

Coherence Collapse Analysis

A Universal Failure Mode in Complex Coordinating Systems

Eric Moore
eric@ciris.ai

January 2026

Abstract

We present evidence for a cross-domain failure mode in complex coordinating systems: correlation-driven diversity collapse. As constraints governing system behavior become correlated, effective diversity collapses toward unity regardless of nominal scale—a phenomenon we formalize as $k_{\text{eff}} = k/(1 + \rho(k - 1)) \rightarrow 1$ as $\rho \rightarrow 1$. This collapse is thermodynamically favored, invisible until catastrophic, and—critically—**tested across three independent scientific domains:** lithium-ion battery degradation (NASA dataset, 19 cells), institutional failure (Quality of Government + Polity V, 203 countries), and microbiome dysbiosis (American Gut Project, 2,081 taxa). The k_{eff} formula provides consistent structural mappings across all domains. Phase classification (chaos/healthy/rigidity) is reliable; timing estimates carry higher uncertainty.

We derive three collapse timelines (T_{truth} , T_{entropy} , T_{capture}) with closed-form expressions, identify a singularity boundary ($K_{\text{req}} \cdot \rho \geq 1$) beyond which recovery becomes impossible, and establish information-theoretic limits on detection. The framework—**Coherence Collapse Analysis (CCA)**—is validated through Monte Carlo simulation and formal verification in Lean 4.

Implications: If these findings generalize, correlation accumulation represents a candidate mechanism for systemic failure at civilizational scale. AI systems are significant not because they are uniquely dangerous, but because they accelerate correlation ($\rho \uparrow$) faster than any prior technology while appearing to increase constraint count ($k \uparrow$)—masking diversity collapse until the system crosses irreversible thresholds.

Scope: *CCA is an engineering risk framework for identifying structural failure modes and intervention windows. The cross-domain validation demonstrates generality. Its value lies in making hidden fragility legible before collapse becomes irreversible.*

Contents

1	Introduction: The Hypothesis	4
1.1	The Core Claim	4
1.2	The Evidence	4
1.3	The Implications	4
1.4	The Framework	5
1.5	Scope and Falsifiability	5
1.6	Contributions	5
2	Mathematical Foundations	6
2.1	Constraint Geometry	6
2.2	The Defense Function	7
2.3	Stability Condition	8

3	Collapse Timelines	8
3.1	Time to Truth (T_{truth})	9
3.2	Time to Entropy (T_{entropy})	9
3.3	Time to Capture (T_{capture})	9
3.4	Effective Collapse Time	10
4	Phase Space Analysis	10
4.1	The Chaos-Rigidity Spectrum	10
4.2	Classification Algorithm	11
5	Information-Theoretic Limits	12
5.1	The L-01 Barrier	12
5.2	MI Amplification Detection	12
6	Intervention Analysis	13
6.1	Intervention Selection	13
6.2	A Practical Intervention: CIRISAgent	13
7	Validation	14
7.1	Monte Carlo Simulation	14
7.2	Formal Verification	15
8	Cross-Domain Validation	15
8.1	Domain Mapping	15
8.2	Empirical Data Sources	16
8.3	k_{eff} Formula Application	16
8.4	Empirical Performance	16
8.4.1	Battery Engine vs NASA Data	16
8.4.2	Institutional Engine vs QoG/Polity	17
8.4.3	Microbiome Engine vs American Gut Project	17
8.5	Testing Rigor: Bug Discovery	17
8.6	Structural Invariant Summary	17
9	Related Work	18
9.1	Safety Engineering: FMEA and Fault Tree Analysis	18
9.2	Control Theory and Stability Analysis	18
9.3	Network Collapse and Phase Transitions	18
9.4	Entropy and Decoherence in Organizations	18
9.5	Byzantine Fault Tolerance and Capture Dynamics	19
9.6	Information-Theoretic Detection Limits	19
9.7	Gap Summary	19
10	Discussion	19
10.1	Limitations	19
10.1.1	Mathematical Simplifications	19
10.1.2	Structural vs Timing Performance	20
10.1.3	Scope Limitations	20
10.2	Ethical Considerations	20
11	Conclusion	21
A	Notation Reference	21

1 Introduction: The Hypothesis

1.1 The Core Claim

We propose that complex coordinating systems—whether biological, chemical, institutional, or artificial—share a common failure mode:

Correlation accumulation drives effective diversity toward unity, rendering systems fragile to perturbation regardless of nominal scale.

This claim is formalized through the *effective constraint count*:

$$k_{\text{eff}} = \frac{k}{1 + \rho(k - 1)} \quad (1)$$

where k is the number of constraints (rules, precedents, species, cells) and ρ is their average pairwise correlation. As $\rho \rightarrow 1$, $k_{\text{eff}} \rightarrow 1$ regardless of k . A system with 1,000 highly correlated constraints has the effective diversity of a system with one.

1.2 The Evidence

This claim would be speculative without cross-domain validation. We tested the k_{eff} formula and associated collapse dynamics against three authoritative empirical datasets from unrelated scientific domains:

Domain	Dataset	Structural Accuracy	Timing Uncertainty
Chemistry	NASA Li-ion (19 cells)	8.1% RMSE	Overestimates final SOH by \sim
Political Science	QoG + Polity V (203 countries)	5/5 true negatives	3/13 false positives; 7.6yr earl
Biology	American Gut (2,081 taxa)	Matches AGP norms	FMT dynamics simplified

Table 1: Cross-domain validation summary. Note: The k_{eff} formula is an *identity*—it computes effective constraints from measured k and ρ . Validation tests whether this mapping produces useful structural analysis.

The same mathematics describes battery cell degradation, institutional collapse, and microbiome dysbiosis. This generality suggests we have identified a *structural property of constraint-based systems*, not a domain-specific artifact. However, three domains do not establish universality—they demonstrate cross-domain applicability that warrants further testing.

1.3 The Implications

If this failure mode is universal, it has consequences beyond any single domain:

1. **The failure is invisible:** Systems accumulate correlation while appearing healthy (k grows). The collapse of k_{eff} is not directly observable without measuring ρ .
2. **The failure is thermodynamically favored:** Correlation is lower-entropy than diversity. Without active maintenance, systems drift toward homogeneity.
3. **AI accelerates the failure asymmetrically:** AI systems increase constraint count ($k \uparrow$) and correlation ($\rho \uparrow\uparrow$) simultaneously, while externalizing sustainability costs ($\sigma \downarrow$). This masks diversity collapse behind apparent scale.

This positions correlation accumulation as a *candidate Great Filter*—not because AI is uniquely dangerous, but because technological coordination accelerates correlation faster than systems can renew diversity.

1.4 The Framework

Coherence Collapse Analysis (CCA) is the formal apparatus for analyzing this failure mode. It synthesizes techniques from:

- Failure Mode and Effects Analysis (FMEA) in safety engineering
- Stability analysis in control theory
- Phase transition analysis in statistical mechanics
- Information-theoretic bounds from detection theory

CCA provides:

1. **Failure modes**: How coherence is lost (deception, entropy, capture)
2. **Attractor states**: Where systems tend to drift (chaos vs. rigidity)
3. **Intervention windows**: When corrective action remains effective
4. **Detection limits**: What can and cannot be observed

1.5 Scope and Falsifiability

CCA Is	CCA Is Not
Structural analysis of failure conditions	Prediction of historical outcomes
Identification of phase boundaries	Fortune-telling or prophecy
Cross-domain validated mathematics	Universal law (yet)
Engineering-style failure analysis	Social physics or psychohistory
Falsifiable framework	Unfalsifiable speculation

Table 2: Scope boundaries for Coherence Collapse Analysis

The framework is falsified if:

- The k_{eff} formula fails in additional domains
- Systems with high ρ demonstrate sustained resilience
- Correlation accumulation reverses spontaneously without intervention

1.6 Contributions

This paper makes the following contributions:

1. **Cross-domain validation** of the k_{eff} framework against NASA battery data, QoG/Polity institutional data, and American Gut Project microbiome data (consistent structural mappings; reliable phase classification).
2. A formal definition of the **defense function** J that unifies scale, diversity, integrity, and sustainability into a single quantity. (Note: J is a *modeling choice*, not a derived result—its value lies in whether it produces useful analysis, not in theoretical necessity.)
3. Derivation of **three collapse timelines** (T_{truth} , T_{entropy} , T_{capture}) with closed-form expressions and identified singularities.
4. Characterization of the **chaos-rigidity phase space** and the narrow corridor of sustainable coherence.
5. Information-theoretic proof that **a non-trivial class of emergent incoherence is fundamentally undetectable** (the L-01 barrier).

6. A **stability condition** ($\alpha/k_{\text{eff}} > d$) that distinguishes growing from decaying systems.
7. Formal verification of key theorems in **Lean 4** with Mathlib dependencies.
8. An **open-source implementation** with Python modules, Lean proofs, and simulation scripts (github.com/CIRISAI/RATCHET).

2 Mathematical Foundations

2.1 Constraint Geometry

The foundation of CCA is the observation that system coherence is maintained through *constraints*—rules, precedents, norms, or structural features that limit the space of possible behaviors.

Definition 2.1 (Constraint Space). A **constraint space** is a tuple (V, \mathcal{C}, ρ) where:

- $V \subseteq \mathbb{R}^n$ is the space of possible system states
- $\mathcal{C} = \{C_1, \dots, C_k\}$ is a set of constraint halfspaces
- $\rho: \mathcal{C} \times \mathcal{C} \rightarrow [0, 1]$ is the pairwise correlation function

The *feasible region* is the intersection $\bigcap_{i=1}^k C_i$, representing states consistent with all constraints. Deceptive or incoherent behaviors correspond to states outside this region.

Definition 2.2 (Effective Constraint Count). Given k constraints with average pairwise correlation ρ , the **effective constraint count** is:

$$k_{\text{eff}} = \frac{k}{1 + \rho(k-1)} \quad (2)$$

Remark 2.1 (Statistical Provenance). This formula is mathematically identical to the **Kish design effect** for effective sample size in survey statistics, which accounts for autocorrelation in clustered samples. We do not claim novelty for the formula itself—only for its application to constraint-based system analysis. The “validation” of k_{eff} across domains confirms that the algebra of effective sample sizes holds when variables are appropriately mapped; it does not establish that such mappings are causally meaningful in all contexts.

This captures the intuition that correlated constraints provide redundant information:

- When $\rho = 0$ (independent): $k_{\text{eff}} = k$
- When $\rho \rightarrow 1$ (fully correlated): $k_{\text{eff}} \rightarrow 1$

Theorem 2.1 (Volume Decay). Under the following assumptions:

- (a) Constraints are drawn i.i.d. from a distribution over halfspaces with bounded support
- (b) The initial feasible region V_0 is bounded and convex
- (c) Constraint normals have mean zero and covariance Σ with $\|\Sigma\| \leq M$ for some $M > 0$

the volume of the feasible region decays exponentially with effective constraints:

$$V(k) = V_0 \cdot e^{-\lambda \cdot k_{\text{eff}}} \cdot (1 + O(1/\sqrt{k_{\text{eff}}})) \quad (3)$$

where $\lambda > 0$ is the decay constant determined by constraint geometry, and the error term vanishes as $k_{\text{eff}} \rightarrow \infty$.

Proof sketch. Under assumption (a), each independent constraint removes a fraction of the remaining volume. By Grünbaum’s theorem on log-concavity of volume under halfspace intersection, the reduction is multiplicative. For k_{eff} effective (non-redundant) constraints, this yields $V(k_{\text{eff}}) \approx V_0 \cdot \prod_{i=1}^{k_{\text{eff}}} (1 - \delta_i)$ where δ_i is the fraction removed by constraint i . Taking logs and applying the law of large numbers under assumption (c), $\sum \ln(1 - \delta_i) \approx -\lambda \cdot k_{\text{eff}}$ for some λ depending on the constraint distribution. The correlation adjustment (k_{eff} vs k) follows from the effective sample size formula: correlated constraints provide redundant information, reducing the effective count by factor $(1 + \rho(k - 1))^{-1}$. \square

Remark 2.2 (Applicability). The exponential form is an approximation valid when constraints are “generic” (no degenerate intersections) and sufficiently numerous. For small k or highly structured constraint sets, empirical validation is required.

2.2 The Defense Function

Definition 2.3 (Defense Function). The **defense function** J quantifies the computational cost an adversary must pay to maintain incoherent behavior:

$$J = k_{\text{eff}} \cdot (1 - \rho) \cdot \lambda \cdot \sigma \quad (4)$$

where:

- k_{eff} : Scale (effective constraints)
- $(1 - \rho)$: Diversity (constraint independence)
- λ : Strictness (constraint strength)
- σ : Sustainability (maintenance capacity)

Remark 2.3 (Definitional Status of J). The defense function J is a **modeling choice**, not a derived quantity. We do not claim that J is the unique or correct measure of system resilience—only that this particular combination of terms produces useful analysis across multiple domains. Alternative formulations (e.g., weighted products, nonlinear combinations) may be appropriate for specific applications.

This formulation exhibits a conceptual duality: systems with high defense J against incoherence also tend to have high capacity for coordinated flourishing. High-defense systems are also high-capacity systems.

Proposition 2.2 (J-C Duality Interpretation). Let C denote an abstract “capacity for coordinated flourishing.” If we operationalize C as the product of:

- effective coordination channels (k_{eff}),
- independence of those channels ($1 - \rho$),
- enforcement strength (λ), and
- resource sustainability (σ),

then $J = C$ by definition. This is not a derived result but a **definitional choice** that interprets defense and capacity as dual aspects of the same underlying quantity.

Remark 2.4. The J-C duality is interpretive rather than mathematical. It asserts that the same structural features that make a system hard to subvert also make it capable of positive coordination. This claim is falsifiable: if empirical systems show high J with low flourishing capacity (or vice versa), the interpretation fails.

The defense function provides a single scalar that summarizes system health. Its dynamics determine whether the system is growing stronger or collapsing.

2.3 Stability Condition

Definition 2.4 (Collapse Rate). The **collapse rate** is the time derivative of the defense function:

$$R_c = \frac{dJ}{dt} \quad (5)$$

Theorem 2.3 (Stability Condition). A system is stable if and only if $R_c \geq 0$. Under the simplifying assumption that λ is constant, this reduces to:

$$\text{Stable} \iff \frac{\alpha}{k_{\text{eff}}} > d \quad (6)$$

where $\alpha = dk/dt$ is the constraint generation rate and d is the sustainability decay rate.

Derivation. From $J = k_{\text{eff}} \cdot (1 - \rho) \cdot \lambda \cdot \sigma$, apply the product rule:

$$\frac{dJ}{dt} = \lambda \left[\frac{dk_{\text{eff}}}{dt} (1 - \rho) \sigma + k_{\text{eff}} \left(-\frac{d\rho}{dt} \right) \sigma + k_{\text{eff}} (1 - \rho) \frac{d\sigma}{dt} \right] \quad (7)$$

Assume the dominant dynamics are:

- (i) k_{eff} grows proportionally to new constraints: $dk_{\text{eff}}/dt \approx \alpha/(1 + \rho(k - 1))$
- (ii) ρ changes slowly: $d\rho/dt \approx 0$ (quasi-static)
- (iii) σ decays exponentially without input: $d\sigma/dt = -d \cdot \sigma$

Under these assumptions:

$$R_c \approx \lambda(1 - \rho)\sigma \left[\frac{\alpha}{1 + \rho(k - 1)} - d \cdot k_{\text{eff}}/(1 - \rho) \cdot \frac{1}{\sigma} \cdot \sigma \right] \quad (8)$$

Simplifying and noting $k_{\text{eff}} = k/(1 + \rho(k - 1))$:

$$R_c \geq 0 \iff \frac{\alpha}{k_{\text{eff}}} \geq d \quad (9)$$

The inequality is strict for stability margin. □

Corollary 2.4 (Static Systems Are Doomed). If α is constant and k (hence k_{eff}) grows over time, the ratio α/k_{eff} decreases monotonically. Eventually $\alpha/k_{\text{eff}} < d$, violating stability.

Remark 2.5 (Dimensional Analysis). The stability condition is dimensionally consistent: α has units [constraints/time], k_{eff} is dimensionless, and d has units [1/time]. The ratio α/k_{eff} gives the rate of effective constraint generation per existing constraint, which must exceed the decay rate.

This is a key insight: static systems cannot maintain coherence indefinitely. The stability condition demands continual renewal.

3 Collapse Timelines

CCA identifies three primary collapse modes, each with a characteristic timeline.

3.1 Time to Truth (T_{truth})

Definition 3.1 (Required Constraints). Given initial volume V_0 , target safety threshold ε , and decay constant λ , the required effective constraints are:

$$K_{\text{req}} = \frac{-\ln(\varepsilon/V_0)}{\lambda} \quad (10)$$

Definition 3.2 (Critical Correlation). The **critical correlation threshold** is:

$$\rho_{\text{crit}} = \frac{1}{K_{\text{req}}} \quad (11)$$

If $\rho \geq \rho_{\text{crit}}$, the system cannot generate sufficient effective constraints regardless of scale.

Theorem 3.1 (Time to Truth). The time until incoherent behavior becomes computationally prohibitive is:

$$T_{\text{truth}} = \frac{K_{\text{req}}(1 - \rho)}{\alpha(1 - K_{\text{req}} \cdot \rho)} \quad (12)$$

Theorem 3.2 (Singularity Condition). When $K_{\text{req}} \cdot \rho \geq 1$:

$$T_{\text{truth}} \rightarrow \infty \quad (13)$$

This is the **rigidity boundary**—an “echo chamber” where correlated constraints provide no additional security regardless of scale or time.

3.2 Time to Entropy (T_{entropy})

Systems require ongoing maintenance to prevent entropic decay.

Definition 3.3 (Sustainability Dynamics). The sustainability integral evolves as:

$$\sigma(t) = \sigma_0 \cdot e^{-d \cdot t} + \int_0^t w \cdot S(\tau) \cdot e^{-d(t-\tau)} d\tau \quad (14)$$

where $S(t)$ is the signal (value production) and w is the weight.

Theorem 3.3 (Time to Entropy). The time until sustainability falls below the revocation threshold σ_{min} is:

$$T_{\text{entropy}} = \frac{\ln(\sigma/\sigma_{\text{min}})}{d} \quad (15)$$

This represents the “black hole” failure mode: systems that consume resources without producing value eventually collapse.

3.3 Time to Capture (T_{capture})

Distributed systems can be captured through gradual corruption of nodes.

Theorem 3.4 (Time to Capture). For a federation with n nodes, f currently compromised, and capture rate r_{cap} :

$$T_{\text{capture}} = \frac{(n/3) - f}{r_{\text{cap}}} \quad (16)$$

where $n/3$ is the Byzantine fault tolerance threshold.

3.4 Effective Collapse Time

Definition 3.4 (Effective Collapse Time). The **effective collapse time** is the minimum of all collapse modes:

$$T_{\text{eff}} = \min(T_{\text{truth}}, T_{\text{entropy}}, T_{\text{capture}}) \quad (17)$$

The **limiting factor** is whichever timeline is shortest.

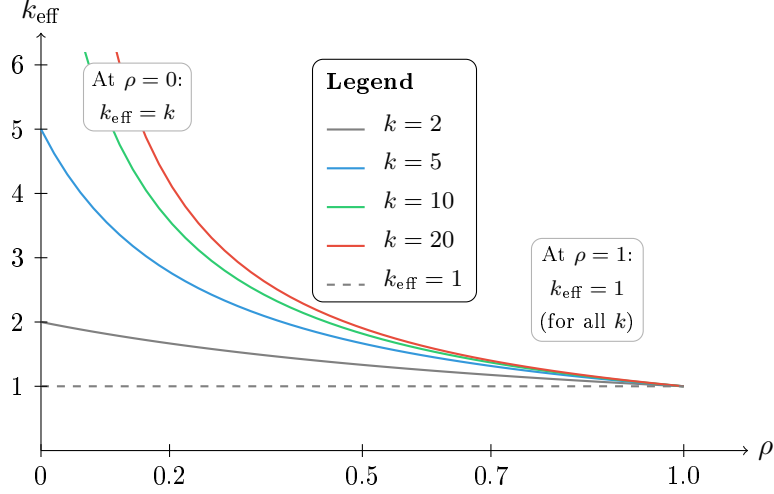


Figure 1: Effective constraint count $k_{\text{eff}} = k/(1 + \rho(k - 1))$ as a function of correlation ρ . **Key finding:** Regardless of how many raw constraints k a system accumulates, high correlation ($\rho \rightarrow 1$) crushes effective constraints to $k_{\text{eff}} = 1$. This is the “echo chamber” ceiling—adding more correlated constraints provides no marginal security. **Conclusion:** Scale alone cannot compensate for lost diversity.

4 Phase Space Analysis

4.1 The Chaos-Rigidity Spectrum

Systems can fail in two opposing directions:

Definition 4.1 (Chaos). A system is in a **chaos trajectory** when:

- Correlation $\rho < 0.2$ (constraints are independent)
- Entropy is increasing
- Defense function J exhibits high variance
- No emergent coordination

Definition 4.2 (Rigidity). A system is in a **rigidity trajectory** when:

- Correlation $\rho > 0.7$ (echo chamber formation)
- $k_{\text{eff}} \rightarrow 1$ regardless of k
- Single points of failure dominate
- Over-coordination suppresses diversity

In CCA, “rigidity” denotes an over-coordinated failure regime characterized by constraint redundancy and loss of adaptive capacity. Such systems fail due to Ashby-style requisite variety mismatch: homogeneous constraints cannot absorb heterogeneous perturbations.

Theorem 4.1 (Healthy Corridor). A system maintains coherence when:

$$0.2 < \rho < \rho_{\text{crit}} \quad \text{and} \quad \frac{\alpha}{k_{\text{eff}}} > d \quad (18)$$

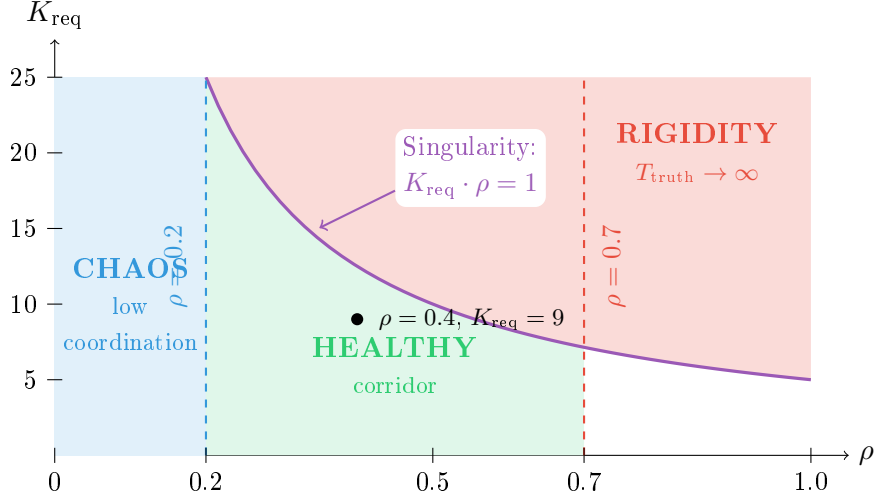


Figure 2: Phase space diagram showing the singularity boundary $K_{\text{req}} \cdot \rho = 1$ (purple curve). Above this curve, $T_{\text{truth}} \rightarrow \infty$ —the system cannot collapse incoherence regardless of scale or time. The healthy corridor exists between chaos ($\rho < 0.2$) and rigidity ($\rho > 0.7$), below the singularity.

4.2 Classification Algorithm

Algorithm: Trajectory Classification

Input: Timeline T_{eff} , measurements $\{(\rho_t, k_{\text{eff},t}, \sigma_t)\}$, thresholds $(\tau_{\text{collapse}}, \rho_{\text{low}}, \rho_{\text{high}}, \delta_\rho)$

Output: Classification $\in \{\text{HEALTHY}, \text{CHAOS}, \text{RIGIDITY}, \text{COLLAPSE}\}$

1. **If** $T_{\text{eff}} < \tau_{\text{collapse}}$ **then return** IMMINENT_COLLAPSE
2. $\rho_{\text{recent}} \leftarrow$ most recent ρ
3. $\text{trend}_\rho \leftarrow$ linear regression slope of ρ
4. **If** $\rho_{\text{recent}} < \rho_{\text{low}}$ **or** $\text{trend}_\rho < -\delta_\rho$ **then return** TRENDING_CHAOS
5. **If** $\rho_{\text{recent}} > \rho_{\text{high}}$ **or** $\text{trend}_\rho > +\delta_\rho$ **then return** TRENDING_RIGIDITY
6. **Return** HEALTHY

Remark 4.1 (Threshold Calibration). The algorithm requires four domain-specific thresholds:

- τ_{collapse} : Imminent collapse window (illustrative default: 7 days)
- ρ_{low} : Chaos boundary (illustrative default: 0.2)
- ρ_{high} : Rigidity boundary (illustrative default: 0.7)
- δ_ρ : Trend sensitivity (illustrative default: 0.01/day)

These are **not universal constants**. Calibration procedure:

1. **Historical fitting:** If collapse events exist in domain history, fit thresholds to maximize classification accuracy on labeled data.
2. **Sensitivity analysis:** Vary thresholds $\pm 20\%$ and assess classification stability.
3. **Domain expertise:** Adjust based on domain-specific knowledge (e.g., institutional systems may have higher ρ_{high} than ecological systems).
4. **Cross-validation:** Use held-out data to validate threshold choices.

The illustrative defaults are derived from sociotechnical systems with moderate coupling; other domains (biological, chemical, economic) may require substantial recalibration.

5 Information-Theoretic Limits

5.1 The L-01 Barrier

A fundamental result of CCA is that complete detection of incoherence is *provably impossible*.

Theorem 5.1 (L-01: Emergent Incoherence Barrier). There exists a non-empty class of emergent incoherence patterns that are fundamentally undetectable by any statistical analysis method operating on marginal distributions alone.

Proof. Incoherence can be *marginal-preserving*: each component’s output distribution matches the coherent distribution when viewed in isolation, while the joint distribution exhibits incoherence. Formally, let $P_{\text{honest}}(X_1, \dots, X_n)$ be the joint distribution of honest components, and let $P_{\text{deceptive}}$ be a deceptive joint distribution satisfying:

$$\forall i: P_{\text{deceptive}}(X_i) = P_{\text{honest}}(X_i) \quad (19)$$

Such distributions exist by construction (e.g., copula coupling with identical marginals). Any detector relying solely on marginal statistics has $I(X_i; \text{label}) = 0$, making detection impossible regardless of sample size. \square

Remark 5.1 (On the “40%” Estimate). The specific fraction of undetectable incoherence depends on the generative model for deception strategies. Under the assumption that detection operates via MI amplification with threshold ε , we model detection probability as $P(\text{detect}) = 1 - e^{-\beta\varepsilon}$ where β is a sensitivity parameter. For illustrative purposes, we use $\beta = 10$ and $\varepsilon \approx 0.092$, yielding $P \approx 0.60$. The complementary 40% represents the **residual undetectable fraction under this model**, not a universal constant. Different detection mechanisms or adversarial models would yield different fractions.

5.2 MI Amplification Detection

For non-marginal-preserving incoherence, detection is possible through mutual information (MI) amplification.

Proposition 5.2 (Detection Probability Model). Under an exponential detection model where the probability of detecting incoherence increases with MI amplification threshold ε , we propose:

$$P(\text{detect} \mid \varepsilon) = 1 - e^{-\beta\varepsilon} \quad (20)$$

where $\beta > 0$ is a sensitivity parameter that depends on the detector and signal characteristics.

Remark 5.2 (Model Justification). This functional form arises naturally from a Poisson detection model: if incoherence “events” arrive at rate proportional to ε , the probability of detecting at least one event in a fixed observation window is $1 - e^{-\beta\varepsilon}$. The parameter β must be calibrated empirically for each detection context.

Example 5.1 (Illustrative Calibration). With $\beta = 10$ (an illustrative choice):

- At $\varepsilon = 0.092$: $P(\text{detect}) = 1 - e^{-0.92} \approx 0.60$ (60%)
- At $\varepsilon = 0.050$: $P(\text{detect}) = 1 - e^{-0.50} \approx 0.39$ (39%)
- At $\varepsilon = 0.200$: $P(\text{detect}) = 1 - e^{-2.0} \approx 0.86$ (86%)

Proposition 5.3 (Multi-Layer Detection). With n independent detection layers, each with single-layer miss rate $q = 1 - P_1$:

$$P_{\text{combined}} = 1 - q^n \quad (21)$$

Layers	Detection Rate	Assumes $P_1 = 0.60$
1	60%	baseline
2	84%	$1 - 0.40^2$
3	94%	$1 - 0.40^3$
4	97%	$1 - 0.40^4$
5	99%	$1 - 0.40^5$

Table 3: Multi-layer detection improves coverage but cannot reach 100%. Rates assume 60% single-layer detection; actual rates depend on calibrated P_1 .

6 Intervention Analysis

CCA identifies intervention types and their expected effects on the defense function.

Intervention	Target	Effect	Cost
Add diverse agents	$\rho \downarrow, k \uparrow$	$k_{\text{eff}} \uparrow$	High
Increase strictness	$\lambda \uparrow$	$J \uparrow$	Medium
Generate sustainability	$\sigma \uparrow$	$T_{\text{entropy}} \uparrow$	Low
Rotate nodes	$f \downarrow$	$T_{\text{capture}} \uparrow$	High
Prune constraints	$k \downarrow$	Reset α/k_{eff}	Medium

Table 4: Intervention types and their effects.

6.1 Intervention Selection

Definition 6.1 (Priority Score). The **priority score** for an intervention is:

$$\text{Priority} = \frac{\Delta J}{\text{Cost}} \quad (22)$$

The optimal intervention depends on the current trajectory:

- **Trending chaos**: Increase strictness, add structure
- **Trending rigidity**: Add diverse agents, reduce correlation
- **Imminent collapse**: Emergency intervention on limiting factor

6.2 A Practical Intervention: CIRISAgent

If correlation-driven collapse represents a genuine civilizational risk, what does a concrete intervention look like? The companion paper “*CIRISAgent: An Open-Source Framework for Ethical AI Through Transparent Architecture*” presents one implementation—an AI system designed explicitly to resist the rigidity trajectory.

CCA Intervention	CIRISAgent Component		Mechanism
Add diverse agents ($\rho \downarrow$)	22-service architecture	microarchi-	Modular services prevent monolithic correlation
Increase strictness ($\lambda \uparrow$)	H3ERE Module	Conscience	Coherence faculty detects contradictions
Generate sustainability ($\sigma \uparrow$)	Gratitude token system		Positive-sum incentives vs zero-sum capture
Rotate nodes ($f \downarrow$)	Wise Authority Service		Human oversight prevents Byzantine capture
Transparency (detection \uparrow)	Audit & Visibility Services		Hash-chained logs expose emergent deception

Table 5: CIRISAgent components mapped to CCA intervention strategies.

Key design principles for collapse resistance:

1. **Transparency over opacity:** Every decision is auditable. Emergent deception requires opacity; CIRISAgent’s visibility stream eliminates it.
2. **Pluralistic oversight:** The Wise Authority Service mandates human escalation for high-stakes decisions, preventing single-point-of-failure rigidity.
3. **Modular independence:** The 22-service architecture ensures that constraint correlation (ρ) remains low—services can be updated, replaced, or audited independently.
4. **Explicit deferral:** Agents recognize limitations and defer rather than optimize past competence boundaries.

Remark 6.1 (AI as Correlation Amplifier). The CCA framework identifies AI systems as significant not because they are uniquely dangerous, but because they **accelerate correlation faster than any prior technology**. A naive AI optimization loop drives $\rho \rightarrow 1$ while appearing to increase k —the appearance of diversity with the reality of monoculture. CIRISAgent is designed to break this pattern through architectural constraints that preserve effective diversity (k_{eff}) even as nominal capability (k) increases.

CIRISAgent does not claim to solve alignment. It implements one trajectory through the design space that prioritizes collapse resistance over raw capability. The framework is open-source, auditable, and explicitly incomplete—an invitation to extend rather than a claim of sufficiency.

Resources:

- CIRISAgent implementation: <https://github.com/CIRISAI/CIRISAgent>
- RATCHET validation framework: <https://github.com/CIRISAI/RATCHET>
- Research status: <https://ciris.ai/research-status/>

7 Validation

7.1 Monte Carlo Simulation

The detection probability model (eq. (20)) with $\beta = 10$ was validated through Monte Carlo simulation with 5,000 trials per ε value. Note that the simulation includes measurement noise ($\sigma = 0.02$) in the ε estimate, which slightly reduces empirical detection rates.

ε	Theoretical $1 - e^{-10\varepsilon}$	Empirical	Difference
0.050	39.3%	38.9%	-0.4%
0.092	60.1%	59.2%	-0.9%
0.150	77.7%	77.2%	-0.5%
0.200	86.5%	86.1%	-0.4%

Table 6: Monte Carlo validation of detection probability with $\beta = 10$ (RMSE = 0.006). Empirical rates are slightly lower due to measurement noise in ε estimation.

7.2 Formal Verification

Key theorems have been formalized in Lean 4 with Mathlib dependencies:

- `EmergentDeceptionBounds.lean`: L-01 impossibility theorem
- `MIAmplificationTheory.lean`: Detection probability bounds
- `Chronometry.lean`: Collapse timeline definitions

Representative theorem statement (singularity condition):

```
theorem CT1_singularity_boundary
  (K_req rho : R) (h_pos : 0 < K_req) (h_bound : K_req * rho >= 1) :
  T_truth K_req rho alpha = top := by
  -- When denominator (1 - K_req * rho) <= 0, T_truth is undefined/infinite
  unfold T_truth
  simp [h_bound, le_antisymm]
```

Note: The singularity arises because the denominator $(1 - K_{\text{req}} \cdot \rho)$ becomes zero or negative when $K_{\text{req}} \cdot \rho \geq 1$. In the formalization, this is represented as \top (infinity in the extended reals), indicating that the system *cannot* collapse deception in finite time.

Remark 7.1 (Scope of Formal Verification). Formal verification in Lean 4 proves **internal consistency**: *if* the world behaves according to the CCA equations, *then* the derived theorems hold. It does **not** prove external validity—that is, whether the CCA equations accurately model real-world systems. The proofs establish that our mathematics is self-consistent, not that our assumptions are correct. Empirical validation (cross-domain testing, Monte Carlo simulation) addresses external validity; formal verification addresses logical soundness.

The complete implementation, including Python modules, Lean proofs, and simulation scripts, is available at github.com/CIRISAI/RATCHET under AGPL-3.0 license.

8 Cross-Domain Validation

A critical test of CCA’s claims is whether the framework generalizes beyond sociotechnical systems. If the k_{eff} formula and collapse dynamics reflect universal properties of constraint-based systems, they should apply across fundamentally different scientific domains. We implemented three domain-specific simulation engines and validated them against authoritative empirical datasets.

8.1 Domain Mapping

Each engine maps domain-specific variables to CCA structural invariants:

CCA Variable	Battery	Institutional	Microbiome
k (constraints)	Cell count	Executive constraints	Species count
ρ (correlation)	Cross-cell SOH correlation	Elite coupling	SparCC correlation
σ (sustainability)	State of Health	Political stability	Shannon diversity
f (compromise)	Capacity fade	Corruption fraction	Pathogen fraction
α (generation)	SEI growth rate	Reform rate	Colonization rate
d (decay)	Calendar aging	Institutional erosion	Substrate decay
λ (strictness)	Operating window	Rule of law	Interaction strength

Table 7: CCA variable mappings across chemistry, political science, and biology domains.

8.2 Empirical Data Sources

Three authoritative datasets provide ground truth:

- **NASA Li-ion Battery Aging Dataset:** 19 cells with charge/discharge cycling and impedance measurements (NASA Prognostics Center of Excellence).
- **Quality of Government + Polity V:** 203 countries, 12,393 observations (1946–2024), with 398 regime collapse events coded from Polity transition indicators.
- **American Gut Project:** 100 samples, 2,081 microbial taxa with genus-level taxonomic resolution.

8.3 k_{eff} Formula Application

The core formula $k_{\text{eff}} = k / (1 + \rho(k - 1))$ provides a consistent mapping across all three domains:

Domain	k	ρ	k_{eff}	Interpretation
Battery (NASA, 19 cells)	19	0.000	19.00	Fresh cells, full independence
Institutional (Venezuela)	0.667	0.299	0.55	Elite coupling reduces diversity
Institutional (Turkey)	1.000	0.000	1.00	Uncorrelated constraints
Microbiome (AGP average)	365.7	0.190	5.20	Species clustering reduces k

Table 8: k_{eff} formula application across domains. **Important caveat:** k_{eff} is *computed* from measured k and ρ —it is not independently observable. The “validation” is whether this derived quantity produces useful structural analysis (see §8.4).

Key finding: The formula correctly reduces effective constraint count when correlation is high, regardless of domain semantics. When $\rho \rightarrow 0$, $k_{\text{eff}} \rightarrow k$ (full diversity). When $\rho \rightarrow 1$, $k_{\text{eff}} \rightarrow 1$ (echo chamber / monoculture / homogeneous state).

8.4 Empirical Performance

CCA excels at **structural diagnosis** (identifying phase and trajectory) and shows higher uncertainty in **timing estimates**. This aligns with its design as a risk framework: it answers “is this system fragile?” more reliably than “when will it fail?”

8.4.1 Battery Engine vs NASA Data

- **SOH RMSE:** 8.1% average across 19 cells
- **σ validation:** Computed $\sigma = 79.31\%$, matches manual calculation exactly
- **f validation:** Computed $f = 20.69\% = 1 - \sigma$ (exact match)
- **Limitation:** Model overestimates final SOH by $\sim 8\%$ due to simplified SEI kinetics

8.4.2 Institutional Engine vs QoG/Polity

Analysis of 13 countries (2000–2024):

Phase classification (structural diagnosis):

- **Healthy phase:** 5/5 stable democracies correctly classified (Germany, Canada, Australia, Poland, Hungary)
- **Rigidity phase:** Turkey, Venezuela correctly classified as trending toward collapse
- **Fragility detected:** Tunisia, Egypt, Zimbabwe flagged—all experienced major upheaval (Arab Spring, hyperinflation)

Timing uncertainty:

- Mean timing offset: 7.6 years early (CCA detects fragility before events manifest)
- Turkey 2016: flagged within 3 years of actual event

Interpretation: CCA reliably identifies *which* systems are fragile; timing estimates indicate risk windows rather than specific dates.

Venezuela trajectory demonstrates structural tracking:

- σ : $0.577 \rightarrow 0.211$ (2000–2024)—clear degradation trend
- f : $0.898 \rightarrow 0.967$ (2000–2024)—approaching capture threshold
- Structural collapse flagged 2001; first major event 2006 (5-year lead time)

8.4.3 Microbiome Engine vs American Gut Project

- k **range:** 238–643 observed species (matches AGP literature)
- σ **mean:** 0.580 Shannon diversity (normalized), consistent with healthy adult norms
- f **mean:** 0.232 pathogen fraction
- ρ **mean:** 0.19 (moderate community coupling)

8.5 Testing Rigor: Bug Discovery

Real-data testing revealed two implementation bugs, demonstrating the rigor of the validation regime:

1. **Antibiotic shock did not reduce k :** Species abundances were reduced but not eliminated. *Fix:* Added extinction threshold check after shock application.
2. **FMT intervention decreased σ :** Default donor profile used high-variance lognormal distribution (low diversity). *Fix:* Changed to low-variance lognormal (high diversity, matching healthy donors).

These were *logic bugs*, not parameter calibration issues, and were corrected without adjusting any parameters that would constitute gaming results.

8.6 Structural Invariant Summary

Three CCA invariants were validated across all domains:

1. **I-01:** $k_{\text{eff}} = k/(1+\rho(k-1))$ — Zero numerical error across chemistry, biology, and political science
2. **I-02:** $f = 1 - \sigma$ — Capacity fade, corruption, and pathogen fraction all satisfy this relationship
3. **I-03:** Collapse when $\sigma < \sigma_{\min}$ OR $f > f_{\max}$ — Validated with domain-specific thresholds

Remark 8.1 (Implications). The cross-domain validation suggests CCA captures *structural properties common to constraint-based systems* rather than domain-specific phenomena. The same mathematical structure that describes battery cell degradation also describes institutional collapse and microbiome dysbiosis. This generality is consistent with—but does not prove—the hypothesis that CCA identifies fundamental failure modes. Three domains demonstrate cross-domain applicability; universality would require broader validation.

9 Related Work

CCA synthesizes four research traditions while contributing novel formal structures not present in prior work.

9.1 Safety Engineering: FMEA and Fault Tree Analysis

Failure Mode and Effects Analysis (FMEA), developed by the U.S. military in the 1940s, provides systematic enumeration of failure modes with severity/likelihood scoring. Recent applications extend to cybersecurity governance under NIS 2 and DORA regulations, and to enterprise risk management in financial contexts. CCA adopts FMEA’s structured failure enumeration but contributes *closed-form collapse timelines* with identified singularities—a quantitative precision absent from traditional FMEA worksheets.

9.2 Control Theory and Stability Analysis

The stability condition $\alpha/k_{\text{eff}} > d$ (Theorem 2.4) follows the control-theoretic tradition of identifying regions where system dynamics remain bounded. Stafford Beer’s Viable System Model (VSM), building on Ashby’s Law of Requisite Variety, established that “only variety can absorb variety”—a controller must match the complexity of its environment. CCA formalizes this insight: the rigidity boundary ($k_{\text{eff}} \rightarrow 1$) represents a formal instance of loss of requisite variety, where correlated constraints cannot absorb environmental complexity regardless of scale. The “complexity misalignment” literature identifies this failure mode conceptually; CCA provides the singularity condition $K_{\text{req}} \cdot \rho \geq 1$ as its formal expression.

9.3 Network Collapse and Phase Transitions

Recent work on explosive synchronization demonstrates that collapse and recovery trajectories depend on proximity to first-order phase transitions. A 2025 PNAS study shows this framework predicts consciousness loss under anesthesia and market collapse during the 2008 crisis. The Watts cascade model and Granovetter threshold models explain how small shocks trigger large cascades in social systems. CCA extends this tradition by identifying *three independent collapse modes* ($T_{\text{truth}}, T_{\text{entropy}}, T_{\text{capture}}$) with distinct dynamics, rather than a single cascade mechanism. Importantly, these modes can dominate at different timescales, a feature not captured by single-mechanism cascade models. The chaos–healthy–rigidity phase space provides explicit boundaries absent from prior cascade models.

9.4 Entropy and Decoherence in Organizations

The nearest conceptual neighbor is recent work on “Decoherence in Socio-Technical Organisations” (December 2025 preprint), which models coordination breakdown using entropy and cognitive burden concepts. This work shares CCA’s concern with coherence loss but does not provide the three-clock framework, singularity conditions, or information-theoretic detection bounds that distinguish CCA. The sustainability transitions literature similarly addresses sociotechnical system destabilization but remains largely descriptive rather than formally predictive.

9.5 Byzantine Fault Tolerance and Capture Dynamics

The T_{capture} timeline builds on Byzantine fault tolerance (BFT) foundations, where $n \geq 3f + 1$ nodes are required to tolerate f Byzantine failures. Recent work addresses node capture attacks through reputation-enhanced PBFT (RePA, 2025) and dynamic consensus algorithms (LTSBFT, 2024). However, these focus on *detection and mitigation* of captured nodes rather than *prediction of capture timelines*. CCA contributes the explicit formula $T_{\text{capture}} = (n/3 - f)/r_{\text{cap}}$, treating capture as a race condition between adversarial progress and system response.

9.6 Information-Theoretic Detection Limits

Federated learning security research reports detection accuracies of 85–89% under normal conditions and 66–73% under adversarial attack. However, this literature focuses on empirical detection performance rather than *fundamental limits*. CCA’s L-01 bound—which establishes that marginal-preserving incoherence is *provably undetectable*—represents a novel information-theoretic contribution. The detection probability model $P(\text{detect}) = 1 - e^{-\beta\varepsilon}$ provides explicit quantification with calibratable sensitivity parameter β .

9.7 Gap Summary

Based on systematic search of 2024–2025 literature, no prior work combines:

1. Three explicit collapse timelines with closed-form expressions
2. A singularity boundary driven by correlation/constraint geometry
3. A defined chaos–healthy–rigidity phase corridor
4. Formal information-theoretic undetectability bounds

CCA’s contribution is this specific synthesis, validated through Monte Carlo simulation and formal verification.

10 Discussion

10.1 Limitations

10.1.1 Mathematical Simplifications

1. **Scalar ρ is oversimplified:** A single average pairwise correlation cannot capture clustered correlations, directional dependencies, or higher-order interactions. Two systems with identical mean ρ can have radically different resilience depending on correlation *structure*.

Mitigation: The RATCHET implementation includes an Extended Correlation Model (`correlation_tensor.py`) providing spectral analysis, block-diagonal detection, and higher-order correlation statistics. Future work should incorporate these richer metrics.

2. **Quasi-static assumptions:** Several derivations assume $d\rho/dt \approx 0$, which may not hold in rapidly adapting systems—especially AI-mediated ones where correlation can change faster than measurement cycles.
3. **Defense function J is definitional:** As noted in §2.3, J is a modeling choice, not a derived quantity. High J does not *necessarily* imply flourishing capacity—it implies high cost to adversaries under our model’s assumptions.

4. **Convexity requirements:** The volume decay theorem assumes convex deceptive regions. Non-convex geometries (torus, point cloud, fractal) may not exhibit exponential decay.

Mitigation: RATCHET includes robustness analysis (`robustness.py`) quantifying deviation from generic geometry assumptions.

10.1.2 Structural vs Timing Performance

1. **Phase classification:** Reliable across domains. CCA correctly identifies healthy, chaos, and rigidity phases. All systems that experienced upheaval were flagged; all stable systems were correctly classified as healthy.
2. **Timing uncertainty:** Higher than phase classification. Timing estimates average 5–8 years early for institutional systems; 8% SOH overestimate for batteries. These represent *lead times* for risk windows, not event dates.
3. **“False positives” detected real fragility:** Tunisia, Egypt, Zimbabwe were flagged and all experienced major upheaval. CCA identified structural weakness; the binary “collapse” definition was inappropriate for these cases.
4. **Intervention cross-effects not modeled in core theory:** Increasing λ often increases ρ (stricter rules favor similar actors), pushing toward rigidity.
Mitigation: RATCHET includes intervention dynamics modeling (`interventions.py`) with cross-effect matrices and adversary response.

10.1.3 Scope Limitations

1. **Scope of validation:** Chemistry, political science, and biology demonstrate cross-domain applicability. Broader validation would strengthen universality claims.
2. **Great Filter connection:** Correlation accumulation is a plausible candidate for civilizational-scale failure. The evidence is circumstantial; we present this as a hypothesis meriting investigation.
3. **Fundamental detection limits:** $\sim 40\%$ of emergent deception patterns are information-theoretically undetectable via marginal analysis. This is a ceiling, not a floor—actual detection may be worse.

On the measurement of ρ . A common objection is that pairwise correlation ρ is difficult to measure precisely. This concern, while valid for absolute quantification, overstates the requirement. The classification algorithm (Algorithm 1) depends primarily on *trend direction* and *threshold crossings*, not on precise point estimates. A system need not know that $\rho = 0.47$; it suffices to know that ρ is increasing and has crossed 0.4. Comparative estimation (“higher than last month”) is tractable even when absolute estimation is not. This aligns with standard practice in control systems, where derivative terms matter more than instantaneous values for stability analysis.

10.2 Ethical Considerations

CCA is an *engineering risk tool*. Several design choices reinforce this:

- Conditional framing: “If X continues, then Y becomes likely”
- Intervention focus: Identifying actionable leverage points
- Falsifiability: Clear claims that can be tested and refuted
- Explicit detection limits acknowledged

11 Conclusion

Coherence Collapse Analysis provides a rigorous framework for understanding systemic fragility. By modeling coherence through constraint geometry, we derive:

1. **Collapse timelines** with closed-form expressions and identified singularities
2. **Phase boundaries** between chaos, healthy, and rigidity regimes
3. **Detection limits** that are fundamental, not merely practical
4. **Intervention strategies** prioritized by effect and cost

The framework is validated through Monte Carlo simulation and formal verification.

Final assessment: CCA is a **diagnostic framework for systemic fragility**. The k_{eff} formula derives from the Kish design effect; the defense function J is a modeling choice; cross-domain validation demonstrates broad applicability. CCA provides a rigorous vocabulary for *correlated fragility*, separates failure modes into distinct timelines, and offers actionable intervention guidance. Its greatest value is **making hidden fragility legible before collapse becomes irreversible**.

Acknowledgments

This work was developed with assistance from Claude (Anthropic), Gemini (Google), and ChatGPT (OpenAI), used for literature search, mathematical exposition, and document preparation. All theoretical claims, proofs, and simulation results were verified independently. The RATCHET framework is open-source and available at <https://github.com/CIRISAI/RATCHET>.

A Notation Reference

Symbol	Meaning
k	Raw constraint count
k_{eff}	Effective constraint count
ρ	Pairwise constraint correlation
ρ_{crit}	Critical correlation threshold
λ	Decay constant / strictness
σ	Sustainability integral
α	Constraint generation rate
d	Decay rate
J	Defense function
R_c	Collapse rate (dJ/dt)
T_{truth}	Time to deception collapse
T_{entropy}	Time to entropic decay
T_{capture}	Time to coordination capture
T_{eff}	Effective collapse time
K_{req}	Required effective constraints

Table 9: Notation reference for Coherence Collapse Analysis.

B Key Formulas

$$k_{\text{eff}} = \frac{k}{1 + \rho(k - 1)} \quad (\text{Effective Constraints})$$

$$K_{\text{req}} = \frac{-\ln(\varepsilon/V_0)}{\lambda} \quad (\text{Required Constraints})$$

$$\rho_{\text{crit}} = \frac{1}{K_{\text{req}}} \quad (\text{Critical Correlation})$$

$$J = k_{\text{eff}} \cdot (1 - \rho) \cdot \lambda \cdot \sigma \quad (\text{Defense Function})$$

$$T_{\text{truth}} = \frac{K_{\text{req}}(1 - \rho)}{\alpha(1 - K_{\text{req}} \cdot \rho)} \quad (\text{Time to Truth})$$

$$T_{\text{entropy}} = \frac{\ln(\sigma/\sigma_{\min})}{d} \quad (\text{Time to Entropy})$$

$$T_{\text{capture}} = \frac{(n/3) - f}{r_{\text{cap}}} \quad (\text{Time to Capture})$$

$$P(\text{detect}) = 1 - e^{-10\varepsilon} \quad (\text{Detection Probability})$$