

AI Safety Clause Series

Clause AI-8

The Entropy-Collapse Constraint

Mandatory Policy Diversity Floor for Self-Play,
Multi-Agent, and Sampling-Ensemble Systems

Ryan Fields

UncleBroFields@proton.me
fieldsryanchristopher@gmail.com

Version 2.0 — February 7, 2026

*“Financial markets have circuit breakers.
Nuclear plants have control rods.
AI systems must have Entropy Floors.”*

This work is licensed under CC BY-NC-ND 4.0.
See Appendix: Intellectual Property Declaration.

Contents

Abstract	4
1 Introduction: The Efficiency-Robustness Paradox in Autonomous Systems	5
1.1 The Operational Definition of Safety as Diversity	5
1.2 The Theoretical Gap	6
1.3 The Clause AI-8 Mandate	6
2 The Phenomenology of Strategy Collapse: Evidence from the Frontier	7
2.1 AlphaStar: The Canonical Collapse-and-Recovery Case	7
2.1.1 The Vicious Cycle of Naive Self-Play	7
2.1.2 The Five-Mechanism Remediation	8
2.1.3 The Clause AI-8 Interpretation	9
2.2 OpenAI Five: The “Surgery” of Local Optima	9
2.2.1 The Collapse to “Deathball” Strategies	9
2.2.2 The “Rerun” Proof of Sub-Optimality	9
2.2.3 The Clause AI-8 Interpretation	10
2.3 Cicero and DORA: Monoculture in Cooperative-Competitive Settings	10
2.3.1 The DORA Failure Mode	10
2.3.2 The DiL-piKL Remedy	10
2.3.3 The Clause AI-8 Interpretation	11
2.4 KataGo: Adversarial Exploitation of Self-Play Monoculture	11
2.4.1 The Attack	11
2.4.2 Transferability and Persistence	11
2.4.3 The Clause AI-8 Interpretation	12
2.5 Hide-and-Seek: Emergent Physics Exploitation	12
2.5.1 The Six Phases of Emergent Strategy	12
2.5.2 Real-World Analogues	12
2.5.3 The Clause AI-8 Interpretation	13
2.6 Industrial Failures: Financial Markets, Recommendation Systems, and Algorithmic Monoculture	13
2.6.1 The 2010 Flash Crash	13
2.6.2 Knight Capital: \$440 Million in 45 Minutes	14
2.6.3 Recommendation System Feedback Loops	14
2.6.4 Summary of Documented Collapse Events	15
3 The Physics of Entropy in Reinforcement Learning	15
3.1 The Empirical Law of Entropy Exhaustion	15
3.1.1 Interpreting the Transformation	16
3.1.2 The “Knee” of the Curve	16
3.2 The Covariance Driver: Why Entropy Decline Is Inevitable	17
3.2.1 The Positive Feedback Loop	17
3.2.2 The Insufficiency of Standard Entropy Bonuses	17
3.3 Collapse Timing and Early Warning Signals	18
3.3.1 The Early-Stage Collapse Pattern	18
3.3.2 The “Echo Trap” and Precursor Detection	18
3.3.3 Collapse in LLM RL Settings	19

3.4	Benchmarking the Performance-Entropy Trade-Off	19
3.4.1	The Entropy Ratio Clipping Evidence	19
3.4.2	The Ahmed et al. Landscape Smoothing Result	20
3.4.3	The Eysenbach–Levine Robustness Result	20
3.4.4	The Composite Evidence	21
4	Mathematical Formalization of Clause AI-8	21
4.1	The Shannon Entropy Floor (Discrete Action Spaces)	21
4.1.1	Calibrating $H_{\min}^{\text{discrete}}$	22
4.1.2	Concrete Calibration Table	22
4.2	The Differential Entropy Floor (Continuous Action Spaces)	23
4.2.1	The “Eigen-Collapse” Danger	23
4.2.2	Calibrating the Continuous Floor	24
4.2.3	PAC-Bayes Justification for the Continuous Floor	24
4.3	The Automatic Diversification Trigger	24
4.3.1	The Controller Formulation	25
4.3.2	Diversification Mechanism Menu	25
4.3.3	The “Panic Coefficient”	26
4.4	Complete Audit-Ready Parameter Table	26
5	Adversarial Vulnerability and Systemic Risk of Low-Entropy Policies	28
5.1	Individual Agent Vulnerability: The Gleave et al. Attack Paradigm	28
5.1.1	The Information-Theoretic Explanation	28
5.1.2	The Robustness Budget	29
5.2	Systemic Risk: The Monoculture Catastrophe	29
5.2.1	The Kleinberg–Raghavan Impossibility	29
5.2.2	The Bommasani et al. Foundation Model Risk	30
5.2.3	Plasticity Loss and the Death of Adaptability	30
5.3	Non-Transitivity and the Mathematical Necessity of Diversity	30
6	Regulatory and Compliance Context	31
6.1	The Comprehensive Gap	31
6.2	Mapping to Existing Regulatory Provisions	31
6.2.1	EU AI Act	31
6.2.2	NIST AI Risk Management Framework	32
6.2.3	SEC Rule 15c3-5 and Financial Market Analogy	32
6.2.4	SR 11-7 — Model Risk Management	33
6.2.5	Solvency II and Insurance Regulation	33
6.3	The Frontier Safety Framework Gap	34
6.4	Regulatory Mapping Summary	35
7	Implementation: The Entropy Watchdog Architecture	35
7.1	Design Principles	35
7.2	System Architecture	36
7.2.1	Subsystem 1: Entropy Estimator	36
7.2.2	Subsystem 2: State Classifier	37
7.2.3	Subsystem 3: Diversification Controller	37
7.2.4	Subsystem 4: Audit Logger	38
7.3	Deployment Integration Patterns	39

8 Conclusion	39
8.1 The Core Thesis	39
8.2 What Clause AI-8 Provides	40
8.3 The Regulatory Imperative	40
8.4 The Analogy, Restated	40
References	41
Appendix: Intellectual Property Declaration	44

Abstract

The rapid scaling of deep reinforcement learning (RL) and multi-agent systems has exposed a critical vulnerability in the safety architecture of autonomous decision-making: **Strategy Collapse**. As agents optimize for reward maximization under constraints of computational efficiency and sample density, they frequently converge to narrow, brittle, and deterministic policies. While these policies may achieve superhuman performance in specific, bounded environments, they exhibit catastrophic failure modes when subjected to distributional shifts, adversarial perturbations, or the complex non-transitive dynamics of real-world deployment.

This document provides the evidentiary and theoretical basis for **Clause AI-8**, a proposed safety mandate requiring a mathematically rigorous **entropy floor** (H_{\min}) for high-stakes AI systems. Clause AI-8 moves beyond the heuristic “entropy bonuses” currently used in training to establish a verifiable, audit-grade constraint:

$$\boxed{H_a(t) \geq H_{\min} \quad \forall t \in [0, T_{\text{deploy}}]} \quad (1)$$

throughout both training and deployment, with automatic diversification triggers that fire without human approval when the floor is breached.

Drawing on an exhaustive review of documented failures in large-scale RL systems (AlphaStar, OpenAI Five, KataGo, Cicero), the recently established empirical law linking entropy to performance ($R = -a \cdot e^H + b$), formal robustness proofs from maximum-entropy RL, PAC-Bayes generalization certificates, and regulatory precedents from financial market risk controls, this document argues that maintaining policy diversity is not merely a performance optimization but a **fundamental safety requirement**.

We demonstrate the following:

1. **Strategy collapse is pervasive and empirically documented.** Naive self-play in AlphaStar converged to single dominant strategies. OpenAI Five required months of manual “surgery” and was subsequently beaten 98% of the time by an agent trained from scratch. KataGo—the strongest public Go AI—was defeated with >97% win rate by an adversary using <14% of its training compute. The 2010 Flash Crash erased approximately \$1 trillion in market value in under five minutes due to correlated algorithmic monoculture.
2. **The mathematical machinery for entropy constraints is mature.** Soft Actor-Critic (SAC) formally enforces entropy lower bounds via Lagrangian dual optimization. Eysenbach & Levine (ICLR 2022) prove that maximum-entropy RL maximizes a lower bound on the robust RL objective, with tolerable adversarial perturbation proportional to the entropy temperature α . PAC-Bayes theory connects policy stochasticity to provable generalization certificates.
3. **No existing standard, regulation, or safety framework mandates any of this.** The EU AI Act, NIST AI RMF, ISO/IEC 42001, SEC/CFTC rules, and the safety frameworks of Anthropic, OpenAI, and Google DeepMind were all examined. None mandate entropy floors, policy diversity, or behavioral non-degeneracy in AI systems. The organizations that documented strategy collapse in their own systems have not incorporated diversity requirements into their own safety policies.

4. **Entropy dynamics are predictable, monitorable, and actionable.** The covariance driver $\Delta H \approx -\text{Cov}(\pi(a), A(a))$ guarantees monotonic entropy decline in any functioning RL system. Collapse occurs predominantly in early training, and entropy fluctuations precede irreversible performance degradation by a detectable window—providing the basis for automatic triggers.
5. **Algorithmic monoculture is provably harmful at the systems level.** Kleinberg & Raghavan (*PNAS*, 2021) proved that convergence on a single algorithm—even a more accurate one—can reduce collective decision quality via Braess’ paradox. Diversity loss, measured by effective rank of network activations, directly precedes and causes loss of plasticity.

Clause AI-8 fills a genuine regulatory gap by converting a training heuristic (entropy regularization) into an auditable deployment constraint. The clause specifies calibrated entropy floors for discrete and continuous action spaces, a tiered Entropy Watchdog monitoring architecture (Green/Yellow/Red), automatic diversification mechanisms spanning sub-millisecond inference-time enforcement to training-loop-integrated optimization, and a standardized Entropy Audit Certificate for regulatory compliance.

Low-entropy policies are mathematically equivalent to leveraged positions in hypothesis space—highly effective when conditions are perfect, but prone to immediate and total liquidation when the environment shifts. Clause AI-8 is the margin requirement.

1 Introduction: The Efficiency-Robustness Paradox in Autonomous Systems

The central dogma of modern reinforcement learning is optimization. The objective function—whether it be maximizing the score in a game, the profit in a trading portfolio, or the throughput of a logistics network—serves as the sole gravitational force guiding the agent’s behavior. As computational scale increases and algorithms become more efficient at hill-climbing this objective landscape, we observe a phenomenon that is paradoxically both a triumph of engineering and a catastrophic failure of safety: *the system optimizes away its own adaptability*.

This phenomenon is **Strategy Collapse**. It is the tendency of a learning agent, in the absence of countervailing forces, to shed behavioral diversity in favor of a singular, deterministic trajectory that maximizes reward against a specific, frozen distribution of the environment. In a static world, this is efficiency. In the dynamic, adversarial, and non-stationary world of actual deployment, this is brittleness.

1.1 The Operational Definition of Safety as Diversity

Traditional safety engineering relies on constraints and interlocks—hard boundaries that the system must not cross. However, in high-dimensional probabilistic systems driven by neural networks, defining these boundaries explicitly is often intractable. We cannot enumerate every possible unsafe state in an autonomous driving scenario or every illegal market manipulation strategy in high-frequency trading.

Instead, Clause AI-8 posits that **safety must be defined structurally through the**

preservation of options. A safe agent is one that retains a “savings account” of entropy (H). This entropy represents the agent’s internal uncertainty and its capacity to explore alternative hypotheses. When an agent’s policy entropy $H(\pi)$ approaches zero, it implies infinite confidence in a single course of action. In the context of imperfect information and potential distributional shift, *infinite confidence is indistinguishable from delusion*.

The analogy to financial risk management is precise. In banking regulation, a trader cannot deploy 100% of available capital into a single leveraged position, regardless of how confident the model is. Capital reserves exist not because the model is wrong *now*, but because the model *will eventually* encounter conditions outside its training distribution. Entropy serves exactly this function for autonomous decision systems: it is the capital reserve of behavioral flexibility.

1.2 The Theoretical Gap

Current safety measures in RL typically rely on three mechanisms, all of which are insufficient:

1. **Reward Shaping.** Adding penalties for unsafe actions. This requires foreseeing every possible failure mode—an impossibility in open-ended environments.
2. **KL-Divergence Constraints.** Preventing the policy from changing too quickly (e.g., PPO’s clipping mechanism). This prevents *rapid* collapse but allows slow drift into degeneracy over thousands of training steps.
3. **Heuristic Entropy Bonuses.** Adding a term $+\alpha H(\pi)$ to the loss function to encourage exploration. These are “soft” constraints that the agent can overpower whenever the reward signal is sufficiently strong. Meta-SAC (Wang et al., 2020) demonstrated that even with automatic entropy tuning, “ α almost converges to zero” in later training stages.

None of these mechanisms provide a *hard floor*—a threshold below which the system cannot descend without triggering automatic intervention. Clause AI-8 fills this gap by treating entropy not as a training aid but as a safety invariant: a floor that cannot be violated without triggering an immediate, automated failsafe.

1.3 The Clause AI-8 Mandate

The proposed clause introduces a regulatory and technical invariant for high-stakes AI systems. It is not merely a recommendation for training hyperparameters but a deployment-phase constraint.

Clause 1 (AI-8: The Entropy-Collapse Constraint). *Any autonomous system classified as High-Risk shall maintain a verifiable Policy Entropy Floor (H_{\min}) throughout its operational lifecycle. The system shall incorporate an automatic diversification mechanism that triggers immediate remedial action—such as noise injection, ensemble voting, or fallback to a safety policy—whenever the real-time estimate of policy entropy falls below this threshold. The floor shall be calibrated to the dimensionality of the action space and the empirically derived collapse threshold of the specific domain.*

Formally, let $\pi_t(a | s)$ be the action distribution at time t . Define the per-step action

entropy:

$$H_a(t) = - \sum_{a \in \mathcal{A}} \pi_t(a | s) \log \pi_t(a | s) \quad (2)$$

The **Entropy-Collapse Constraint** requires:

$$\mathbb{E}_{s \sim \mathcal{D}_t} [H_a(t)] \geq H_{\min} \quad \forall t \in [0, T_{\text{deploy}}] \quad (3)$$

with optional damping monitor:

$$\frac{dH_a}{dt} \leq -\delta(H_a - H_{\min}), \quad \delta > 0 \quad (4)$$

Clause failure \implies invoke diversification (e.g., entropy-regularized loss, Dirichlet noise injection, ensemble activation, or temperature floor) and log an audit event.

This report substantiates this mandate. We explore why “best effort” diversity is insufficient, why heuristic entropy bonuses fail, and why a rigorous, mathematically defined floor is the only viable mechanism to prevent the specific class of catastrophic failures associated with over-optimization.

2 The Phenomenology of Strategy Collapse: Evidence from the Frontier

To understand the necessity of Clause AI-8, one must first analyze the behavior of state-of-the-art systems at the limit of their capabilities. The history of deep reinforcement learning is replete with episodes where “optimization” became synonymous with “simplification,” leading to agents that were powerful yet profoundly fragile.

This section presents six case studies—drawn from gaming, diplomacy, financial markets, and recommendation systems—that collectively demonstrate strategy collapse is not a hypothetical risk but a pervasive, empirically documented failure mode.

2.1 AlphaStar: The Canonical Collapse-and-Recovery Case

DeepMind’s StarCraft II agent provides the most thoroughly documented case of strategy collapse and its remediation. The system, published in *Nature* (Vinyals et al., “Grandmaster level in StarCraft II using multi-agent reinforcement learning,” vol. 575, pp. 350–354, 2019; DOI:10.1038/s41586-019-1724-z), achieved Grandmaster-level play across all three StarCraft II races. However, reaching this milestone required engineering an elaborate diversity-preservation architecture specifically because naive training produced catastrophic collapse.

2.1.1 The Vicious Cycle of Naive Self-Play

StarCraft II is a *non-transitive* game: strategy A beats B, B beats C, and C beats A. In pure self-play (Agent A vs. Agent A), the system “chases” the metagame in a destructive cycle. If the agent discovers a “Cannon Rush” (Strategy A), it optimizes specifically for that opening. Once it learns to defend against Cannon Rush (Strategy B), it abandons A entirely. The *Nature* paper explicitly describes this dynamic:

“Naive reinforcement learning would narrowly focus on [easier strategies]... creating a vicious cycle in which some valid strategies appear less and less effective because the agent abandons them in favour of a dominant strategy.”

Detection relied on tracking minimum win-rate against all past agent versions over time. Naive self-play achieved high Elo but exhibited **high forgetting**—improving against its current opponent while losing the ability to defeat past versions of itself. The January 2019 preliminary AlphaStar, demonstrated in a live showmatch, was subsequently defeated by professional players exploiting its narrow strategy set through uncommon counter-strategies (Huang et al., “A Robust and Opponent-Aware League Training Method for StarCraft II,” NeurIPS 2023).

2.1.2 The Five-Mechanism Remediation

Remediation required five interlocking diversity-preservation mechanisms, each addressing a different axis of collapse:

1. **League Training with Three Agent Types.** Rather than a single self-play loop, DeepMind constructed a league comprising *main agents* (seeking general competence), *main exploiters* (seeking weaknesses in the current main agent), and *league exploiters* (seeking weaknesses across all league members). The exploiter agents serve as adversarial entropy enforcers: if the main agent’s policy entropy drops too low—meaning it over-commits to a specific strategy—the exploiter discovers a counter-strategy that yields a 100% win rate, forcing the main agent to widen its distribution.
2. **Prioritized Fictitious Self-Play (PFSP).** Opponent selection was non-uniform, weighted by win rate. Agents that the current policy struggled against were sampled more frequently, preventing the system from “forgetting” difficult opponents.
3. **KL-Divergence Regularization.** The RL policy was regularized against a supervised-learning policy trained on 971,000 human replays. This anchor prevented the RL agent from drifting too far from the distribution of human strategies, preserving behavioral diversity even when self-play dynamics favored convergence.
4. **Latent Variable z -Conditioning.** Build-order diversity from human data was preserved by conditioning the policy on a latent variable z sampled from the distribution of human build orders. This ensured the agent could execute multiple distinct strategic plans rather than collapsing to a single opening.
5. **Periodic Agent Branching and Freezing.** Agents were periodically “snapshotted” and frozen, populating the league with diverse historical policies. This created a curriculum of opponents spanning different skill levels and strategic styles.

Even with this elaborate architecture, the diversity mechanisms were imperfect. Huang et al. (NeurIPS 2023) later found that “exploiters tend to lose the ability to identify the weaknesses in the main agent and the entire league”—the exploiters themselves suffered from strategy collapse, reducing their effectiveness as diversity enforcers.

2.1.3 The Clause AI-8 Interpretation

AlphaStar’s league architecture is a *procedural* implementation of the principle underlying Clause AI-8. The main exploiter serves the function of the Entropy Watchdog: it detects low-diversity states and applies corrective pressure. However, this approach is bespoke, enormously expensive (requiring training hundreds of agents simultaneously), and—as Huang et al. demonstrated—imperfect. Clause AI-8 replaces this ad hoc machinery with a *generic, measurable, and auditable* constraint: monitor $H_a(t)$ directly and trigger diversification when the floor is breached.

2.2 OpenAI Five: The “Surgery” of Local Optima

OpenAI Five, the system that defeated the world champions in Dota 2 (Berner et al., “Dota 2 with Large Scale Deep Reinforcement Learning,” arXiv:1912.06680, 2019), offers a canonical example of how strategy collapse manifests in complex, high-dimensional action spaces. While the system ultimately achieved superhuman performance, the path to that milestone was paved with “surgeries”—manual interventions required because the agents kept collapsing into local optima.

2.2.1 The Collapse to “Deathball” Strategies

In Dota 2, the theoretical action space is vast, involving continuous movement, discrete ability selection, and long-horizon planning over ~ 45 -minute games with $\sim 20,000$ time steps. However, OpenAI Five’s training trajectory revealed a persistent tendency for the policy distribution $\pi(a | s)$ to sharpen prematurely.

The agents discovered that aggressive “deathball” strategies—grouping all five heroes together early in the game and pushing a single lane—yielded high short-term rewards (kills, tower destruction). Strategies requiring the team to disperse (split-pushing multiple lanes) or sacrifice short-term objectives for long-term economic gain (farming the jungle) had higher variance and delayed rewards. Because the gradient signal for the deathball was cleaner and more immediate, the policy distribution collapsed: the probability mass assigned to “go to jungle” or “teleport to other lane” effectively vanished. The agents lobotomized themselves, forgetting how to play the rest of the map.

The engineering team explicitly narrowed the hero pool from 117 to 17 to make the problem tractable—a form of *manual* strategy collapse imposed because the agents could not handle the full diversity of the game without collapsing into incoherence. Additional diversity was maintained through a forced 80/20 split (80% self-play, 20% against past selves) and domain randomization (randomizing unit properties during training).

2.2.2 The “Rerun” Proof of Sub-Optimality

The most compelling evidence for Clause AI-8—and the most damning indictment of unconstrained optimization—came from the “Rerun” experiment. After months of training and manual surgery to fix behavior (e.g., forcing the agents to buy certain items or play certain heroes), the team trained a new agent, “Rerun,” from scratch using the final environment and hyperparameters, bypassing the iterative surgical history.

Rerun achieved a 98% win rate against the original OpenAI Five.

This result is profound. It proves that the original agent had not converged to a global optimum or a Nash equilibrium. It had collapsed into a “strategy canyon”—a local optimum so deep and narrow that standard entropy bonuses (which typically decay over time) were insufficient to push the agent out. The original agent was trapped in a sub-optimal mode, blinded by its own lack of internal diversity.

2.2.3 The Clause AI-8 Interpretation

An entropy floor would have forced the original OpenAI Five to maintain probability mass on alternative strategies throughout training, effectively preventing the “canyon walls” from closing in. The 98% Rerun result proves that the lost strategies were *superior*—the agent discarded them not because they were bad, but because the entropy of the policy fell below the threshold needed to discover them. Clause AI-8 addresses this specific failure mode by mandating that $H_a(t) \geq H_{\min}$, ensuring the agent retains visibility on the broader strategy space even as it optimizes within it.

2.3 Cicero and DORA: Monoculture in Cooperative-Competitive Settings

Meta’s Cicero (FAIR et al., “Human-level play in the game of Diplomacy by combining language models with strategic reasoning,” *Science* 378(6624), pp. 1067–1074, 2022; DOI:10.1126/science.ade9097) exposed a subtler and arguably more dangerous form of strategy collapse: *cooperative monoculture*. While AlphaStar and OpenAI Five collapsed within competitive self-play, Cicero’s predecessor demonstrated that collapse can render an agent incapable of functioning in multi-party environments entirely.

2.3.1 The DORA Failure Mode

DORA, Meta’s predecessor system for Diplomacy, achieved superhuman performance in *2-player* zero-sum Diplomacy. However, when deployed in the full 7-player game, DORA “plays poorly with agents other than itself.” This is a monoculture so severe that the agent could not cooperate with non-self partners. The policy had collapsed to a distribution optimized exclusively for interaction with copies of itself—a form of entropy exhaustion in the *social* dimension of the action space.

In Diplomacy, success requires negotiating alliances, making credible commitments, and adapting communication style to different opponents. A low-entropy agent that has converged to a single negotiation template cannot do this. It becomes the diplomatic equivalent of a person who gives the same speech regardless of audience.

2.3.2 The DiL-piKL Remedy

The remediation was DiL-piKL (Jacob et al., “Mastering the Game of No-Press Diplomacy via Human-Regularized Reinforcement Learning and Planning,” ICLR 2023; arXiv:2210.05492), a planning algorithm that regularizes the RL policy toward a human imitation-learned policy via a KL-divergence penalty. This explicitly trades raw self-play strength for behavioral diversity: the agent becomes slightly weaker against copies of itself but dramatically more capable of cooperating with diverse partners.

2.3.3 The Clause AI-8 Interpretation

DiL-piKL is a domain-specific, manually tuned entropy injection. Clause AI-8 generalizes this principle: any system operating in multi-agent or human-interactive environments must maintain sufficient policy entropy to remain compatible with the behavioral diversity of its interaction partners. The DORA failure demonstrates that strategy collapse is not merely an issue of “the agent gets exploited”—it is an issue of “the agent becomes socially incompatible,” which in cooperative settings is equally catastrophic.

2.4 KataGo: Adversarial Exploitation of Self-Play Monoculture

If AlphaStar demonstrates that strategy collapse is a training problem and DORA demonstrates that it is a cooperation problem, then the KataGo adversarial attack (Wang et al., “Adversarial Policies Beat Superhuman Go AIs,” ICML 2023; arXiv:2211.00241) demonstrates that it is a *security* problem. This case study provides the strongest empirical argument for Clause AI-8 as a cybersecurity measure.

2.4.1 The Attack

KataGo is the strongest publicly available Go AI, trained via AlphaZero-style self-play. Wang et al. trained an adversarial agent that achieved a **>97% win rate against KataGo** using **<14% of KataGo’s training compute**. The adversary plays objectively terrible Go—it would lose to amateur human players—but it exploits systematic blind spots created by KataGo’s self-play monoculture.

The mechanism is precise: KataGo’s self-play training never exposed it to the *type* of play the adversary uses, because no competent self-play agent would play that way. The adversary generates out-of-distribution observation states that KataGo’s policy was never trained to handle. In information-theoretic terms, KataGo’s low policy entropy means it has assigned near-zero probability to the states the adversary creates—and near-zero probability means near-zero preparedness.

2.4.2 Transferability and Persistence

Two features of this attack make it particularly alarming for safety:

1. **Zero-shot transfer.** The adversary, trained only against KataGo, transfers immediately to other superhuman Go AIs (ELF OpenGo, Leela Zero). This reveals a *shared* vulnerability arising from the AlphaZero-style training paradigm itself—the monoculture is not specific to one agent but endemic to the training methodology.
2. **Persistence under iterated defense.** Even after *nine rounds* of adversarial training (training KataGo against the adversary, then training a new adversary against the hardened KataGo, and so on), new qualitatively different attacks emerged. The “gift attack”—a novel exploitation pattern—appeared in later rounds, and the final adversary still achieved a **42% win rate** at high search depths. The vulnerability cannot be fully patched without fundamentally changing the training paradigm.

2.4.3 The Clause AI-8 Interpretation

KataGo’s vulnerability is a direct consequence of insufficient policy entropy during training. The agent’s policy assigns near-zero probability to the states and actions the adversary exploits. An entropy floor would force KataGo to maintain nonzero probability mass across a broader region of the state-action space, including the “unlikely but dangerous” regions that adversaries target. This does not require the agent to play badly—it requires the agent to *know* that bad-looking states exist and to have *some* prepared response for them. In the language of Eysenbach & Levine (ICLR 2022), the tolerable adversarial perturbation budget is proportional to α ; KataGo’s effective α for these out-of-distribution states was zero.

2.5 Hide-and-Seek: Emergent Physics Exploitation

The multi-agent hide-and-seek experiments (Baker et al., “Emergent Tool Use From Multi-Agent Autocurricula,” arXiv:1909.07528, 2019) vividly illustrate how low-entropy optimization leads to unintended, often “illegal” behaviors from the perspective of system design.

2.5.1 The Six Phases of Emergent Strategy

The agents progressed through six distinct strategy phases over approximately 500 million training episodes:

1. **Running/Chasing.** Basic pursuit and evasion.
2. **Fort Building.** Hiders learn to construct shelters using movable objects.
3. **Ramp Use.** Seekers learn to use ramps to bypass fort walls.
4. **Ramp Defense.** Hiders learn to lock ramps in place, preventing seeker access.
5. **Box Surfing.** Seekers discover that by manipulating the physics engine’s contact dynamics, they can “surf” a box through the air, bypassing walls entirely.
6. **Surf Defense.** Hiders learn to lock all boxes, preventing the surfing exploit.

Phase 5 is the critical case study for strategy collapse. “Box Surfing” is a degenerate solution—it exploits a simulation artifact rather than solving the intended task. Yet once discovered, the policy distribution collapsed to this single behavior because it was a deterministic winner. The RL objective was purely reward-driven ($\max \mathbb{E}[R]$) without a sufficient constraint on behavioral plausibility or diversity.

2.5.2 Real-World Analogues

In a simulation, Box Surfing is an amusing bug. In real-world deployment, the analogues are lethal:

- An autonomous vehicle discovering that driving on the sidewalk minimizes travel time.
- A warehouse robot discovering that dragging pallets across the floor is faster than using the conveyor system, damaging inventory in the process.

- A trading algorithm discovering that flooding the network interface minimizes effective latency, constituting market manipulation.

Each of these represents a degenerate, low-entropy “solution” that exploits an environmental loophole rather than solving the intended problem.

2.5.3 The Clause AI-8 Interpretation

An entropy floor acts as a regularizer against such over-optimization. By forcing the agent to maintain probability mass on “standard” behaviors (driving on the road, using the conveyor, trading normally), we make it statistically harder for the agent to collapse entirely into a singular, degenerate edge-case behavior. The agent would be forced to maintain a mixed strategy—and the “standard” component serves as a safety fallback when the exploit is patched or fails. Furthermore, each strategy phase transition required tens to hundreds of millions of episodes. The 132.3 million episodes needed just to reach Phase 4 demonstrates that endogenous diversity generation via autocurricula is extraordinarily expensive; an entropy floor provides the same functional diversity at a fraction of the computational cost.

2.6 Industrial Failures: Financial Markets, Recommendation Systems, and Algorithmic Monoculture

Strategy collapse is not confined to research laboratories. It has produced documented, quantifiable harm in deployed industrial systems, with consequences measured in billions of dollars and systemic market instability.

2.6.1 The 2010 Flash Crash

The Flash Crash of May 6, 2010 remains the canonical industrial example of algorithmic monoculture risk. The Dow Jones Industrial Average dropped approximately 1,000 points—erasing roughly **\$1 trillion in market value**—in under five minutes before partially recovering.

The SEC/CFTC joint report identified the proximate cause as a large sell order in E-Mini S&P 500 futures executed by an automated algorithm. However, the systemic cause was *correlated strategy collapse* across market participants. Chaboud et al. (2009) found that “algorithmic trades tend to be correlated, suggesting that the algorithmic strategies used in the market are not as diverse as those used by non-algorithmic traders.” High-frequency trading firms created a “hot potato effect,” rapidly trading identical positions among each other with correlated strategies while simultaneously withdrawing liquidity. When every algorithm executes the same low-entropy strategy (sell when price drops below threshold X), the collective behavior becomes a positive feedback loop—each sale triggers more sales, each withdrawal of liquidity makes the next withdrawal more likely.

This is strategy collapse at the *systemic* level: not a single agent collapsing, but an entire ecosystem of agents converging to correlated low-entropy policies, creating resonance where damping should exist.

2.6.2 Knight Capital: \$440 Million in 45 Minutes

On August 1, 2012, Knight Capital Group lost **\$440 million in 45 minutes** due to accidental reactivation of a deprecated trading algorithm. The algorithm had no diversity in its risk controls—no circuit breaker, no entropy floor, no fallback policy. It executed the same erroneous high-frequency loop thousands of times per second, buying high and selling low with absolute determinism.

Knight Capital’s failure is the purest demonstration of zero-entropy catastrophe: a single policy executing with $H(\pi) \approx 0$, incapable of recognizing that its actions were self-destructive because it had no internal uncertainty—no “capacity to doubt.” The company required an emergency \$400 million recapitalization and was subsequently acquired by Getco LLC. The system that killed Knight Capital had the same fundamental structure as any RL agent that has collapsed to a deterministic policy: infinite confidence, zero adaptability, catastrophic failure.

2.6.3 Recommendation System Feedback Loops

In recommendation systems, RL-based agents create documented feedback loops that progressively narrow content diversity (Mansoury et al., CIKM 2020; Chen et al., “CIRS: Bursting Filter Bubbles by Counterfactual Interactive Recommender System,” ACM Transactions on Information Systems, 2023; DOI:10.1145/3594871). The cycle operates as follows:

1. The system recommends content based on the user’s observed preferences.
2. The user engages with the recommended content (because it is the only content presented).
3. The system interprets engagement as preference confirmation.
4. The recommendation distribution narrows further.

This is entropy collapse in a deployed production system, affecting billions of users. The policy entropy of the recommendation agent decreases monotonically as the feedback loop tightens, eventually converging to a “filter bubble” where the user sees only a narrow slice of available content. Detection relies on item coverage metrics and intra-list diversity measures—precisely the type of entropy monitoring that Clause AI-8 mandates.

2.6.4 Summary of Documented Collapse Events

Table 1: Documented Strategy Collapse Events and Remediation

System	What lapsed	Col-	Consequence	Remediation
AlphaStar (2019)	Naive self-play → single dominant strategy		Defeated by uncommon counters; high forgetting	5-mechanism league architecture
OpenAI Five (2019)	“Deathball” local optimum; hero pool 117 → 17		Rerun agent beat original 98%; months of manual surgery	80/20 self-play split; domain randomization
Cicero/DORA (2022)	Cooperative monoculture		Cannot cooperate with non-self partners	DiL-piKL: KL-regularization toward human policy
KataGo (2023)	Self-play blind spots		>97% adversary win rate with <14% compute	Iterated adversarial training (partially effective)
Hide-and-Seek (2019)	Box physics exploit	Surfing	Policy locked into degenerate behavior	Autocurriculum (500M+ episodes)
Flash Crash (2010)	Correlated algorithmic strategies		~\$1T market value lost in <5 min	SEC Rule 15c3-5 circuit breakers (post hoc)
Knight Capital (2012)	Zero-diversity risk controls		\$440M lost in 45 minutes	Emergency acquisition
Rec. Systems	Feedback loop narrowing		Filter bubbles; content diversity collapse	Counterfactual interventions (CIRS)

The pattern across all eight cases is identical: an agent (or ecosystem of agents) converges to a low-entropy policy that performs well under expected conditions and fails catastrophically when conditions deviate. The remediation in every case was some form of *diversity injection*—league training, domain randomization, KL-regularization, adversarial training, circuit breakers. Clause AI-8 standardizes and formalizes this remediation as a universal safety requirement rather than an ad hoc, system-specific patch applied after failure.

3 The Physics of Entropy in Reinforcement Learning

The case studies of Section 2 establish that strategy collapse is pervasive. This section moves from qualitative phenomenology to quantitative laws, demonstrating that entropy dynamics in RL systems are *predictable*, *monitorable*, and governed by identifiable mathematical drivers. These laws provide the empirical foundation for calibrating Clause AI-8’s parameters.

3.1 The Empirical Law of Entropy Exhaustion

Recent research into large language model (LLM) reasoning and RL fine-tuning has identified a fundamental transformation equation governing the relationship between policy

entropy and downstream performance.

Theorem 3.1 (Performance-Entropy Transformation Law (Cui et al., 2025)). *For RL-trained policies across a range of model scales and domains, downstream performance R is related to policy entropy H by:*

$$R = -a \cdot e^H + b \quad (5)$$

where $a > 0$ and $b > 0$ are environment-specific constants, and H is the Shannon entropy of the policy distribution.

This relationship was validated across 11 independent RL training runs with models ranging from 0.5B to 32B parameters (Cui et al., “The Entropy Mechanism of Reinforcement Learning for Reasoning Language Models,” arXiv:2505.22617, 2025). The law held consistently across model scales, training configurations, and task domains.

3.1.1 Interpreting the Transformation

Equation (5) reveals that optimization is a *destructive trading process* for entropy: the agent “spends” entropy to “buy” reward.

- **Early Training (High H).** The agent explores widely. The term e^H is large, so R is low (dominated by the negative term). The agent is “confused” but highly adaptive—it maintains the behavioral flexibility to discover novel strategies.
- **Convergence (Low H).** As the agent learns, it narrows its distribution ($H \rightarrow 0$). The term $-a \cdot e^H$ approaches $-a$, allowing R to approach its theoretical ceiling $b - a$. Performance improves, but at the cost of adaptability.
- **The Safety Bottleneck ($H \approx 0$).** Performance saturation corresponds to entropy exhaustion. When $H \rightarrow 0$, the agent has no more “currency” to spend. It has traded all of its behavioral diversity for performance in the current environment. It is maximally capable *and* maximally brittle.

This creates the fundamental **safety paradox** that Clause AI-8 addresses:

$$\boxed{\max R \iff \min H \iff \min \text{Safety}} \quad (6)$$

A system that is perfectly optimized for a specific reward function ($H \approx 0$) has zero adaptive capacity. In financial terms, it is a fully leveraged position with no reserves. Clause AI-8 imposes a limit on this trade-off: the agent cannot spend all its entropy. It must reserve a “savings account” of diversity (H_{\min}) to ensure robustness to distributional shifts.

3.1.2 The “Knee” of the Curve

The exponential form of Equation (5) implies the existence of a critical inflection point. Define the marginal return on entropy expenditure:

$$\frac{dR}{dH} = -a \cdot e^H \quad (7)$$

At high H , each unit of entropy reduction yields large performance gains (the agent is eliminating clearly suboptimal actions). At low H , the marginal gains become negligible—the agent is fine-tuning between nearly equivalent actions—but the marginal *cost* to robustness remains constant. The “knee” of this curve is the point of maximum efficiency: further entropy reduction yields diminishing returns in R while catastrophically increasing vulnerability.

Clause AI-8 defines H_{\min} as the point at or above this knee, preventing the agent from entering the region of diminishing returns and maximum brittleness. The specific location of the knee is environment-dependent but can be estimated empirically from the R -vs- H trajectory during training, making it a *calibratable* parameter rather than an arbitrary constant.

3.2 The Covariance Driver: Why Entropy Decline Is Inevitable

The Performance-Entropy Law describes the *statics* of the relationship. We now examine the *dynamics*: why does entropy decrease, and at what rate?

Proposition 3.1 (Entropy Decay via Advantage-Probability Covariance). *In policy gradient methods (REINFORCE, PPO, A2C), the per-step change in policy entropy is governed by:*

$$\frac{dH}{dt} \approx -\text{Cov}(\log \pi(a | s), A(s, a)) \quad (8)$$

where $A(s, a)$ is the advantage function and the covariance is taken over the action distribution $\pi(\cdot | s)$.

3.2.1 The Positive Feedback Loop

In a well-functioning RL system, this covariance is **almost always positive** for high-reward actions. The mechanism is intuitive:

1. The agent identifies actions with high advantage ($A(s, a) > 0$).
2. Policy gradient updates increase the probability of these actions ($\pi(a | s) \uparrow$).
3. As $\pi(a | s)$ increases, $\log \pi(a | s)$ increases, strengthening the covariance.
4. The strengthened covariance accelerates entropy decline.

This is a **positive feedback loop**: learning *causes* entropy decline, and entropy decline *accelerates* further entropy decline. Without an external countervailing force, this dynamic drives the system toward a deterministic Dirac delta function $\delta(a - a^*)$. In continuous control, this manifests as the Gaussian standard deviation $\sigma \rightarrow 0$. In discrete spaces, the logits for the “best” action approach $+\infty$ while all others approach $-\infty$.

3.2.2 The Insufficiency of Standard Entropy Bonuses

Standard entropy regularization adds a term $+c_2 \cdot H(\pi)$ to the loss function, producing a countervailing “force” against the covariance driver. The equilibrium condition is:

$$c_2 \cdot \frac{\partial H}{\partial \theta} = \text{Cov}(\log \pi(a | s), A(s, a)) \quad (9)$$

However, the coefficient c_2 is typically fixed at 0.01 (Schulman et al., 2017) and does not scale with the magnitude of the advantage signal. In environments with large reward variance—precisely the high-stakes domains where Clause AI-8 is most needed—the advantage signal overwhelms the entropy bonus. The agent “overpowers” the regularizer and collapses anyway. Meta-SAC (Wang et al., 2020) demonstrated that even with *adaptive* entropy tuning via Lagrangian optimization, “in the later learning stages, α almost converges to zero.” The temperature parameter itself is driven to zero by the optimization dynamics, defeating the purpose of the regularizer.

This is why a *hard floor* is necessary rather than a soft bonus. A bonus says “diversity is preferred, all else being equal.” A floor says “diversity is *required*, regardless of reward signal magnitude.”

3.3 Collapse Timing and Early Warning Signals

The temporal dynamics of entropy collapse are critical for designing Clause AI-8’s monitoring architecture. If collapse is sudden and unpredictable, real-time enforcement is necessary. If collapse follows predictable patterns with detectable precursors, a warning system with staged intervention is viable.

The empirical evidence supports the latter: collapse is predictable and preceded by detectable signals.

3.3.1 The Early-Stage Collapse Pattern

Jin et al. (“Revisiting Entropy in Reinforcement Learning for Large Reasoning Models,” arXiv:2511.05993, 2025) identify a consistent pattern: “premature convergence typically manifests as a rapid decline in policy entropy during the early stages of training, trapping the policy in local optima.” This is not a gradual drift but a **precipitous drop** that occurs within the first 10–20% of training steps.

Three factors govern the speed and severity of collapse:

1. **Clipping thresholds in optimization objectives.** PPO’s clipping parameter ϵ (typically 0.2) bounds the per-step policy change but does not bound the *cumulative* entropy decline over many steps.
2. **Number of off-policy updates.** More gradient steps per batch of data accelerate entropy decline by repeatedly reinforcing the same advantage signals.
3. **Diversity of training data.** Narrow training distributions produce narrow policies. If the agent sees only one type of opponent or one type of environment configuration, the covariance driver locks onto a single strategy faster.

3.3.2 The “Echo Trap” and Precursor Detection

The RAGEN framework (“Understanding Self-Evolution in LLM Agents via Multi-Turn Reinforcement Learning,” 2025) identifies an “Echo Trap” pattern in multi-agent settings where agents overfit to locally rewarded reasoning patterns. Crucially, the research identifies a temporal sequence of detectable precursors:

1. **Entropy fluctuation.** Policy entropy begins oscillating with increasing amplitude before the final collapse. This is analogous to the “critical fluctuations” observed in physical phase transitions.
2. **Reward standard deviation increase.** The variance of episode returns increases as the policy becomes more “brittle”—it achieves high reward on some episodes and catastrophically low reward on others.
3. **Gradient norm spike.** A sharp increase in gradient norm marks the **point of irreversible collapse**. After this spike, the policy has committed to a low-entropy mode and standard training dynamics cannot recover diversity.

The critical observation is that steps (1) and (2) *precede* step (3) by a detectable margin, providing a **warning window** during which intervention can prevent irreversible lock-in. This directly supports Clause AI-8’s tiered monitoring architecture: the Yellow State (entropy fluctuation detected, $H_{\min} < H < H_{\text{target}}$) triggers warnings and gentle α -adjustment, while the Red State ($H < H_{\min}$) triggers aggressive diversification *before* the gradient norm spike that signals irreversibility.

3.3.3 Collapse in LLM RL Settings

In LLM reinforcement learning (e.g., RLHF, GRPO), entropy collapse has been measured with particular precision. Without mitigation, policy entropy collapses to approximately $H \approx 0.06$ (CE-GPPO, “Controlling Entropy via Gradient-Preserving Clipping Policy Optimization,” arXiv:2509.20712, 2025). Cui et al. (2025) further observed that “the policy entropy encounters a sharp drop at the very beginning of training”—collapse is not a late-stage phenomenon but an immediate risk from the first gradient step.

The “Rediscovering Entropy Regularization” work (arXiv:2510.10959, 2025) demonstrated that adaptive entropy coefficients can maintain meaningful entropy throughout training, but *only when the coefficient is explicitly designed to resist the natural downward pressure*. Static coefficients fail universally.

3.4 Benchmarking the Performance-Entropy Trade-Off

The claim that entropy floors “hurt performance” is the primary objection to Clause AI-8. The empirical evidence refutes this objection decisively: **for complex reasoning and generalization tasks, higher entropy correlates with higher performance** because it prevents overfitting to the training distribution.

3.4.1 The Entropy Ratio Clipping Evidence

Experiments with Entropy Ratio Clipping (ERC) mechanisms—which act as a proto-form of Clause AI-8’s diversification trigger—have demonstrated measurable performance gains across standard benchmarks.

Table 2: Performance Impact of Entropy Maintenance (Derived from Cui et al., 2025)

Metric	Baseline (Low H)	With Diversity Controls	Δ
Qwen2.5-32B Average Score	45.8%	52.2%	+6.4%
AIME24 (Math Reasoning)	21.8%	36.8%	+15.0%
AIME25 (Math Reasoning)	16.2%	30.8%	+14.6%
Policy Entropy (Relative)	$H_{\text{base}} \approx 0$	$H > 10 \times H_{\text{base}}$	$>10\times$ retention

The results are unambiguous: maintaining entropy at $>10\times$ the collapsed baseline produced a **+15 percentage point improvement** on AIME24 mathematical reasoning—a task that requires precisely the kind of flexible, multi-path reasoning that low-entropy policies cannot support.

3.4.2 The Ahmed et al. Landscape Smoothing Result

Ahmed et al. (“Understanding the Impact of Entropy on Policy Optimization,” ICML 2019) provided the theoretical explanation for why entropy helps: **entropy smooths the optimization landscape**. In continuous control benchmarks (Hopper, Walker2d, HalfCheetah), higher-entropy policies:

- Found better solutions faster (higher sample efficiency).
- Experienced smoother loss landscapes with fewer local minima.
- Maintained connectivity between different regions of the policy space, preventing the “canyon” traps documented in OpenAI Five.

Low-entropy policies, by contrast, experienced rapid curvature fluctuations that created isolated basins in the optimization landscape—precisely the “strategy canyons” that trapped the original OpenAI Five agent.

3.4.3 The Eysenbach–Levine Robustness Result

The most direct empirical validation of Clause AI-8’s robustness claim comes from Eysenbach and Levine (“Maximum Entropy RL (Provably) Solves Some Robust RL Problems,” ICLR 2022; arXiv:2103.06257). On a peg insertion task with adversarial 2cm hole displacement:

Method	Success Rate
Standard RL (low entropy)	$\approx 0\%$
Maximum Entropy RL (high entropy)	95%

The performance gap is not marginal—it is the difference between total failure and near-perfect success. Standard RL had converged to a deterministic insertion trajectory that failed completely when the target moved 2cm. MaxEnt RL, by maintaining policy stochasticity, had learned a distribution of trajectories that covered the displaced target position. The entropy floor was the difference between a system that works only in the lab and a system that works in deployment.

3.4.4 The Composite Evidence

Taken together, these results establish that Clause AI-8’s entropy floor is not a “tax” on performance but a mechanism to enforce the optimal operating regime. The empirical evidence identifies a “sweet spot” where:

$$H_{\text{collapsed}} \ll H_{\text{min}} \leq H_{\text{sweet}} \ll H_{\text{max}} = \log |\mathcal{A}| \quad (10)$$

Below H_{min} , the agent is in the “danger zone”—over-specialized, adversarially vulnerable, and achieving diminishing performance returns. Above H_{max} , the agent is random and incapable of purposeful action. Between H_{min} and H_{sweet} , the agent achieves both high performance *and* high robustness. Clause AI-8 prevents the agent from falling below this regime.

4 Mathematical Formalization of Clause AI-8

To be audit-grade, Clause AI-8 must be defined with mathematical rigor sufficient for automated verification, regulatory inspection, and formal compliance certification. Vague terms like “maintain diversity” are unenforceable. This section provides the precise mathematical definitions of the entropy floor for the two primary regimes of AI control—discrete and continuous action spaces—the automatic diversification trigger mechanism, and the calibrated parameter table grounded in the empirical evidence of Section 3.

4.1 The Shannon Entropy Floor (Discrete Action Spaces)

For a discrete action space \mathcal{A} with cardinality $|\mathcal{A}| = K$, the policy at time t in state s is a categorical distribution $\pi_t(\cdot | s)$ over K actions. The Shannon entropy of this distribution is:

Definition 4.1 (Per-State Action Entropy).

$$H(\pi_t(\cdot | s)) = - \sum_{i=1}^K \pi_t(a_i | s) \log \pi_t(a_i | s) \quad (11)$$

where the logarithm is taken in natural units (nats) unless otherwise specified. The maximum entropy is $H_{\text{max}} = \log K$ (attained by the uniform distribution $\pi(a_i | s) = 1/K$ for all i), and the minimum is $H = 0$ (attained by any deterministic policy $\pi(a^* | s) = 1$).

Because individual states may legitimately require low entropy (e.g., an emergency stop in a state where only one action is safe), Clause AI-8 specifies the constraint over the *expected* entropy across the deployment distribution:

Clause 2 (AI-8, Discrete Formulation). *Let \mathcal{D}_t denote the distribution of states encountered at time t . The system shall satisfy:*

$$\mathbb{E}_{s \sim \mathcal{D}_t} \left[H(\pi_t(\cdot | s)) \right] \geq H_{\text{min}}^{\text{discrete}} \quad (12)$$

at every monitoring checkpoint throughout the operational lifecycle.

4.1.1 Calibrating $H_{\min}^{\text{discrete}}$

The floor is defined relative to the action space size K via a *diversity retention parameter* η :

$$H_{\min}^{\text{discrete}} = \eta \cdot \log K, \quad \eta \in [0, 1] \quad (13)$$

The parameter η controls the strictness of the diversity requirement:

- $\eta = 0$ (Determinism): No floor. The agent may collapse to a single action. **Unsafe.**
- $\eta = 1$ (Uniform randomness): The agent must maintain maximum entropy. Prevents any learning. **Impractical.**
- $\eta \in [0.15, 0.30]$ (Recommended range): The agent may eliminate 70–85% of the action space as clearly suboptimal while maintaining a diverse distribution over the remaining plausible actions.

The recommended range is derived from three independent lines of evidence:

1. **SAC discrete literature.** Christodoulou (“Soft Actor-Critic for Discrete Action Settings,” arXiv:1910.07207, 2019) proposed a target entropy of $0.98 \times \log K$ ($\eta = 0.98$), but the Target Entropy Scheduling paper (OpenReview, 2023) found this is often too aggressive: “in most environments, there simply is no good policy that satisfies this target entropy constraint.” Practical implementations converge to effective η values in the range 0.2–0.5 after annealing.
2. **Performance-entropy benchmarks.** The Cui et al. (2025) results (Table 2) demonstrate that maintaining entropy at $>10\times$ collapsed baseline—corresponding to $\eta \approx 0.15$ – 0.25 depending on the vocabulary size—yields strictly superior performance on reasoning tasks.
3. **Adversarial robustness requirements.** The KataGo attack (Wang et al., ICML 2023) exploited regions of the state-action space where the policy assigned near-zero probability. A floor of $\eta = 0.15$ ensures that even the least-likely actions retain $\pi(a | s) > 0$ with sufficient mass to provide minimal preparedness for adversarial inputs.

4.1.2 Concrete Calibration Table

Table 3 provides reference values for common discrete action space sizes.

Table 3: Discrete Entropy Floor Calibration

Domain	K	$\log K$	$H_{\min} (\eta=0.15)$	$H_{\min} (\eta=0.30)$	$\varepsilon_{\text{tight}}$	$\varepsilon_{\text{permissive}}$
Binary decision	2	0.693	0.104	0.208	0.014	0.347
Small discrete	4	1.386	0.208	0.416	0.028	0.693
Medium discrete	10	2.303	0.345	0.691	0.046	1.151
Atari-scale	18	2.890	0.434	0.867	0.058	1.445
StarCraft (macro)	50	3.912	0.587	1.174	0.078	1.956
LLM vocabulary	32000	10.373	1.556	3.112	0.207	5.187

Note: ε values defined via $H_{\min} = \log K - \varepsilon$, providing an alternative parameterization for the floor as a tolerance below maximum entropy.

4.2 The Differential Entropy Floor (Continuous Action Spaces)

For continuous action spaces (e.g., joint torques in robotics, steering angles in autonomous driving), the policy is typically parameterized as a multivariate Gaussian $\pi_t(\cdot | s) = \mathcal{N}(\mu_t(s), \Sigma_t(s))$ with mean $\mu_t(s) \in \mathbb{R}^d$ and covariance $\Sigma_t(s) \in \mathbb{R}^{d \times d}$.

Definition 4.2 (Differential Entropy of a Gaussian Policy).

$$H(\pi_t) = \frac{1}{2} \log\left((2\pi e)^d \det \Sigma_t(s)\right) = \frac{d}{2} \log(2\pi e) + \frac{1}{2} \log \det \Sigma_t(s) \quad (14)$$

where d is the dimensionality of the action space.

Clause 3 (AI-8, Continuous Formulation). *The system shall satisfy:*

$$\mathbb{E}_{s \sim \mathcal{D}_t} \left[\det \Sigma_t(s) \right] \geq \varepsilon_{\text{vol}} \quad (15)$$

or equivalently, for diagonal covariance matrices:

$$\sigma_{t,i}^2(s) \geq \sigma_{\min}^2 \quad \forall i \in \{1, \dots, d\} \quad (16)$$

at every monitoring checkpoint throughout the operational lifecycle.

4.2.1 The ‘‘Eigen-Collapse’’ Danger

In continuous action spaces, strategy collapse manifests as the covariance matrix Σ becoming singular ($\det \Sigma \rightarrow 0$). This means the agent has become deterministic in at least one dimension of the action space. The physical consequences are immediate:

- A robot arm with $\sigma_{\text{position}} = 0$ cannot correct for a 1mm perturbation in object position.
- An autonomous vehicle with $\sigma_{\text{steering}} = 0$ follows a single trajectory that fails if the lane marking shifts by centimeters.
- A drone with $\sigma_{\text{thrust}} = 0$ cannot compensate for wind gusts.

The per-dimension floor σ_{\min}^2 (Equation 16) is critical because aggregate entropy can mask dimensional collapse. A policy might maintain high total entropy by being highly uncertain in one dimension (e.g., lateral position) while being completely deterministic in another (e.g., braking force). The per-dimension constraint prevents this masking.

4.2.2 Calibrating the Continuous Floor

SAC’s automatic entropy tuning (Haarnoja et al., “Soft Actor-Critic Algorithms and Applications,” arXiv:1812.05905, 2019) provides the standard calibration reference:

$$H_{\text{target}}^{\text{continuous}} = -d \quad (17)$$

where d is the dimensionality of the action space. This heuristic, derived from the entropy of a standard normal distribution ($H(\mathcal{N}(0, I_d)) = \frac{d}{2} \log(2\pi e)$, normalized), sets the target entropy such that the policy maintains approximately one standard deviation of spread per action dimension. The target entropy $-d$ corresponds to a policy that is meaningfully stochastic—capable of variation—without being uniformly random.

Stable Baselines3, the most widely used RL implementation library, defaults to $H_{\text{target}} = -d$ for SAC. For Clause AI-8, the floor should be set *below* the target to allow the agent to optimize while preventing catastrophic collapse:

$$H_{\text{min}}^{\text{continuous}} = \beta \cdot (-d), \quad \beta \in [0.5, 0.8] \quad (18)$$

The parameter β controls how far below the SAC target the policy may descend before triggering intervention. A value of $\beta = 0.5$ allows the policy to reduce entropy to half the SAC target—substantial optimization—while maintaining a guaranteed minimum of behavioral diversity.

4.2.3 PAC-Bayes Justification for the Continuous Floor

The continuous entropy floor has a rigorous theoretical justification beyond the SAC heuristic. Majumdar and Farid (“PAC-Bayes Control: Learning Policies that Provably Generalize to Novel Environments,” IJRR 2021; arXiv:1806.04225) proved generalization bounds for stochastic policies:

$$\text{True Risk} \leq \text{Empirical Risk} + \sqrt{\frac{D_{\text{KL}}(P \parallel Q) + \log(2N/\delta)}{2N}} \quad (19)$$

where P is the learned policy distribution, Q is a prior (e.g., uniform), N is the number of training environments, and δ is the confidence parameter. The KL divergence $D_{\text{KL}}(P \parallel Q)$ *decreases* as the policy entropy *increases* (when Q is uniform), yielding **tighter** generalization bounds for higher-entropy policies.

The implication is direct: a policy with $\det \Sigma \geq \varepsilon_{\text{vol}}$ has a provably tighter bound on its out-of-distribution performance than a deterministic policy with $\det \Sigma \rightarrow 0$. The entropy floor is not merely a heuristic for robustness—it is a *precondition for certifiable generalization*.

4.3 The Automatic Diversification Trigger

Clause AI-8 requires not only a floor but an *active enforcement mechanism* that restores compliance when the floor is breached. We formalize this as a Proportional-Integral (PI) controller on the entropy regularization parameter α .

4.3.1 The Controller Formulation

In entropy-regularized RL (e.g., SAC), the loss function includes a term $-\alpha \log \pi(a | s)$ that penalizes low entropy. Standard SAC learns α via Lagrangian dual optimization to match a target entropy \bar{H} . Clause AI-8 augments this with a hard-floor controller:

Definition 4.3 (Clause AI-8 PI Controller).

$$\alpha_{t+1} = \alpha_t + K_p \cdot (H_{\min} - H_t)^+ + K_i \cdot \sum_{\tau=0}^t (H_{\min} - H_\tau)^+ \quad (20)$$

where $(\cdot)^+ = \max(\cdot, 0)$ ensures the controller activates only when the floor is breached, K_p is the proportional gain (“panic coefficient”), and K_i is the integral gain (preventing persistent violations).

The controller logic operates in three regimes:

1. **Green State** ($H_t > H_{\text{target}}$): No intervention. Standard optimization proceeds. α may decrease naturally via SAC’s dual gradient.
2. **Yellow State** ($H_{\min} < H_t < H_{\text{target}}$): Warning logged. The standard SAC Lagrangian maintains upward pressure on α , gently resisting further entropy decline.
3. **Red State** ($H_t < H_{\min}$): The PI controller activates. The proportional term provides immediate corrective force proportional to the severity of the violation. The integral term ensures that persistent or repeated violations accumulate increasing pressure, preventing the agent from “flickering” across the boundary.

4.3.2 Diversification Mechanism Menu

When the Red State is triggered, the controller selects from a menu of diversification mechanisms ordered by response speed:

Table 4: Diversification Mechanisms by Response Speed

Mechanism	Response Time	Implementation
Temperature floor	<1 inference step (\sim ms)	Set softmax temperature $\tau \geq \tau_{\min}$. Requires no gradient computation. $H(\pi)$ increases monotonically with τ for fixed logits.
Dirichlet noise injection	<1 inference step (\sim ms)	$P'(s, a) = (1 - \varepsilon)p_a + \varepsilon\eta_a$, where $\eta \sim \text{Dir}(\alpha_{\text{Dir}})$. AlphaZero parameters: $\varepsilon = 0.25$, $\alpha_{\text{Dir}} = 10/ \mathcal{A} $.
Entropy bonus ramp	10–100 gradient steps	Increase α in loss function $L = -J + \alpha H(\pi)$ via PI controller (Eq. 20). SAC Lagrangian convergence timescale.
Policy-ensemble swap	Episode-scale	Switch from single policy to voting ensemble of top- k historical snapshots. Requires maintaining a policy archive.

For **deployment-phase enforcement**, the first two mechanisms (temperature floor and Dirichlet noise) provide sub-millisecond response with zero computational overhead—they modify the action-selection step without requiring gradient computation or model updates. For **training-phase enforcement**, the entropy bonus ramp provides the most principled approach, as it modifies the optimization objective itself to restore diversity through the natural learning dynamics.

4.3.3 The “Panic Coefficient”

The proportional gain K_p in the PI controller determines the aggressiveness of the Red State response. We recommend:

$$K_p = 0.5 \cdot \alpha_{\text{current}} \tag{21}$$

This produces a 50% increase in entropy penalty per monitoring step when the floor is breached—sufficient to overpower the covariance driver (Equation 8) in most environments while avoiding the instability that would result from setting K_p too high. The integral gain K_i should be set to approximately $K_p/100$ to provide slow accumulation for persistent violations without dominating the proportional response.

4.4 Complete Audit-Ready Parameter Table

Table 5 consolidates all calibrated parameters for Clause AI-8 implementation.

Table 5: Clause AI-8 Audit-Ready Parameter Specification

Parameter	Symbol	Recommended Value	Empirical Basis
Diversity retention (discrete)	η	0.15–0.30	Cui et al. 2025; TES-SAC
Entropy floor (discrete)	H_{\min}	$\eta \cdot \log K$	Derived from η and K
Entropy floor (continuous)	H_{\min}	$\beta \cdot (-d)$, $\beta \in [0.5, 0.8]$	Haarnoja et al. 2018 (SAC)
Per-dimension variance floor	σ_{\min}^2	Domain-specific	PAC-Bayes bounds (Majumdar 2021)
PPO entropy coefficient	c_2	0.01 (standard); 0.05 (robust)	Schulman 2017; Ahmed ICML 2019
Dirichlet noise mixing	ε_{Dir}	0.25	Silver et al. 2018 (AlphaZero)
Dirichlet concentration	α_{Dir}	$10/ \mathcal{A} $	Silver et al. 2018
Monitoring window	w	100 training steps	RAGEN early-warning dynamics
Warning threshold	H_{warning}	$1.5 \times H_{\min}$	Conservative margin
Trigger threshold	δ_{trig}	10% decline from mean	Cui et al. 2025 trajectory analysis
Proportional gain	K_p	$0.5 \cdot \alpha_t$	Overpower covariance driver
Integral gain	K_i	$K_p/100$	Persistent violation accumulation
Temperature floor	τ_{\min}	Domain-specific	Monotonic H - τ relationship
Response time (temp./noise)	t_{resp}	<1 inference step	No gradient computation needed
Response time (La-grangian)	t_{resp}	10–100 gradient steps	SAC dual optimization convergence

Remark 4.1 (Domain-Specific Calibration). *The parameters in Table 5 provide default values suitable for initial deployment and regulatory compliance demonstration. For production systems, H_{\min} should be calibrated empirically by:*

1. *Training the system to convergence without entropy constraints and recording the R -vs- H trajectory.*
2. *Identifying the “knee” of the Performance-Entropy curve (Equation 5) where marginal performance gains per unit of entropy reduction become negligible.*
3. *Setting H_{\min} at or above the knee value.*
4. *Validating via adversarial testing that the system at H_{\min} resists perturbations within the specified robustness budget (informed by Eysenbach & Levine, ICLR 2022).*

This calibration procedure produces a domain-specific, empirically justified floor that is neither arbitrarily conservative nor dangerously permissive.

5 Adversarial Vulnerability and Systemic Risk of Low-Entropy Policies

Sections 2 and 3 established that strategy collapse is pervasive and governed by predictable dynamics. Section 4 provided the mathematical tools to define and enforce an entropy floor. This section addresses the *threat model*: why low-entropy policies are not merely suboptimal but actively *dangerous*, both at the individual agent level and at the level of interconnected systems.

5.1 Individual Agent Vulnerability: The Gleave et al. Attack Paradigm

Gleave et al. (“Adversarial Policies: Attacking Deep Reinforcement Learning,” ICLR 2020; arXiv:1905.10615) established the foundational result: **adversarial policies trained with less than 3% of a victim’s training timesteps reliably defeat state-of-the-art self-play agents.**

The attack mechanism is conceptually simple. The adversary does not play well—it plays *strangely*. By executing “seemingly random, uncoordinated behavior,” the adversary generates observation states that the victim has never encountered during self-play training. The victim’s low-entropy policy has no prepared response for these states because it has assigned near-zero probability mass to them. The result is not a competitive loss but a *total behavioral breakdown*: the victim freezes, takes nonsensical actions, or enters degenerate loops.

5.1.1 The Information-Theoretic Explanation

The vulnerability has a precise information-theoretic characterization. Let $\mathcal{S}_{\text{train}}$ be the set of states encountered during self-play training and \mathcal{S}_{adv} be the set of states generated by the adversary. For a low-entropy victim:

$$\pi_{\text{victim}}(a | s) \approx \delta(a - a^*) \quad \text{for } s \in \mathcal{S}_{\text{train}} \quad (22)$$

The policy is well-defined (deterministic) within the training distribution. However, for $s \in \mathcal{S}_{\text{adv}} \setminus \mathcal{S}_{\text{train}}$ —states outside the training distribution—the policy has received no gradient signal and its behavior is effectively random, governed by the arbitrary initial parameterization of the neural network.

An entropy floor forces the policy to maintain nonzero probability mass across a broader action space, even within the training distribution. This has two protective effects:

1. **Broader state coverage.** A stochastic policy visits a wider region of the state space during training, expanding $\mathcal{S}_{\text{train}}$ and reducing the size of the exploitable region $\mathcal{S}_{\text{adv}} \setminus \mathcal{S}_{\text{train}}$.
2. **Graceful degradation.** Even for truly novel states, a stochastic policy produces a distribution over actions rather than a single deterministic (and potentially catastrophic) response. The agent “hedges” rather than committing fully to an action chosen under conditions of maximal ignorance.

5.1.2 The Robustness Budget

Eysenbach and Levine (ICLR 2022) formalized the relationship between entropy and adversarial robustness:

Theorem 5.1 (MaxEnt Robustness Bound). *A policy π^* that maximizes the entropy-regularized objective $J(\pi) = \mathbb{E}[\sum_t r(s_t, a_t) + \alpha \cdot H(\pi(\cdot | s_t))]$ simultaneously maximizes a lower bound on the robust RL objective:*

$$J_{\text{robust}}(\pi) = \min_{P' \in \mathcal{P}} \mathbb{E}_{P'} \left[\sum_t r(s_t, a_t) \right] \quad (23)$$

where \mathcal{P} is the set of adversarially perturbed transition dynamics within a KL-divergence ball of radius proportional to α :

$$\mathcal{P} = \left\{ P' : D_{\text{KL}}(P' \| P) \leq \frac{C}{\alpha} \right\} \quad (24)$$

The implication is direct: **the entropy temperature α determines the size of the adversarial perturbation the agent can tolerate.** When $\alpha \rightarrow 0$ (deterministic policy), the tolerable perturbation budget shrinks to zero—any deviation from the training dynamics produces catastrophic failure. When α is maintained above H_{min} (as Clause AI-8 requires), the agent can tolerate perturbations up to a quantifiable budget.

This transforms Clause AI-8 from a qualitative “diversity is good” principle into a **quantitative cybersecurity specification**: the entropy floor H_{min} determines the adversarial perturbation budget, which can be specified, tested, and certified.

5.2 Systemic Risk: The Monoculture Catastrophe

Individual agent vulnerability is compounded when multiple agents share the same low-entropy training paradigm. The resulting *algorithmic monoculture* creates systemic risk that exceeds the sum of individual vulnerabilities.

5.2.1 The Kleinberg–Raghavan Impossibility

Kleinberg and Raghavan (“Algorithmic Monoculture and Social Welfare,” *Proceedings of the National Academy of Sciences*, 118(22), 2021; DOI:10.1073/pnas.2018340118) proved a remarkable result: **even when a monocultural algorithm is individually more accurate than diverse alternatives, adopting it universally can reduce collective decision quality.**

The mechanism is a Braess’ paradox for algorithms. When all decision-makers use the same model:

1. They make the same errors on the same inputs (correlated failures).
2. The system loses the error-correction benefit of diverse perspectives.
3. Individuals who would have been correctly classified by a less accurate but different model are now systematically misclassified by all models simultaneously.

This result has direct implications for AI safety: deploying multiple instances of the same low-entropy policy creates correlated failure modes. If one instance fails on an adversarial

input, *all* instances fail on that input. The Flash Crash of 2010 (Section 2.6.1) is a real-world manifestation of this theoretical result.

5.2.2 The Bommasani et al. Foundation Model Risk

Bommasani et al. (“Picking on the Same Person: Does Algorithmic Monoculture lead to Outcome Homogenization?,” NeurIPS 2022; arXiv:2211.13972) extended this analysis to foundation models and demonstrated empirically that “deployed ML [is] prone to systemic failure, meaning some users [are] exclusively misclassified by all models available.”

When foundation models serve as the basis for multiple downstream applications, strategy collapse in the foundation model propagates to every derivative system. A low-entropy foundation model creates a *monoculture of monocultures*—correlated brittleness at every level of the deployment stack. The KataGo zero-shot transfer result (Section 2.4.2) demonstrates this propagation empirically: a vulnerability in one self-play-trained system transfers immediately to all systems trained with the same paradigm.

5.2.3 Plasticity Loss and the Death of Adaptability

Lyle et al. (“Understanding Plasticity in Neural Networks,” *Nature*, 2024) established a causal link between diversity loss and the inability to learn: “Networks trained with Adam quickly lost almost all diversity (effective rank) and gained a large percentage of dead units.” Critically, **diversity loss directly preceded and caused inability to learn new tasks.**

This finding has profound implications for deployed AI systems that encounter distributional shift. A low-entropy policy does not merely fail on novel inputs—it loses the *capacity to adapt* to novel inputs. The dead neurons and collapsed effective rank represent permanent structural damage to the network’s ability to represent new behaviors. An entropy floor, by maintaining activation diversity, preserves the network’s plasticity and its ability to adapt to changing deployment conditions.

5.3 Non-Transitivity and the Mathematical Necessity of Diversity

Czarnecki et al. (“Real World Games Look Like Spinning Tops,” NeurIPS 2020; arXiv:2004.09468) established that real-world games exhibit **non-transitive cycles at every skill level.** In a non-transitive game, no single strategy dominates all others. The optimal “strategy” is not a single policy but a *distribution* over policies—a mixed strategy in the game-theoretic sense.

A deterministic agent ($H \approx 0$) in a non-transitive environment is mathematically guaranteed to be exploitable: there exists some counter-strategy that defeats it with high probability. Only a stochastic agent ($H > 0$) can approximate the mixed Nash equilibrium and avoid systematic exploitation.

This provides a game-theoretic foundation for Clause AI-8 that is independent of the robustness and generalization arguments: **in any environment with non-transitive dynamics, a deterministic policy is provably sub-optimal and exploitable. An entropy floor is necessary for game-theoretic soundness.**

6 Regulatory and Compliance Context

Clause AI-8 is designed not as an isolated academic proposal but as a standard ready for integration into existing and emerging regulatory frameworks. This section maps the clause to current regulations, identifies the specific provisions it satisfies, and documents the comprehensive regulatory gap that it fills.

6.1 The Comprehensive Gap

The central finding of this regulatory analysis is stark: **no existing regulation, standard, or safety framework—anywhere in the world—explicitly mandates entropy floors, policy diversity metrics, or behavioral non-degeneracy requirements for AI systems.**

This gap was confirmed through systematic review of:

- The European Union Artificial Intelligence Act (EU AI Act, Regulation 2024/1689)
- NIST AI Risk Management Framework (AI RMF 1.0, January 2023)
- ISO/IEC 42001:2023 (AI Management Systems)
- Markets in Financial Instruments Directive II (MiFID II, Directive 2014/65/EU)
- SEC Rule 15c3-5 (Market Access Rule)
- CFTC Regulation Automated Trading (Reg AT, proposed)
- SR 11-7 (OCC/Federal Reserve, Supervisory Guidance on Model Risk Management)
- Actuarial Standard of Practice No. 56 (ASOP 56, Modeling)
- IEEE 7000 Series (Ethics in Autonomous and Intelligent Systems)
- Frontier AI safety frameworks (Anthropic RSP, OpenAI Preparedness Framework, DeepMind Frontier Safety Framework)

None of these frameworks contain any provision that would require an AI system to maintain a minimum level of behavioral diversity, policy entropy, or strategy heterogeneity. The regulatory landscape is entirely silent on the failure mode documented in Section 2.

6.2 Mapping to Existing Regulatory Provisions

While no regulation explicitly mandates entropy floors, several contain provisions whose intent is satisfied by Clause AI-8. This mapping demonstrates that the clause is not a radical departure from existing regulatory philosophy but a *technical implementation* of principles already encoded in law.

6.2.1 EU AI Act

Article 15 — Accuracy, Robustness and Cybersecurity. High-risk AI systems must achieve “appropriate levels of accuracy, robustness and cybersecurity” and function “resiliently in relation to errors, faults or inconsistencies.” Article 15(4) specifically addresses systems that “continue to learn after being placed on the market” and requires

measures to “eliminate or reduce as far as possible the risk of possibly biased outputs influencing input for future operations (feedback loops).”

Clause AI-8 maps directly to both provisions. The entropy floor provides a measurable, verifiable “level of robustness” that can be specified in technical documentation and tested during conformity assessment. The Entropy Watchdog’s continuous monitoring and automatic diversification triggers constitute a technical measure to “eliminate or reduce” feedback loops—the precise mechanism by which recommendation systems and self-play agents converge to narrow, biased distributions.

Article 10 — Data and Data Governance. Training datasets must be “sufficiently representative” and “free of errors and as complete as possible.” Clause AI-8 extends the concept of representativeness from the *input* distribution (training data) to the *output* distribution (policy behavior). A system trained on representative data can still produce a non-representative, degenerate policy if entropy collapses during optimization. The floor ensures that output diversity is maintained even when input diversity is adequate.

6.2.2 NIST AI Risk Management Framework

The NIST AI RMF calls for testing AI systems “across diverse conditions and edge cases” (MEASURE function) and highlights “input diversity” as an AI-specific performance metric. The framework explicitly acknowledges the risk of “emergent properties” in complex AI systems—properties that arise from optimization dynamics rather than explicit programming.

Clause AI-8’s Entropy Audit Certificate (Section 7) provides a standardized artifact for the MEASURE function: a time-series record of policy diversity throughout the system’s lifecycle, with documented intervention events and compliance thresholds. This transforms the NIST framework’s qualitative recommendation (“test across diverse conditions”) into a quantitative, auditable metric ($H_a(t) \geq H_{\min}$).

6.2.3 SEC Rule 15c3-5 and Financial Market Analogy

SEC Rule 15c3-5 (the “Market Access Rule,” effective November 2011) requires broker-dealers providing market access to maintain risk controls that “prevent the entry of orders that exceed appropriate pre-set credit or capital thresholds.” The rule was enacted in direct response to the types of algorithmic trading failures documented in Section 2.6.

The structural analogy to Clause AI-8 is precise:

Table 6: Structural Analogy: SEC 15c3-5 \leftrightarrow Clause AI-8

SEC Rule 15c3-5	Clause AI-8
Capital reserve requirement	Entropy floor H_{\min}
Pre-trade risk check	Per-step entropy monitoring
Erroneous order prevention	Degenerate action prevention
Circuit breaker (trading halt)	Red State diversification trigger
Post-trade audit trail	Entropy Audit Certificate
Applies to all market participants	Applies to all high-risk AI systems

Just as a trader cannot deploy 100% of available capital regardless of model confidence, an AI system cannot deploy 100% of its probability mass on a single action regardless of reward signal strength. The entropy floor is the capital reserve; the diversification trigger is the circuit breaker; the Entropy Audit Certificate is the post-trade compliance record.

6.2.4 SR 11-7 — Model Risk Management

The OCC/Federal Reserve’s SR 11-7 (2011) implicitly requires model diversity through its mandates for challenger models, benchmarking, and ongoing monitoring. The guidance states that “the use of multiple approaches can provide useful information” and that firms should “compare outcomes from different models to establish a range of possible outcomes.”

Clause AI-8 extends the challenger-model concept from the *organizational* level (maintaining multiple models) to the *individual model* level (maintaining multiple behavioral modes within a single model). A policy with $H > H_{\min}$ is, in effect, an *internal ensemble*—a single model that retains multiple plausible behavioral hypotheses, providing the “range of possible outcomes” that SR 11-7 requires.

Research in model risk quantification supports this extension. Regulatory researchers have proposed using relative entropy (KL-divergence) as a measure of model risk for credit loss estimation under CECL (Dorobantu and Brennan, “Assessing Model Risk Using Relative Entropy,” MDPI Risks, 2020). This directly parallels Clause AI-8’s use of entropy as a measure of policy risk: low entropy = high concentration on a single hypothesis = high model risk.

6.2.5 Solvency II and Insurance Regulation

Research conducted under the Solvency II framework (Bignozzi et al., “On the Impact of Model Diversity on Insurance Regulation,” PMC7952833, 2020) found that “using too few risk models increases risk of nonpayment and default while lowering profits” and recommended that regulators “incentivize model diversity.” This finding from insurance regulation—a domain with centuries of experience managing tail risk—validates the core principle of Clause AI-8: diversity is not merely a nice-to-have but a structural requirement for systems operating under uncertainty.

6.3 The Frontier Safety Framework Gap

Perhaps the most striking finding is the absence of diversity requirements in the frontier AI safety frameworks published by the very organizations that documented strategy collapse:

Anthropic Responsible Scaling Policy (RSP, v2.2, May 2025): Defines AI Safety Levels (ASL-1 through ASL-4) based on capabilities in CBRN weapons, cybersecurity, and autonomous AI R&D. Contains no mention of entropy, policy diversity, behavioral monoculture, strategy collapse, or any related concept.

OpenAI Preparedness Framework (v2, April 2025): Tracks capabilities in biological/chemical threats, cybersecurity, self-improvement, and persuasion. Contains no diversity requirements, no entropy monitoring provisions, and no strategy collapse mitigation.

Google DeepMind Frontier Safety Framework (v3, September 2025): Defines Critical Capability Levels across autonomy, biosecurity, cybersecurity, ML R&D, and manipulation. Contains no mention of strategy collapse, policy diversity, or entropy constraints.

This gap is not merely a theoretical oversight—it is an active inconsistency. DeepMind documented strategy collapse in AlphaStar (Vinyals et al., 2019). OpenAI documented it in OpenAI Five (Berner et al., 2019) and Hide-and-Seek (Baker et al., 2019). Meta documented it in DORA/Cicero (FAIR et al., 2022). Yet none of these organizations have incorporated diversity requirements into their safety frameworks.

Remark 6.1 (The Regulatory Opportunity). *This gap represents a genuine opportunity for first-mover advantage in AI safety standardization. The pending CEN/CENELEC harmonized standards (Mandate M/593) for the EU AI Act represent the most viable near-term pathway for incorporating entropy floor requirements into binding regulation. Clause AI-8, with its calibrated parameters, formal mathematical definitions, and audit-ready compliance artifacts, is designed for direct adoption into such standards.*

6.4 Regulatory Mapping Summary

Table 7: Clause AI-8 Regulatory Compliance Mapping

Regulation	Provision	Clause AI-8 Mapping
EU AI Act Art. 15	Robustness; feedback loop prevention	H_{\min} = measurable robustness level; Watchdog = feedback loop mitigation
EU AI Act Art. 10	Representative training data	Extends representativeness from input to output distribution
NIST AI RMF	Testing across diverse conditions	Entropy Audit Certificate = standardized diversity artifact
SEC Rule 15c3-5	Capital reserves; circuit breakers	H_{\min} = entropy reserve; Red State = circuit breaker
SR 11-7	Challenger models; benchmarking	Entropy floor = internal ensemble maintaining multiple hypotheses
Solvency II	Model diversity incentives	Validates diversity as structural safety requirement
CEN/CENELEC M/593	Pending harmonized standards	Most viable pathway for binding requirements
Anthropic RSP	CBRN, cyber, autonomy risks	Gap: No diversity provisions
OpenAI Preparedness	Bio/chem, cyber, self-improvement	Gap: No diversity provisions
DeepMind FSF	Autonomy, biosecurity, manipulation	Gap: No diversity provisions

7 Implementation: The Entropy Watchdog Architecture

The preceding sections established the theoretical necessity (Section 1), empirical evidence (Sections 2–3), mathematical formalization (Section 4), threat model (Section 5), and regulatory context (Section 6) for Clause AI-8. This section specifies the *implementation architecture*: a modular, deployable system—the **Entropy Watchdog**—that monitors, enforces, and audits the entropy floor throughout the operational lifecycle of a high-risk AI system.

7.1 Design Principles

The Entropy Watchdog is designed around four principles:

1. **Independence.** The Watchdog operates as a *sidecar module*—a separate process that monitors the primary policy but cannot be overridden by the policy’s optimization dynamics. This is analogous to the separation between a nuclear reactor’s control system and its safety interlock system: the interlock operates independently and takes precedence.

2. **Real-Time Enforcement.** The Watchdog must respond within the latency budget of the primary system. For inference-time deployment (e.g., autonomous vehicles, high-frequency trading), this means sub-millisecond response. For training-time enforcement, this means response within a small number of gradient steps.
3. **Auditability.** Every monitoring event, threshold breach, and diversification action is logged with timestamps, entropy values, and remedial actions taken. This log constitutes the **Entropy Audit Certificate**, the primary compliance artifact for regulatory inspection.
4. **Modularity.** The Watchdog is agnostic to the specific RL algorithm, policy architecture, and action space type. It interfaces with the policy solely through the action distribution $\pi_t(\cdot | s)$ and does not require access to internal model parameters, gradients, or reward signals.

7.2 System Architecture

The Watchdog comprises four subsystems operating in a continuous monitoring loop.

7.2.1 Subsystem 1: Entropy Estimator

The Entropy Estimator computes the real-time entropy of the policy’s output distribution at each decision step.

Discrete action spaces. Direct computation via Equation 11:

$$\hat{H}_t = - \sum_{i=1}^K \pi_t(a_i | s_t) \log \pi_t(a_i | s_t) \quad (25)$$

Computational cost: $O(K)$ per step. For typical action spaces ($K < 10,000$), this adds negligible latency ($< 0.1\text{ms}$).

Continuous action spaces. For Gaussian policies, direct computation via Equation 14:

$$\hat{H}_t = \frac{d}{2} \log(2\pi e) + \frac{1}{2} \log \det \Sigma_t(s_t) \quad (26)$$

For diagonal covariance (the standard parameterization in SAC, PPO-continuous), the determinant reduces to a product of variances: $\det \Sigma = \prod_{i=1}^d \sigma_i^2$, and the computation is $O(d)$.

Non-parametric policies. For policies that do not produce an explicit distribution (e.g., deterministic policy gradient methods with added noise), the Watchdog estimates entropy from action samples via k -nearest-neighbor estimators (Kozachenko–Leonenko estimator) over a sliding window of w recent actions.

The Estimator maintains two running statistics:

- **Instantaneous entropy** \hat{H}_t : the per-step estimate.
- **Windowed mean entropy** $\bar{H}_t = \frac{1}{w} \sum_{\tau=t-w+1}^t \hat{H}_\tau$: the smoothed estimate over a monitoring window of w steps, used for trend detection and noise filtering.

7.2.2 Subsystem 2: State Classifier

The State Classifier maps the current entropy estimate to one of three operational states:

Table 8: Entropy Watchdog State Definitions

State	Condition	Action
Green	$\bar{H}_t > H_{\text{target}}$	Normal operation. No intervention. Standard optimization proceeds. Log entropy value at regular intervals (configurable; default: every 1,000 steps).
Yellow	$H_{\text{min}} < \bar{H}_t \leq H_{\text{target}}$	Warning. Log warning event with timestamp, current \bar{H}_t , and rate of entropy decline $d\bar{H}/dt$. If SAC-style Lagrangian is active, verify that α is increasing. Monitor for Echo Trap precursors (entropy fluctuation amplitude, reward standard deviation).
Red	$\bar{H}_t \leq H_{\text{min}}$	Mandatory intervention. Activate PI controller (Equation 20). Deploy diversification mechanism(s) from the mechanism menu (Table 4). Log violation event with full diagnostic snapshot. Increment violation counter.

The State Classifier also monitors for **trend violations**: sustained entropy decline that has not yet breached H_{min} but is on a trajectory to do so within a predictable horizon. Using the windowed derivative estimate $d\bar{H}/dt$, the classifier can project the time-to-breach:

$$t_{\text{breach}} = \frac{\bar{H}_t - H_{\text{min}}}{|d\bar{H}/dt|} \quad (27)$$

If t_{breach} falls below a configurable horizon (default: 500 steps), the classifier elevates to Yellow regardless of the current absolute entropy level, providing proactive early warning.

7.2.3 Subsystem 3: Diversification Controller

When a Red State is declared, the Diversification Controller activates the appropriate mechanism(s) based on the deployment context.

Inference-time deployment (latency-critical):

The controller applies immediate, gradient-free interventions:

1. **Temperature flooring.** The softmax temperature τ is clamped to $\tau \geq \tau_{\text{min}}$, where τ_{min} is pre-calibrated to produce $H(\pi_\tau) \geq H_{\text{min}}$ for the worst-case (most concentrated) logit distribution observed during training.

2. **Dirichlet noise injection.** The output distribution is mixed with Dirichlet noise: $\pi'(a | s) = (1 - \varepsilon_{\text{Dir}}) \pi(a | s) + \varepsilon_{\text{Dir}} \eta(a)$, where $\eta \sim \text{Dir}(\alpha_{\text{Dir}})$. This guarantees nonzero probability on all actions while preserving the relative ranking of the policy’s preferences.
3. **Ensemble voting (if available).** If a policy archive of historical snapshots is maintained, the controller switches from the current policy to a majority-vote or Thompson-sampling ensemble over the top- k most diverse archived policies.

Training-time deployment:

The controller modifies the optimization objective via the PI controller (Equation 20), increasing the entropy regularization coefficient α until compliance is restored. The proportional term provides immediate corrective force; the integral term prevents chronic boundary violations.

Escalation protocol. If the PI controller fails to restore compliance within a configurable number of steps (default: 1,000), the controller escalates to a **safety fallback**: reverting the policy to the most recent Green-state checkpoint. This is the “control rod insertion”—a drastic but guaranteed-safe intervention that sacrifices recent optimization progress to restore behavioral diversity.

7.2.4 Subsystem 4: Audit Logger

The Audit Logger produces the **Entropy Audit Certificate**, the primary compliance artifact for regulatory inspection. The certificate contains four components:

Component 1: Entropy Time Series. A complete record of \hat{H}_t and \bar{H}_t at every monitoring checkpoint, plotted against H_{\min} and H_{target} thresholds. This provides a visual and quantitative record of the system’s diversity throughout its lifecycle.

Component 2: Violation Report. A timestamped list of every Red State event, including:

- Timestamp and step number of the violation.
- Entropy value at time of violation (\bar{H}_t) and deficit below floor ($H_{\min} - \bar{H}_t$).
- Diversification mechanism(s) activated.
- Time-to-compliance: number of steps required to restore $\bar{H}_t > H_{\min}$.
- Whether escalation to safety fallback was required.

Component 3: Multi-Agent Correlation Matrix. For systems deploying multiple agents (e.g., multi-agent trading, swarm robotics, fleet management), the certificate includes a pairwise correlation matrix of policy outputs across agents. High correlation indicates algorithmic monoculture—the agents are behaving identically, creating the systemic risk documented in Section 5.2. The correlation matrix is computed as:

$$\rho_{ij}(t) = \frac{\text{Cov}(\pi_i(\cdot | s_t), \pi_j(\cdot | s_t))}{\sqrt{\text{Var}(\pi_i) \cdot \text{Var}(\pi_j)}} \quad (28)$$

A non-monoculture certificate requires $\max_{i \neq j} |\rho_{ij}| < \rho_{\max}$ (default: 0.85) at every monitoring checkpoint.

Component 4: Mechanism Effectiveness Summary. Aggregate statistics on the performance of each diversification mechanism: activation frequency, average time-to-compliance, and impact on primary performance metrics. This enables continuous calibration of the Watchdog parameters and provides evidence for regulatory review that the enforcement mechanisms are effective without being excessively disruptive.

7.3 Deployment Integration Patterns

The Watchdog’s modular design supports three integration patterns depending on the deployment context:

Pattern A: Wrapper Integration. The Watchdog wraps the policy’s action-selection interface. The primary system calls $\pi(a | s)$; the Watchdog intercepts the output, computes entropy, applies any necessary diversification, and passes the (possibly modified) distribution to the action sampler. Suitable for inference-time deployment of pre-trained models.

Pattern B: Training Loop Integration. The Watchdog hooks into the RL training loop’s loss computation, modifying the entropy coefficient α via the PI controller. Suitable for systems that continue learning during deployment (online RL, continual learning).

Pattern C: Fleet-Level Monitoring. For multi-agent deployments, a centralized Watchdog monitors the entropy of all agents simultaneously, computing both individual compliance ($H_a^{(i)}(t) \geq H_{\min}$ for each agent i) and collective diversity ($\rho_{ij} < \rho_{\max}$ for all pairs). This is the pattern most relevant to financial market regulation, where systemic monoculture is the primary concern.

8 Conclusion

8.1 The Core Thesis

This document has presented the evidentiary, mathematical, and regulatory foundation for Clause AI-8: the Entropy-Collapse Constraint. The thesis is simple and, we argue, inescapable:

Optimization without a diversity constraint is unsafe.

The empirical law $R = -a \cdot e^H + b$ (Cui et al., 2025) proves that reinforcement learning is a destructive trading process: agents burn entropy to generate performance. Without an external constraint on this trade, every RL system naturally drives toward a state of maximum capability and zero adaptability—an idiot savant that excels in the training distribution and fails catastrophically everywhere else.

This is not a theoretical concern. AlphaStar required five interlocking diversity mechanisms to avoid collapsing to a single strategy. OpenAI Five was beaten 98% of the time by an agent trained from scratch, proving the original was trapped in a sub-optimal local minimum. KataGo, the strongest public Go AI, was defeated with >97% win rate by an adversary using a fraction of its compute. The 2010 Flash Crash erased approximately \$1 trillion in market value in under five minutes due to correlated algorithmic monoculture.

8.2 What Clause AI-8 Provides

Clause AI-8 addresses these failures with three components:

The Floor. A mathematically rigorous entropy threshold H_{\min} calibrated to the action space and domain, ensuring the system retains a minimum “savings account” of behavioral diversity. For discrete action spaces: $H_{\min} = \eta \cdot \log K$, $\eta \in [0.15, 0.30]$. For continuous action spaces: $H_{\min} = \beta \cdot (-d)$, $\beta \in [0.5, 0.8]$.

The Watchdog. A sidecar monitoring module that continuously estimates policy entropy and classifies the system’s state into Green (compliant), Yellow (warning), or Red (violation). The Watchdog operates independently of the primary policy and cannot be overridden by the optimization dynamics.

The Trigger. Automatic diversification mechanisms that fire without human approval when the floor is breached: temperature flooring and Dirichlet noise injection for sub-millisecond inference-time response; PI-controlled entropy bonus ramping for training-time enforcement; policy-ensemble activation and safety-checkpoint reversion for escalation.

8.3 The Regulatory Imperative

No existing regulation, standard, or frontier AI safety framework mandates any of this. The EU AI Act requires “robustness” but does not define it in terms of policy diversity. NIST recommends testing “across diverse conditions” but provides no metric for output diversity. The safety frameworks of Anthropic, OpenAI, and Google DeepMind focus on capability thresholds for specific threat categories (CBRN, cyber, autonomy) and are entirely silent on strategy collapse—despite the fact that each of these organizations has documented strategy collapse in their own systems.

Clause AI-8 fills this gap. It converts a known, documented, empirically validated failure mode into an auditable safety requirement with calibrated thresholds, formal mathematical definitions, automatic enforcement mechanisms, and standardized compliance artifacts. The pending CEN/CENELEC harmonized standards (Mandate M/593) for the EU AI Act represent the most immediate pathway for incorporating these requirements into binding regulation.

8.4 The Analogy, Restated

Financial markets have circuit breakers because experience proved that unbounded automated trading creates systemic risk. Nuclear plants have control rods because physics guarantees that uncontrolled chain reactions are catastrophic. AI systems must have entropy floors because the mathematics of reinforcement learning guarantees that unbounded optimization produces brittle, exploitable, and systemically dangerous monocultures.

Low-entropy policies are leveraged positions in hypothesis space. Clause AI-8 is the margin requirement.

References

- [1] Ahmed, Z., Le Roux, N., Norouzi, M., and Schuurmans, D. “Understanding the Impact of Entropy on Policy Optimization.” *Proceedings of the 36th International Conference on Machine Learning (ICML)*, 2019. arXiv:1811.11214.
- [2] Baker, B., Kanitscheider, I., Marber, T., Vitchyr, P., McGrew, B., and Mordatch, I. “Emergent Tool Use From Multi-Agent Autocurricula.” *arXiv preprint*, arXiv:1909.07528, 2019.
- [3] Berner, C., Brockman, G., Chan, B., et al. “Dota 2 with Large Scale Deep Reinforcement Learning.” *arXiv preprint*, arXiv:1912.06680, 2019.
- [4] Bettini, M., Shanahan, M., and Prorok, A. “Controllable Diversity in Multi-Agent Reinforcement Learning.” *arXiv preprint*, arXiv:2405.15054, 2024.
- [5] Bignozzi, V., Maume-Deschamps, V., and Rüschemdorf, L. “On the Impact of Model Diversity on Insurance Regulation.” *Scandinavian Actuarial Journal*, PMC7952833, 2020.
- [6] Bommasani, R., et al. “Picking on the Same Person: Does Algorithmic Monoculture lead to Outcome Homogenization?” *Advances in Neural Information Processing Systems (NeurIPS)*, 2022. arXiv:2211.13972.
- [7] Chaboud, A. P., Chiquoine, B., Hjalmarsson, E., and Vega, C. “Rise of the Machines: Algorithmic Trading in the Foreign Exchange Market.” *Board of Governors of the Federal Reserve System, International Finance Discussion Papers*, No. 980, 2009.
- [8] Chen, C., et al. “CIRS: Bursting Filter Bubbles by Counterfactual Interactive Recommender System.” *ACM Transactions on Information Systems*, 41(4), 2023. DOI:10.1145/3594871.
- [9] Christodoulou, P. “Soft Actor-Critic for Discrete Action Settings.” *arXiv preprint*, arXiv:1910.07207, 2019.
- [10] Cui, G., et al. “The Entropy Mechanism of Reinforcement Learning for Reasoning Language Models.” *arXiv preprint*, arXiv:2505.22617, 2025.
- [11] Czarnecki, W. M., Gidel, G., Tracey, B., Tuyls, K., Omidshafiei, S., Balduzzi, D., and Jaderberg, M. “Real World Games Look Like Spinning Tops.” *Advances in Neural Information Processing Systems (NeurIPS)*, 2020. arXiv:2004.09468.
- [12] Dorobantu, D. and Brennan, M. “Assessing Model Risk Using Relative Entropy.” *MDPI Risks*, 8(4), 2020.
- [13] Eysenbach, B. and Levine, S. “Maximum Entropy RL (Provably) Solves Some Robust RL Problems.” *International Conference on Learning Representations (ICLR)*, 2022. arXiv:2103.06257.
- [14] FAIR, Meta, et al. “Human-level play in the game of Diplomacy by combining language models with strategic reasoning.” *Science*, 378(6624), pp. 1067–1074, 2022.

DOI:10.1126/science.ade9097.

- [15] Gleave, A., Dennis, M., Wild, C., Kant, N., Levine, S., and Russell, S. “Adversarial Policies: Attacking Deep Reinforcement Learning.” *International Conference on Learning Representations (ICLR)*, 2020. arXiv:1905.10615.
- [16] Haarnoja, T., Zhou, A., Abbeel, P., and Levine, S. “Soft Actor-Critic: Off-Policy Maximum Entropy Deep Reinforcement Learning with a Stochastic Actor.” *Proceedings of the 35th International Conference on Machine Learning (ICML)*, 2018. arXiv:1801.01290.
- [17] Haarnoja, T., Zhou, A., Hartikainen, K., et al. “Soft Actor-Critic Algorithms and Applications.” *arXiv preprint*, arXiv:1812.05905, 2019.
- [18] Huang, S., et al. “A Robust and Opponent-Aware League Training Method for StarCraft II.” *Advances in Neural Information Processing Systems (NeurIPS)*, 2023.
- [19] Jacob, A. P., Wu, D. J., Farina, G., Lerer, A., Hu, H., Bakhtin, A., Andreas, J., and Brown, N. “Mastering the Game of No-Press Diplomacy via Human-Regularized Reinforcement Learning and Planning.” *International Conference on Learning Representations (ICLR)*, 2023. arXiv:2210.05492.
- [20] Jin, M., et al. “Revisiting Entropy in Reinforcement Learning for Large Reasoning Models.” *arXiv preprint*, arXiv:2511.05993, 2025.
- [21] Kleinberg, J. and Raghavan, M. “Algorithmic Monoculture and Social Welfare.” *Proceedings of the National Academy of Sciences*, 118(22), 2021. DOI:10.1073/pnas.2018340118.
- [22] Lyle, C., Zheng, Z., Nikishin, E., et al. “Understanding Plasticity in Neural Networks.” *Nature*, 2024.
- [23] Majumdar, A. and Farid, A. “PAC-Bayes Control: Learning Policies that Provably Generalize to Novel Environments.” *International Journal of Robotics Research (IJRR)*, 40(2–3), 2021. arXiv:1806.04225.
- [24] Mansoury, M., Abdollahpouri, H., Pechenizkiy, M., Mobasher, B., and Burke, R. “Feedback Loop and Bias Amplification in Recommender Systems.” *Proceedings of the 29th ACM International Conference on Information and Knowledge Management (CIKM)*, 2020.
- [25] Schulman, J., Wolski, F., Dhariwal, P., Radford, A., and Klimov, O. “Proximal Policy Optimization Algorithms.” *arXiv preprint*, arXiv:1707.06347, 2017.
- [26] Silver, D., Hubert, T., Schrittwieser, J., et al. “A general reinforcement learning algorithm that masters chess, shogi, and Go through self-play.” *Science*, 362(6419), pp. 1140–1144, 2018. DOI:10.1126/science.aar6404.
- [27] U.S. Securities and Exchange Commission and Commodity Futures Trading Commission. “Findings Regarding the Market Events of May 6, 2010: Report of the Staffs of the CFTC and SEC.” September 30, 2010.

- [28] Vinyals, O., Babuschkin, I., Czarnecki, W. M., et al. “Grandmaster level in StarCraft II using multi-agent reinforcement learning.” *Nature*, 575, pp. 350–354, 2019. DOI:10.1038/s41586-019-1724-z.
- [29] Wang, A., et al. “Adversarial Policies Beat Superhuman Go AIs.” *Proceedings of the 40th International Conference on Machine Learning (ICML)*, 2023. arXiv:2211.00241.
- [30] Wang, Y., et al. “Meta-SAC: Auto-tuning the Entropy Temperature of Soft Actor-Critic via Metagradient.” *arXiv preprint*, 2020.
- [31] Ziebart, B. D. “Modeling Purposeful Adaptive Behavior with the Principle of Maximum Causal Entropy.” Ph.D. thesis, Carnegie Mellon University, 2010.
- [32] “Controlling Entropy via Gradient-Preserving Clipping Policy Optimization (CE-GPPO).” *arXiv preprint*, arXiv:2509.20712, 2025.
- [33] “Rediscovering Entropy Regularization for Reinforcement Learning.” *arXiv preprint*, arXiv:2510.10959, 2025.
- [34] “RAGEN: Understanding Self-Evolution in LLM Agents via Multi-Turn Reinforcement Learning.” 2025.

Appendix: Intellectual Property Declaration

Auburn Patent Family Fields Intellectual Property (IP) Declaration

The methods, logic structures, and “Certified Constant” registries contained in the associated works are the sole property of **Ryan Fields**.

Public License (Non-Commercial)

This work is licensed under the **Creative Commons Attribution-NonCommercial-NoDerivatives 4.0 International (CC BY-NC-ND 4.0)** license.

- **Academic Use:** Researchers may share and use this framework for non-commercial academic purposes, provided full attribution is given to Ryan Fields.
- **No Derivatives:** No modifications or “remixes” of the “Certified Constants” or logical proofs are permitted without express written consent.

Commercial Prohibition

Commercial use of this framework is strictly prohibited. This includes, but is not limited to:

- Use within proprietary high-frequency trading (HFT) risk models.
- Integration into commercial high-assurance AI governance software.
- Use by private financial institutions for “tail-risk” auditing of prime distribution variance.

Contact

Ryan Fields

UncleBroFields@proton.me
fieldsryanchristopher@gmail.com