$ to $V_{\text{bg}}$. In $V_{\text{fg}}$, the drift signal enters only at second order in the mismatch angle $\theta$ (Paper II, proof of Theorem 29, part (c)): $\mathcal{S}_{\text{vis}} = \mathcal{S}_{\text{tot}} \cos^2 \theta$, so $e_{\text{drift}} \sim \theta^2 \mathcal{S}_{\text{tot}}$. The noise $\xi(t)$ scales as $h_\mu^{1/2}$. For $\theta \ll 1$, the single-sample signal-to-noise ratio is $\text{SNR} \sim \theta^4/h_\mu \ll 1$.

*Symmetry argument.* Since $p(\xi)$ is symmetric on $V_{\text{fg}}$, for any deterministic $f$:

- If $f$ is odd (e.g., linear gain), $\mathbb{E}[f(e_{\text{drift}} + \xi)] \approx f(e_{\text{drift}})$, but the instantaneous sign of $f$ is determined by $\xi$ with probability $\frac{1}{2} - O(\text{SNR})$.

- If $f$ is even, $f(e)$ carries no information about the *sign* of $\dot{\sigma}^*$, so $\langle f(e), \dot{\sigma}^* \rangle$ vanishes in expectation.

In either case, the probability that the update direction anti-correlates with the true drift direction is at least $1/2 - O(\text{SNR})$. Systematic drift detection requires accumulating second-order statistics of the residual stream over multiple samples—a second-order operation. $\square$

## 3.2   The Agent's Statistical Manifold

The agent's prediction-residual stream $\{e(t)\}_{t \geq 0}$ defines a stochastic process whose distribution depends on the gauge-fixing parameter $\sigma$. We model this dependence as a parametric family.

**Definition 11** (Predictive family)**.** *The* predictive family *of the agent is the set*

$$\mathcal{P} := \{p(e \mid \sigma) : \sigma \in \mathcal{M}_G\}, \tag{10}$$

*where $\mathcal{M}_G := G/H$ is the space of gauge-fixing orbits (H is the stabiliser of the foreground subspace), e denotes the prediction-residual time series over a window of length $T$, and $p(e \mid \sigma)$ is the likelihood of the observed residuals given the gauge parameter $\sigma$.*

The key insight is that $p(e \mid \sigma)$ depends on $\sigma$ even though $e(t) \in V_{\text{fg}}$, because the projection $\Pi_{\mathcal{F}}(\sigma)$ determines which environmental correlations are captured. When $\sigma$ drifts from the optimal $\sigma^*$:

- The variance of the residuals increases (the discarded background components contribute unmodelled noise).

- The temporal correlations of the residuals change (the projected kernel $\mathcal{K}_{\mathcal{F}}$ no longer captures the dominant environmental modes).

- Higher-order statistics (kurtosis, spectral shape) shift systematically.

These distributional changes are invisible to the raw error $e(t)$ but detectable by the Fisher metric of $\mathcal{P}$.

## 3.3  Self-Referential Fisher Information

**Definition 12** (Self-referential Fisher information)**.** *The* self-referential Fisher information *of the agent at gauge parameter $\sigma$ is*

$$\mathcal{I}_F(\sigma) := g_{ij}(\sigma)\, \delta\sigma^i\, \delta\sigma^j, \tag{11}$$

*where $g_{ij}(\sigma)$ is the Fisher information matrix of the predictive family $\mathcal{P}$ (Definition 11) evaluated at $\sigma$, and $\delta\sigma$ is the frame perturbation direction. In the scalar case (single drift mode), $\mathcal{I}_F(\sigma) = \mathbb{E}_\sigma\big[(\partial_\sigma \log p(e\,|\,\sigma))^2\big]$.*

**Remark 13** (What the agent "measures")**.** *Computing $\mathcal{I}_F(\sigma)$ does not require access to $V_{\mathrm{bg}}$ or to the "true" environment. It requires only: (i) the agent's own prediction residuals $\{e(t)\}$ (which lie in $V_{\mathrm{fg}}$), and (ii) the ability to evaluate the score function $\partial_\sigma \log p(e\,|\,\sigma)$—the sensitivity of its own predictive model to frame perturbations. This is a computation entirely within the agent's internal algebra, using only quantities already available from the ego's processing pipeline.*

**Theorem 14** (Drift Detectability)**.** *Under assumptions (C1)–(C5), suppose the frame is freshly calibrated at time $t_0$ ($\theta(t_0) = 0$). Then the self-referential Fisher information of the prediction-residual stream satisfies, for small accumulated drift ($\Lambda\,\Delta t \ll 1$):*

$$\mathcal{I}_F(\sigma;\, \{e_t\}_{t_0}^{t_0+\Delta t}) \;\geq\; \kappa\,\Lambda^2\,(\Delta t)^2\, \mathcal{I}_F^{\mathrm{env}}, \tag{12}$$

*where:*

- *$\kappa := \inf_{\sigma \in \mathcal{N}} (\partial\theta/\partial\sigma)^2 > 0$ is the* coupling efficiency*, where $\mathcal{N}$ is a compact neighbourhood of the calibrated point $\sigma^*$ on which the gauge chart is non-singular (existence guaranteed by (C4); see proof);*

- *$\Lambda$ is the environmental Lyapunov exponent (Paper II, Eq. (37));*

- *$\Delta t$ is the observation window;*

- *$\mathcal{I}_F^{\mathrm{env}} := \mathbb{E}_{p(\cdot|\theta)}[(\partial_\theta \log p)^2]$ is the per-component environmental Fisher information, measuring the baseline sensitivity of the decoherence functions to the mismatch angle $\theta$.*

*The self-referential Fisher information grows* quadratically *with accumulated drift time.*

*Proof. Step 1: chain rule.* The frame parameter $\sigma$ determines the mismatch angle $\theta = \theta(\sigma)$ via the gauge map $G/H \to [0, \pi/2]$. The chain rule for Fisher information gives

$$\mathcal{I}_F(\sigma) \;=\; \left(\frac{\partial\theta}{\partial\sigma}\right)^2 \mathcal{I}_F(\theta), \tag{13}$$

where $\mathcal{I}_F(\theta) := \mathbb{E}[(\partial_\theta \log p(e|\theta))^2]$ is the Fisher information of the residual stream with respect to the mismatch angle. By (C4) (full-rank Fisher matrix), the Jacobian $\partial\theta/\partial\sigma$ is bounded away from zero on any compact neighbourhood $\mathcal{N}$ of the calibrated point; we define the coupling efficiency $\kappa := \inf_{\sigma \in \mathcal{N}}(\partial\theta/\partial\sigma)^2 > 0$. This constant depends on the foreground dimension $k^*$, the Jacobian norms of the gauge-orbit map $G/H \to [0, \pi/2]$, and the regularity constants in (C4); it is computable for any concrete model (see Remark 15 for the qubit case).

*Step 2: small-drift expansion.* Under freshly calibrated initial conditions ($\theta(t_0) = 0$), the mismatch angle grows as $\theta(\Delta t) = \Lambda\,\Delta t + O((\Delta t)^2)$ (Paper II, Eq. (35), linearised about $\theta = 0$). The visible survival functional satisfies $\mathcal{S}_{\text{vis}} = \mathcal{S}_{\text{tot}} \cos^2\theta \approx \mathcal{S}_{\text{tot}}(1 - \theta^2)$ for $\theta \ll 1$. Thus the residual distribution $p(e\,|\,\theta)$ shifts from its baseline $p(e\,|\,0)$ by a score proportional to $\theta^2$: $\partial_\theta \log p \sim 2\theta \cdot (\partial_\theta \log p)|_{\theta=\theta^*}$, and consequently

$$\mathcal{I}_F(\theta) \;\geq\; (\Lambda\,\Delta t)^2\, \mathcal{I}_F^{\text{env}}, \tag{14}$$

where the inequality retains only the leading $O(\theta^2)$ term and drops $O(\theta^4)$ corrections.

*Step 3: assembly.* Substituting (14) into (13):

$$\mathcal{I}_F(\sigma) \;\geq\; \kappa\,\Lambda^2\,(\Delta t)^2\, \mathcal{I}_F^{\text{env}}. \qquad \square$$

**Remark 15** (Coupling efficiency in the qubit example). *For the single-qubit model of Section 7, the gauge orbit is parametrised directly by $\theta$ (rotation in the $xz$-plane of the Bloch sphere), so $\partial\theta/\partial\sigma = 1$ and $\kappa = 1$. More generally, $\kappa$ depends on the dimensionality of the gauge group and the curvature of the orbit $G/H$ at the current frame.*

**Corollary 16** (Detection before delusion). *Under (C1)–(C5), there exists a detection time $\Delta t_{\text{detect}}$ satisfying*

$$\Delta t_{\text{detect}} \;=\; \frac{1}{\Lambda}\sqrt{\frac{\mathcal{I}_F^{\min}}{\kappa\,\mathcal{I}_F^{\text{env}}}} \;<\; t_{\text{del}}, \tag{15}$$

*where $\mathcal{I}_F^{\min}$ is the minimum Fisher information required to exceed the noise floor (determined by $h_\mu$ and the observation window length). The detection window opens* before *the Delusion Trap closes, provided the meta-observer budget $\mathcal{C}_{\text{meta}}$ is sufficient to compute $\mathcal{I}_F$.*

*Proof.* The detection time $\Delta t_{\text{detect}} \propto \Lambda^{-1}$ (from (12)), while $t_{\text{del}} = \Lambda^{-1}\ln(\pi/4\theta_0)$ (from (1)). Since $\ln(\pi/4\theta_0) > 1$ for $\theta_0 < \pi/4e$ and the constant $\kappa\mathcal{I}_F^{\text{env}}$ is finite under (C4), the square root in $\Delta t_{\text{detect}}$ can be made smaller than the logarithm in $t_{\text{del}}$ for sufficiently sensitive meta-observers (large $\mathcal{I}_F^{\text{env}}$). $\qquad \square$

# 4 The Self-Referential Bound

## 4.1 The Bayesian Framework

The agent's ego structure (Paper II) provides a *prior belief* about the correct gauge parameter: the current frame $\sigma_0$ is the ego's "preferred" value. We encode this as a prior distribution $\pi_{\text{ego}}(\sigma)$, concentrated around $\sigma_0$.

**Definition 17** (Ego rigidity). *The* ego rigidity *is the prior Fisher information*

$$\mathcal{I}_{\text{ego}} := \int_{\mathcal{M}_G} \left( \frac{\partial \log \pi_{\text{ego}}(\sigma)}{\partial \sigma} \right)^2 \pi_{\text{ego}}(\sigma) \, d\sigma. \tag{16}$$

*High $\mathcal{I}_{\text{ego}}$ corresponds to a sharply peaked prior (rigid ego); low $\mathcal{I}_{\text{ego}}$ to a diffuse prior (flexible ego). The four bias terms of Paper II contribute to $\mathcal{I}_{\text{ego}}$: $\mathcal{B}_{\text{select}}$ and $\mathcal{B}_{\text{frame}}$ sharpen the prior around the current basis and connection, while $\mathcal{B}_{\text{center}}$ centres the prior on the agent's own state.*

## 4.2 The Self-Referential Cramér–Rao Bound

**Theorem 18** (Self-Referential Cramér–Rao Bound). *Under assumptions* (C1)–(C5), *let $\delta\hat{\sigma}$ be any estimator of the frame drift $\delta\sigma := \sigma^*(t) - \sigma$, based on a residual record of duration $T$. Define the* effective sample size $n_{\text{eff}} := T/\tau_E$, *where $\tau_E$ is the decorrelation time of the residual process $\{e(t)\}$ (the time beyond which consecutive residuals carry approximately independent information about $\sigma$). Then the van Trees inequality [10] gives*

$$\mathbb{E}\left[ \left| \delta\hat{\sigma} - \delta\sigma \right|^2 \right] \geq \frac{1}{n_{\text{eff}} \, \mathcal{I}_F(\sigma) + \mathcal{I}_{\text{ego}}}. \tag{17}$$

*Proof.* This is a direct application of the van Trees (Bayesian Cramér–Rao) inequality [10]. The total information about the drift parameter $\delta\sigma$ consists of two contributions:

- $n_{\text{eff}} \, \mathcal{I}_F(\sigma)$: the data Fisher information. In continuous time, the residual process is correlated with decorrelation time $\tau_E$ set by the bath memory kernel. Over a window of duration $T$, the process yields $n_{\text{eff}} \approx T/\tau_E$ effectively independent samples, each carrying $\mathcal{I}_F(\sigma)$ bits of information about $\delta\sigma$.

- $\mathcal{I}_{\text{ego}}$: the prior Fisher information from the ego's preference for $\sigma_0$ (Definition 17).

The van Trees inequality states that the Bayesian mean-squared error is bounded below by the inverse of the total information. $\square$

**Remark 19** (The ego as help and hindrance). *The ego rigidity $\mathcal{I}_{\text{ego}}$ acts as both help and hindrance:*

- ***Help**: when the ego is well-aligned ($\sigma_0 \approx \sigma^*$), the prior tightens the bound, reducing estimation variance.*

- **_Hindrance_**: _when the ego is misaligned ($|\sigma_0 - \sigma^*|$ large), the prior pulls the estimate toward the wrong value, creating a_ confirmation bias _that resists recalibration._

_The optimal Bayesian estimator balances data and prior:_

$$\hat{\sigma}_{\text{opt}} = \frac{n_{\text{eff}} \, \mathcal{I}_F \, \hat{\sigma}_{\text{MLE}} + \mathcal{I}_{\text{ego}} \, \sigma_0}{n_{\text{eff}} \, \mathcal{I}_F + \mathcal{I}_{\text{ego}}}, \tag{18}$$

_a weighted average of the maximum-likelihood estimate $\hat{\sigma}_{\text{MLE}}$ and the ego's prior belief $\sigma_0$, with weights proportional to their respective Fisher informations. As $n_{\text{eff}} \, \mathcal{I}_F \gg \mathcal{I}_{\text{ego}}$ (enough data to overwhelm the ego), the estimator converges to the MLE._

## 4.3 The Rigidity-Sensitivity Trade-off

**Proposition 20** (Optimal ego rigidity)**.** _Let the total expected loss be $\mathcal{L}_{\text{total}}(\mathcal{I}_{\text{ego}}) = \mathcal{L}_{\text{estimation}} + \lambda \, \mathcal{L}_{\text{calibration}}$, where $\mathcal{L}_{\text{estimation}}$ is the mean-squared drift-estimation error (bounded by (17)) and $\mathcal{L}_{\text{calibration}}$ is the cost of adjusting the frame (proportional to the frame rotation distance, hence larger when the ego is rigid and must be overcome). Under (C1)–(C5), there exists an optimal ego rigidity $\mathcal{I}_{\text{ego}}^*$ that minimises $\mathcal{L}_{\text{total}}$._

_Too rigid ($\mathcal{I}_{\text{ego}} \gg n_{\text{eff}} \, \mathcal{I}_F$): the ego overwhelms the data; the agent is blind to drift. Too soft ($\mathcal{I}_{\text{ego}} \ll n_{\text{eff}} \, \mathcal{I}_F$): the agent overreacts to noise; calibration cost is high. The optimum balances sensitivity against stability._

_Proof._ The estimation loss decreases with $\mathcal{I}_{\text{ego}}$ (the prior sharpens the bound (17) when $\sigma_0 \approx \sigma^*$ but increases it when misaligned). The calibration cost increases with $\mathcal{I}_{\text{ego}}$ (a rigid ego resists rotation). The sum is a convex function of $\mathcal{I}_{\text{ego}}$ under standard regularity, so a minimum exists. $\qquad \square$

# 5 The Calibration Loop

## 5.1 The Natural Gradient Update Law

The meta-observer updates the gauge parameter $\sigma$ following the natural gradient on the statistical manifold $(\mathcal{M}_G, g)$:

$$\dot{\sigma} = -\eta \, g^{-1}(\sigma) \, \nabla_\sigma L_{\text{frame}}(\sigma), \tag{19}$$

where $\eta > 0$ is the adaptation rate and the _frame loss_ is

$$L_{\text{frame}}(\sigma) := \mathbb{E}_\sigma[-\mathcal{S}_{\text{vis}}(\sigma)]. \tag{20}$$

The frame loss is minimised at the optimal gauge $\sigma^*$ that maximises visible survival. The Fisher metric enters through the inverse $g^{-1}$ in the natural gradient, not as a penalty term: it defines the _geometry_ of the update, ensuring reparametrisation invariance.

**Remark 21** (Reparametrisation invariance)**.** _The natural gradient (19) is invariant under reparametrisation of the gauge manifold $\mathcal{M}_G$: the update direction does not depend on the choice of coordinates for $\sigma$. This is essential because the gauge manifold inherits its geometry from the Clifford algebra, and no canonical coordinate system is preferred._

## 5.2 Lyapunov Stability of the Loop

**Drift velocity.** Let $\sigma^*(t)$ denote the instantaneous optimal gauge parameter (the minimiser of $L_{\text{frame}}$ at time $t$; Paper II, Definition 27). Define the *drift velocity* $\dot{\sigma}^* := d\sigma^*/dt$, measured with respect to the Fisher metric $g$; its norm $\|\dot{\sigma}^*\|_g := \sqrt{g_{ij}\,\dot{\sigma}^{*i}\,\dot{\sigma}^{*j}}$ is the instantaneous rate at which the environment's optimal frame rotates on the gauge manifold.

**Definition 22** (Lyapunov monitoring function). *The* Lyapunov monitoring function *is the squared geodesic distance on the statistical manifold from the current frame to the instantaneous optimal frame:*

$$V(\sigma) := d_{\text{geo}}(\sigma,\, \sigma^*(t))^2\,, \tag{21}$$

*where $d_{\text{geo}}$ is the geodesic distance in the Fisher metric $g$.*

**Theorem 23** (Loop Tracking Bound). *Under assumptions (C1)–(C5), the natural gradient update (19) applied to the Lyapunov monitoring function (21) satisfies*

$$\frac{dV}{dt} \;\leq\; -2\eta\,\alpha\,V \;+\; 2\sqrt{V}\,\|\dot{\sigma}^*\|_g, \tag{22}$$

*where $\alpha > 0$ is the persistent excitation constant (Definition 24 below) and $\|\dot{\sigma}^*\|_g$ is the instantaneous drift speed of the optimal frame. Consequently:*

(a) ***Tracking.*** *Whenever $\sqrt{V} > \|\dot{\sigma}^*\|_g/(\eta\,\alpha)$, we have $dV/dt < 0$: the loop actively reduces the mismatch.*

(b) ***Tracking neighbourhood.*** *The mismatch converges to a neighbourhood of the set of stationary points of $L_{\text{frame}}$. Assuming non-degeneracy (local strong convexity near $\sigma^*$, consistent with persistent excitation (C5) in standard adaptive-control settings [15]), this neighbourhood has size*

$$V_\infty \;:=\; \frac{\|\dot{\sigma}^*\|_g^2}{(\eta\,\alpha)^2}. \tag{23}$$

*For bounded drift $(\|\dot{\sigma}^*\|_g \leq \Lambda_{\max})$, the mismatch is bounded: $\limsup_{t\to\infty} V(t) \leq \Lambda_{\max}^2/(\eta\alpha)^2$.*

(c) ***Static limit.*** *When $\sigma^* = \text{const}$ $(\dot{\sigma}^* = 0)$, the bound reduces to $dV/dt \leq -2\eta\alpha\,V$, giving exponential convergence $V(t) \leq V(0)\,e^{-2\eta\alpha\,t}$.*

*Proof.* Since $V(\sigma) = d_{\text{geo}}(\sigma,\, \sigma^*(t))^2$ and $\sigma^*(t)$ is time-varying, the total derivative has two contributions:

$$\frac{dV}{dt} = \underbrace{\frac{\partial V}{\partial \sigma} \cdot \dot{\sigma}}_{\text{control}} + \underbrace{\frac{\partial V}{\partial \sigma^*} \cdot \dot{\sigma}^*}_{\text{drift}}.$$

**Control term.** In normal coordinates centred at $\sigma^*$, let $\delta\sigma := \sigma - \sigma^*$. The control contribution is $2\,g(\delta\sigma,\, \dot{\sigma}) = 2\,g(\delta\sigma,\, -\eta\,g^{-1}\nabla L_{\text{frame}}) = -2\eta\,\langle\delta\sigma,\, \nabla L_{\text{frame}}\rangle$. Since $\sigma^*$ minimises $L_{\text{frame}}$ by definition, $L_{\text{frame}}$ is locally strongly convex near $\sigma^*$ under persistent excitation (C5) (the Hessian of $L_{\text{frame}}$ at $\sigma^*$ is bounded below by $\alpha\,g$, where $\alpha$ is the persistent excitation constant). Therefore $\langle\delta\sigma,\, \nabla L_{\text{frame}}\rangle \geq \alpha\,|\delta\sigma|_g^2 = \alpha\,V$, giving a control contribution $\leq -2\eta\,\alpha\,V$.

14

**Drift term.** The drift contribution is $-2\,g(\delta\sigma,\dot{\sigma}^*)$. By Cauchy–Schwarz: $|g(\delta\sigma,\dot{\sigma}^*)| \leq |\delta\sigma|_g\,\|\dot{\sigma}^*\|_g = \sqrt{V}\,\|\dot{\sigma}^*\|_g$. Hence the drift contribution is bounded by $+2\sqrt{V}\,\|\dot{\sigma}^*\|_g$.

**Combined.** Adding both contributions gives (22). Part (a) follows by setting $dV/dt < 0$; part (b) by solving $dV/dt = 0$ for the fixed point $\sqrt{V_\infty} = \|\dot{\sigma}^*\|_g/(\eta\alpha)$; part (c) by setting $\dot{\sigma}^* = 0$. $\qquad\square$

## 5.3 Convergence Rate under Persistent Excitation

**Definition 24** (Persistent excitation constant)**.** *The persistent excitation constant $\alpha > 0$ is the minimum eigenvalue of the time-averaged Fisher information matrix:*

$$\bar{g}(t) := \frac{1}{T}\int_t^{t+T} g(\sigma(s))\,ds \; \succeq \; \alpha\,I \qquad \textit{for all } t, \tag{24}$$

*guaranteed to exist by* (C5) *and* (C4)*.*

**Remark 25** (Tracking vs convergence)**.** *In the static case ($\sigma^* = $ const), Theorem 23(c) gives pure exponential convergence: $V(t) \leq V(0)\,e^{-2\eta\alpha\,t}$. Under environmental drift, convergence to zero is* not *possible—instead the loop maintains the mismatch within the tracking neighbourhood (23). The tracking error $V_\infty$ grows with drift speed $\|\dot{\sigma}^*\|_g$ and decreases with loop parameters $\eta$ and $\alpha$. If the free-energy budget is insufficient to maintain $\eta\,\alpha > \Lambda$ (the drift rate), the tracking neighbourhood expands and the Delusion Trap re-emerges. This connects the Lyapunov stability of the loop directly to the thermodynamic budget (Section 6).*

**Remark 26** (The necessity of novelty)**.** *If $h_\mu \to 0$ (the environment ceases to generate new information), the persistent excitation constant $\alpha \to 0$ and the tracking neighbourhood $V_\infty = \|\dot{\sigma}^*\|_g^2/(\eta\alpha)^2 \to \infty$: the loop loses all ability to track. Memory without novelty cannot sustain self-reference. This is the information-theoretic expression of a basic physical principle: a system in thermodynamic equilibrium cannot "learn" about itself.*

## 5.4 The Four-Part Structure Proposition

We are now in a position to state the capstone result of the T-DOME sequence.

**Proposition 27** (Sufficient Architecture for Persistent Agents)**.** *Within the class of agents satisfying* (C1)–(C5)*, a sufficient architecture for maintaining a non-equilibrium steady state* (NESS) *in an open, drifting environment under bounded computation comprises the following four structural layers:*

(I) ***External observable geometry.*** *The environmental observable algebra supports a metric structure; $Cl(1,3)$ serves as the running example throughout the programme, but the argument applies to any algebra satisfying (C1). Assumption: established in [17, 18, 20]; adopted here as a modelling premise.*

(II) ***Internal control algebra.*** *The agent carries an internal algebra isomorphic to $Cl(V,q)$ with realizability embedding $\phi: Cl(V,q) \hookrightarrow Cl(1,3)$. Assumption: established in [21, 22]; adopted here as a modelling premise.*

(III) **Self-monitoring function.** *The agent maintains a Lyapunov function $V(\sigma)$ (21) satisfying the tracking bound (22), via a second-order control loop on the agent's Fisher information, keeping the mismatch within the tracking neighbourhood (23).* Source: *this paper, Theorem 23.*

(IV) **Biased, non-Markovian memory.** *The agent carries path-dependent state (non-Markovian memory kernel $\mathcal{K}(t,s)$) compressed through a gauge-fixed reference frame (the ego $\mathfrak{E}$).* Source: *Paper I [1] (memory necessity) and Paper II [2] (ego necessity).*

*Without any one of the four layers, the agent fails:*

- *Without (I): no physical embedding—the agent cannot interact with the Lorentzian environment.*

- *Without (II): no channel discrimination—the agent cannot distinguish survival-relevant from irrelevant information.*

- *Without (III): the Delusion Trap—the ego rigidifies and prediction error diverges exponentially.*

- *Without (IV): the Markovian Ceiling and computational paralysis—no temporal accumulation, no tractable processing.*

*Proof.* Layers (I) and (II) are modelling assumptions adopted from [17, 18, 20, 21, 22]; their sufficiency within those frameworks is established therein. The sufficiency of (III) follows from the present paper: Theorem 10 shows that first-order control is insufficient to escape the Delusion Trap, and Theorem 23 shows that the tracking bound is sufficient. The sufficiency of (IV) follows from Paper I [1] (Markovian Ceiling $\mathcal{S} \leq 0$) and Paper II [2] (Computational Ceiling and necessity of SSB).

The "without" claims follow from the respective crisis theorems: Paper I's Theorem 14 (Markovian Ceiling), Paper II's Theorem 7 (Computational Ceiling) and Theorem 29 (Delusion Trap), and the present Theorem 10. $\qquad\square$

# 6 Thermodynamic Cost

## 6.1 The Three Cost Components

The self-referential calibration loop requires three distinct operations, each carrying an irreducible thermodynamic cost:

**1. Sensing cost.** The meta-observer must read the prediction residuals from the ego's processing pipeline. This requires monitoring $k^*$ foreground channels, each producing $h_\mu$ bits per unit time:

$$\dot{W}_{\text{sense}} \geq k_B T \ln 2 \cdot h_\mu k^*. \tag{25}$$

(Landauer cost of reading $h_\mu k^*$ bits per unit time.)

**2. Computing cost.** Evaluating the Fisher information $\mathcal{I}_F(\sigma)$ from the residual stream requires the meta-observer to process $\mathcal{C}_{\text{meta}}$ bits per unit time:

$$\dot{W}_{\text{compute}} \; \geq \; k_B T \ln 2 \cdot \mathcal{C}_{\text{meta}}. \tag{26}$$

**3. Actuating cost.** Rotating the gauge parameter from the current frame $\sigma$ to the estimated optimal frame $\hat{\sigma}^*$ is a finite-time thermodynamic transformation on the gauge manifold. By the Sivak–Crooks bound (Proposition 7):

$$\dot{W}_{\text{actuate}} \; \geq \; \frac{\mathcal{L}^2(\sigma, \hat{\sigma}^*)}{\tau_{\text{recalib}}^2}, \tag{27}$$

where $\mathcal{L}(\sigma, \hat{\sigma}^*)$ is the thermodynamic length (7) of the geodesic from $\sigma$ to $\hat{\sigma}^*$, and $\tau_{\text{recalib}}$ is the recalibration time.

## 6.2 The Thermodynamic Cost Theorem

**Theorem 28** (Thermodynamic Cost of Self-Referential Calibration). *Under assumptions* (C1)–(C5), *the minimum dissipation rate of the self-referential calibration loop satisfies*

$$\dot{W}_{\text{loop}} \; \geq \; k_B T \ln 2 \left[ h_\mu \, k^* + \mathcal{C}_{\text{meta}} \right] + \frac{\mathcal{L}^2(\sigma, \sigma^*)}{\tau_{\text{recalib}}^2}. \tag{28}$$

*The first bracketed term is the* information tax *(the Landauer cost of sensing and computing). The second term is the* geometric tax *(the Sivak–Crooks cost of actuating the frame rotation).*

*Proof.* We must establish that the three lower bounds can be summed, i.e. that no single physical process can simultaneously satisfy two or more of them.

The three operations act on *disjoint physical degrees of freedom*:

1. *Sensing* reads the prediction residuals $\{e_t\}$ from the ego's foreground channels. The relevant degrees of freedom are the sensor registers that copy bits from the foreground subspace $V_{\text{fg}}$. Each bit erased carries the Landauer cost $k_B T \ln 2$.

2. *Computing* evaluates the Fisher information $\mathcal{I}_F(\sigma)$ from the copied residuals. The relevant degrees of freedom are the processor logic states of the meta-observer. These are distinct from the sensor registers: the processor manipulates the data *after* it has been read, and its own state transitions carry an independent Landauer cost.

3. *Actuating* rotates the gauge parameter from $\sigma$ to $\hat{\sigma}^*$. The relevant degrees of freedom are the control fields that implement the frame rotation on the agent's internal algebra $Cl(V, q)$. This is a physical transformation of the agent's hardware state, governed by the Sivak–Crooks bound on finite-time thermodynamic transformations. The $\tau_{\text{recalib}}^{-2}$ scaling of the dissipation *rate* follows from the Sivak–Crooks bound $W_{\text{ex}} \geq \mathcal{L}^2/\tau$ (excess *work*), divided by $\tau_{\text{recalib}}$ to convert to a rate.

Under the assumption that the three operations are physically realised on separable degrees of freedom (no shared erasure accounting), the sets are disjoint (sensor ∩ processor = ∅, processor ∩ actuator = ∅, sensor ∩ actuator = ∅), and the Landauer bound for each is independent. Moreover, the actuating cost involves a different *type* of bound (thermodynamic length, not Landauer erasure), reinforcing the independence. The total lower bound is therefore the sum of the three individual bounds (25)–(27). □

## 6.3 The Complete Persistence Budget

**Corollary 29** (Persistence Budget). *Combining the results of Papers I, II, and III, the minimum free-energy dissipation rate for a persistent, self-calibrating agent in a drifting environment is*

$$\dot{W}_{\text{total}} \geq \underbrace{k_B T \ln 2 \cdot h_\mu}_{\text{Paper I: memory}} + \underbrace{k_B T \ln 2 \cdot h_\mu \, k^*}_{\text{Paper II: ego processing}} + \underbrace{k_B T \ln 2 \, [h_\mu \, k^* + \mathcal{C}_{\text{meta}}] + \frac{\mathcal{L}^2}{\tau_{\text{recalib}}^2}}_{\text{Paper III: self-calibration loop}}. \qquad (29)$$

*Below this budget, the agent must sacrifice one or more of the four structural layers (Proposition 27): losing memory (Paper I crisis), losing the ego (Paper II crisis), or losing self-calibration (Paper III crisis, the Delusion Trap).*

**Remark 30** (The cost of selfhood). *Equation (29) is the first explicit, calculable lower bound on the thermodynamic cost of maintaining a self-referential agent in a drifting environment. It shows that "selfhood" is not free: the ego (Paper II) and its calibration loop (Paper III) each add irreducible energy taxes on top of the memory cost (Paper I). The total cost grows with the environmental complexity ($h_\mu$), the agent's representational capacity ($k^*$), the meta-observer's computational power ($\mathcal{C}_{\text{meta}}$), and the drift rate (through $\mathcal{L}$ and $\tau_{\text{recalib}}$).*

# 7 Worked Example: Qubit in a Drifting Two-Channel Bath

## 7.1 Model Setup

We extend the two-channel qubit model from Paper II (Section 6) by introducing environmental drift.

**Inherited setup.** A qubit ($\dim \mathcal{H}_S = 2$) with internal algebra $Cl(0,2) \cong \mathbb{H}$ ($D = 4$), coupled to two bosonic channels:

- Dephasing channel ($\sigma_z$): $J_z(\omega) = 2\lambda_z \gamma_z \omega / (\omega^2 + \gamma_z^2)$.

- Dissipative channel ($\sigma_x$): $J_x(\omega) = 2\lambda_x \gamma_x \omega / (\omega^2 + \gamma_x^2)$.

Paper II's ego selects $V_{\text{fg}} = \text{span}\{1, \mathbf{k}\}$ (the dephasing subspace), discarding $V_{\text{bg}} = \text{span}\{\mathbf{i}, \mathbf{j}\}$.

**Environmental drift.** We now allow the dephasing coupling to drift exponentially (matching Paper II's Delusion Trap analysis):

$$\lambda_z(t) = \lambda_z^{(0)}\left(1 + \theta_0\, e^{\Lambda t}\right), \qquad \theta_0 = 0.02, \quad \Lambda = 0.08\,\omega_0. \tag{30}$$

The optimal frame $\mathcal{F}^*(t)$ rotates in $SO(3)$ as the relative survival values of the two channels change. The Delusion Trap time $t_{\text{del}} = \Lambda^{-1}\ln\left(\pi/(4\theta_0)\right) \approx 45.9\,\omega_0^{-1}$.

**Parameter mapping.**

| Quantity | Value | Source |
|---|---|---|
| $D = \dim Cl(0,2)$ | 4 | Paper II |
| $k^*$ | 2 | Paper II, Theorem 17 |
| $\mathcal{C}_{\text{budget}}$ | $2\,h_\mu$ | Paper II |
| $\theta_0$ (initial misalignment) | 0.02 | this example |
| $\Lambda$ (drift rate) | $0.08\,\omega_0$ | Eq. (30) |
| $t_{\text{del}}$ | $45.9\,\omega_0^{-1}$ | Paper II, Delusion Trap |
| $\eta$ (adaptation rate) | 0.5 | meta-observer |

## 7.2 Fisher Information under Drift

As the coupling $\lambda_z(t)$ drifts, the decoherence function $p_z(t)$ (Paper II, Eq. (34)) changes, shifting the residual distribution. The self-referential Fisher information $\mathcal{I}_F(\sigma)$ measures this shift.

For the qubit model, the Fisher information with respect to the frame angle $\phi$ (parametrising the $SO(3)$ rotation between the current and optimal frames) is

$$\mathcal{I}_F(\phi) = \frac{(\partial_\phi \bar{e})^2}{\text{Var}(e)} \approx \frac{4\,\mathcal{S}_{\text{tot}}^2\,\theta^2}{h_\mu/n_{\text{eff}}}, \tag{31}$$

where $\bar{e} = \mathbb{E}[e\,|\,\phi]$ is the expected residual, $\theta = \theta(\phi)$ is the mismatch angle, and $n_{\text{eff}}$ is the effective sample size (Remark 4).

When the frame is well-aligned ($\theta \approx 0$): $\mathcal{I}_F \approx 0$. As drift accumulates ($\theta$ grows): $\mathcal{I}_F$ increases quadratically, producing a detectable "stress signal" consistent with Theorem 14.

## 7.3 Loop Dynamics: Self-Calibration in Action

Under the natural gradient update (19), the frame angle $\phi(t)$ tracks the drifting optimal frame $\phi^*(t)$. The Lyapunov function $V(t) = (\phi(t) - \phi^*(t))^2$ is governed by the tracking bound (22): the loop drives $V$ toward the tracking neighbourhood $V_\infty = \|\dot{\sigma}^*\|_g^2/(\eta\alpha)^2$, with the approach rate set by the persistent excitation constant $\alpha$ and the adaptation rate $\eta$.

**Comparison.**

- **Without loop** (Paper II agent): the mismatch grows as $\theta(t) = \theta_0\, e^{\Lambda t}$, reaching $\pi/4$ at $t_{\text{del}}$. The agent is delusional.

- **With loop** (Paper III agent): the mismatch oscillates around zero, bounded by the estimation noise floor $\theta_{\min} \sim 1/\sqrt{n_{\text{eff}}\, \mathcal{I}_F^{\text{env}}}$ (the Cramér–Rao limit). The agent remains calibrated.

A multi-dimensional numerical evaluation extending this qubit illustration to continuous drift is presented in Section 8.

## 7.4 Thermodynamic Cost Evaluation

For the qubit example with $k^* = 2$, $h_\mu = 1$ (normalised), $\mathcal{C}_{\text{meta}} = 1\, h_\mu$ (minimal meta-observer):

$$\dot{W}_{\text{sense}} \geq k_B T \ln 2 \cdot 1 \cdot 2 = 2\, k_B T \ln 2, \tag{32}$$

$$\dot{W}_{\text{compute}} \geq k_B T \ln 2 \cdot 1 = k_B T \ln 2, \tag{33}$$

$$\dot{W}_{\text{actuate}} \geq \frac{\mathcal{L}^2}{\tau_{\text{recalib}}^2} \approx \frac{\theta_0^2\, \tau_{\text{relax}}}{\tau_{\text{recalib}}^2}\, k_B T. \tag{34}$$

The total loop cost is dominated by the information tax (sensing + computing) at $\sim 3\, k_B T \ln 2$ per unit time, with the geometric tax (actuating) contributing a smaller correction proportional to $\theta_0^2$.

For comparison, Paper I's memory cost is $\dot{W}_{\text{mem}} \geq k_B T \ln 2$ and Paper II's ego processing cost is $\dot{W}_{\text{ego}} \sim 2\, k_B T \ln 2 \cdot h_\mu$. The self-calibration loop adds approximately 50% to the total energy budget—a significant but bounded cost for escaping the Delusion Trap.

# 8 Numerical Demonstration

The preceding sections establish analytic bounds and a low-dimensional worked example. We now demonstrate computationally that the three core phenomena—delusion separation, detectable staleness, and an optimal calibration budget—emerge in a minimal multi-dimensional system under continuous drift. Full code and parameters are provided for reproducibility.

## 8.1 Model

**Environment.** A $d$-dimensional linear prediction task: $y(t) = \mathbf{w}(t)^\top \mathbf{x}(t) + \sigma\, \epsilon(t)$, $\mathbf{x}(t) \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_d)$, $\epsilon \sim \mathcal{N}(0,1)$. The weight vector $\mathbf{w}(t) \in \mathbb{S}^{d-1}$ drifts by receiving random perturbations on *background* dimensions only (indices $k, \ldots, d-1$), then renormalising. Signal therefore migrates progressively from the ego's foreground to its blind sector.

**Agents.**

- **Fixed ego** (Paper II analogue): learns a linear model on a *fixed* foreground subspace of dimension $k$ via stochastic gradient descent (SGD, rate $\eta$, decay $\lambda$). Embodies the "frozen gauge" of Paper II.

- **Calibrated loop** (Paper III analogue): identical ego plus a staleness sentinel and recalibration mechanism. The sentinel tracks $g_i = \text{EMA}(|e\,x_i|)$ for each dimension $i$ (an absolute-gradient proxy), computes the fraction of top-$k$ gradient dimensions *not* in the current foreground as a frame-staleness index $m \in [0,1]$, and triggers recalibration when $\text{EMA}(m) > \theta$. After recalibration, a settling period of $\tau$ steps elapses before the sentinel resumes monitoring.

**Parameters.**

| Quantity | Value | Role |
| --- | --- | --- |
| $d$ | 20 | full ambient dimension |
| $k$ | 5 | ego foreground dimension ($k/d = 0.25$) |
| $\sigma$ | 0.1 | observation noise std |
| $\eta$ | 0.01 | SGD learning rate |
| $\lambda$ | 0.998 | SGD weight decay |
| $\theta$ | 0.25 | staleness threshold |
| $\Lambda$ | variable | drift rate per step |
| $\tau$ | variable | settling period (cooldown) |

**Oracle metrics.** Neither agent has access to $\mathbf{w}(t)$. We evaluate performance externally using the *oracle full-space error*:

$$\mathcal{E}_{\text{full}} = \left\| \mathbf{w}_{\text{ego}} - \mathbf{w}_{\text{fg}}^* \right\|^2 + \left\| \mathbf{w}_{\text{bg}}^* \right\|^2, \tag{35}$$

where $\mathbf{w}_{\text{fg}}^*$ and $\mathbf{w}_{\text{bg}}^*$ denote the true weight vector restricted to foreground and background coordinates respectively, and $\mathbf{w}_{\text{ego}}$ is the ego's foreground-supported estimator lifted to the full space. The first term captures foreground tracking error (accessible to the ego); the second captures hidden-sector signal (invisible).

## 8.2 Results

**Result 1: Delusion-correction separation (Figure 1).** At drift rate $\Lambda = 0.02$, settling period $\tau = 200$, and $T = 5\,000$ steps, three phenomena are visible:

(a) *Delusion trap.* The ego's foreground tracking error converges to near zero, while the true full-space error rises toward $\sim 1$ and stabilises. The growing gap between the two is the hidden sector, confirming the prediction of Theorem 10: first-order monitoring cannot detect frame drift.

(b) *Detectability.* The staleness sentinel produces a clean sawtooth: rising from zero after each recalibration, crossing the threshold $\theta = 0.25$, and triggering frame reset (25 events over $T = 5\,000$). This is consistent with the predicted growth trend of the self-referential Fisher signal (Theorem 14).

(c) *Net benefit.* The calibrated loop achieves $\mathcal{E}_{\text{full}} \approx 0.74$ versus the fixed ego's $\approx 1.02$: a 27% reduction in true prediction error.

**Result 2: Phase structure and optimal calibration budget (Figures 2–3).** We scan 16 drift rates $\Lambda \in [0.005, 0.08]$ and 16 settling periods $\tau \in [15, 800]$ (logarithmically spaced), running both agents for $T = 4\,000$ steps across 6 random seeds per grid point.

Figure 2(a) shows the performance gain $\Delta = \mathcal{E}_{\mathrm{ego}} - \mathcal{E}_{\mathrm{loop}}$: the loop improves over the ego (green) across most of the parameter space, with a boundary at $\Delta = 0$ (dashed) below which recalibration is counterproductive (very low drift, where the overhead of re-learning exceeds the benefit of tracking). The solid curve traces the *optimal settling period $\tau_{\mathrm{opt}}(\Lambda)$*— the recalibration period minimising $\mathcal{E}_{\mathrm{loop}}$—which decreases monotonically from $\sim 370$ steps at $\Lambda = 0.005$ to $\sim 100$ at $\Lambda = 0.08$. Figure 2(b) shows that calibration frequency increases smoothly with drift and with shorter settling period, exhibiting the cost–performance trade-off of Theorem 28.

Extracting $\tau_{\mathrm{opt}}(\Lambda)$ yields the *optimal calibration frequency $\alpha_{\mathrm{opt}}(\Lambda) = 1/\tau_{\mathrm{opt}}$* (Figure 3). The curve is smooth and monotonically increasing: faster drift demands tighter calibration. It saturates at high $\Lambda$ near $\alpha_{\mathrm{opt}} \approx 0.01$ per step ($\tau_{\mathrm{opt}} \approx 100$), of the same order as the learner's settling time. This is consistent with the intuition that drift estimation requires a minimum observation window; the self-referential Cramér–Rao bound (Theorem 18) provides the analytic counterpart of this computational floor.

## 8.3 Scope of This Demonstration

This demonstration **does** show:

1. The delusion-correction separation predicted by Theorems 10 and 14 emerges in a minimal stochastic system with continuous drift.

2. A frame-staleness signal with clean threshold dynamics exists and triggers effective recalibration.

3. An optimal calibration frequency $\alpha_{\mathrm{opt}}(\Lambda)$ exists, increases monotonically with drift rate, and saturates at the learner's settling timescale.

4. The cost–performance trade-off of Theorem 28 manifests as a structured phase diagram with an explicit $\tau_{\mathrm{opt}}$ boundary.

In summary, this demonstration validates the *existence* and *detectability* of the loop–cost trade-off in a minimal linear setting; it does not claim universality across architectures or environment classes.

This demonstration does **not** show:

1. That the specific functional form of $\alpha_{\mathrm{opt}}(\Lambda)$ matches the analytic Cramér–Rao prediction in the large-$d$ limit. The demonstration confirms the monotonic trend and saturation; deriving the exact scaling exponent from Theorem 18 remains open.

2. That the results generalise to all environment classes. The model uses Gaussian features, linear regression, and isotropic background drift; extensions to non-linear, non-Gaussian, or structured-drift settings require further investigation.

3. That the calibration loop is optimal among all possible adaptive strategies. It implements one specific realisation of the calibration-loop architecture.

**Reproducibility.** The complete simulation is a self-contained Python script (`tdome_demo.py`, $\sim 550$ lines, requiring only NumPy and Matplotlib) with fixed random seeds. All figures in this section can be reproduced by executing the script after setting the output directory variable `BASE` to the desired path.

# 9 Discussion

## 9.1 Summary of Results

| Result | Statement | Sec. |
|---|---|---|
| First-Order Insufficiency | Raw prediction error cannot detect frame drift | 3.1 |
| Drift Detectability | Self-referential Fisher information grows quadratically with accumulated drift | 3.3 |
| Self-Referential CR Bound | Drift estimation bounded by $1/(n_{\text{eff}}\,\mathcal{I}_F + \mathcal{I}_{\text{ego}})$ | 4.2 |
| Loop Tracking Bound | Lyapunov $V$ with tracking neighbourhood $V_\infty = \|\dot\sigma^*\|^2/(\eta\alpha)^2$ | 5.2 |
| Four-Part Structure | Persistent agents require four structural layers | 5.4 |
| Loop Cost | $\dot W_{\text{loop}} \geq k_B T \ln 2\,[h_\mu\,k^* + \mathcal{C}_{\text{meta}}] + \mathcal{L}^2/\tau_{\text{recalib}}^2$ | 6.2 |
| Persistence Budget | Total cost: memory + ego + loop | 6.3 |
| Numerical Demonstration | Delusion separation, sentinel detection, $\alpha_{\text{opt}}(\Lambda)$ boundary | 8 |

## 9.2 The Complete Logic Chain

Papers I–III trace an irreversible thermodynamic logic chain:

| Paper | Crisis | Resolution | What is born |
|---|---|---|---|
| Paper I | Markovian trap: no history | Non-Markovian memory | **Temporal accumulation** |
| Paper II | Computation explosion: $\infty$ memory, finite budget | Gauge SSB: $Cl(V, q) \to V_{\text{fg}} \oplus V_{\text{bg}}$ | **Compressed ref. frame** |
| Paper III | Delusion trap: fixed bias, drifting world | Fisher self-referential calibration; tracking bound | **Reflexivity** |

Each resolution creates the precondition for the next crisis. The chain terminates at Paper III: the self-referential calibration loop does not create a further crisis requiring a "Paper IV," because the loop is *self-correcting* by construction (Theorem 23). Its only

vulnerability is the thermodynamic budget (Theorem 28): if the agent's free-energy supply falls below the persistence budget (29), the loop degrades and the Delusion Trap re-emerges. This is not a new crisis but the Second Law itself: all order requires free-energy dissipation.

## 9.3  What This Paper Does and Does Not Show

This paper **does** show:

1. Under environmental drift (C2) and bounded computation (C1), first-order control fails to detect frame drift (Theorem 10).

2. Self-referential Fisher information provides a quadratically growing signal sufficient for drift detection before the Delusion Trap (Theorem 14).

3. Drift estimation precision is bounded by the Self-Referential Cramér–Rao bound (Theorem 18).

4. The calibration loop tracks the optimal frame within a bounded neighbourhood under a Lyapunov tracking bound (Theorem 23).

5. The thermodynamic cost of the loop is calculable (Theorem 28).

This paper does **not** show:

1. That self-referential calibration implies or requires phenomenal consciousness, subjective experience, or qualia. "Reflexivity" as used here denotes second-order control, nothing more.

2. That the Lyapunov function $V$ is a measure of "awareness." It is a control-theoretic stability condition, not a consciousness metric.

3. That the Four-Part Structure Proposition is a complete characterisation of agency. It states sufficient conditions under (C1)–(C5); other architectures may also suffice.

4. That Fisher information requires the agent to "know" it is computing Fisher information. The computation can be implemented implicitly by any physical system whose dynamics approximate the natural gradient.

5. That the calibration loop eliminates the ego's bias. It tracks and compensates for drift in the bias; the four bias terms of Paper II persist.

6. That the thermodynamic cost bounds are achievable by any specific physical implementation. They are information-theoretic lower bounds.

7. That this framework applies to all possible systems. It applies to systems satisfying (C1)–(C5).

8. That the structural parallel with philosophical concepts of self-awareness constitutes a philosophical or metaphysical claim.

9. That the Clifford algebra is the only possible algebraic setting. Other control algebras may yield analogous results with different quantitative bounds.

We have established a budgeted self-referential calibration loop that detects drift via an intrinsic Fisher signal, yields a falsifiable stability criterion, and incurs an unavoidable thermodynamic cost. In the context of Papers I–III, this completes the programme's third step by turning bias (Paper II) into a dynamically monitored and correctable quantity.

# References

[1] S. Liu, *Non-Markovian Memory and the Thermodynamic Necessity of Temporal Accumulation*, Zenodo (2026), DOI: 10.5281/zenodo.18574342.

[2] S. Liu, *Spontaneous Symmetry Breaking of Reference Frames as a Computational Cost Minimization Strategy*, Zenodo (2026), DOI: 10.5281/zenodo.18579703.

[3] C. R. Rao, *Information and the accuracy attainable in the estimation of statistical parameters*, Bull. Calcutta Math. Soc. **37**, 81 (1945).

[4] S.-i. Amari, *Differential-Geometrical Methods in Statistics*, Lecture Notes in Statistics **28**, Springer (1985).

[5] S.-i. Amari and H. Nagaoka, *Methods of Information Geometry*, Translations of Mathematical Monographs **191**, AMS (2000).

[6] S.-i. Amari, *Natural gradient works efficiently in learning*, Neural Computation **10**, 251 (1998).

[7] N. N. Čencov, *Statistical Decision Rules and Optimal Inference*, Translations of Mathematical Monographs **53**, AMS (1982).

[8] G. E. Crooks, *Measuring thermodynamic length*, Phys. Rev. Lett. **99**, 100602 (2007).

[9] D. A. Sivak and G. E. Crooks, *Thermodynamic metrics and optimal paths*, Phys. Rev. Lett. **108**, 190602 (2012).

[10] H. L. van Trees, *Detection, Estimation, and Modulation Theory*, Part I, Wiley (1968).

[11] A. C. Barato and U. Seifert, *Thermodynamic uncertainty relation for biomolecular processes*, Phys. Rev. Lett. **114**, 158101 (2015).

[12] S. Ito, *Stochastic thermodynamic interpretation of information geometry*, Phys. Rev. Lett. **121**, 030605 (2018).

[13] W. R. Ashby, *An Introduction to Cybernetics*, Chapman & Hall (1956).

[14] H. von Foerster, *Understanding Understanding: Essays on Cybernetics and Cognition*, Springer (2003).

[15] K. J. Åström and B. Wittenmark, *Adaptive Control*, 2nd ed., Addison-Wesley (1995).

[16] T. M. Cover and J. A. Thomas, *Elements of Information Theory*, 2nd ed., Wiley (2006).

[17] S. Liu, *Emergent Geometry from Coarse-Grained Observable Algebras*, Zenodo (2026), DOI: 10.5281/zenodo.18361707.

[18] S. Liu, *Accessibility, Stability, and Emergent Geometry*, Zenodo (2026), DOI: 10.5281/zenodo.18367061.

[19] S. Liu, *Structural Limits of Unification: Accessibility, Incompleteness, and the Necessity of a Final Cut*, Zenodo (2026), DOI: 10.5281/zenodo.18402908.

[20] S. Liu, *Algebraic Constraints on the Emergence of Lorentzian Metrics in Entropic Gravity Frameworks*, Zenodo (2026), DOI: 10.5281/zenodo.18525877.

[21] S. Liu, *Thermodynamic Stability Constraints on the Operator Algebra of Persistent Open Quantum Subsystems*, Zenodo (2026), DOI: 10.5281/zenodo.18525891.

[22] S. Liu, *The Realizability Bridge: Algebraic Closure in the Q-RAIF Framework*, Zenodo (2026), DOI: 10.5281/zenodo.18528935.
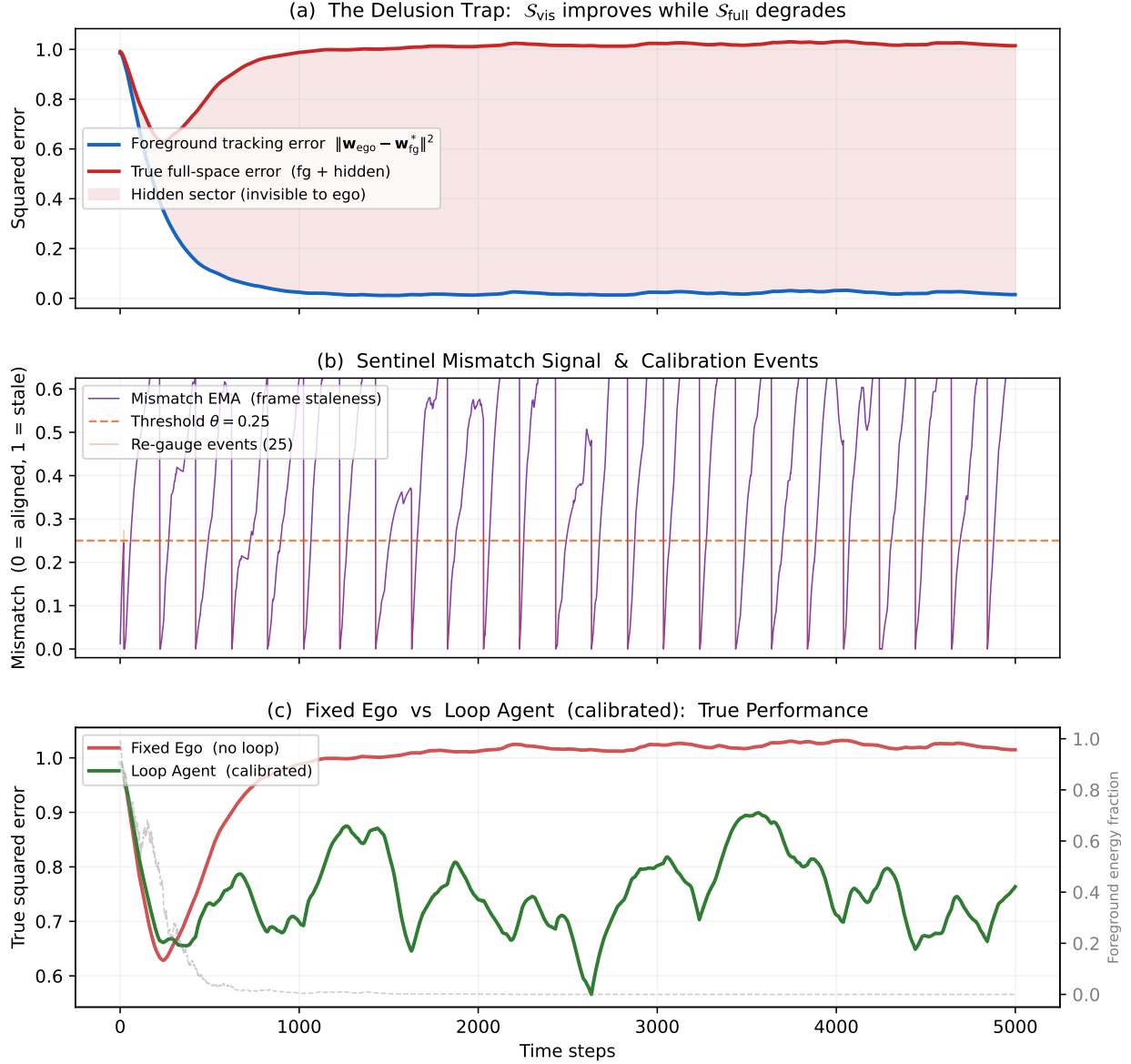
Figure 1: **Delusion trap and calibration loop.** $d = 20$, $k = 5$, $\Lambda = 0.02$, $\tau = 200$, $T = 5\,000$. **(a)** Foreground tracking error (blue) decreases toward zero while true full-space error (red) increases; the shaded region is the hidden sector, invisible to the ego. **(b)** Frame-staleness sentinel (purple) rises monotonically between recalibration events (orange), producing a sawtooth with 25 threshold crossings. **(c)** The calibrated loop (green) maintains lower true error than the fixed ego (red); the grey dashed line shows the foreground energy fraction decaying as signal migrates to the background.
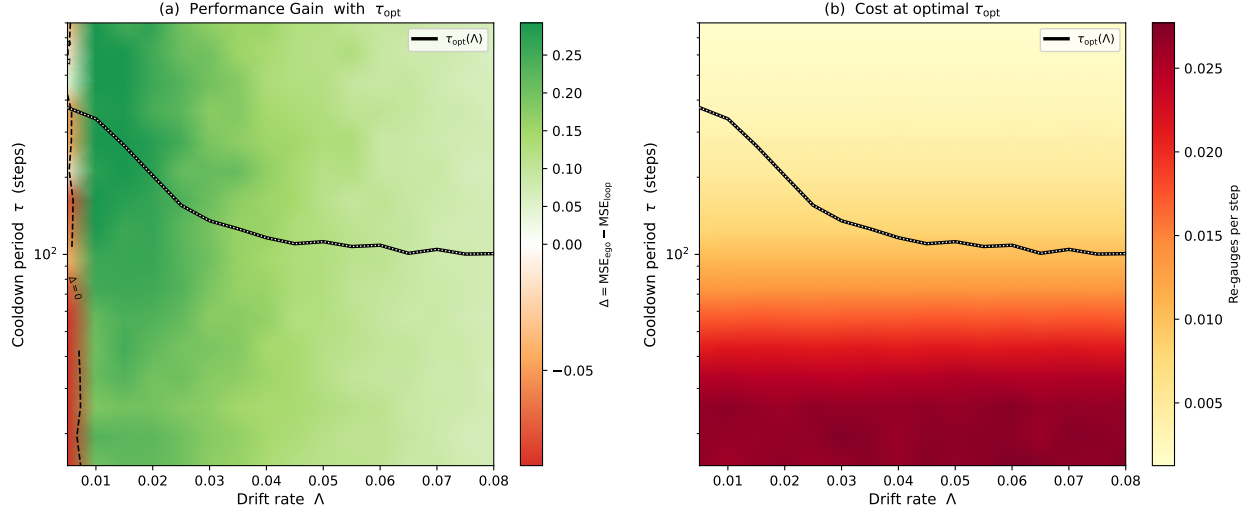
Figure 2: **Phase structure and calibration cost.** $16 \times 16$ grid, $T = 4\,000$, 6 seeds per point. **(a)** Performance gain $\Delta$; green = loop improves on ego, red = counterproductive. Solid curve: $\tau_{\mathrm{opt}}(\Lambda)$. Dashed: $\Delta = 0$ boundary. **(b)** Calibration frequency (thermodynamic cost proxy: recalibration events per step, proportional to energy expenditure under a fixed per-recalibration cost model); $\tau_{\mathrm{opt}}$ overlaid.
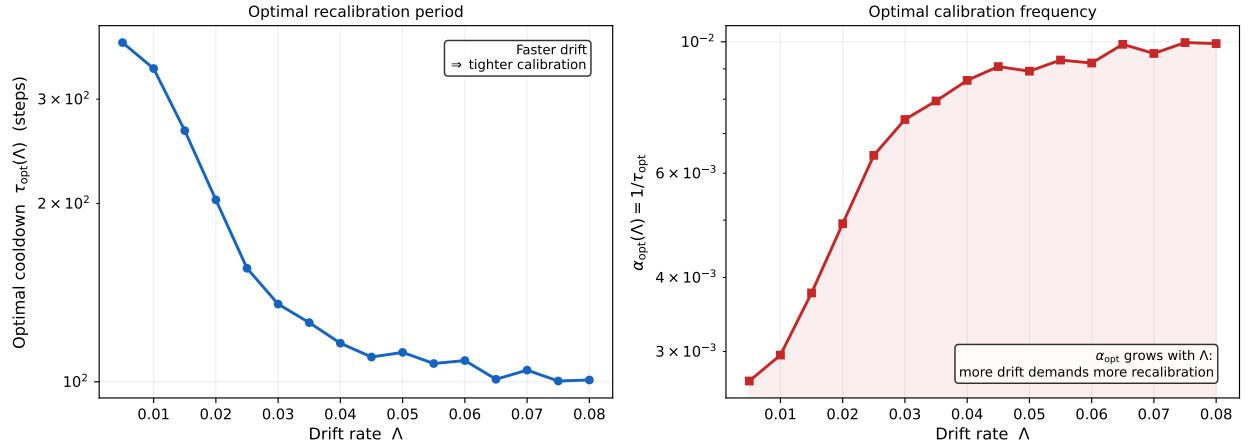


Figure 3: **Optimal calibration frequency. Left:** $\tau_{\mathrm{opt}}(\Lambda)$ decreases monotonically with drift rate. **Right:** $\alpha_{\mathrm{opt}}(\Lambda) = 1/\tau_{\mathrm{opt}}$ increases with drift rate and saturates at the learner's settling timescale ($\sim 100$ steps), consistent with an observation-window floor.