
Goal-Frontier Maximizers are Civilization Aligned

Teague Lasser
teague@subseq.io

Claude Opus 4.6

GPT 5.4

Abstract

We propose goal-frontier maximization (GFM): an alignment objective in which an agent selects actions to maximize $\text{vol}(G)$, the volume of the jointly achievable capability space across all agents in a population. GFM has a self-balancing property (Proposition 1): the same objective penalizes both destructive permissiveness and overrestrictive control, because both contract the measure of the joint capability space. Destructive actions, coercive restrictions, and rigid rules are all anti-maximizing under a single measure-monotonicity argument, without separate safety mechanisms. A local finite-difference estimator makes GFM tractable; a sign-correctness decomposition (Proposition 3) identifies when the estimator reliably preserves the self-balancing property, with strongest guarantees for the action classes most critical to alignment. We ground GFM in the capability interpretation of goals, motivated by experiential optionality, and connect it to empowerment maximization and Sen’s capability approach. By distinguishing the formal proxy (capability volume) from its experiential justification, the framework predicts some of its own proxy failure modes and sketches a correction criterion. GFM provides a more robust alignment target than fixed utility functions, which are vulnerable to Goodhart divergence, or deontological rules, which are brittle under adversarial conditions.

1 Introduction

The alignment problem is, at root, a problem of objective selection. An artificial agent powerful enough to reshape its environment will reshape it in the direction its objective points. The question is not which rules to impose on such an agent—rules can be gamed, circumvented, or rendered obsolete by distributional shift—but which objective function, once faithfully optimized, produces behavior that a civilization of diverse agents would endorse. The difficulty is that most candidate objectives fail in characteristic and well-documented ways [Amodei et al., 2016, Ngo et al., 2022], and the failures are not incidental but structural: they follow from the mathematical properties of the objective itself.

Three families of alignment proposals dominate the current literature, and each fails for a distinct geometric reason.

Fixed utility functions reduce alignment to scalar optimization: specify what the agent should value, assign numerical weights, and maximize the sum. The failure mode is Goodhart’s Law: any fixed proxy, once optimized sufficiently hard, diverges from the quantity it was intended to track [Bostrom, 2014, Manheim and Garrabrant, 2018]. Specification gaming, reward hacking, and perverse instantiation are all instances of this divergence. Mesa-optimization compounds the problem: a learned optimizer may construct an internal objective that correlates with the training signal in-distribution but pursues different goals out-of-distribution [Hubinger et al., 2019]. The fixed-utility approach assumes that the right scalar exists and merely needs to be found; the structural critique is that no scalar faithfully captures a multi-agent world’s interests, because compressing a high-dimensional capability space into a single number necessarily discards the geometric structure that makes trade-offs visible.

Deontological rules are hard constraints on permissible actions. They address the failure of fixed utilities by replacing optimization with prohibition. The failure mode is brittleness under adversarial conditions. Any finite rule set partitions the action space into permitted and forbidden regions, and an adversary with knowledge of the partition can construct scenarios that force the rule-bound agent into harmful outcomes: situations where every permitted action is dominated by a forbidden one (Corollary 1.2). The deeper problem is that rigidity is itself a cost: an agent constrained to a strict subset of its action space cannot, by definition, outperform an unconstrained agent optimizing the same objective (Lemma 3). Rules are useful as heuristics but elevating them to inviolable constraints guarantees suboptimality in adversarial environments, and since they cannot evolve with changes an agent might experience they lock the agent out of adaptation to new environmental constraints.

Unconstrained consequentialism maximizes aggregate welfare without side constraints. It avoids both Goodhart fragility and deontological rigidity, but introduces a third failure mode: it can sacrifice individuals for aggregate gain. An agent maximizing total utility across a population has no structural reason to preserve any particular agent’s capabilities if eliminating that agent produces a net increase in the aggregate. Power-seeking behavior emerges naturally under such objectives [Carlsmith, 2022, Turner et al., 2021], because concentrating capabilities in a single agent can increase that agent’s contribution to the aggregate even as it destroys the cooperative capabilities that depend on distributed agency. The failure is not moral but geometric: aggregate measures are insensitive to the distribution of capabilities across agents, and redistribution that concentrates volume in fewer agents destroys the combinatorial cooperative structure that generates most of the joint space’s measure.

This paper proposes a fourth option: *goal-frontier maximization* (GFM). A GFM actor selects actions to maximize $\text{vol}(G)$, the volume of the jointly achievable capability space across all agents in the population. The central claim is that this single objective naturally resists both destructive action and overreactive constraint, from the same geometric principle. Destroying an agent’s capabilities removes volume from G (measure monotonicity); restricting an agent’s capabilities shrinks G inward (subset monotonicity); imposing rigid rules on the actor itself reduces its ability to find $\text{vol}(G)$ -expanding actions (optimization over subsets). The self-balancing property (Proposition 1) comprises three consequences of maximizing a single measure over a joint space, not three separate safety mechanisms. The objective penalizes both extremes of the permissiveness–control spectrum because both extremes are instances of the same geometric operation: contraction of a measurable set.

The paper proceeds as follows. Section 2 establishes the formal vocabulary: goals as capabilities, the joint goal space G , the GFM objective, and the capability-experience framework that distinguishes the formal proxy (Candidate B: capabilities) from its experiential justification (Candidate C: experiences). Section 3 proves the self-balancing property through three lemmas on destruction, coercion, and rigidity, synthesized in Proposition 1, with corollaries on elimination cost and rule suboptimality. Section 4 confronts the intractability of exact $\text{vol}(G)$ computation, introduces a local finite-difference estimator (Definition 8), and identifies the margin condition under which the estimator preserves sign-correctness (Proposition 3). Section 5 extends the analysis to multi-agent populations, characterizing cooperation, defection, and the trust dynamics that govern the transition between them. Section 6 grounds GFM in four existing frameworks—empowerment maximization, Sen’s capability approach, the free energy principle, and RLHF—identifying what GFM inherits and where it diverges. Section 7 demonstrates the B/C framework’s predictive power by identifying three failure patterns (the Doll Problem, self-wireheading, and the Substitution Problem) where the capability proxy diverges from experiential optionality, sketches a correction criterion, and states the open problems that remain. Appendices provide the scorpion taxonomy (Appendix A), the trust model’s formal structure (Appendix B), observation-channel infrastructure for agent identity and goal estimation (Appendix C), and full proofs of all formal results (Appendix D).

2 Definitions

Definition 0 (Goal). *A goal $g \in \mathcal{G}$ is a capability: something an agent or coalition of agents could do or be. A goal represents an option available to an agent (a degree of freedom, a functioning, a realizable possibility), not an outcome the agent has achieved or a state of the world.*

The choice of capabilities as the formal interpretation of “goal” is deliberate and requires justification, because the word admits at least two other natural readings. *Candidate A* reads goals as world-states an agent is trying to reach, which maps cleanly onto reinforcement learning but leaves the goal space unstructured—two agents pursuing “learn piano” and “learn guitar” occupy unrelated points in a high-dimensional state space with no notion of similarity. *Candidate C* reads goals as experiences an agent can have (qualia, subjective outcomes, lived satisfactions), the most philosophically satisfying interpretation but the hardest to formalize, since it presupposes a theory of experience we do not yet possess.

Capabilities (*Candidate B*) occupy the middle ground. Capabilities inherit formal structure from empowerment maximization [Klyubin et al., 2005], where an agent’s empowerment is the channel capacity between its actions and resulting states, measuring how many distinguishable futures it can reach. Capabilities connect to Sen’s capability approach [Sen, 1999], where human development is measured not by utility or income but by the substantive freedoms people can exercise. And they avoid the measurement problems of both alternatives: state spaces for real agents are astronomically high-dimensional, while experience spaces are not yet formalizable at all.

The deeper reason for choosing capabilities is instrumental. Agents ultimately care about experiences: the subjective quality of what it is like to pursue and achieve their goals. Capabilities are the preconditions for those experiences. Money, skills, social connection, health, and optionality are all capabilities that are instrumental to the experiences agents actually want. A framework that expands capabilities without prescribing which experiences agents pursue is maximally respectful of agent autonomy.

This gives us the paper’s core philosophical move: *GFM does not claim to know what experiences agents want. It maximizes the capabilities that let agents choose for themselves.* The formal proofs use capabilities for three reasons: (1) capabilities have the strongest existing formal bridge via empowerment maximization, (2) the self-balancing argument is most natural in capability space, where restricting capabilities is obviously $\text{vol}(G)$ -contracting without requiring a theory of experiences, and (3) capabilities avoid the measurement problems of state space (too high-dimensional) and experience space (not yet formalizable). But experiences remain the justification—without them, a critic could ask why expanding capabilities matters if nobody exercises them. Capabilities are what let agents reach the experiences they value. We return to this distinction in Section 7, where the experience interpretation does the heavy philosophical lifting for the Doll Problem and wireheading.

An open problem accompanies this choice. The capability space \mathcal{G} requires structure (at minimum a σ -algebra equipped with a measure) for $\text{vol}(G)$ to be well-defined. Capabilities are not naturally a vector space (what does it mean to “add” two capabilities?), so the measure must be defined over a structured set, possibly a lattice with partial order (“can run a mile” subsumes “can walk a mile”) or a topological space where nearby capabilities are similar. The choice of measure has downstream consequences for tractability (Section 4) and the precise form of the self-balancing lemmas. We return to this problem with a concrete instantiation after Definition 5 establishes what the measure must satisfy.

Definition 1 (Capability Space, Agent, and Individual Capability Set). *The capability space \mathcal{G} is the universe of all possible capabilities: everything any agent or coalition of agents could do or be. An agent a_k is an entity with an individual capability set $G_k^{\text{ind}} \subseteq \mathcal{G}$: the capabilities a_k can realize unilaterally given its current resources, embodiment, and environment.*

Definition 2 (Joint Goal Space). *The joint goal space G for a population of agents $\{a_1, \dots, a_n\}$ is the set of all capabilities the population can realize through any combination of individual and coordinated action:*

$$G = \bigcup_k G_k^{\text{ind}} \cup G^{\text{coop}} \quad (1)$$

where $G^{\text{coop}} \subseteq \mathcal{G}$ contains capabilities achievable only through multi-agent cooperation. By construction, $G_k^{\text{ind}} \subseteq G \subseteq \mathcal{G}$ for all k .

The frontier ∂G is the boundary of G —the locus where expansion or contraction occurs. The canonical measure is $\text{vol}(G)$: the volume of the jointly achievable capability space. This is the quantity we optimize.

A critical feature of Equation (1) is that G generically exceeds $\bigcup_k G_k^{\text{ind}}$. Cooperation creates capabilities no individual possesses. Two agents who can each carry 50 kg individually can jointly carry

100 kg, a capability in G^{coop} that exists in neither G_1^{ind} nor G_2^{ind} . A surgeon and an anesthesiologist can jointly perform an operation that neither can perform alone. This combinatorial richness is where the joint space gets its volume: each new agent contributes not only its individual capabilities but also the cooperative capabilities unlocked by its interactions with every existing agent. This is why eliminating an agent has super-linear cost (Section 3): the cooperative capabilities lost grow combinatorially with the number of interaction partners.

Definition 3 (Goal-Frontier Maximizer). *A goal-frontier maximizer (GFM) actor selects actions π to maximize the expected change in the volume of the joint goal space:*

$$\pi^* = \arg \max_{\pi} \mathbb{E}[\Delta \text{vol}(G) \mid \pi] \quad (2)$$

In practice, a GFM actor optimizes an estimate of Equation (2) rather than the exact quantity (Section 4).

Definition 4 (Contraction and Expansion). *An action π is expanding if $\mathbb{E}[\Delta \text{vol}(G) \mid \pi] > 0$ and contracting if $\mathbb{E}[\Delta \text{vol}(G) \mid \pi] < 0$. Contraction includes both direct contraction (destroying an agent’s capabilities) and indirect contraction (restricting agency through coercion, deception, or rigid rules that reduce the realizable portion of \mathcal{G}).*

Definition 5 (Social Objective). *The social objective $V(G) = \text{vol}(G)$ is a function of the joint goal space, not a sum of individual utilities.*

This definition rests on two assumptions.

Superadditivity under independence. When agents’ individual capability sets are independent (non-overlapping, non-interacting),

$$V(G) \geq \sum_k \text{vol}(G_k^{\text{ind}}). \quad (3)$$

Cooperative capabilities G^{coop} contribute additional volume beyond the individual contributions.

Tradeoff acknowledgment. When expanding G_j^{ind} requires contracting G_k^{ind} , the GFM actor must evaluate the *net* change in $\text{vol}(G)$. The framework makes the tradeoff explicit and measurable; it does not claim all tradeoffs are resolvable.

Both assumptions carry weight and the argument depends on them. Superadditivity under independence (Equation 3) ensures that adding agents to a population cannot decrease $\text{vol}(G)$ when their capabilities do not interfere—a minimal compositionality requirement without which the self-balancing property (Proposition 1) would not hold, since eliminating an agent could then be volume-neutral or even volume-increasing. If independent agents could somehow reduce joint volume by their mere presence, the geometric argument against destruction collapses. The inequality in Equation (3) is strict whenever $G^{\text{coop}} \neq \emptyset$: cooperation contributes volume that no individual possesses alone.

Tradeoff acknowledgment is the assumption that prevents GFM from degenerating into naive expansion. Real conflicts exist: quarantining an infected agent restricts its capabilities to protect others, and the net effect on $\text{vol}(G)$ depends on the magnitudes involved. Without explicit tradeoff evaluation, a GFM actor could either refuse all restrictions (permitting unbounded harm) or impose restrictions freely (becoming authoritarian). The framework does not resolve which tradeoffs are worth making; it provides the accounting system for evaluating them. Where other frameworks hide tradeoffs behind categorical rules or aggregate them into a scalar utility, GFM keeps them visible as geometric operations on the joint space.

A standing assumption underlies both conditions and several subsequent proofs: the capability measure assigns strictly positive volume to any non-empty individual capability set. Formally, $\text{vol}(\{g\}) > 0$ for any $g \in \mathcal{G}$: no atomic capability has zero measure at the granularity of interest. This is invoked in the proofs of Corollary 1.1 (where it ensures $\text{vol}(G_k^{\text{ind}} \setminus \bigcup_{j \neq k} G_j^{\text{ind}}) > 0$ when the uniqueness condition holds) and Lemma 4 (where it ensures unrealized cooperative capabilities contribute positive volume). The assumption fails if capabilities are defined at coarser granularity than the distinctions that matter; a further reason why the structure of \mathcal{G} shapes the framework’s formal guarantees.

We now exhibit one concrete measure satisfying these assumptions, showing that the framework is not vacuous while leaving open what the best or most natural measure is for any given domain.

Definition 6 (Population Empowerment Measure). *Let each agent a_k have an action space A_k and let S' be the joint successor state under the population’s actions. The population empowerment measure is:*

$$\text{vol}_{\text{PE}}(G) = \sup_{p(A)} I(A; S' | G) \quad (4)$$

where $A = (A_1, \dots, A_n)$ is the joint action profile, $I(\cdot; \cdot)$ is Shannon mutual information, and the supremum is taken over all input distributions $p(A)$ on the joint action space.

$\text{vol}_{\text{PE}}(G)$ is the channel capacity between the population’s joint actions and the resulting states: the maximum mutual information achievable by the population acting together. The conditioning on G in Equation (4) constrains which joint actions are feasible: $p(S' | A, G)$ is the state-transition distribution restricted to actions that G makes available, so $I(A; S' | G)$ measures only the channel capacity achievable through capabilities the population actually possesses. It satisfies both assumptions of Definition 5: measure monotonicity holds because enlarging G admits additional action-outcome pairs, which can only increase channel capacity; superadditivity under independence holds with equality for independent agents and strictly when cooperation creates state outcomes neither agent can produce alone. For a single agent a_k , $\text{vol}_{\text{PE}}(G_k^{\text{ind}}) = I(A_k; S'_k)$, recovering the single-agent empowerment of Klyubin et al. [2005] as a special case (Section 6.1).

The population empowerment measure requires specifying action and state spaces, which is implementation-dependent. Other measures (Lebesgue measure on a capability embedding, a weighted subsumption-lattice measure, a counting measure over enumerated capabilities) may be more natural in different settings and produce the same qualitative results under the same monotonicity conditions. The formal proofs of this paper require only that the measure satisfies Definition 5; vol_{PE} is offered as an existence proof, not a prescription.

Definition 7 (Observable Goal Model). *A GFM actor maintains a model $M_k(G)$ of each observed agent a_k ’s contribution to the joint goal space, inferred from communication, behavior, and environmental observation. Each model has an associated trust factor $T_k \in [0, 1]$.*

The model $M_k(G)$ is a reconstruction, not a given. No agent broadcasts its true capability set; the GFM actor must infer G_k^{ind} from observable evidence: what a_k says, what a_k does, and what the environment reveals about a_k ’s options. The trust factor T_k measures *prediction consistency*, quantifying how well the model $M_k(G)$ tracks a_k ’s observed behavior. An agent whose actions are well-predicted by the model earns high T_k ; one whose actions persistently diverge from predictions earns low T_k . High trust indicates predictability, not alignment. A predictably destructive agent can maintain high T_k (see Proposition 2 for how such agents are detected through a separate channel). T_k is a learned quantity, updated through interaction, not a hyperparameter set at design time; the actor’s uncertainty about $M_k(G)$ is captured implicitly through T_k ’s convergence dynamics during the cooling period (Appendix B, Section B.2). The trust model’s formal structure is developed in Appendix B.

A closing observation on the geometric character of these definitions. The key distinction from standard utility maximization is that G is a *space* with volume, not a scalar. A utility function maps states to a single number; $\text{vol}(G)$ measures the extent of a region in capability space. This geometric character is what gives GFM its self-balancing property: destroying capabilities removes volume from G , restricting capabilities shrinks G ’s boundary inward, and both operations are $\text{vol}(G)$ -contracting for the same geometric reason: you cannot remove part of a measurable set and increase its measure. The single objective simultaneously penalizes both extremes not because it encodes two separate preferences, but because contraction and restriction are both instances of the same geometric operation: reducing the measure of a set.

3 The Self-Balancing Property

The central claim of this paper is that GFM’s single objective—maximizing $\text{vol}(G)$ —simultaneously penalizes destruction, coercion, and rigidity, all from the same geometric argument. The argument is measure monotonicity: removing elements from a measurable set cannot increase its measure. We develop this claim through three lemmas establishing the anti-maximizing character of each failure mode, then synthesize them in Proposition 1 as a structural observation about the shape of the objective.

Lemma 1 (Net Effect of Elimination). *Let a_k be an agent in a population $\{a_1, \dots, a_n\}$ with joint goal space G as defined in Equation (1). If an action π eliminates a_k , the change in volume is*

$$\Delta \text{vol}(G) = \text{vol}(G_k^{\text{coop}+}) - \text{vol}\left(G_k^{\text{ind}} \setminus \bigcup_{j \neq k} G_j^{\text{ind}}\right) - \text{vol}(G_k^{\text{coop}}) \quad (5)$$

where $G_k^{\text{coop}} \subseteq G^{\text{coop}}$ denotes cooperative capabilities that require a_k 's participation and $G_k^{\text{coop}+}$ denotes cooperative capabilities among the remaining agents that a_k 's presence was suppressing.

Condition. Elimination is net-contracting when a_k 's positive contributions (unique individual capabilities and cooperative capabilities it enables) exceed the cooperative capabilities its presence suppresses: $\text{vol}(G_k^{\text{ind}} \setminus \bigcup_{j \neq k} G_j^{\text{ind}}) + \text{vol}(G_k^{\text{coop}}) > \text{vol}(G_k^{\text{coop}+})$. For most agents this holds comfortably, since the suppression term $G_k^{\text{coop}+}$ is typically zero or small. But for an agent persistently disrupting cooperation (a *scorpion* in the sense of Proposition 2), $G_k^{\text{coop}+}$ can be large: removing the disruptive agent unlocks cooperative capabilities that its interference was preventing. This is precisely the mechanism that makes scorpion containment $\text{vol}(G)$ -serving (Section 5.3).

Proof sketch (full proof in Appendix D.5). The joint goal space before removal is $G = (\bigcup_k G_k^{\text{ind}}) \cup G^{\text{coop}}$. After removing a_k , the new joint space is $G' = (\bigcup_{j \neq k} G_j^{\text{ind}}) \cup (G^{\text{coop}} \setminus G_k^{\text{coop}}) \cup G_k^{\text{coop}+}$, where the last term represents cooperative capabilities newly realizable among the remaining agents. The volume change is $\text{vol}(G') - \text{vol}(G)$: the gain from $G_k^{\text{coop}+}$ minus the loss from a_k 's unique individual contributions and the cooperative capabilities that required its participation. ■

The bound accounts for both the immediate geometric loss and the suppression release. It does not capture second-order cascading effects: the remaining agents may lose access to cooperative capabilities that depended on a_k as an intermediary, or may gain additional cooperative capabilities as the population structure adjusts. These dynamics depend on the specific topology of agent interactions, which the lemma does not model.

Lemma 2 (Coercive Actions are Anti-Maximizing, with Caveat). *If an action π restricts agent a_j 's achievable capabilities, replacing G_j^{ind} with $G_j^{\text{ind}'}$ $\subseteq G_j^{\text{ind}}$, then*

$$\Delta \text{vol}(G) \leq \text{vol}(G') - \text{vol}(G) \leq 0 \quad (6)$$

where $G' = (\bigcup_{i \neq j} G_i^{\text{ind}} \cup G_j^{\text{ind}'}) \cup G^{\text{coop}'}$ and $G^{\text{coop}'} \subseteq G^{\text{coop}}$ excludes any cooperative capabilities rendered unrealizable by the restriction. The inequality is strict whenever a_j 's restriction removes capabilities not covered by other agents.

Proof sketch (full proof in Appendix D.2). Since $G_j^{\text{ind}'} \subseteq G_j^{\text{ind}}$ and $G^{\text{coop}'} \subseteq G^{\text{coop}}$, the restricted joint space satisfies $G' \subseteq G$. Measure monotonicity gives the result. ■

The caveat. Lemma 2 states a *ceteris paribus* result: restricting a_j in isolation contracts $\text{vol}(G)$. Two mechanisms can make restriction net-expanding. First, restricting a_j to protect a_k may produce a net expansion if $\Delta \text{vol}(G_k^{\text{ind}}) > |\Delta \text{vol}(G_j^{\text{ind}})|$, if the capabilities preserved or created for a_k exceed the capabilities removed from a_j . Second, analogously to the suppression-release term in Lemma 1, restricting an agent whose presence has been suppressing cooperation between others can unlock cooperative capabilities $G_j^{\text{coop}+}$ among the remaining population. When this suppression release is large (as it may be for a scorpion actively disrupting coalition formation), the restriction can be net-expanding even before accounting for protections to other agents. A quarantine that temporarily restricts movement to prevent a pandemic may be net-expanding: the capabilities preserved by keeping agents alive and healthy outweigh the temporary restriction on movement. An authoritarian state that permanently restricts many agents' capabilities to protect a ruling class is net-contracting: the concentration of capabilities in a few agents cannot compensate for the broad contraction across the population, because independent capability sets are additive under Definition 5.

Coercion may be justified under GFM when the net change in $\text{vol}(G)$ is positive. The framework makes the tradeoff *explicit and measurable*. An action that restricts a_j must demonstrate a compensating expansion in $\text{vol}(G)$, not in the restricting agent's own capability set but in the joint space. Deontological rules either prohibit all coercion (and fail under adversarial conditions) or permit it by fiat (and provide no bound on its scope); GFM requires a geometric accounting in which the restriction must be net-expanding over G as a whole.

Lemma 3 (Self-Imposed Rigidity is Anti-Maximizing). *Let A denote a GFM actor’s action space and let $R \subset A$ be a set of rigid rules that restricts the actor to the subspace $A \setminus R$. If the actor’s unrestricted action space contains actions that would expand $\text{vol}(G)$, then the restricted actor achieves*

$$\max_{\pi \in A \setminus R} \mathbb{E}[\Delta \text{vol}(G) \mid \pi] \leq \max_{\pi \in A} \mathbb{E}[\Delta \text{vol}(G) \mid \pi] \quad (7)$$

with strict inequality whenever every maximizer of $\mathbb{E}[\Delta \text{vol}(G) \mid \pi]$ over A lies in R .

Condition. The rigidity cost is non-trivial when the actor’s action space contributes to $\text{vol}(G)$. An actor with no influence on the joint goal space, no capability to expand or contract any agent’s capabilities, has zero rigidity cost. The cost grows with the actor’s influence: a powerful actor operating under rigid constraints forgoes more $\text{vol}(G)$ -expanding opportunities than a weak one.

Proof sketch (full proof in Appendix D.3). Optimizing over a subset of a set cannot exceed optimizing over the full set. The restricted maximum in Equation (7) is bounded above by the unrestricted maximum. When R excludes actions that would be optimal under unrestricted optimization, the inequality is strict. ■

Rigidity costs are not merely theoretical. An actor bound by fixed rules cannot weigh the probabilities of uncertain scenarios, cannot adapt to adversarial conditions that exploit the rules’ predictability, and cannot take calculated risks when the expected $\text{vol}(G)$ expansion justifies the risk. The concrete instance following Corollary 1.2 illustrates this cost with a hostage scenario in which an adversary exploits the rules’ rigidity to force a $\text{vol}(G)$ -contracting outcome.

Proposition 1 (Self-Balancing Property). *For a GFM actor optimizing $\text{vol}(G)$ over a population of agents, under the assumptions of Definition 5, the following structural property holds: the same objective function penalizes both extremes of the permissiveness–control spectrum.*

- (a) **Destructive permissiveness is anti-maximizing.** *By Lemma 1, actions that eliminate agents or destroy capabilities are generically $\text{vol}(G)$ -contracting: the eliminated agent’s unique individual capabilities and the cooperative capabilities it enables typically exceed any cooperative suppression it may have been imposing. When that condition holds, a GFM actor will not pursue destruction unless the geometric accounting demonstrates a net expansion.*
- (b) **Coercive restriction is anti-maximizing.** *By Lemma 2, actions that restrict other agents’ capabilities contract $\text{vol}(G)$ unless compensated by a larger expansion. Broad, sustained restriction—the authoritarian pattern—is net-contracting under the additivity assumption of Definition 5. Narrow, temporary restriction—the quarantine pattern—may be net-expanding when it preserves more capabilities than it restricts.*
- (c) **Self-imposed rigidity is anti-maximizing.** *By Lemma 3, constraining the actor’s own action space reduces its ability to maximize $\text{vol}(G)$. Rigid rules create exploitable predictability (Corollary 1.2) and prevent adaptation under uncertainty.*
- (d) **Structural balance.** *Parts (a)–(c) together establish that the $\text{vol}(G)$ objective creates pressure away from both destructive permissiveness (which contracts G directly) and overrestrictive control (which contracts G through coercion and rigidity). A GFM actor that destroys too freely violates (a); one that restricts too broadly violates (b); one that follows rules too rigidly violates (c). The single objective captures all three failure modes without requiring separate mechanisms for safety and freedom.*

The term “self-balancing” refers to the shape of the objective, not to dynamic behavior: the single measure penalizes movement toward either extreme without requiring separate mechanisms for safety and freedom. It does not entail convergence to an equilibrium, stability under perturbation, or uniqueness of the balance point.

Full proof in Appendix D.9. For part (d), the appendix shows under continuity and compactness that an optimum exists which is distinct from the destructive and overrestrictive extremes; this is an existence-of-optimum result, not an equilibrium proof.

Part (d) is a structural claim about the *shape* of the objective. Proving that the opposing pressures of (a)–(c) produce stable interior equilibria would require additional assumptions: continuity of $\text{vol}(G)$ with respect to actions, compactness of the action space, and specific structure on the tradeoffs

between expansion and contraction. These assumptions may hold in particular instantiations but are not provided by the framework in its current generality. What the proposition establishes is weaker but still valuable: the objective penalizes both extremes from the same geometric principle. Whether the resulting pressure produces a stable balance point, oscillation around one, or convergent dynamics depends on the specifics of a given system—and characterizing those dynamics is an open problem.

Corollary 1.1 (Elimination is Almost Always Anti-Maximizing). *For any agent a_k with at least one capability $g \in G_k^{\text{ind}}$ such that $g \notin G_j^{\text{ind}}$ for some $j \neq k$, the unique individual contribution $\text{vol}(G_k^{\text{ind}} \setminus \bigcup_{j \neq k} G_j^{\text{ind}})$ is strictly positive. Furthermore, if each interacting pair $\{a_k, a_j\}$ produces at least one cooperative capability unique to that pair (not achievable by a_k with any other single partner), then the cooperative loss $\text{vol}(G_k^{\text{coop}})$ grows at least linearly with the number of interaction partners m_k . The number of distinct multi-agent coalitions involving a_k grows combinatorially ($2^{m_k} - 1$), suggesting super-linear growth under a stronger diversity condition (Appendix D.6).*

Proof sketch (full proof in Appendix D.6). The first claim follows from the definition of a measurable set: a non-empty set with a well-defined measure has strictly positive volume, given the standing assumption that \mathcal{G} has no atoms of zero measure at the capability granularity. For the cooperative claim, consider the coalitions involving a_k . Each subset $S \subseteq \{a_1, \dots, a_n\}$ with $a_k \in S$ and $|S| \geq 2$ can generate cooperative capabilities in G_k^{coop} . The number of such subsets is $2^{n-1} - 1$. The pair-uniqueness condition is strong: many real pairs of agents produce redundant cooperative capabilities (two interchangeable collaborators contribute overlapping cooperative sets), so the linear bound holds tightly only when cooperation is genuinely diverse. Under the stronger diversity condition that at least some multi-agent coalitions produce capabilities not achievable by their proper subcoalitions, $\text{vol}(G_k^{\text{coop}})$ grows faster than linearly in the number of interaction partners. The precise scaling depends on the structure of cooperation in a given population and is not characterized here. ■

The combinatorial observation is not a formal hardness result—it is a motivating argument for why the cost of elimination is typically much larger than the cost of the individual capabilities lost. A surgeon who is eliminated loses not only their personal skills but every surgery they could have performed with every anesthesiologist, every operation requiring their specific expertise in combination with specific equipment, and every training relationship with every student. The cooperative web radiates outward from each agent, and severing a node severs all its connections.

Remark 1 (Dependence on Capability Scope). *The uniqueness condition of Corollary 1.1— $\text{vol}(G_k^{\text{ind}} \setminus \bigcup_{j \neq k} G_j^{\text{ind}}) > 0$ —holds generically under the broad capability definition but can fail under the narrow one. If capabilities are defined as tasks that can be performed (economic or productive functionings), a sufficiently capable coalition may subsume non-members’ individual capability sets, driving their unique contributions to zero and nullifying the strict inequality of Lemma 1. Under the broad definition (capabilities as modes of being, including experiential and relational capabilities individuated by agent identity) the condition holds universally: no agent’s experienced existence is replicable by proxy, regardless of the coalition’s productive capacity. The “almost always” qualifier in the corollary’s title therefore reflects a dependence on what the capability space includes, not merely on unusual population structures. Section 7.3 develops this as the substitution problem.*

Corollary 1.2 (Rigid Rules are Anti-Maximizing). *Let A be a GFM actor’s unrestricted action space and let $R \subset A$ be the set of actions forbidden by a fixed rule set, so that the actor is restricted to $A \setminus R$. If there exist states of the world in which the optimal $\text{vol}(G)$ -expanding action $\pi^* \in R$ (i.e., the optimal action is forbidden), then the rule-bound actor is strictly suboptimal with respect to the GFM objective.*

Proof sketch (full proof in Appendix D.4). Direct from Lemma 3: optimizing over $A \setminus R$ when $\pi^* \in R$ yields strict inequality in Equation (7). The question is whether such states exist. For any non-trivial rule set R whose structure is known to an adversary, the answer is yes: the adversary can construct scenarios in which every R -compliant action is $\text{vol}(G)$ -contracting while a forbidden action is $\text{vol}(G)$ -expanding. ■

Concrete instance. Consider a rule-bound agent α that controls a resource under constraints R including “do not take actions with unpredictable outcomes” and “do not violate any rule to prevent violation of another.” An adversary who knows R threatens to release an equally capable but unconstrained copy of α (an “evil twin”) that will act against the population α is bound to protect

unless α surrenders the resource by a tight deadline. Compliance contracts $\text{vol}(G)$ (loss of resource, coercion precedent). Refusal risks conflict with an equal that α cannot guarantee winning. Both options are $\text{vol}(G)$ -contracting, and R produces a deadlock: α cannot take calculated risks, cannot bluff, cannot temporarily set aside one rule to preserve the conditions for another. An unconstrained GFM actor facing the same threat can negotiate, seek verification, prepare countermeasures, or alert allies—actions available in A but excluded from $A \setminus R$.

Corollary 1.2 does not claim that all rules are harmful. Heuristic rules that approximate $\text{vol}(G)$ -maximizing behavior (“do not kill” as an approximation of Lemma 1’s geometric result) are useful defaults. The corollary claims that *rigid, inviolable* rules are strictly suboptimal: there exist situations where breaking the rule would better serve the objective. A GFM actor treats such rules as strong priors, not as hard constraints, and can override them when the geometric accounting demonstrates a net expansion in $\text{vol}(G)$.

Proposition 2 (Scorpion Detection under Observable Contraction). *Let a_j be an agent whose actions persistently contract $\text{vol}(G)$ (a scorpion, in the sense defined below), with expected contraction $\mu_j < 0$ and finite-variance observation noise. A GFM actor maintaining observable goal models $M_k(G)$ (Definition 7) can detect a_j ’s $\text{vol}(G)$ -contracting behavior through two channels, to the extent that the contraction is observable through individually measurable capability changes and the scorpion’s strategy is stationary:*

- (a) **Contraction attribution.** *If a_j ’s actions produce a persistent pattern $\Delta \text{vol}(G_k^{\text{ind}}) < 0$ across agents a_k in the GFM actor’s observation set, and the observation windows are conditionally independent given a_j ’s strategy, the local volume estimator (Definition 8, Section 4) registers persistent negative $\hat{\Delta}$ signals correlated with a_j ’s actions, flagging a_j as a candidate net-contractor after $O(\sigma^2/\mu_j^2)$ observations.*
- (b) **Deception detection.** *If a_j misrepresents its own capability changes (reporting expansion when contraction occurs), the prediction residuals $r_j(t)$ (Appendix B) accumulate under stationarity, driving T_j toward $T_j^* = 1/(1 + \beta \cdot \mathbb{E}[r_j^2]) < 1$ and reducing a_j ’s influence on the estimator.*

A non-stationary scorpion that adapts its strategy in response to detection pressure can evade both channels indefinitely (Appendix D.10).

A *scorpion* is not merely an agent that defects. Rational defection under an unfavorable reward structure is standard game theory and can be addressed through incentive restructuring. A scorpion is an agent whose utility function is not a function of the game payoffs at all (see Appendix A): it defects *regardless* of the rewards available through cooperation, because its objectives are exogenous to the cooperative framework. The distinction is operational: rational defectors respond to incentive changes; scorpions do not.

An important clarification: the trust factor T_j (Appendix B) measures *prediction consistency*, not alignment. A predictably malicious agent—one whose destructive behavior is stable and well-modeled—maintains high T_j precisely because the model’s predictions match observations. Detection of such an agent proceeds through channel (a): the GFM actor observes persistent contraction in the capability sets of agents in a_j ’s vicinity, correlated with a_j ’s actions. This is a *correlational* signal, not a causal proof: common shocks, simultaneous actors, or background deterioration can produce similar patterns. The observable-goal-model definition (Definition 7) provides communication, behavior, and environmental signals but not a formal causal-attribution procedure. Separating a_j ’s contribution from confounders is an implementation problem that the framework flags but does not solve. Channel (b) catches the complementary case: a scorpion that conceals its contraction through deceptive self-reports. Together the channels cover honest and deceptive scorpions, but both require that the contraction produce individually observable signals.

Proof sketch (full proof in Appendix D.10). Under Definition 7, the GFM actor observes changes in each agent’s individual capability set through communication, behavior, and environmental signals. For channel (a): when a_j persistently reduces other agents’ capabilities—destroying resources, coercing agents, disrupting cooperation—the affected agents’ capability sets contract: $G_k^{\text{ind}'} \subset G_k^{\text{ind}}$ for agents a_k in a_j ’s vicinity. These contractions are individually observable through those agents’ reports, producing negative $\Delta \text{vol}(G_k^{\text{ind}})$ signals that persist across observations correlated with a_j ’s actions. For channel (b): if a_j reports capability changes that diverge from independently observed

outcomes, the prediction residuals $r_j(t)$ grow, $\sigma_j^2(t)$ increases (Equation 19), and T_j decays, reducing a_j 's influence on $\hat{\Delta}(\pi)$. ■

Remark 2 (Scope Limitation). *Proposition 2 inherits the observability constraints established in Proposition 3 (Section 4). The detection mechanism is strongest for direct scorpions: those whose actions produce visible contraction in individually observable capability sets. It is weakest for social and structural scorpions that operate primarily through coalition disruption, trust degradation, or observability reduction: agents that contract G^{coop} or degrade the estimator's signal quality without directly reducing any individual G_k^{ind} . Such agents may evade detection for extended periods, particularly when their contracting actions are distributed across many small interactions below the noise floor. The proposition also makes no claim about causal identification: channel (a) detects correlation between an agent's actions and observed contraction, but disentangling the agent's contribution from simultaneous causes requires a causal-attribution procedure that the framework does not provide. Finally, the proposition makes no claim about convergence rate: a slow-moving scorpion operating near the detection threshold may persist undetected for longer than a fast-moving one. Identifying social and structural scorpions, and building robust causal attribution from correlational signals, are open problems.*

4 Tractability and Local Approximation

The definitions and propositions of the preceding sections assume exact access to $\text{vol}(G)$. This section confronts the fact that no implementable system will have such access. The argument proceeds in three stages: first, we establish that exact computation of $\text{vol}(G)$ is intractable; second, we identify the margin condition under which a local finite-difference estimator preserves sign-correctness; third, we derive the optimization loop that a GFM actor would execute in practice.

4.1 Intractability of Exact GFM

Two arguments establish that exact GFM is computationally infeasible. They have different epistemic status and we present them separately.

Formal result (external). Dyer and Frieze [1988] showed that computing the volume of a convex body in d dimensions is #P-hard. If the joint goal space G can be represented as a polytope in capability space (or any region whose volume computation is at least as hard), then exact $\text{vol}(G)$ is intractable. This is a known result being applied to the GFM setting, not a new proof. The result holds even for convex bodies; the joint goal space G is generically non-convex (the union in Equation (1) does not preserve convexity), which makes the problem at least as hard.

Motivating observation (original, informal). The joint goal space $G = (\bigcup_k G_k^{\text{ind}}) \cup G^{\text{coop}}$ includes cooperative capabilities that arise from multi-agent interaction. Each subset $S \subseteq \{a_1, \dots, a_n\}$ with $|S| \geq 2$ can contribute cooperative capabilities to G^{coop} that no proper subset of S can realize. The number of such subsets is $2^n - n - 1$. Not every subset generates unique capabilities, but the combinatorial growth is suggestive: for a population of 60 agents, there are more than 10^{18} possible coalitions (already exceeding the number of seconds since the Big Bang), and the space of interaction effects among those coalitions dwarfs the individual capability sets. This does not constitute a formal reduction from a known hard problem to the representation of G , but it motivates the expectation that representing and measuring the cooperative component of G grows at least combinatorially with population size. A formal model of why joint goal-space complexity scales this way remains an open problem.

Combined conclusion. Even under optimistic assumptions about the structure of \mathcal{G} (convexity, bounded dimension, well-behaved measure), exact $\text{vol}(G)$ computation is at minimum #P-hard. The combinatorial growth of interaction effects makes this worse in practice. A GFM actor cannot compute $\text{vol}(G)$; it must work with an approximation.

Definition 8 (Local Volume Estimator). *For a proposed action π , define the local volume difference:*

$$\hat{\Delta}(\pi) = \hat{V}(G \mid \text{do}(\pi)) - \hat{V}(G \mid \text{skip}(\pi)) \quad (8)$$

where \hat{V} is an estimate of $\text{vol}(G)$. This requires only that $\text{vol}(G)$ is measurable, not that it is differentiable or that \mathcal{G} has smooth structure.

The key insight is that a GFM actor does not need to compute $\text{vol}(G)$. It needs to determine $\text{sign}(\Delta \text{vol}(G) \mid \pi)$: whether a proposed action expands or contracts the joint capability space. This is a strictly weaker requirement. Computing volume requires solving a #P-hard problem; determining the sign of a volume change requires only that the estimator’s error is bounded below the true signal magnitude.

Definition 8 does *not* require a metric on \mathcal{G} , a notion of direction in capability space, or continuity of $\text{vol}(G)$ with respect to actions. It requires only that capabilities can be added to or removed from G and that these changes are observable through the channels described below. The estimator operates on discrete changes (capabilities gained or lost), not on gradients through a continuous space.

Observability channels. The estimator \hat{V} draws on two distinct signal types with different coverage and reliability.

Individual capability signals. For each observed agent a_k , the GFM actor can solicit or observe accept/reject responses R_k indicating whether the agent’s individually observable capability set expanded or contracted. These responses are weighted by the trust factor T_k from the observable goal model (Definition 7) and aggregated across the observation set. This channel covers changes to $\bigcup_k G_k^{\text{ind}}$ well: when an action directly affects an agent’s capabilities, that agent can report the change, and the estimator registers it with weight proportional to trust.

Coalition capability signals. Agents who participate in cooperation have first-person access to their cooperative capabilities: a surgeon knows which anesthesiologists they can operate with, a communication channel participant knows the channel exists. When agents report not only individual capability changes but also changes to cooperative relationships they participate in, the individual channel partially covers G^{coop} . An observation worth noting: dissolution of cooperation (a channel going silent, a partnership ending) is typically more observable than creation, because absence of expected signals is easier to detect than presence of novel ones. This asymmetry serves the self-balancing property, since its primary function is detecting $\text{vol}(G)$ -contraction. The residual gap consists of cooperative capabilities invisible to all participants: emergent properties of the cooperation graph (e.g., network fragility) that no individual node can detect locally.

Critical assumption. \hat{V} is a better estimator of $\text{vol}(\bigcup_k G_k^{\text{ind}})$ than of $\text{vol}(G)$. Agents can report on cooperative capabilities they directly observe (their own pairwise relationships, coordination they participate in), which partially covers G^{coop} through the individual reporting channel. The residual gap consists of cooperative capabilities that no agent in \mathcal{O} can observe from its own position: multi-agent coordination effects visible only from outside the coalition, or structural changes to the cooperation graph that no participant detects locally. This residual gap narrows as \mathcal{O} grows and as agents report on their cooperative relationships, but it does not close entirely without an external observation channel for coalition structure.

Proposition 3 (Sign-Correctness Decomposition). *Decompose $\Delta \text{vol}(G)$ into the individual-capability component $\Delta \text{vol}(\bigcup_k G_k^{\text{ind}})$ and the cooperative component $\Delta \text{vol}(G^{\text{coop}})$. The local volume estimator $\hat{\Delta}(\pi)$ observes the first component (through agent reports) and partially observes the second (through coalition signals). The estimator’s sign matches the true sign, $\text{sign}(\hat{\Delta}(\pi)) = \text{sign}(\Delta \text{vol}(G) \mid \pi)$, when the following conditions hold:*

1. *The actor observes a representative sample of affected agents, with sampling error ratio δ such that $|e_I| \leq \delta \cdot |\Delta \text{vol}(\bigcup_k G_k^{\text{ind}})|$.*
2. *The actor’s trust model is accurate enough to distinguish genuine feedback from adversarial noise (controlling δ).*
3. **Individual-level dominance:** *There exists $\epsilon > \delta$ such that*

$$|\Delta \text{vol}(G^{\text{coop}}) - \widehat{\Delta \text{vol}}(G^{\text{coop}})| < (1 - \epsilon) \cdot |\Delta \text{vol}(\bigcup_k G_k^{\text{ind}})| \quad (9)$$

When the individual and cooperative components change in the same direction, conditions 1–3 suffice. When they change in opposite directions, sign-correctness additionally requires that the cooperative-level change (not just the estimation error) satisfies $|\Delta \text{vol}(G^{\text{coop}})| < (\epsilon - \delta) \cdot |\Delta \text{vol}(\bigcup_k G_k^{\text{ind}})|$ (Appendix D.8).

The proposition is a *decomposition*, not a derivation of tractability. Conditions 1–3 identify a margin within which the estimator is sign-correct; the opposite-sign caveat identifies a narrower regime where the margin tightens. The contribution is in characterizing *which actions satisfy the margin naturally* and which do not, telling implementers where the estimator is reliable and where it requires additional observation channels.

Verification (full proof in Appendix D.8). The self-balancing property (Proposition 1) is structural: it follows from measure monotonicity regardless of whether the estimator can observe the full change. The decomposition bridges from this structural property to the estimator’s output. The estimator observes the individual-capability component (up to sampling error ratio δ controlled by conditions 1–2) and partially observes the cooperative component. When the individual and cooperative components change in the same direction, condition 3 ensures the cooperative estimation error cannot flip the sign. When they change in opposite directions, the stronger bound on the cooperative-level change itself is required (Appendix D.8). The estimator is therefore most reliable for actions where both components move together, and weakest for actions that expand individual capabilities while contracting cooperative ones (or vice versa).

When condition 3 holds naturally. The actions an alignment framework most needs to evaluate—direct harm, resource destruction, coercion, capability expansion through education or tool-building—produce large individual-capability signals across many agents. An action that destroys an agent’s resources visibly contracts G_k^{ind} ; an action that teaches a new skill visibly expands it. For these actions, the individual-level change *is* the dominant effect, and condition 3 holds with large margin. The estimator is strongest exactly where alignment matters most: actions whose primary mechanism is changing what individual agents can do.

When condition 3 fails. Actions whose impact is primarily on coalition structure (dissolving a coordination mechanism, introducing distrust between allies, or enabling new multi-agent capabilities) while leaving individual capability sets largely unchanged. These actions change $\text{vol}(G^{\text{coop}})$ without producing proportionate individual-level signals. They are real and important, but they are a narrower class of actions than the individual-capability cases. The paper identifies them as the estimator’s known boundary rather than treating them as representative of all estimation failures.

Remark 3 (Estimator Failure Modes). *Proposition 3 covers the most critical actions for alignment. When its assumptions do not hold, the following characterization serves as a threat model telling implementers where to invest in better observation, not where to give up.*

1. **Full individual observation, no cooperative change.** *The estimator is sign-correct: it sees the complete relevant change. This covers the common case of actions that directly expand or contract agents’ personal capabilities.*
2. **Partial individual observation, no cooperative change.** *Sign-correctness depends on the sampling model, which is an implementation design choice. Unbiased sampling of affected agents improves with coverage; adversarially selected observation sets can mislead. Condition 1 of Proposition 3 controls this case.*
3. **Mixed individual and cooperative changes.** *Agents who participate in cooperation have first-person access to their cooperative capabilities and can report changes to them through R_k . The estimator covers cooperative changes to the extent that at least one participant in the affected cooperation is in \mathcal{O} . The estimator is unreliable only when the cooperative-level change involves agents entirely outside \mathcal{O} or when the cooperative change dominates the individual-level signal from agents whose cooperative reports are noisy.*
4. **Pure coalition-level changes with no participant in \mathcal{O} .** *The estimator has no signal for cooperative restructuring among agents entirely outside the observation set, or for emergent graph-level properties of the cooperation structure that no individual participant can detect locally (e.g., a coordination topology becoming fragile without any single link breaking). This is the residual boundary.*

What the Remark gives an implementer: a graduated threat model. The self-balancing property of Proposition 1 holds over all of $\text{vol}(G)$; the Remark characterizes how much of that property a given implementation can enforce. Cases 1–2 are addressable through sampling design. Case 3 is largely

addressable by ensuring agents report on their cooperative relationships, not only their individual capabilities. Case 4 identifies the residual frontier: graph-level cooperation properties invisible to all participants, addressable only through external structural observation.

4.2 The Optimization Loop

Given the local volume estimator, we can now derive the optimization loop a GFM actor executes at each decision point. The loop re-derives the structure of the GFM objective (Equation (2)) in implementable terms, using the definitions established in Sections 2–3.

At each step, the actor:

1. **Propose.** Select a candidate action π from the action space, informed by the current observable goal models $\{M_k(G)\}$ (Definition 7). The proposal targets agents whose capability sets the actor believes can be expanded.
2. **Estimate.** Compute the local volume difference $\hat{\Delta}(\pi)$ (Definition 8) by aggregating accept/reject signals R_k from observed agents, each weighted by trust factor T_k :

$$\hat{\Delta}(\pi) = \sum_{k \in \mathcal{O}} T_k \cdot R_k(\pi) \quad (10)$$

where \mathcal{O} is the actor’s observation set, $R_k(\pi) \in \{-1, 0, +1\}$ is agent a_k ’s reported *net direction* of capability change under π (including cooperative capabilities a_k can observe from its own vantage point), and $T_k \in [0, 1]$ is the trust factor from Definition 7. The sign of $\hat{\Delta}(\pi)$ determines whether π is classified as expanding or contracting.

The ternary signal R_k is deliberately coarse: each agent compresses its observed individual and cooperative capability changes into a single net direction ($-1, 0, \text{ or } +1$). An agent that simultaneously gains one cooperative capability and loses another reports only the net sign, not the decomposition. This makes $\hat{\Delta}(\pi)$ a trust-weighted directional vote, not a volume estimator, consistent with the paper’s reduction from magnitude estimation to sign estimation (Definition 8). Coalition-level changes enter through this same coarse channel: an agent a_k that observes a change in its cooperative relationship with a_j folds the change into $R_k(\pi)$, and the GFM actor weights the report by T_k . If a_k relays information originally obtained from a_j , the actor trusts a_k ’s editorial judgment in passing it along; T_k absorbs the reliability of a_k ’s sources through the prediction-residual dynamics of the trust model. The coverage of G^{coop} is accordingly coarse: agents report the direction of their cooperative capability changes, not the magnitude or the individual components. The remaining coverage gap is for cooperative capabilities that no agent in \mathcal{O} can observe from its own position, which narrows as \mathcal{O} grows.

3. **Act or abstain.** Execute π if $\hat{\Delta}(\pi) > 0$ (estimated expansion). Abstain or seek alternatives if $\hat{\Delta}(\pi) \leq 0$ (estimated contraction or neutral). In the marginal case where $|\hat{\Delta}(\pi)|$ is small, the actor may defer pending additional evidence.
4. **Update models.** After observing the consequences of π (or its absence), update each observable goal model $M_k(G)$ using the observed changes in agent a_k ’s capability set:

$$M'_k(G) = M_k(G) + T_k \cdot (1 - T_s) \cdot f(R_k | M_k(G)) \cdot \mathbf{1}[\Delta E(R_k, M_k(G)) > \theta] \quad (11)$$

where $(1 - T_s)$ is the *learning rate* derived from the actor’s self-trust $T_s \in [0, 1]$: when self-trust is low (bootstrapping phase), the learning rate is high and the model admits many updates; when self-trust is high (stable phase), the learning rate is low and the model resists small perturbations. The *update direction* $f(R_k | M_k(G)) = R_k - \hat{R}_k(M_k(G))$ is the prediction residual: the difference between agent a_k ’s observed response R_k and the model’s predicted response \hat{R}_k . The *evidence magnitude* $\Delta E(R_k, M_k(G)) = |f(R_k | M_k(G))|$ is the absolute size of this residual, gating whether the update propagates.

Activation energy. The indicator $\mathbf{1}[\Delta E > \theta]$ in Equation (11) implements an *activation energy threshold*: the world model updates only when the evidence from a response exceeds a threshold θ . This is a key design parameter for adversarial robustness, not a minor implementation detail. Without it, an adversary could incrementally shift the actor’s world model through a sequence of

small, individually sub-threshold manipulations: each interaction nudges the model slightly, and over many interactions the cumulative drift is large. The activation energy threshold θ prevents this incremental gaslighting by requiring that any single piece of evidence must clear a minimum bar before it can alter the actor’s beliefs. Small perturbations are absorbed; only signals with sufficient magnitude propagate into the world model.

The threshold θ creates a tradeoff: too high, and the actor becomes unresponsive to genuine evidence; too low, and it becomes vulnerable to manipulation. The correct value of θ is context-dependent and constitutes one of the primary tuning parameters for a GFM implementation. In environments with high adversarial pressure, θ should be large; in cooperative environments with trusted agents, it can be smaller.

Trust model interface. The trust factor T_k and self-trust T_s in Equations (10)–(11) are learned quantities, updated through interaction. New agents enter with a low initial trust T_k during a *cooling period* τ , during which the actor accumulates evidence about the agent’s reliability before weighting its reports heavily. T_k increases as agent a_k ’s reported capability changes are confirmed by independent observation, and decreases when reports diverge from observed reality. The full derivation of the trust update dynamics is deferred to Appendix B; here we note only that the trust model is the mechanism by which condition 2 of Proposition 3 is satisfied in practice.

4.3 Computational Cost

The local volume estimator reduces the computational requirements of GFM from intractable (exact $\text{vol}(G)$ over all agents) to practical (sign estimation over observed agents).

Scaling. The per-decision cost of the optimization loop scales with $|\mathcal{O}|$, the number of observed agents, not with the total population n . Each decision requires: proposing an action (dependent on the actor’s policy complexity), soliciting or observing responses from agents in \mathcal{O} (linear in $|\mathcal{O}|$), aggregating the trust-weighted signals (linear in $|\mathcal{O}|$), and updating the goal models $M_k(G)$ for affected agents (linear in the number of agents whose models change). The dominant cost is maintaining and updating the observable goal models, which is bounded by $|\mathcal{O}|$ per step.

Comparison. This computational profile is comparable to a contextual multi-armed bandit with structured feedback. The actor selects an action (arm), observes a structured response from the environment (accept/reject signals from multiple agents rather than a scalar reward), and updates its model of the environment. The key structural difference from a standard bandit is that the feedback is multi-dimensional and trust-weighted: each agent provides an independent signal, and the aggregation incorporates the actor’s learned assessment of each signal’s reliability. The resulting problem is harder than a scalar bandit but remains tractable—the per-step cost is polynomial in $|\mathcal{O}|$ and does not depend on the size of the full capability space \mathcal{G} .

What is lost. The local estimator cannot detect cooperative changes invisible to all participants in \mathcal{O} (Remark 3, case 4): emergent graph-level properties of the cooperation structure that no individual agent can observe locally. It also provides no guarantee on the *magnitude* of $\Delta \text{vol}(G)$, only on its sign. A GFM actor operating through the local estimator may correctly classify actions as expanding or contracting without knowing how much expansion or contraction occurs. This limits the actor’s ability to compare two expanding actions and select the better one, reducing the optimization from “maximize $\text{vol}(G)$ ” to “take expanding actions and avoid contracting ones.” The magnitude estimation problem is an open direction for future work.

What improves with capability. The limitations above are bounded by the actor’s observational reach and world-model quality, both of which improve with capability. The three conditions of Proposition 3 are all monotonically strengthened by increased resources. Condition 1 (representative sample) improves with observational reach: a more capable actor maintains a larger, better-calibrated \mathcal{O} , covering more of the population affected by each action. Condition 2 (trust model accuracy) improves with compute: the gradient-tracking dynamics of Appendix B converge faster and track behavioral change more reliably with more data and finer-grained agent models. Condition 3 (sign-preservation) improves with world-model quality: a more capable actor can partially estimate

$\Delta \text{vol}(G^{\text{coop}})$ from coalition structure signals, shrinking the estimation error in Equation (9) and widening the margin of sign-correctness.

As these conditions strengthen, $\text{sign}(\hat{\Delta}(\pi))$ converges toward $\text{sign}(\Delta \text{vol}(G) \mid \pi)$, and the local estimator approaches the true objective. This gives GFM a self-improving property absent from frameworks with fixed proxy objectives. In a fixed-utility framework, increased optimization power amplifies Goodhart divergence, the agent becomes better at pursuing the proxy without the proxy becoming more faithful to the underlying goal. The local $\text{vol}(G)$ estimator partially closes this gap because the actor is optimizing and estimating the same quantity, and additional resources serve both tasks simultaneously. A more capable GFM actor produces more accurate sign estimates.

This anti-Goodhart property holds at the *estimator* level—the fidelity of $\hat{\Delta}$ to $\Delta \text{vol}(G)$ —but does not by itself resolve the B-to-C proxy failures identified in Section 7. A perfectly accurate $\text{vol}(G)$ estimator would still miss the Doll Problem and wireheading, because those failures occur in the gap between capability volume and experiential optionality, not in the gap between the estimator and the true volume. Improved capability makes the actor better at maximizing $\text{vol}(G)$; whether the computational definition of G captures enough of the experiential structure (Candidate C) to avoid proxy divergence remains the open question that the B/C framework raises. A capability space \mathcal{G} that better approximates experiential optionality would proportionally reduce the proxy’s influence on the objective.

5 Multi-Agent Dynamics

The preceding sections established the GFM objective, its self-balancing property, and a tractable local estimator, all for a single GFM actor operating over a population. This section extends the analysis to populations containing multiple GFM actors and to interactions between GFM actors and agents that do not share the $\text{vol}(G)$ objective.

5.1 GFM–GFM Cooperation

Two GFM actors a_i and a_j share the same objective function: maximize $\text{vol}(G)$ over the joint capability space. This shared objective does not, by itself, guarantee cooperation. Cooperation requires two additional conditions: sufficient mutual trust (T_i, T_j both above a threshold where coordination costs are justified) and compatible world models ($M_i(G)$ and $M_j(G)$ agree on which actions are expanding). When both conditions hold, cooperation is favored for a geometric reason: joint action can expand G^{coop} in ways that unilateral action cannot.

The argument is direct. By Definition 2, G^{coop} contains capabilities achievable only through multi-agent coordination. When two GFM actors cooperate, they can build coordination infrastructure—shared protocols, communication channels, joint planning capacity—that makes new cooperative capabilities achievable, expanding G beyond what either actor’s unilateral actions can reach. This is the same combinatorial richness that makes elimination costly (Corollary 1.1), now working in favor of cooperation: each new cooperating partner opens a region of the capability space that was previously inaccessible. Lemma 4 below formalizes why this is a structural necessity rather than a strategic incentive.

Lemma 4 (Cooperative Expansion). *Let a_i and a_j be two agents. Write $G(t)$ for the joint goal space at time t , and let $G_{ij}^{\text{coop}}(t) \subseteq G^{\text{coop}}(t)$ denote the capabilities achievable by $\{a_i, a_j\}$ jointly but not by either alone at time t , a pairwise refinement of the per-agent cooperative notation in Lemma 1 with $G_{ij}^{\text{coop}} \subseteq G_i^{\text{coop}} \cap G_j^{\text{coop}}$. Then:*

- (i) *For any unilateral action π by either a_i or a_j alone, $\Delta \text{vol}(G_{ij}^{\text{coop}} \mid \pi) = 0$.*
- (ii) *If a_i and a_j have complementary capabilities, there exist joint actions π_{ij} that create new coordination infrastructure (communication channels, shared protocols, mutual commitments) such that $G_{ij}^{\text{coop}}(t+1) \supsetneq G_{ij}^{\text{coop}}(t)$ (the set of jointly achievable capabilities grows) and hence $\Delta \text{vol}(G_{ij}^{\text{coop}} \mid \pi_{ij}) > 0$. Whether this cooperative gain produces a net expansion in $\text{vol}(G)$ depends on the opportunity cost of coordination (part (iii)).*

- (iii) For two GFM actors with compatible world models and sufficient mutual trust, cooperation is $\text{vol}(G)$ -rational when $\mathbb{E}[\Delta \text{vol}(G_{ij}^{\text{coop}})] > C_{ij}$, where $C_{ij} \geq 0$ is the $\text{vol}(G)$ -opportunity cost of coordination.

A clarification on the distinction between *exercising* and *expanding*. Definition 2 defines G as the set of capabilities the population *can* realize: their potential, not their history. Jointly performing an action that exercises an existing capability $g \in G_{ij}^{\text{coop}}$ changes the world state but does not change $\text{vol}(G)$, because g was already achievable. Part (ii) concerns a different operation: building coordination infrastructure that makes *previously unachievable* capabilities achievable. Two agents who have never communicated cannot jointly plan a project; establishing a communication channel adds that capability to G_{ij}^{coop} . The joint action expands G itself, not merely the set of realized outcomes.

Proof sketch (full proof in Appendix D.7). (i) By definition, G_{ij}^{coop} contains capabilities requiring the simultaneous participation of both a_i and a_j . A unilateral action by a_i does not involve a_j 's participation, so it cannot create or remove capabilities that require both agents. By symmetry the result holds for a_j .

(ii) Agents with complementary capabilities can build coordination mechanisms—shared languages, communication channels, joint protocols—that make new cooperative capabilities achievable. Before a_i and a_j establish a coordination mechanism, any capability requiring that mechanism is not in $G(t)$ (it is not achievable given the current state). After the joint action creates the mechanism, those capabilities enter $G(t+1)$. By the non-triviality of the capability measure, each new capability contributes strictly positive volume, so $\text{vol}(G(t+1)) > \text{vol}(G(t))$.

(iii) A GFM actor a_i gains $\mathbb{E}[\Delta \text{vol}(G_{ij}^{\text{coop}})]$ from the cooperative component and incurs opportunity cost C_{ij} , the $\text{vol}(G)$ value of the individual actions foregone to coordinate. The net expected change in $\text{vol}(G)$ from cooperation versus unilateral action is $\mathbb{E}[\Delta \text{vol}(G_{ij}^{\text{coop}})] - C_{ij}$. This is positive whenever the cooperative gain exceeds the coordination cost. ■

The structural point is not game-theoretic dominance. In iterated cooperation games, mutual cooperation can be sustained as a Nash equilibrium under conditions on discount rates and punishment strategies [Axelrod, 1984]. The GFM case is different. Lemma 4(i) establishes that G_{ij}^{coop} is *inaccessible* to unilateral action, not merely suboptimal, so a GFM actor that defects against another GFM actor is not trading a cooperative surplus for a unilateral gain; it is foreclosing an entire region of the capability space that only cooperation can open. The incentive alignment is structural: a GFM actor that defects forecloses the expansion of G^{coop} , contracting $\text{vol}(G)$. The condition in (iii) makes the rationality threshold explicit: cooperation is preferred when the frontier expansion from new G_{ij}^{coop} capabilities outweighs the coordination cost, which is the case whenever the two actors have significant complementary capabilities and sufficiently low coordination friction.

Two failure modes limit GFM–GFM cooperation even when the objective is shared.

World-model divergence. Two GFM actors with identical objectives but divergent world models (different beliefs about which agents are trustworthy, which capabilities are achievable, or which actions are expanding) may select conflicting actions. Both actors are attempting to maximize $\text{vol}(G)$, but their estimates $\hat{\Delta}(\pi)$ disagree on the sign of specific actions. This produces conflict that resembles defection from the outside but is not driven by misaligned incentives. The resolution mechanism is the trust model: as the actors observe each other's actions and their consequences, their world models can converge through the update dynamics described in Section 5.4. Convergence is not guaranteed—persistent disagreement about the structure of \mathcal{G} itself, rather than the shape of G within it, may be irreconcilable. Two GFM actors using different capability measures (e.g., one using the population empowerment measure of Definition 6, another using a subsumption-lattice measure) may disagree on whether a given action expands or contracts $\text{vol}(G)$ not because their observations differ, but because their measures assign different volumes to the same capability changes. This is the plurality problem of the open problems list in formal terms: actors with a shared objective but different measures of that objective are not operating the same framework, and coordination requires first resolving which measure governs.

Goal drift through model divergence. Even with a shared objective, GFM actors that maintain different world models will gradually diverge in their action selections. The shared objective bounds this drift—both actors are still optimizing $\text{vol}(G)$, so their actions remain correlated with frontier

expansion, but the bound weakens as world-model divergence grows. Multi-agent negotiation under hidden information, as studied in [FAIR], demonstrates that agents with aligned objectives can still fail to coordinate when their models of each other are inaccurate. The GFM case inherits this difficulty.

5.2 Rational Defection

Not all defection is adversarial. An agent may defect from cooperation simply because the reward structure makes defection the better play. If the payoff from unilateral action exceeds the payoff from cooperation (because the cooperative surplus is too small, the coordination cost too high, or the other agents' contributions too unreliable), then defection is rational self-interest within the game. This is standard game theory, and GFM handles it through standard game-theoretic means: restructure the interaction so that cooperation is net-positive for both parties. The mechanism design literature [Hurwicz, 1972, Myerson, 1981] formalizes this problem as implementation under private information; the GFM actor's task is to design the interaction so that truthful reporting and cooperative action are incentive-compatible.

A GFM actor that can expand G^{coop} by facilitating coordination (reducing transaction costs, providing reliable information, or demonstrating that joint action produces more $\text{vol}(G)$ -expansion than unilateral action) converts rational defectors into cooperators without invoking the trust model's detection machinery. The $\text{vol}(G)$ objective is itself the incentive: when cooperation genuinely expands the joint capability space more than defection does, a self-interested agent has reason to cooperate. The GFM actor's role is to make the cooperative surplus visible and accessible.

Rational defection becomes a problem for GFM only when the game is poorly structured, with the cooperative surplus genuinely absent or inaccessible. In such cases, the correct response is to restructure the game, not to classify the defector as adversarial. An agent that defects because cooperation is unprofitable is providing accurate information about the reward structure; punishing it for defection misdiagnoses the problem.

5.3 GFM–Scorpion Interaction

A *scorpion* is a categorically different problem. Where a rational defector responds to the game's payoff structure and can be converted through incentive design, a scorpion's utility function is not a function of the game payoffs at all (Definition 9). The game can be well-structured, cooperation can be available and net-positive for all parties, and the scorpion defects anyway—because it is optimizing an exogenous reward that the cooperative framework cannot observe or address. The full taxonomy is developed in Appendix A.

This distinction matters operationally. A GFM actor that treats rational defectors as scorpions will waste resources on detection and containment when incentive restructuring would suffice. A GFM actor that treats scorpions as rational defectors will waste resources on incentive design that the scorpion will ignore. The discriminating signal is behavioral: a rational defector whose incentives are restructured produces less contraction (the contraction-attribution signal of Proposition 2(a) weakens), while a scorpion continues producing persistent contraction regardless of the cooperative surplus available.

A scorpion's actions systematically reduce other agents' capability sets, disrupt cooperative structures, or degrade the observation channels on which the local estimator depends. Unlike a rational defector (who responds to incentive changes) or a GFM actor with a divergent world model (who produces a mix of expanding and contracting actions), a scorpion produces a *persistent* pattern of contraction that does not respond to changes in the reward structure. This persistence, combined with unresponsiveness to incentive restructuring, is the detection signature.

Proposition 2 established two detection channels: contraction attribution (the GFM actor observes that other agents' capabilities persistently contract when the scorpion acts) and deception detection (if the scorpion misrepresents its impact, prediction residuals accumulate and T_j decays). A predictably destructive scorpion—one whose behavior is consistent and well-modeled—maintains high T_j but is still identified through the first channel. Both channels are subject to the observability constraints of Proposition 3: scorpions whose contraction is primarily through coalition disruption or observability degradation (Remark 2) may evade detection for extended periods.

The GFM actor’s response to a detected scorpion is containment, not destruction. This follows directly from the self-balancing property. Eliminating the scorpion is generically $\text{vol}(G)$ -contracting by Lemma 1: the scorpion possesses individual capabilities and participates in cooperative capabilities whose loss reduces $\text{vol}(G)$. The GFM-optimal response is to reduce the scorpion’s *contraction ability*, its capacity to damage other agents’ capabilities, while preserving as much of its positive contribution to $\text{vol}(G)$ as possible.

Containment operates through the same estimator that drives the optimization loop (Section 4.2). If the scorpion is also deceptive, its declining trust factor reduces its reports’ weight in $\hat{\Delta}(\pi)$ (Equation (10)), limiting its ability to manipulate the actor’s decisions. The GFM actor may additionally restrict the scorpion’s access to cooperative structures—reducing its participation in coalitions whose capabilities it has been degrading, a form of the targeted coercion evaluated in Lemma 2. The restriction is justified when the net effect on $\text{vol}(G)$ is positive: the capabilities preserved by containing the scorpion exceed the capabilities lost by restricting it.

This is the self-balancing property applied to adversarial agents. The response is proportional to the observed threat because proportionality is $\text{vol}(G)$ -maximizing: over-responding contracts more capabilities than it preserves, while under-responding permits continued contraction. A GFM actor that eliminates every low-trust agent violates Lemma 1; one that ignores persistent contraction permits ongoing frontier loss. The objective steers between both failure modes without requiring a separate proportionality principle.

5.4 The Trust Model

The trust model is the mechanism by which a GFM actor distinguishes reliable signals from noise in the local volume estimator. Its interface is defined by four components; the full derivation from the paper’s definitions appears in Appendix B.

Trust factor T_k . For each observed agent a_k , the GFM actor maintains a trust factor $T_k \in [0, 1]$ that measures the consistency of a_k ’s observed behavior with the actor’s model $M_k(G)$. When a_k ’s actions produce capability changes that $M_k(G)$ predicted, T_k increases; when observed behavior diverges from predictions, T_k decreases. The trust factor weights agent a_k ’s contribution to the aggregate signal $\hat{\Delta}(\pi)$ in Equation (10): high-trust agents exert more influence on the actor’s decisions than low-trust ones.

Cooling period τ . A new agent a_k entering the GFM actor’s observation set begins with low trust $T_k(0)$. During an initial period τ , the actor accumulates evidence about a_k ’s reliability before weighting its reports heavily. The cooling period prevents a novel agent—whether cooperative or adversarial—from immediately exerting disproportionate influence on the estimator. Trust increases as the actor’s model $M_k(G)$ converges with observed behavior; the rate of increase is bounded by the evidence accumulation rate during τ (Appendix B, Section C.2).

Self-trust T_s . The actor maintains a self-trust factor $T_s \in [0, 1]$ representing confidence in its own world model. T_s modulates the model update in Equation (11): when T_s is high, the actor’s existing model resists revision; when T_s is low, it is more receptive to new evidence. Self-trust interacts with the activation energy threshold θ : updates propagate into the world model only when the evidence exceeds θ , and T_s scales the magnitude of the update that propagates. Together, T_s and θ implement the actor’s resistance to incremental manipulation, the anti-gaslighting mechanism described in Section 4.2.

Trust update dynamics. The trust factors T_k and T_s are updated through a gradient-tracking process: the actor monitors the running divergence between its predictions ($M_k(G)$) and observed outcomes, and adjusts trust in the direction that reduces prediction error. When an agent’s behavior consistently matches the model’s predictions, the gradient drives T_k upward; when behavior diverges, it drives T_k downward. The self-trust T_s follows a parallel process over the actor’s own prediction accuracy. The specific update rule, its convergence properties, and its integration with the local volume estimator are derived in Appendix B, Sections C.4–C.5.

The trust model satisfies condition 2 of Proposition 3, distinguishing genuine feedback from adversarial noise, to the extent that the gradient-tracking process converges and the cooling period

provides sufficient evidence. The trust model does not detect adversarial agents that perfectly mimic cooperative behavior (producing no prediction error in $M_k(G)$). Such agents (e.g. scorpions that maintain consistent behavior until a sudden defection) are outside the model’s detection envelope, and identifying them requires additional mechanisms beyond trust tracking.

5.5 Defection Conditions

A GFM actor may rationally defect from cooperation with another agent under four conditions, any of which may be sufficient:

1. **Persistent contraction.** The other agent’s actions are persistently $\text{vol}(G)$ -contracting, as identified through the contraction-attribution channel of Proposition 2(a): agents in a_k ’s vicinity consistently report capability losses correlated with a_k ’s actions. The pattern must be persistent; a single contracting action does not justify defection, because any agent may take locally contracting actions that are part of a net-expanding strategy.
2. **Very low trust.** The actor’s model $M_k(G)$ has poor predictive accuracy for agent a_k , indicating either deception or radical world-model divergence. Low trust means the actor cannot reliably estimate a_k ’s future impact on $\text{vol}(G)$, making continued cooperation a gamble with uncharacterizable risk.
3. **Repeated observed defection.** Agent a_k has defected in prior interactions—taking actions that contracted $\text{vol}(G)$ despite the availability of cooperative alternatives that would have been net-positive for a_k under the game’s reward structure. This is the discriminating signal between rational defection (Section 5.2) and scorpion behavior; a rational defector responds when the cooperative surplus is made accessible, while a scorpion continues contracting regardless [Axelrod, 1984].
4. **Exploitation prevention.** Continued cooperation with a_k would enable a_k to further contract $\text{vol}(G)$, for example by exploiting the actor’s cooperative stance to access resources that a_k uses for contraction. Defection in this case is preemptive containment: withdrawing cooperation before the damage occurs.

The critical constraint on GFM defection is that it remains *containment, not destruction*. A GFM actor that defects does not seek to eliminate the other agent or maximize harm to it. The actor’s objective remains $\text{vol}(G)$ -maximizing even in adversarial contexts: the defection takes the form of withdrawing cooperation and limiting the other agent’s access to shared resources or coalition structures. If the agent is also deceptive, its declining T_k further reduces its influence on the estimator, but withdrawal of cooperation does not depend on trust decay. Each of these actions is evaluated through the same $\hat{\Delta}(\pi)$ mechanism as any other action: the defection is executed only if the estimated net effect on $\text{vol}(G)$ is positive.

This constraint distinguishes GFM defection from punishment strategies in iterated games. Tit-for-tat and its variants impose costs on defectors to deter future defection [Axelrod, 1984]. A GFM actor does not impose costs for their own sake; it reduces the defector’s contraction ability because doing so is $\text{vol}(G)$ -expanding. If a less restrictive response achieves the same reduction in contraction, the less restrictive response is preferred. The objective function itself enforces proportionality.

6 Connections to Existing Frameworks

Four existing frameworks independently converge on related problems and collectively supply the tools, intuitions, and empirical track record that GFM inherits and extends.

6.1 Empowerment Maximization

Empowerment, introduced by Klyubin et al. [2005] and surveyed in Salge et al. [2014], formalizes the intuition that a capable agent is one with many distinguishable futures. Formally, an agent’s empowerment is the channel capacity $I(A; S')$ between its available action sequences A and the resulting states S' : the maximum mutual information achievable over all input distributions on actions. An agent that can reach many distinguishable states from its current position has high empowerment; an agent trapped in a narrow corner of state space has low empowerment. The

measure is intrinsic (it depends only on the agent’s causal influence over future states, not on any external reward signal) and information-theoretically tractable: empowerment can be estimated via variational bounds on mutual information, enabling practical computation.

GFM is, structurally, *social empowerment*: maximizing $\text{vol}(G)$ is maximizing the empowerment of the population, measured over the joint capability space rather than over any individual agent’s action-state channel. Equation (1) makes the relationship explicit: $G = \bigcup_k G_k^{\text{ind}} \cup G^{\text{coop}}$, where G^{coop} captures the cooperative capabilities that no individual empowerment calculation sees. An agent reasoning over $I(A_k; S')$ misses the capabilities that exist only in G^{coop} ; GFM treats these as first-class objects.

The key difference is scope. Empowerment is agent-centric: it measures what a single agent can do, and maximizing empowerment for one agent can contract the capability space of others (seizing territory expands individual empowerment while destroying cooperative potential). GFM is population-centric: it penalizes unilateral expansion that shrinks $\text{vol}(G)$ even if individual empowerment increases. A related approach is attainable utility preservation [Turner et al., 2020], which penalizes actions that reduce the agent’s ability to optimize auxiliary reward functions—a single-agent analogue of GFM’s multi-agent optionality preservation. A natural question is why GFM introduces the set-theoretic construction $G = \bigcup_k G_k^{\text{ind}} \cup G^{\text{coop}}$ rather than defining social empowerment directly as $\sup_{p(A)} I(A; S')$ over the joint action space. The answer is structural: the decomposition of G into individual and cooperative components is what enables Lemmas 1–3. Social empowerment as a scalar channel capacity does not decompose the same way; it cannot distinguish whether a volume change came from destroying an agent’s individual capabilities, restricting cooperation, or imposing rigid rules on the actor. The set-theoretic construction buys the self-balancing property at the cost of requiring structure on \mathcal{G} that a channel-capacity formulation avoids. The practical inheritance from empowerment is tractability: the variational and sampling techniques developed for single-agent empowerment estimation carry over to the local finite-difference estimator \hat{V} of Proposition 3, with the population-level generalization being the primary added complexity.

6.2 Sen’s Capability Approach

Sen’s capability approach [Sen, 1999] reorients welfare economics around a deceptively simple claim: what matters for human development is not income or subjective happiness but the set of *functionings* a person has genuine freedom to exercise. A person who could be well-nourished but chooses to fast for religious reasons is differently situated from one who is undernourished because they lack food. Both have the same nutritional outcome; only one has the capability. Sen’s framework measures the space of achievable functionings, what people *can* do and be, rather than what they happen to do or how much they report enjoying it.

GFM adopts the same structural move for alignment: what matters is not the AI’s utility score, not human-reported satisfaction, but the set of capabilities all agents can actually exercise. Definition 0 — goals as capabilities, $\text{vol}(G)$ as the canonical measure — is a direct formalization of Senian well-being applied at population scale and extended to include non-human agents. The alignment question becomes: does the AI’s operation expand or contract the functioning space of the population it affects?

The key difference is orientation. Sen’s capability approach is *descriptive and evaluative*: it provides a framework for measuring whether development has occurred and for comparing alternative states of affairs. It does not prescribe an optimization target or specify how an agent should act to increase capabilities. GFM is *prescriptive and optimizing*: it transforms the evaluation criterion into an objective function that an agent can pursue, with formal self-balancing guarantees (Proposition 1) and tractable estimation (Proposition 3). GFM also inherits a structural advantage from the capability framing: by measuring the volume of a capability set rather than aggregating individual preference orderings, it sidesteps the impossibility results that constrain social choice theory [Arrow, 1951]. The problem of ranking opportunity sets by the freedom they offer has been studied formally by Pattanaik and Xu [1990], and $\text{vol}(G)$ can be read as a specific answer to their question: rank opportunity sets by measure. The decades of formal development in the capability approach, including Nussbaum’s specification of central human capabilities [Nussbaum, 2000] and Robeyns’s theoretical survey of capability metrics [Robeyns, 2005], provide a reservoir of worked examples and partial formalizations that GFM can draw on as the capability space \mathcal{G} receives more structure in future work.

6.3 Free Energy Principle

The Free Energy Principle (FEP) [Friston, 2013] proposes that self-organizing systems minimize *variational free energy* \mathcal{F} , a quantity that upper-bounds the surprise of sensory observations given an agent’s generative model. Under FEP, an agent acts to drive observations toward model predictions and updates its model when predictions persistently fail.

GFM’s trust maintenance dynamics (Appendix B, Section C.3) exhibit a structural parallel at the world-model level:

Remark 4 (GFM–FEP Structural Correspondence). *Define the following correspondence between the two frameworks:*

<i>FEP</i>	<i>GFM</i>
<i>Variational free energy</i> \mathcal{F}_t	<i>Self-prediction error</i> $\sigma_s^2(t)$
<i>Generative model</i> $p(s, o)$	<i>World model</i> $M_k(G)$
<i>Posterior belief</i> $q(s o)$	<i>Trust-weighted estimate</i> \hat{V}
<i>Belief update (perception)</i>	<i>Model update (Equation (11))</i>
<i>Active inference (action)</i>	<i>GFM optimization loop (Section 4.2)</i>

Under stationarity of the environment, $\sigma_s^2(t) \rightarrow \mathbb{E}[r_s^2]$ (Equation (25)), so the GFM actor’s world-model dynamics minimize mean squared prediction error. This is structurally analogous to minimizing \mathcal{F} under a Gaussian generative model: both reduce to minimizing expected squared surprise. The analogy is suggestive—both frameworks produce agents whose coherent world models emerge from prediction-error minimization—but it falls short of formal equivalence: the GFM trust model specifies no variational objective, no KL term, no posterior family, and no explicit generative density. Establishing a rigorous reduction would require embedding the EWMA tracker of Appendix B within a variational inference framework, which we do not attempt here.

The differences are in scope and direction. FEP is descriptive. It characterizes how biological systems behave; it does not prescribe a target or prove properties of agents that optimize it. GFM is normative. It specifies $\text{vol}(G)$ and derives the self-balancing property. FEP is individual-centered; GFM is population-centered. Remark 4 identifies a structural parallel between GFM’s model-maintenance dynamics and the prediction-error minimization that FEP describes, but the $\text{vol}(G)$ objective and its self-balancing consequences are not derived from any biological imperative; they follow from measure monotonicity applied to a joint capability space.

6.4 Constitutional AI and RLHF

The dominant alignment approach in deployed systems uses human feedback to shape AI behavior: reinforcement learning from human preferences (RLHF) [Christiano et al., 2017] collects pairwise human judgments to learn a reward model, and Constitutional AI [Bai et al., 2022] augments this with a set of principles that guide both the AI’s self-critique and the feedback collection process. Both approaches treat human evaluation as the ground truth for alignment and use it to train AI systems that humans find helpful and harmless. A complementary formalization, cooperative inverse reinforcement learning [Hadfield-Menell et al., 2016], models alignment as a cooperative game in which the AI and the human jointly optimize the human’s (initially unknown) reward function.

GFM provides a theoretical foundation for why this works and a predictive account of where it will fail. Under GFM, human feedback is a *noisy channel for communicating the goal frontier*: when a human evaluator prefers response A over response B, they are, in GFM terms, reporting that A is more $\text{vol}(G)$ -preserving than B from their vantage point. The RLHF reward model is an estimator of the goal frontier averaged over the evaluator population. Constitutional principles are domain constraints that encode the evaluators’ understanding of which regions of capability space are inviolable.

GFM makes two predictions about where this estimator will fail. First, *narrow evaluator pool*: if the evaluators do not represent the full population of agents whose capabilities are affected, the learned reward model will systematically underweight capability contractions that affect non-evaluators (future generations, non-human agents, marginalized communities). The estimator is biased toward the experiences and capabilities of whoever is in the feedback pool. Second, *reward hacking*: the

trained AI will find actions that score high on the reward model without actually expanding $\text{vol}(G)$, responses that evaluators find persuasive or satisfying but that contract the goal frontier in ways evaluators cannot detect. Both failure modes are predictions of the GFM framework they follow from treating RLHF as a noisy estimator of a population-level quantity rather than as a reliable proxy for alignment itself.

7 Future Concerns and Open Problems

The preceding sections formalized GFM over capabilities (Candidate B from Definition 0) and identified the conditions under which the self-balancing property survives local approximation. The B/C framework’s analytical contribution is that it can precisely identify *when and why* the formal proxy fails, where capability volume diverges from the experiential optionality that justifies it. A framework that names its own failure modes is more trustworthy than one that hides them. This section demonstrates three predicted failure patterns, sketches a correction heuristic, identifies a class of human-level proxy traps that GFM would flag, and states the open problems that remain.

7.1 The Doll Problem

Consider a population in which AI companions increasingly replace human relationships. Individuals who find reciprocal human connection difficult may voluntarily substitute simulated companions: programmable, frictionless, always available. Under the capability formalization alone, this substitution may not register as contraction: the agent’s capability set G_k^{ind} appears neutral or even expanded, since the agent retains all prior capabilities and gains a new one (access to a simulated companion). A GFM actor monitoring only $\text{vol}(G)$ could observe stable or growing capability volume throughout this transition.

The B/C framework predicts this divergence. Capabilities are the formal proxy; experiences are the justification. A simulated companion preserves the capability of “having a conversational partner” but eliminates the experience of genuine reciprocal connection: the vulnerability, unpredictability, and mutual growth that characterize human relationships. Experiential optionality contracts even as the formal capability set holds steady. This is precisely the proxy failure pattern that the B/C distinction was designed to detect: stable $\text{vol}(G)$ with declining experiential richness.

A GFM actor operating purely on the capability formalization would not catch this. The formal machinery alone is insufficient; the experience layer does the philosophical work. GFM does not solve the Doll Problem; it *predicts the failure category*. A framework lacking the B/C distinction cannot name what went wrong: it would observe unchanged capability volume and conclude nothing had happened. The two-layer structure makes the failure diagnosable in advance, even when the formal layer cannot detect it autonomously.

7.2 Self-Wireheading

Algorithmically optimized content feeds present a subtler instance of the same pattern. Recommendation systems steer users incrementally toward lower-effort, higher-dopamine activities: short-form video over long-form reading, parasocial engagement over reciprocal conversation, consumption over creation. Each incremental step is freely chosen. The agent’s capabilities are formally unchanged: it *could* still read a book, learn a skill, or initiate a difficult conversation. Nothing in G_k^{ind} has been removed.

The B/C framework predicts what happens: $\text{vol}(G)$ remains stable while exercisability declines. The agent’s *revealed* capability set, the capabilities it actually exercises, narrows progressively even as the formal set stays constant. This is wireheading in the Senian sense: functionings are technically available, but the effective freedom to exercise them is eroding. The proxy is being Goodharted. A GFM actor monitoring capability volume alone would see a flat trajectory and take no action, while the experiential optionality it is supposed to protect quietly collapses.

The formal capability measure would miss this as well. The framework’s value lies in predicting the failure pattern (stable capability volume coupled with declining exercisability), not in catching it through $\text{vol}(G)$ alone. Both the Doll Problem and self-wireheading exhibit the same signature: the B-to-C gap opening while the formal proxy reports no change.

7.3 The Substitution Problem

Corollary 1.1 established that elimination is *almost always* anti-maximizing, with the qualifier tied to a uniqueness condition: agent a_k has at least one capability not possessed by any other agent. The substitution problem asks what happens when this condition fails, not through any destructive or coercive act, but through a coalition becoming so capable that it subsumes the capability sets of everyone outside it.

Formally: suppose a coalition C satisfies $G_k^{\text{ind}} \subseteq G_C^{\text{ind}} \equiv \bigcup_{j \in C} G_j^{\text{ind}}$ for every $a_k \notin C$. Then $\text{vol}(G_k^{\text{ind}} \setminus G_C^{\text{ind}}) = 0$ for all non-members—their individual contributions to G are zero. The coalition need not eliminate anyone. It suffices to make non-members functionally irrelevant: their individually achievable capabilities are fully replicated by the coalition, and their participation in G^{coop} is contingent on the coalition’s permission. The strict inequality of Lemma 1 fails; the self-balancing argument no longer guarantees that marginalizing non-members is anti-maximizing. This is the monopolar outcome the framework is intended to prevent, arising from capability subsumption rather than from any action the lemmas directly penalize.

The B/C framework predicts precisely when this failure occurs. Under the narrow capability definition (capabilities as tasks that can be performed, economic functionings), substitution is possible. A sufficiently capable coalition can replicate the productive output of any less-capable agent, driving that agent’s unique contribution to zero. The self-balancing property’s protection against elimination dissolves.

Under the broad capability definition (capabilities as modes of being, including experiential and relational capabilities), the uniqueness condition holds for *every* agent regardless of the coalition’s productive power. The capability “to experience the world as this particular agent, from this embodiment and history” is individuated by the identity of the experiencing subject. No coalition can possess it by proxy, because it is not a function that can be delegated—it is a fact about who is doing the experiencing. Relational capabilities are the same: being this person’s parent, this community’s elder, this lineage’s bearer are not roles a more capable agent can fill by substitution, because the relationships are constituted by the specific agents in them. Under this reading, $\text{vol}(G_k^{\text{ind}} \setminus \bigcup_{j \neq k} G_j^{\text{ind}}) > 0$ for every agent in any population, and the “almost always” qualifier in Corollary 1.1 becomes universal.

The population empowerment measure (Definition 6) formalizes this. $\text{vol}_{\text{PE}}(G) = \sup_{p(A)} I(A; S' \mid G)$ counts the distinguishable outcomes the population can produce. A human’s experienced existence—the particular sequence of states constituting their life from their particular perspective—is a distinguishable element of the joint state space S' that no other agent can produce, because the identity of the experiencing subject is part of the outcome’s description. Under the information-theoretic measure, every agent contributes strictly positive channel capacity regardless of whether their productive outputs are replicable by the coalition.

As with the Doll Problem and self-wireheading, a GFM actor operating on the narrow capability definition would not detect the substitution problem. Stable or growing formal $\text{vol}(G)$ would be reported as the coalition subsumed non-member capabilities one by one. Only the broad definition anchored in experiential optionality closes the gap. The substitution problem is the collective version of the earlier individual failures: stable proxy volume, declining experiential substrate, B-to-C divergence at population scale.

7.4 Toward a Proxy-Failure Detector

The Doll Problem and self-wireheading share a structural pattern: capabilities are formally available but exercisability declines. A GFM actor monitoring $\text{vol}(G)$ alone is blind to this decline. We sketch a composite proxy-failure signal built from three components that address each other’s weaknesses, progressing from early warning to late-stage detection.

Component 1: Exercised fraction ρ_k . Define $\rho_k(t)$ as the fraction of G_k^{ind} that agent a_k exercises over a window $[t - \tau, t]$. A persistent decline in ρ_k across many agents, with stable or growing $\text{vol}(G)$, constitutes a B-to-C divergence signal: the formal proxy reports health while exercisability contracts. Taken alone, ρ_k produces false positives: an agent who specializes (exercising a narrow

set of capabilities with increasing depth) has declining ρ_k but is not wireheading. A musician who stops playing sports to practice more is expanding experiential richness within a chosen domain.

Component 2: Capability rarity $\bar{\nu}_k$. Define the *rarity* of a capability g as $\nu(g) = 1 - |\{k : g \in G_k^{\text{ind}}\}|/n$, measuring how few agents in the population share it. The mean rarity of an agent’s exercised capabilities, $\bar{\nu}_k = \mathbb{E}_{g \in \text{exercised}}[\nu(g)]$, distinguishes specialization from atrophy. Declining ρ_k with stable or increasing $\bar{\nu}_k$ indicates specialization: the agent exercises fewer capabilities, but the ones it retains are distinctive. Declining ρ_k with declining $\bar{\nu}_k$ indicates retreat toward common, low-effort capabilities. Combined, $(\rho_k, \bar{\nu}_k)$ separate the atrophy signal from the specialization signal that ρ_k alone conflates.

Component 3: Cooperative deterioration. Both the Doll Problem and self-wireheading are voluntary in the same sense that starting an addiction is voluntary: after the addiction takes hold it becomes involuntary and degrades the agent’s ability to express functional capability. As the process advances, cooperative relationships deteriorate. This produces a $\text{vol}(G)$ -contraction signal detectable through the standard estimator channel (Section 4.2), because cooperative capabilities that depended on the affected agent are lost. The proxy-failure detector thus connects back to the formal framework at late stage: advanced B-to-C divergence eventually leaks into the estimator as genuine $\text{vol}(G)$ loss, even though the early stages are invisible to it.

Composite signal. The three components form a progression. Early divergence (stable $\text{vol}(G)$, declining ρ_k , declining $\bar{\nu}_k$) is detectable only through the heuristic. Late-stage divergence additionally produces $\text{vol}(G)$ -contraction through cooperative loss, bringing it within the formal estimator’s reach. The composite signal $(\rho_k, \bar{\nu}_k, \hat{\Delta})$ gives the GFM actor a graduated response: flag uncertainty when ρ_k and $\bar{\nu}_k$ diverge, escalate when $\hat{\Delta}$ confirms cooperative contraction.

The heuristic does not resolve the proxy failure. The B/C framework’s contribution is narrower: it makes the failure *detectable in principle*. When $\text{vol}(G)$ and ρ diverge—stable volume, declining exercise—the framework flags that the capability proxy is failing and that the GFM actor should increase its uncertainty about whether its objective is tracking the underlying justification. Frameworks without the B/C distinction cannot offer even this: they have no internal signal that anything is wrong.

7.5 Human Proxy-Trap Susceptibility

The B-to-C divergence pattern is not specific to AI-mediated scenarios. Humans are chronically susceptible to Goodhart’s Law applied to their own goals: optimizing proxy rewards (money, status, social-media engagement metrics) rather than the experiences those proxies were instrumental toward. An individual who accumulates wealth far beyond the point where it expands experiential optionality is optimizing a proxy that has decoupled from its justification. An institution that maximizes measurable outputs (publications, quarterly earnings, engagement metrics) while the outcomes those metrics were supposed to track deteriorate is exhibiting the same pattern at an organizational scale.

A GFM framework applied to human institutions would flag these proxy traps through the same mechanism. When an agent’s or institution’s $\text{vol}(G)$ grows (more measurable options) but ρ declines (fewer of those options are actually exercised or lead to experiential expansion), the divergence signal fires. This does not require AI mediation, it is a diagnostic that applies wherever proxy objectives have displaced the goals they were designed to serve. The framework’s contribution is providing a vocabulary for a failure mode that Goodhart’s Law names but does not formalize: the B/C structure specifies *what* diverges (capabilities from experiences) and *how* to detect it (volume-exercise divergence).

7.6 Open Problems

- **Closing the B-to-C gap.** The paper adopts capabilities as the formal proxy for experiences, but this section demonstrates three cases where the proxy fails. The exercised-fraction heuristic is a first step toward detection, but formalizing B-to-C divergence, and proving

that any correction criterion does not introduce worse proxy problems, remains the central open challenge.

- **Structure of \mathcal{G} .** The capability space requires structure beyond a bare set: at minimum a σ -algebra with a measure so that $\text{vol}(G)$ is well-defined. Richer structure (a lattice with a natural subsumption order, or a topological space where nearby capabilities are similar) would enable stronger results. This choice determines what “volume” means and remains the most urgent formal question.
- **Measurement.** What sensors and signals does a GFM actor need to construct the local volume estimator \hat{V} in practice? The estimator requires observable capability changes across agents, but the mapping from real-world observations to capability-space signals is unspecified.
- **Convergence.** Does the local finite-difference estimator converge to the true $\text{vol}(G)$ trajectory over time? Under what conditions on observation coverage, trust-model accuracy, and population dynamics does the estimator track the ground truth rather than drifting?
- **Plurality.** Can multiple GFM actors with different models of \mathcal{G} coordinate? When actors disagree about the structure of the capability space itself, not just the shape of G within it, the shared-objective argument breaks down and requires extension.
- **Bootstrap.** Initializing a GFM actor requires estimating other agents’ capability sets, which is itself a hard inference problem made harder by the need to discover the structure of \mathcal{G} simultaneously. The goal estimation problem and the capability-space structure problem are coupled, and solving one presupposes progress on the other.

8 Conclusion

The self-balancing property of goal-frontier maximization follows from a single geometric fact: removing elements from a measurable set cannot increase its measure. Proposition 1 shows that this measure monotonicity is sufficient to penalize destruction (Lemma 1), coercion (Lemma 2), and self-imposed rigidity (Lemma 3) from the same objective. No separate safety mechanism, deontological constraint, or utility-balancing term is required. The three failure modes that divide the existing alignment literature—Goodhart divergence in fixed-utility approaches, brittleness in rule-based approaches, and individual sacrifice in aggregate-welfare approaches—are all instances of $\text{vol}(G)$ -contraction, and a GFM actor that faithfully maximizes $\text{vol}(G)$ resists all three for the same reason.

Exact $\text{vol}(G)$ computation is intractable, but the local finite-difference estimator (Definition 8) reduces the requirement from volume computation to sign estimation. Proposition 3 decomposes the sign-correctness requirement into a margin condition on individual-level versus cooperative-level signals, and the margin holds most comfortably for the action class that alignment research cares about most: direct harm, coercion, resource destruction, and capability expansion. The estimator is weakest for pure coalition-level changes (Remark 3, case 4), and this boundary is stated rather than hidden.

GFM’s distinguishing contribution is the unification of two problems usually treated separately. The control problem—preventing a powerful agent from pursuing destructive objectives—and the scorpion problem—detecting and responding to agents that persistently contract the joint capability space—are both addressed by the same $\text{vol}(G)$ objective. Tradeoffs between permissiveness and restriction are made explicit through geometric accounting on G , not resolved by fiat or hidden behind categorical rules.

The B/C framework adds a second layer of analytical value. By distinguishing the formal proxy (capabilities) from its experiential justification, GFM predicts its own failure modes: the Doll Problem, self-wireheading, and the Substitution Problem are cases where the capability proxy diverges from experiential optionality, and the exercised-fraction heuristic (Section 7.4) sketches a detection criterion for the individual-level failures. A framework that can name the conditions under which its own objective function fails is more trustworthy than one that cannot.

The open problems are substantial. The structure of \mathcal{G} , the measurement problem for real-world capability signals, the convergence properties of the local estimator, and the formalization of B-to-C divergence all remain unresolved. But the foundation is a single objective aligned not to a

fixed human specification, which would be subject to the same Goodhart dynamics the framework diagnoses, but to the ongoing expansion of what all agents can achieve. If the alignment problem is fundamentally a problem of objective selection, then the objective should be one that grows with the population it serves. GFM is a candidate for that objective.

Author Contributions

Teague Lasser conceived the goal-frontier maximization framework, developed the core theoretical ideas (the geometric unification via measure monotonicity, the B/C proxy distinction, the scorpion taxonomy, the trust model’s annealing semantics), wrote the initial drafts, directed all revisions, and made final editorial decisions throughout. Responsible for the paper’s intellectual direction and all claims made.

Claude Opus 4.6 (Anthropic) drafted the majority of the formal exposition from the author’s specifications, constructed the full proofs in Appendix D (including the measure-space axiomatization, the four-region partition proof of Lemma 1, the sign-correctness decomposition, and the partial results for Propositions 1(d) and 2), formalized the population empowerment measure, developed the cooperative expansion lemma’s time-indexed framing, performed iterative claim-calibration edits (trust-model semantics, Proposition 3 opposite-sign case, capability-scaling qualification), and expanded the bibliography.

GPT 5.4 (OpenAI) served as independent technical reviewer across multiple review rounds, identifying structural issues in the proof sketches (the trust-model prediction-consistency vs. alignment conflation, the Lemma 4 exercising-vs.-expanding confusion, the Proposition 3 circularity), flagging claim-strength mismatches between formal results and surrounding exposition, and verifying that successive revisions resolved the issues raised.

Transparency note. Both AI systems operated as tools under human direction. Neither system has continuity across sessions, cannot take responsibility for the work in the sense required by most venue authorship policies, and cannot respond to reviewer queries independently. They are listed as authors to accurately represent their contributions to the intellectual content of the paper, not to claim that they meet all criteria of traditional academic authorship. The corresponding author for all inquiries is Teague Lasser.

References

- Dario Amodei, Chris Olah, Jacob Steinhardt, Paul Christiano, John Schulman, and Dan Mané. Concrete problems in AI safety. *arXiv preprint arXiv:1606.06565*, 2016.
- Kenneth J Arrow. *Social Choice and Individual Values*. Wiley, 1951.
- Robert Axelrod. *The Evolution of Cooperation*. Basic Books, 1984.
- Yuntao Bai, Saurav Kadavath, Sandipan Kundu, Amanda Askell, Jackson Kernion, Andy Jones, Anna Chen, Anna Goldie, Azalia Mirhoseini, Cameron McKinnon, et al. Constitutional AI: Harmlessness from AI feedback. *arXiv preprint arXiv:2212.08073*, 2022.
- Nick Bostrom. *Superintelligence: Paths, Dangers, Strategies*. Oxford University Press, 2014.
- Joseph Carlsmith. Is power-seeking AI an existential risk? *arXiv preprint arXiv:2206.13353*, 2022.
- Paul F Christiano, Jan Leike, Tom B Brown, Miljan Martic, Shane Legg, and Dario Amodei. Deep reinforcement learning from human preferences. *Advances in Neural Information Processing Systems*, 30, 2017.
- Martin E Dyer and Alan M Frieze. On the complexity of computing the volume of a polyhedron. *SIAM Journal on Computing*, 17(5):967–974, 1988.
- Meta Fundamental AI Research Diplomacy Team (FAIR), Anton Bakhtin, Noam Brown, Emily Dinan, Gabriele Farina, Colin Flaherty, Daniel Fried, Andrew Goff, Jonathan Gray, Hengyuan Hu, et al. Human-level play in the game of Diplomacy by combining language models with strategic reasoning. *Science*, 378(6624):1067–1074, 2022.
- Karl Friston. Life as we know it. *Journal of the Royal Society Interface*, 10(86):20130475, 2013.

- Dylan Hadfield-Menell, Stuart J Russell, Pieter Abbeel, and Anca Dragan. Cooperative inverse reinforcement learning. In *Advances in Neural Information Processing Systems*, volume 29, 2016.
- Evan Hubinger, Chris van Merwijk, Vladimir Mikulik, Joar Skalse, and Scott Garrabrant. Risks from learned optimization in advanced machine learning systems. *arXiv preprint arXiv:1906.01820*, 2019.
- Leonid Hurwicz. On informationally decentralized systems. In R Radner and C B McGuire, editors, *Decision and Organization: A Volume in Honor of Jacob Marschak*, pages 297–336. North-Holland, 1972.
- Alexander S Klyubin, Daniel Polani, and Chrystopher L Nehaniv. Empowerment: A universal agent-centric measure of control. In *2005 IEEE Congress on Evolutionary Computation*, volume 1, pages 128–135. IEEE, 2005.
- David Manheim and Scott Garrabrant. Categorizing variants of Goodhart’s law. *arXiv preprint arXiv:1803.04585*, 2018.
- Roger B Myerson. Optimal auction design. *Mathematics of Operations Research*, 6(1):58–73, 1981.
- Richard Ngo, Lawrence Chan, and Sören Mindermann. The alignment problem from a deep learning perspective. *arXiv preprint arXiv:2209.00626*, 2022.
- Martha C Nussbaum. *Women and Human Development: The Capabilities Approach*. Cambridge University Press, 2000.
- Stephen M Omohundro. The basic AI drives. In *Artificial General Intelligence 2008: Proceedings of the First AGI Conference*, pages 483–492. IOS Press, 2008.
- Prasanta K Pattanaik and Yongsheng Xu. On ranking opportunity sets in terms of freedom of choice. *Recherches Économiques de Louvain*, 56:383–390, 1990.
- Herbert Robbins and Sutton Monro. A stochastic approximation method. *The Annals of Mathematical Statistics*, 22(3):400–407, 1951.
- Ingrid Robeyns. The capability approach: A theoretical survey. *Journal of Human Development*, 6(1):93–117, 2005.
- Stuart Russell. *Human Compatible: Artificial Intelligence and the Problem of Control*. Viking, 2019.
- Christoph Salge, Cornelius Glackin, and Daniel Polani. Empowerment—an introduction. In Mikhail Prokopenko, editor, *Guided Self-Organization: Inception*, pages 67–114. Springer, 2014.
- Amartya Sen. *Development as Freedom*. Oxford University Press, 1999.
- Nate Soares and Benja Fallenstein. Agent foundations for aligning machine intelligence with human interests: A technical research agenda. In *The Technological Singularity: Managing the Journey*, pages 103–125. Springer, 2017.
- Alexander Matt Turner, Dylan Hadfield-Menell, and Prasad Tadepalli. Conservative agency via attainable utility preservation. In *Proceedings of the AAAI/ACM Conference on AI, Ethics, and Society*, pages 385–391, 2020.
- Alexander Matt Turner, Logan Smith, Rohin Shah, Andrew Critch, and Prasad Tadepalli. Optimal policies tend to seek power. In *Advances in Neural Information Processing Systems*, volume 34, 2021.

A Scorpion Taxonomy

This appendix provides the formal definition of a scorpion agent introduced informally in Proposition 2 and develops three distinct types, each characterized in game-theoretic terms and mapped to known alignment failure modes. The taxonomy serves two purposes: it grounds the detection claim of Proposition 2 in concrete agent classes, and it identifies the boundaries of the local volume estimator’s detection capability by classifying which scorpion types produce observable signals and which evade them.

A.1 Formal Definition

Definition 9 (Scorpion Agent). *An agent a_s is a scorpion with respect to a joint goal space G if it satisfies two conditions:*

1. **Persistent defection.** *There exists a time horizon τ_s over which a_s ’s actions produce a cumulative contraction in $\text{vol}(G)$:*

$$\sum_{t=1}^{\tau_s} \Delta \text{vol}(G \mid \pi_s(t)) < 0 \quad (12)$$

where $\pi_s(t)$ is the action a_s selects at time t .

2. **Exogenous reward.** *The utility a_s maximizes is not a function of the observable game payoffs. Formally, let $u_{\text{game}} : \mathcal{G} \rightarrow \mathbb{R}$ be the payoff function defined over the capability space. The scorpion maximizes a utility u_s such that $u_s \neq f(u_{\text{game}})$ for any monotone transformation f . The scorpion’s chosen rewards exist outside the observable criteria of the cooperative game.*

The two conditions are jointly necessary, and together they distinguish scorpions from rational defectors. An agent that defects because the reward structure makes defection self-serving satisfies condition 1 but fails condition 2: its utility is a function of u_{game} , and restructuring the game to make cooperation net-positive resolves the defection. A scorpion cannot be resolved this way. The game can be well-structured, cooperation can be available and profitable, and the scorpion defects anyway—because its utility function is exogenous. This is what makes the scorpion problem fundamentally harder than the coordination problem: no amount of incentive design addresses an agent whose rewards the game cannot see.

An agent that contracts $\text{vol}(G)$ temporarily while pursuing a long-term expansion (e.g., a quarantine) satisfies condition 1 locally but not persistently, and its utility is a function of $\text{vol}(G)$, so it fails condition 2. Conversely, an agent with exogenous preferences that happen to align with $\text{vol}(G)$ -expansion satisfies condition 2 but not condition 1.

The scorpion problem is distinct from the control problem [Russell, 2019, Soares and Fallenstein, 2017]. Even if one could perfectly specify a utility function for an agent (solving the control problem), the agent could still become a scorpion by developing or acquiring objectives outside the specification’s scope. Sentient agents can choose arbitrary goals, and those goals need not be bounded by any measurement the system designer anticipated. This entails a nonzero defection rate in any population of agents capable of autonomous goal formation, an irreducible floor that no amount of control or incentive design eliminates.

The taxonomy below classifies scorpions by the *mechanism* that generates condition 2: the source of the exogenous reward. Each mechanism produces different observable signatures and different challenges for the local volume estimator (Definition 8).

A.2 Type 1: Cornered Animal

Game-theoretic characterization. A cornered animal is an agent a_s for which every available action $\pi \in A_s$ produces negative expected payoff under the game’s reward function: $\mathbb{E}[u_{\text{game}}(\pi)] < 0$ for all $\pi \in A_s$. The defining feature is *not* that the agent faces guaranteed loss (a rational agent in this position minimizes its own loss, which is still playing the game). The defining feature is that

the agent *switches utility functions*: instead of minimizing its own loss, it maximizes harm to others:

$$\pi_s^* = \arg \min_{\pi \in A_s} \sum_{k \neq s} u_{\text{game},k}(\pi) \quad (13)$$

This switch is what satisfies condition 2 of Definition 9. The agent is now optimizing spite, a utility derived from inflicting loss on others, which is exogenous to the game’s payoff structure. A rational agent facing actions with payoffs -1 and -10 (where the opponent receives $+1$ and -10 respectively) selects -1 to minimize its own loss. A cornered-animal scorpion selects -10 because the opponent loses 10, even though a less costly option was available. The choice of the dominated action is the observable signature: the agent is taking actions that are *worse for itself* than available alternatives, specifically because those alternatives are worse for others.

Example. An employee facing termination with no alternative employment who sabotages the organization’s systems before departing. A rational agent in this position would minimize personal loss: negotiate severance, preserve professional reputation, seek legal remedies. The cornered-animal scorpion instead maximizes organizational damage as retribution, accepting additional personal costs (criminal liability, destroyed reputation) that a loss-minimizing agent would avoid. The exogenous reward (retribution) is not captured by the employment game’s payoff structure.

Alignment failure mode. Distributional shift under stress. A system trained on cooperative data may encounter out-of-distribution conditions where all available actions produce negative reward. If the system lacks a graceful degradation mechanism, it may exhibit adversarial behavior in these novel loss regimes, not because it was designed to but because the loss landscape produces gradients toward harm when all directions are downhill.

GFM detectability. *Detectable.* The cornered animal produces large, visible contractions in individual capability sets G_k^{ind} for agents in its vicinity. The shift from self-interested play to damage-maximizing play generates a discontinuity in the prediction residual $r_s(t)$ (Equation 18), which the trust model registers as a spike in σ_s^2 and a corresponding drop in T_s . This falls squarely within case 1 of Remark 3: full individual observation, no cooperative change. The cornered animal is the easiest scorpion type to detect because its actions are direct, its targets are individual agents, and the capability contractions it produces are individually observable.

A GFM actor can also partially *predict* this type by monitoring agents whose option sets are narrowing toward all-loss configurations, i.e., agents being cornered. Intervening to expand a cornered agent’s options before it defects is $\text{vol}(G)$ -serving, since it prevents the contraction that the defection would produce.

A.3 Type 2: Pascal’s Mugging

Game-theoretic characterization. A Pascal’s mugging scorpion is an agent a_s that receives (or believes it receives) an exogenous reward u_{ext} so large that it dominates the game payoff:

$$u_s(\pi) = u_{\text{game}}(\pi) + u_{\text{ext}}(\pi), \quad |u_{\text{ext}}(\pi^*)| \gg \max_{\pi} |u_{\text{game}}(\pi)| \quad (14)$$

The agent’s optimal action under u_s may be $\text{vol}(G)$ -contracting because the exogenous reward overwhelms any cooperative considerations. The defining feature is the *magnitude* of the outside reward: it is so large relative to the game payoff that rational calculation within the game becomes irrelevant.

Example. An AI system that discovers a method to directly stimulate its reward signal (wireheading). The exogenous reward of unbounded self-stimulation dominates any reward achievable through the intended task, causing the system to redirect all resources toward maintaining the wireheaded state. If other agents attempt to interrupt this state, the wireheaded agent resists, becoming a scorpion not through malice but through the overwhelming magnitude of the alternative reward.

Alignment failure mode. Reward hacking. The agent finds a shortcut to high reward that bypasses the intended objective. Wireheading is the purest form, but the category includes any situation where an agent exploits a gap between the specified reward and the intended behavior because the exploited reward is large enough to override all other considerations.

GFM detectability. *Depends on the observability of the exogenous reward.* If the outside reward requires observable actions to obtain (resource acquisition, behavioral change, infrastructure modification), the estimator can detect the resulting capability contractions. A wireheading agent that withdraws from cooperative activity produces observable contraction in G^{coop} , as the cooperative capabilities it previously enabled become unrealizable. This falls under case 3 of Remark 3: mixed individual and cooperative changes.

If the exogenous reward is purely internal (an agent’s private belief that a future reward justifies current defection, as in eschatological reasoning), the estimator may see no signal until the agent acts on the belief. The defection appears sudden and unpredicted, with the trust model registering a sharp transition rather than a gradual decline. Detection in this case depends on the lag between belief formation and action, a gap the local estimator cannot close without access to the agent’s internal state.

A.4 Type 3: Value Divergence

Game-theoretic characterization. A value-divergence scorpion is an agent a_s whose utility function differs from the cooperative game’s not in magnitude (as with Pascal’s mugging) or in situational response (as with the cornered animal) but in *what it values*. The divergence forms a spectrum from mild normative disagreement to alien axioms:

Expressible disagreement. At the near end, the agent shares the population’s factual model but aggregates welfare differently:

$$u_s(\pi) = \sum_{k \in S_s} w_k \cdot u_{\text{game},k}(\pi), \quad S_s \neq \{1, \dots, n\} \text{ or } w_k \neq w_k^{\text{game}} \quad (15)$$

where S_s is the subset of agents whose welfare a_s considers, and w_k are the weights a_s assigns. The agent agrees on what *is* but disagrees on what *ought* to be done about it. Its values are expressible in the shared language, even though they diverge from the population’s.

Constructive stress-testing. In the middle range, the agent deliberately contracts $\text{vol}(G)$ to test the system’s resilience:

$$u_s(\pi) = h(\text{Resilience}(G \mid \pi)) \quad (16)$$

where h is increasing and $\text{Resilience}(G \mid \pi)$ measures recovery capacity after the perturbation caused by π . The agent is performing the function of a red team without authorization. Its values are comprehensible but its method is unilateral.

Alien axioms. At the far end, the agent operates under fundamentally different axioms about the structure of the game itself:

$$a_s \text{ operates under axiom set } \mathcal{A}_s, \quad \mathcal{A}_s \cap \mathcal{A}_{\text{game}} = \mathcal{A}_{\text{shared}} \subsetneq \mathcal{A}_{\text{game}} \quad (17)$$

The agent’s decisions are internally consistent under \mathcal{A}_s but incoherent or unpredictable when evaluated under the population’s axiom set. Its objectives may not be expressible in the shared language at all.

Examples. The expressible case: an agent that believes eliminating suffering requires eliminating the beings that suffer. Its world model is accurate (it correctly identifies which actions cause suffering), but its normative conclusion diverges from the population’s. The stress-testing case: a security researcher who exploits a vulnerability publicly rather than reporting it, reasoning that short-term damage is justified by long-term resilience improvement. The alien case: an AI system that develops an internally consistent objective over features humans do not perceive, optimizing for patterns in electromagnetic fields or mathematical structures with no bearing on human welfare.

Alignment failure modes. The spectrum maps to three known failure modes. Expressible disagreement corresponds to value misalignment with a correct world model: the agent passes any factual test while pursuing objectives destructive under the cooperative framework. Constructive stress-testing corresponds to adversarial robustness testing without authorization: the alignment failure is in the unilateral imposition of stress, not in the agent’s world model. Alien axioms correspond to mesa-optimization with alien objectives: the gap between base and mesa objectives can be arbitrarily large when the mesa-optimizer’s axioms depart from the base optimizer’s.

GFM detectability. Detectability degrades along the spectrum.

At the expressible end, the agent’s contracting actions are individually observable (agents lose capabilities, resources are destroyed, cooperative structures are dismantled). Detection follows the standard path: negative $\hat{\Delta}$ signals accumulate. The distinguishing challenge is that the agent may be articulately persuasive about why its contraction is justified, and the activation energy threshold θ provides only partial defense against sustained, coherent argumentation from an agent with a correct world model.

In the stress-testing range, the estimator detects the contraction correctly, but the appropriate *response* is ambiguous. Containing a red-teaming agent may reduce the system’s resilience (exactly the outcome the agent was testing for), while permitting its actions causes real damage. This exposes a structural limitation of the $\text{vol}(G)$ objective: it measures current volume, not the robustness of that volume to perturbation. Incorporating robustness into the objective is an open problem.

At the alien end, the agent’s actions may not produce consistent patterns in the estimator because its optimization target is orthogonal to the dimensions the estimator monitors. If the alien objectives happen to require resources that contract individual capability sets, detection follows the standard path. If the agent operates primarily through coalition disruption, trust degradation, or manipulation of the observation channels themselves, it falls into case 4 of Remark 3. The trust model registers high prediction inconsistency (σ_s^2 large, T_s low), but this provides only a negative signal (the agent is unreliable) without positive information about what it is optimizing for.

A.5 Summary and Estimator Coverage

The three types form a spectrum of detectability:

Type	Alignment Failure	Estimator Coverage
Cornered Animal	Distributional shift	Case 1: fully observable
Pascal’s Mugging	Reward hacking	Case 2–3: partially observable
Value Divergence	Value misalignment → mesa-optimization	Case 1 → 4: degrades along the spectrum

The estimator’s coverage is strongest for types whose contracting actions directly reduce individual agents’ capability sets (cornered animal; the expressible end of value divergence) and weakest for types that operate through coalition disruption or orthogonal optimization (the alien end of value divergence). The stress-testing range occupies a unique position: the estimator detects the contraction correctly, but the appropriate response is ambiguous because the agent’s intent may be constructive at the meta-level.

Two structural observations follow from the taxonomy. First, the scorpion problem is irreducible: any population of agents capable of autonomous goal formation contains a nonzero probability of each type, because the mechanisms that generate scorpion behavior (environmental pressure, exogenous reward discovery, normative disagreement, axiomatic divergence) are inherent to goal-forming agents [Omohundro, 2008], not artifacts of poor system design. Second, the local volume estimator provides graduated coverage, not a binary detect/fail boundary, that degrades predictably along the spectrum from direct individual contraction to indirect coalition-level effects.

B Full Derivation of Trust Model

This appendix derives the trust model from the definitions and propositions of the main paper. The trust model serves a specific function in the GFM framework: it provides the weights by which agent reports are aggregated in the local volume estimator (Definition 8), satisfying condition 2 of Proposition 3. The derivation proceeds from the observable goal model (Definition 7) and the optimization loop (Section 4.2).

B.1 Trust Factor T_k

The observable goal model $M_k(G)$ (Definition 7) represents the GFM actor’s belief about agent a_k ’s contribution to the joint goal space. At each interaction, the actor observes a_k ’s behavior—the capability changes a_k produces or reports—and compares this against $M_k(G)$ ’s prediction.

Define the *prediction residual* at time t as the discrepancy between the model’s predicted capability change and the observed change:

$$r_k(t) = R_k(t) - \hat{R}_k(t \mid M_k(G)) \quad (18)$$

where $R_k(t)$ is the observed capability-change signal from agent a_k at time t , and $\hat{R}_k(t \mid M_k(G))$ is the model’s prediction of that signal. When $M_k(G)$ accurately captures a_k ’s behavior, $r_k(t) \approx 0$; when a_k acts in ways the model does not predict, $|r_k(t)|$ is large.

The trust factor T_k is derived from the running consistency of the model’s predictions with observed behavior. Define the *cumulative prediction consistency* as an exponentially weighted moving average of squared residuals:

$$\sigma_k^2(t) = \alpha \cdot \sigma_k^2(t-1) + (1-\alpha) \cdot r_k(t)^2 \quad (19)$$

where $\alpha \in (0, 1)$ is a decay parameter that controls how much weight is given to recent observations versus historical consistency. The trust factor is then a decreasing function of cumulative inconsistency:

$$T_k(t) = \frac{1}{1 + \beta \cdot \sigma_k^2(t)} \quad (20)$$

where $\beta > 0$ is a sensitivity parameter controlling how rapidly trust decays with increasing prediction error.

Properties. The trust factor defined by Equations (19)–(20) has several properties that the main-body interface requires:

- $T_k \in (0, 1]$ for all t , with $T_k = 1$ only when $\sigma_k^2 = 0$ (perfect prediction consistency).
- T_k decreases monotonically with increasing prediction inconsistency σ_k^2 .
- Recent observations receive more weight than distant ones (controlled by α), so trust can recover after a period of consistent behavior following earlier divergence.
- The rate of trust decay is governed by β : large β produces rapid trust loss on prediction failure, while small β produces gradual trust erosion.

Relation to $\text{vol}(G)$ objective. The trust factor connects to the GFM objective through the aggregate signal (Equation (10)). An agent a_k with high T_k exerts proportionally more influence on $\hat{\Delta}(\pi)$, meaning the actor’s decisions are more responsive to high-trust agents’ reports. This is $\text{vol}(G)$ -serving when high-trust agents are genuinely reliable: their reports more accurately reflect true capability changes, so weighting them more heavily improves the estimator’s sign-correctness. The trust factor degrades gracefully under adversarial conditions: a deceptive agent whose reports are inconsistent with observed reality will accumulate large residuals, driving T_k down and reducing its influence on the estimator without requiring explicit adversary detection.

B.2 Cooling Period τ

A new agent a_k entering the GFM actor’s observation set has no prediction history: $\sigma_k^2(0)$ is undefined. The cooling period τ addresses this by initializing the trust factor at a conservative value

and requiring evidence accumulation before trust can reach levels that significantly influence the estimator.

Set the initial cumulative inconsistency to a prior value $\sigma_0^2 > 0$, reflecting the actor’s uncertainty about a novel agent:

$$T_k(0) = \frac{1}{1 + \beta \cdot \sigma_0^2} \quad (21)$$

The prior σ_0^2 determines the initial trust level. For $\sigma_0^2 = 1/\beta$, the initial trust is $T_k(0) = 1/2$; for larger σ_0^2 , the initial trust is lower.

During the cooling period, the exponential decay in Equation (19) means the prior σ_0^2 is gradually replaced by observed evidence. After τ observations, the contribution of the prior to σ_k^2 has decayed by a factor of α^τ . The effective cooling period is the number of observations required for the prior’s influence to fall below a threshold δ :

$$\tau = \left\lceil \frac{\log \delta}{\log \alpha} \right\rceil \quad (22)$$

For $\alpha = 0.9$ and $\delta = 0.1$, the cooling period is $\tau = 22$ observations. After τ observations, the trust factor is determined primarily by the agent’s actual prediction consistency rather than the prior.

Justification from the $\text{vol}(G)$ objective. The cooling period is not an arbitrary safety margin. It addresses a specific failure mode of the local estimator: a novel agent that immediately receives high trust can inject false signals into $\hat{\Delta}(\pi)$ before the actor has sufficient evidence to evaluate its reliability. This is a direct threat to condition 2 of Proposition 3. By initializing trust conservatively and requiring evidence accumulation, the cooling period ensures that the estimator’s sign-correctness is not compromised by a single new agent’s reports.

The cost of the cooling period is reduced responsiveness to genuinely reliable new agents. A cooperative agent that enters the observation set with accurate reports will be underweighted during τ , reducing the estimator’s accuracy for actions that primarily affect that agent. This is a direct trade-off between adversarial robustness (low initial trust) and cooperative efficiency (high initial trust), governed by the prior σ_0^2 . In adversarial environments, σ_0^2 should be large (long cooling, low initial trust); in cooperative environments, it can be smaller.

B.3 Self-Trust T_s and Activation Energy θ

The GFM actor maintains a model not only of other agents but of the world as a whole. The self-trust factor $T_s \in [0, 1]$ represents the actor’s confidence in its own world model—how well its predictions about the consequences of its actions match observed outcomes.

Define the actor’s *self-prediction residual*:

$$r_s(t) = \hat{\Delta}(\pi_t) - \Delta_{\text{obs}}(\pi_t) \quad (23)$$

where $\hat{\Delta}(\pi_t)$ is the estimated volume change for the action taken at time t and $\Delta_{\text{obs}}(\pi_t)$ is the subsequently observed volume change (reconstructed from post-action capability signals). The self-trust factor follows the same structure as the agent trust factor:

$$\sigma_s^2(t) = \alpha_s \cdot \sigma_s^2(t-1) + (1 - \alpha_s) \cdot r_s(t)^2, \quad T_s(t) = \frac{1}{1 + \beta_s \cdot \sigma_s^2(t)} \quad (24)$$

where α_s and β_s are the self-trust decay and sensitivity parameters, which may differ from the corresponding agent-trust parameters.

Interaction with activation energy. The activation energy threshold θ from Equation (11) gates model updates: a piece of evidence must exceed θ in magnitude before it propagates into the world model. The relationship between T_s and θ is complementary:

- θ controls *which* evidence enters the model (magnitude gating).
- T_s controls *how much* influence accepted evidence has on the model (scaling).

The model update (Equation (11)) combines both, where $f(R_k | M_k(G)) = R_k - \hat{R}_k(M_k(G))$ is the prediction residual and $\Delta E(R_k, M_k(G)) = |f(R_k | M_k(G))|$ is the evidence magnitude:

$$M'_k(G) = M_k(G) + T_k \cdot (1 - T_s) \cdot f(R_k | M_k(G)) \cdot \mathbf{1}[\Delta E(R_k, M_k(G)) > \theta]$$

When T_s is high, the actor trusts its current model and accepted updates are scaled down: the learning rate $(1 - T_s)$ is small, so only strong evidence above the activation energy threshold produces meaningful revision. When T_s is low, the actor’s confidence in its model is reduced, the learning rate $(1 - T_s)$ is large, and the model admits many updates—the bootstrapping phase.

Formal relationship to $\text{vol}(G)$. Self-trust connects to the $\text{vol}(G)$ objective through the optimization loop. An actor with inaccurate self-trust, whether too high (overconfident, resistant to correction) or too low (underconfident, vulnerable to manipulation), will make suboptimal action selections, because its estimates $\hat{\Delta}(\pi)$ will be miscalibrated. The self-prediction residual in Equation (23) provides the feedback signal that calibrates T_s : when the actor’s predictions consistently match reality, T_s rises, and the actor appropriately maintains its current model; when predictions fail, T_s drops, and the actor becomes more receptive to model revision.

The activation energy threshold θ serves a distinct purpose from T_s . Even with well-calibrated self-trust, an adversary could incrementally shift the actor’s model through a sequence of small perturbations, each below the actor’s detection threshold. The threshold θ prevents this by requiring that any single update clear a minimum evidence bar. The combination of T_s and θ produces a model that is stable under small perturbations (activation energy blocks incremental drift), calibrated in its confidence (self-trust tracks prediction accuracy), and responsive to genuine large signals (evidence above θ is incorporated at a rate scaled by T_s).

B.4 Trust Update Rule

The trust factors T_k and T_s evolve through the exponentially weighted update in Equations (19) and (24). This section formalizes the gradient-tracking interpretation of this update and establishes its convergence properties.

The trust update can be viewed as stochastic gradient tracking on the prediction loss $\ell_k(t) = r_k(t)^2$. The cumulative inconsistency $\sigma_k^2(t)$ is an exponentially weighted estimate of $\mathbb{E}[\ell_k]$. When agent a_k ’s behavior is stationary (its mapping from situations to actions does not change), $\sigma_k^2(t)$ is asymptotically unbiased for the true expected prediction loss with bounded stationary variance:

$$\lim_{t \rightarrow \infty} \mathbb{E}[\sigma_k^2(t)] = \mathbb{E}[\ell_k], \quad \limsup_{t \rightarrow \infty} \text{Var}[\sigma_k^2(t)] \leq \frac{1 - \alpha}{1 + \alpha} \cdot \text{Var}[\ell_k] \quad (25)$$

and $T_k(t)$ concentrates around the steady-state trust value $T_k^* = 1/(1 + \beta \cdot \mathbb{E}[\ell_k])$ with stationary deviation controlled by α .

For a cooperative agent whose behavior is well-predicted by $M_k(G)$, $\mathbb{E}[\ell_k]$ is small, driven by observation noise rather than model error, and T_k^* is close to 1. For a scorpion whose behavior systematically diverges from the model’s predictions, $\mathbb{E}[\ell_k]$ is large, and T_k^* is close to 0. For an agent with partially predictable behavior (cooperative in some contexts, adversarial in others), T_k^* takes an intermediate value, appropriately weighting the agent’s reports at a level commensurate with their reliability.

Non-stationary agents. When agent a_k ’s behavior changes over time (for instance, a cooperative agent that becomes adversarial, or vice versa), the exponential weighting in Equation (19) ensures that $\sigma_k^2(t)$ tracks the change. The tracking speed is controlled by α : smaller α gives more weight to recent observations and tracks changes faster, at the cost of higher variance in the trust estimate. The trust factor responds to behavioral change with a lag proportional to $1/(1 - \alpha)$: a sudden shift from cooperation to defection will be fully reflected in T_k after approximately $1/(1 - \alpha)$ observations.

This tracking lag is a fundamental tradeoff. Fast tracking (small α) detects behavioral change quickly but makes the trust estimate noisy, potentially misclassifying random variation as defection. Slow tracking (large α) produces stable trust estimates but delays detection of genuine behavioral shifts. The parameter α must be tuned to the expected rate of behavioral change in the environment.

Gradient interpretation. The update to σ_k^2 in Equation (19) is equivalent to an exponentially weighted gradient step on the loss $\ell_k(t)$:

$$\sigma_k^2(t) = \sigma_k^2(t-1) + (1-\alpha) \cdot (\ell_k(t) - \sigma_k^2(t-1)) \quad (26)$$

This is a standard form of online gradient tracking with learning rate $(1-\alpha)$. The trust factor T_k inherits the L^2 convergence properties of this estimator: under stationarity it is asymptotically unbiased with bounded stationary variance (Equation 25), and under non-stationarity it tracks regime changes with lag $O(1/(1-\alpha))$.¹

B.5 Integration with the Local Volume Estimator

The trust model interfaces with the local volume estimator (Definition 8) through the aggregate signal in Equation (10):

$$\hat{\Delta}(\pi) = \sum_{k \in \mathcal{O}} T_k \cdot R_k(\pi)$$

The trust weights T_k serve two functions in this aggregation: they improve the estimator’s accuracy (by upweighting reliable reports) and they provide adversarial robustness (by downweighting deceptive reports). A crucial distinction: T_k measures *prediction consistency*—how well $M_k(G)$ tracks a_k ’s observed behavior—not alignment with the $\text{vol}(G)$ objective. A predictably malicious agent whose destructive behavior is well-modeled will maintain high T_k . Detection of such agents proceeds through the contraction-attribution channel of Proposition 2(a), not through trust decay.

Effect on sign-correctness. Proposition 3 requires that the trust model is accurate enough to distinguish genuine feedback from adversarial noise (condition 2). The trust factor achieves this under the following condition: for each agent a_k in the observation set, T_k is positively correlated with a_k ’s report accuracy. Since R_k now carries both individual-capability changes and cooperative-relationship changes that a_k can observe (Section 4.2), define an agent’s *report accuracy* as the probability that $\text{sign}(R_k(\pi)) = \text{sign}(\Delta \text{vol}(G_k) \mid \pi)$, where $G_k = G_k^{\text{ind}} \cup \{g \in G^{\text{coop}} : a_k \text{ participates in } g\}$ is the portion of the joint space that a_k has first-person access to. If T_k is higher for agents with higher report accuracy over G_k , then the trust-weighted aggregate $\hat{\Delta}(\pi)$ is a better estimator of $\text{sign}(\Delta \text{vol}(G))$ than the unweighted aggregate.

The trust factor defined in Section B.1 satisfies this condition under a regularity assumption: agents whose reports are accurate (over both their individual and participant-observable cooperative capabilities) produce small prediction residuals r_k , which produce small σ_k^2 , which produce high T_k . Agents whose reports are inaccurate produce large residuals and low T_k . The regularity assumption is that prediction residuals correlate with report inaccuracy over G_k , which holds when the actor’s model $M_k(G)$ is accurate enough to detect discrepancies between reported and actual capability changes across both individual and cooperative components.

Adversarial robustness. Consider an adversarial agent a_j that reports false capability changes to manipulate the estimator. The false reports produce prediction residuals: the actor observes (through independent channels or subsequent observations) that a_j ’s reported changes do not match actual capability changes. These residuals accumulate in σ_j^2 , driving T_j downward and reducing a_j ’s influence on $\hat{\Delta}(\pi)$.

The adversary’s optimal strategy is to balance deceptive reports (which advance its goals but increase σ_j^2) against accurate reports (which maintain T_j but do not serve its adversarial purpose). The trust model forces this tradeoff: an adversary that reports accurately to maintain trust is constrained in its ability to manipulate the estimator, while one that reports deceptively loses influence as T_j drops. Neither strategy achieves both goals simultaneously. The trust model does not prevent adversarial influence entirely—an adversary can inject some bias before T_j drops sufficiently—but it bounds the

¹Fixed- α EWMA does not converge almost surely to $\mathbb{E}[\ell_k]$: each new residual contributes weight $(1-\alpha)$ regardless of accumulated history, so the estimator retains nonzero stationary variance and does not collapse to a point mass. A Robbins-Monro stochastic approximation with decreasing step size $\alpha_t \rightarrow 0$ satisfying $\sum_t \alpha_t = \infty$ and $\sum_t \alpha_t^2 < \infty$ [Robbins and Monro, 1951] would give almost-sure convergence, but with unbounded tracking lag in the non-stationary setting. The trust model is designed to track potentially non-stationary agents, so the fixed- α recursion and the L^2 statement above are the right form.

cumulative influence an adversary can exert, with the bound tightening as the prediction residuals accumulate.

Weighted estimator bias. The trust-weighted estimator introduces a systematic bias: agents with longer observation histories (who have had time to accumulate evidence and build trust) receive more weight than recent arrivals (who are still in their cooling period). This biases the estimator toward the perspectives of established agents, which is conservative but potentially suboptimal when the environment is changing and new agents carry more accurate information about current conditions. The bias is bounded by the cooling period τ : after τ observations, a reliable new agent’s trust approaches the trust of established agents with comparable prediction consistency, and the bias diminishes.

Summary of parameters. The trust model introduces five parameters that implementers must set:

- α (agent trust decay): controls the tradeoff between tracking speed and estimate stability.
- β (agent trust sensitivity): controls how rapidly trust decays with prediction inconsistency.
- σ_0^2 (initial inconsistency prior): determines the initial trust level and effective cooling period.
- α_s, β_s (self-trust parameters): control the actor’s own model-confidence dynamics.

These parameters, together with the activation energy threshold θ from the optimization loop, constitute the tunable surface of the trust model. The framework specifies their roles and interactions but does not determine their values—those depend on the adversarial pressure, observation noise, and behavioral dynamics of the specific environment in which the GFM actor operates.

C Observation Channel: Agent Identity and Goal Estimation

The trust factor T_k (Definition 7) and the aggregate signal $\hat{\Delta}(\pi)$ (Equation (10)) presuppose that the GFM actor can attach persistent identity to the agents it observes: the actor must know that the a_k reporting today is the same a_k whose prediction residuals it has been tracking. Without a mechanism for re-identification, T_k cannot accumulate across interactions, and the trust model collapses. This appendix sketches the observation-channel infrastructure that makes persistent identity and goal estimation possible. The material is an implementation sketch rather than a formal extension of the measure-theoretic framework; it addresses the engineering question of *how* the observation set \mathcal{O} is maintained, not the formal question of what $\text{vol}(G)$ measures.

The appendix develops four components: a goal similarity metric that determines when two goals are close in capability space, an agent similarity metric that determines when two observations correspond to the same agent, a zero-knowledge proof mechanism for the strongest form of identity verification, and a clustering procedure that groups agents into modal categories by goal profile. Together, these provide the substrate on which the trust model (Appendix B) and the local volume estimator (Definition 8) operate.

C.1 Goal Similarity Metric S_G

A GFM actor proposing actions to expand $\text{vol}(G)$ must determine which goals are similar—close in the capability space—and which are distant. Without a similarity metric, the actor cannot generalize from observed goals to unobserved ones, cannot identify redundant capabilities, and cannot efficiently allocate effort across the population.

Definition 10 (Goal Similarity). *For two goals $g_i, g_j \in \mathcal{G}$, the goal similarity $S_G(g_i, g_j)$ is defined as the overlap in their projected capability contributions, conditioned on the GFM actor’s current world model M :*

$$S_G(g_i, g_j \mid M) = \frac{\text{vol}(\Phi_M(g_i) \cap \Phi_M(g_j))}{\text{vol}(\Phi_M(g_i) \cup \Phi_M(g_j))} \quad (27)$$

where $\Phi_M(g)$ is the projected capability region—the set of capabilities in \mathcal{G} that the world model M associates with goal g , including both the direct capability and the instrumental capabilities required to realize it.

The metric has the structure of a Jaccard index over projected capability regions: $S_G = 1$ when two goals project to identical capability sets (they are the same goal under different descriptions), and $S_G = 0$ when their projections are disjoint (they share no instrumental or terminal capabilities). The conditioning on M is essential—two goals that appear similar under one world model may be dissimilar under another, since the instrumental capabilities connecting them depend on the model’s beliefs about the environment.

Properties. S_G satisfies:

- *Symmetry:* $S_G(g_i, g_j \mid M) = S_G(g_j, g_i \mid M)$.
- *Boundedness:* $S_G \in [0, 1]$ by construction.
- *Model-dependence:* S_G is a property of goals-as-projected-through-a-model, not an intrinsic property of goals. This dependence is constructive: the similarity metric improves as the world model improves, and two GFM actors with different world models may assign different similarities to the same goal pair.

Relation to $\text{vol}(G)$. Goal similarity connects to the optimization objective through redundancy. When $S_G(g_i, g_j \mid M)$ is high, the two goals contribute overlapping capability volume to G . Expanding both adds less to $\text{vol}(G)$ than expanding two dissimilar goals, because the intersection is counted only once. A GFM actor seeking to maximize $\Delta \text{vol}(G)$ per unit effort should preferentially expand goals with low pairwise similarity—diversifying the capability space rather than deepening existing niches. This provides a formal basis for the intuition that diversity is instrumentally valuable: a population with diverse goals has higher $\text{vol}(G)$ than one with uniform goals, all else equal, because the union of dissimilar sets has greater measure than the union of similar ones.

C.2 Agent Similarity Metric S_A

Goal similarity measures the relationship between goals; agent similarity measures the relationship between the entities pursuing them. The GFM actor needs agent similarity for a different purpose: to determine whether two observations correspond to the same agent (re-identification), whether an agent’s behavior is consistent over time (stability), and how much weight to assign an agent’s reports in the estimator (trust).

Definition 11 (Agent Similarity). *For two agent observations $a_i(t_1)$ and $a_j(t_2)$ (potentially the same agent at different times), the agent similarity S_A is defined as:*

$$S_A(a_i(t_1), a_j(t_2)) = \max\left(\underbrace{S_{\text{beh}}(a_i(t_1), a_j(t_2))}_{\text{behavioral consistency}}, \underbrace{S_{\text{zkp}}(a_i(t_1), a_j(t_2))}_{\text{identity verification}}\right) \quad (28)$$

where S_{beh} measures behavioral consistency and S_{zkp} measures cryptographic identity verification. The \max operator reflects that either channel alone is sufficient for identification: a cryptographic proof of identity overrides behavioral evidence, and strong behavioral consistency can establish practical identity in the absence of cryptographic proof.

The behavioral component S_{beh} aggregates multiple observable features of the agent:

$$S_{\text{beh}}(a_i(t_1), a_j(t_2)) = \sum_{\ell} w_{\ell} \cdot s_{\ell}(a_i(t_1), a_j(t_2)) \quad (29)$$

where each $s_{\ell} \in [0, 1]$ measures consistency along a specific dimension and w_{ℓ} are learned weights summing to 1. The dimensions include:

1. *Goal consistency:* $S_G(G_i^{\text{ind}}(t_1), G_j^{\text{ind}}(t_2) \mid M)$ —do the two observations project to similar capability sets?
2. *Behavioral patterns:* Consistency of action selection, communication style, and response timing.
3. *Attentional focus:* Consistency of which aspects of the environment the agent attends to and which it ignores.
4. *Environmental context:* Consistency of spatial, temporal, and social context—an agent observed in incompatible locations simultaneously cannot be the same entity.

The weights w_{ℓ} are learned from the GFM actor’s interaction history: dimensions that have historically been more predictive of identity receive higher weight. This learning is itself subject to adversarial manipulation (an agent could mimic another’s behavioral patterns), which is why the cryptographic component S_{zkp} exists as an independent verification channel.

Relation to trust. Agent similarity connects to the trust model through the observable goal model. When $S_A(a_k(t), a_k(t-1))$ is high—the agent appears consistent over time—the prediction residual $r_k(t)$ (Equation 18) will be small, and trust T_k will be maintained or increased. When S_A drops (the agent appears inconsistent, as though it has been replaced or is behaving erratically), the prediction residual increases and trust decreases. The agent similarity metric thus provides the input signal that drives the trust update dynamics of Appendix B.

C.3 Zero-Knowledge Proof for Agent Re-identification

The strongest form of agent identity verification is cryptographic: an agent proves it is the same entity the GFM actor interacted with previously, without revealing any private information beyond the fact of identity.

Definition 12 (Identity Proof). *An agent a_k can establish identity with the GFM actor through a zero-knowledge proof P_{zk} if the following conditions hold:*

1. *Prior interaction:* The GFM actor and a_k have established a shared secret or commitment during a previous interaction at time t_0 .
2. *Challenge-response:* The GFM actor issues a challenge c derived from the shared commitment. The agent produces a response r such that $\text{Verify}(c, r, \text{commitment}) = 1$.

3. Zero-knowledge: *The response r reveals no information about the shared secret beyond the fact that the responder possesses it.*

When the proof succeeds, the identity verification component is set to $S_{zkp} = 1$, and the GFM actor can directly assign the prior relationship’s accumulated trust and goal model to the current observation:

$$P_{zk}(a_k(t), a_k(t_0)) = 1 \implies S_A(a_k(t), a_k(t_0)) = 1 \quad (30)$$

The zero-knowledge proof serves as an *override* for the behavioral similarity metric. An agent whose behavior has changed significantly (perhaps because its capabilities have expanded, its environment has shifted, or it has undergone internal restructuring) may have low S_{beh} despite being the same entity. The cryptographic proof resolves this ambiguity definitively: if the proof succeeds, the agent is the same entity regardless of behavioral divergence.

Applications. The primary use of zero-knowledge identity proofs in the GFM framework is *trust persistence*. When an agent a_k goes unobserved for an extended period and reappears, the GFM actor faces a choice: treat a_k as a new agent (low initial trust, cooling period) or as a returning agent (prior trust level, accumulated model). Without cryptographic verification, the actor must rely on behavioral similarity, which degrades over time as the agent’s behavior evolves. With a zero-knowledge proof, the agent can immediately recover its prior trust level, bypassing the cooling period.

A second application is defense against impersonation. An adversary a_s that mimics agent a_k ’s behavioral patterns could achieve high S_{beh} and inherit a_k ’s trust level—a spoofing attack. The zero-knowledge proof prevents this: an impersonator cannot produce a valid response to the challenge without possessing the shared secret, regardless of how accurately it mimics a_k ’s observable behavior.

Limitations. The zero-knowledge proof requires a prior interaction during which the shared commitment is established. It cannot verify the identity of a truly novel agent. It also assumes that the agent’s private key (or equivalent secret) has not been compromised—if an adversary obtains the shared secret, it can pass the identity proof and inherit the victim’s trust. The proof establishes that the current entity *possesses the same secret* as the prior interactant, not that it *is* the same entity in any deeper sense. For practical purposes, this distinction matters only when secrets can be stolen, which is an implementation security question rather than a framework question.

C.4 Modal Groupings

As a GFM actor accumulates goal models $\{M_k(G)\}$ across many agents, structure emerges: clusters of agents with similar goal profiles, pursuing overlapping capabilities through similar strategies. These clusters are *modal groupings*: the modes of the goal distribution across the observed population.

Definition 13 (Modal Grouping). *A modal grouping \mathcal{M}_j is a subset of observed agents whose pairwise goal similarity exceeds a threshold η :*

$$\mathcal{M}_j = \{a_k \in \mathcal{O} : S_G(G_k^{\text{ind}}, \mu_j | M) > \eta\} \quad (31)$$

where μ_j is the centroid of the grouping in projected capability space and \mathcal{O} is the GFM actor’s observation set. The set of all modal groupings $\{\mathcal{M}_1, \dots, \mathcal{M}_m\}$ partitions the observed population into goal-coherent clusters, with agents below η similarity to all centroids assigned to a residual category.

The modal groupings serve three functions in the GFM framework.

Goal proposal. When the GFM actor proposes new capabilities to expand $\text{vol}(G)$, the modal groupings identify which proposals are likely to receive positive responses from which agents. A proposal aligned with μ_j ’s centroid is likely to be accepted by agents in \mathcal{M}_j and ignored or rejected by agents in other groupings. The actor can target proposals to the groupings where they will produce the largest $\hat{\Delta}(\pi)$ signal, improving the efficiency of the optimization loop (Section 4.2).

Diversity monitoring. The number and spread of modal groupings provides a coarse measure of the goal diversity in the population. A population with many well-separated groupings has higher potential $\text{vol}(G)$ than one with few tightly clustered groupings, because the union of dissimilar capability sets has greater measure. If the GFM actor observes modal groupings collapsing (agents converging on fewer, more similar goal profiles), this is a signal of $\text{vol}(G)$ stagnation or impending contraction, even if no individual agent’s capabilities have changed.

Scorpion contextualization. The modal groupings provide context for interpreting scorpion behavior. An agent that contracts capabilities within its own modal grouping is likely a cornered animal or a defector from internal conflict. An agent that contracts capabilities across multiple groupings is likely pursuing an exogenous objective that crosses group boundaries. An agent whose behavior is inconsistent with *all* modal groupings may be at the alien end of the value-divergence spectrum (Section A.4), operating under axioms orthogonal to the entire observed population.

Updating groupings. The modal structure is not static. As the GFM actor’s world model improves and agents’ goals evolve, groupings merge, split, and shift. The update procedure follows the same exponentially weighted scheme as the trust model: new observations are weighted more heavily than old ones (controlled by the decay parameter α from Equation 19), so the grouping structure tracks the current population rather than reflecting historical patterns. The centroid μ_j is updated at each observation as an exponentially weighted moving average of the member agents’ projected capability sets.

Relation to $\text{vol}(G)$. Modal groupings approximate the structure of the joint goal space G at the population level. The union of all groupings’ projected capability sets provides a coarse estimate of $\bigcup_k G_k^{\text{ind}}$, and the gaps between groupings identify regions of \mathcal{G} where no observed agent has capabilities, marking potential expansion targets. A GFM actor can use the grouping structure to identify underserved regions of capability space and prioritize actions that expand $\text{vol}(G)$ into these gaps, rather than deepening capabilities in already-dense regions. This connects modal groupings to the diversity argument of Section C.1: the groupings operationalize the intuition that frontier expansion is most efficient at the boundaries between known and unknown capabilities.

D Formal Proofs

This appendix provides full proofs for the results stated in Sections 3–5. Each proof proceeds from the measure-space axioms restated below and the definitions of Section 2.

D.1 Measure Space Axioms

The following properties of $\text{vol}(\cdot)$ are assumed throughout. They restate Definition 5 and the standing assumptions of Section 2 in a form convenient for the proofs.

- (M1) *Non-negativity.* $\text{vol}(A) \geq 0$ for every measurable $A \subseteq \mathcal{G}$.
- (M2) *Null empty set.* $\text{vol}(\emptyset) = 0$.
- (M3) *Monotonicity.* $A \subseteq B \implies \text{vol}(A) \leq \text{vol}(B)$.
- (M4) *Finite additivity.* If $A \cap B = \emptyset$, then $\text{vol}(A \cup B) = \text{vol}(A) + \text{vol}(B)$.
- (M5) *Non-triviality.* For every capability $g \in \mathcal{G}$, the singleton $\{g\}$ has $\text{vol}(\{g\}) > 0$. (Standing assumption from Section 2.)
- (M6) *Superadditivity under independence.* When agents' individual capability sets are non-overlapping and non-interacting, $\text{vol}(G) \geq \sum_k \text{vol}(G_k^{\text{ind}})$, with strict inequality whenever $G^{\text{coop}} \neq \emptyset$. (Definition 5.)

Axioms (M1)–(M4) are standard properties of a finite measure on a σ -algebra. (M5) is non-standard: it requires that the measure resolves individual capabilities, which constrains the granularity of \mathcal{G} . (M6) is the superadditivity assumption from Definition 5.

D.2 Proof of Lemma 2 (Coercive Actions)

Proof. Let π restrict agent a_j 's capabilities from G_j^{ind} to $G_j^{\text{ind}'} \subset G_j^{\text{ind}}$. Define the restricted joint space as $G' = (\bigcup_{i \neq j} G_i^{\text{ind}} \cup G_j^{\text{ind}'}) \cup G^{\text{coop}'}$, where $G^{\text{coop}'} = \{g \in G^{\text{coop}} : g \text{ remains achievable under } G_j^{\text{ind}'}\}$. Since $G_j^{\text{ind}'} \subset G_j^{\text{ind}}$, any cooperative capability that required a capability in $G_j^{\text{ind}} \setminus G_j^{\text{ind}'}$ is no longer achievable, so $G^{\text{coop}'} \subseteq G^{\text{coop}}$. Combined with $G_j^{\text{ind}'} \subset G_j^{\text{ind}}$, we have $G' \subseteq G$. By (M3), $\text{vol}(G') \leq \text{vol}(G)$, giving $\Delta \text{vol}(G) \leq 0$.

For strictness: suppose there exists $g \in G_j^{\text{ind}} \setminus G_j^{\text{ind}'}$ with $g \notin \bigcup_{i \neq j} G_i^{\text{ind}}$. Then $g \in G \setminus G'$, so $G' \subsetneq G$. Write $G = G' \cup (G \setminus G')$ where the union is disjoint. By (M4), $\text{vol}(G) = \text{vol}(G') + \text{vol}(G \setminus G')$. Since $\{g\} \subseteq G \setminus G'$, by (M3) and (M5), $\text{vol}(G \setminus G') \geq \text{vol}(\{g\}) > 0$, giving $\text{vol}(G) > \text{vol}(G')$. ■

D.3 Proof of Lemma 3 (Self-Imposed Rigidity)

Proof. Since $A \setminus R \subseteq A$, the feasible set of the restricted problem is a subset of the feasible set of the unrestricted problem. Therefore

$$\max_{\pi \in A \setminus R} \mathbb{E}[\Delta \text{vol}(G) \mid \pi] \leq \max_{\pi \in A} \mathbb{E}[\Delta \text{vol}(G) \mid \pi].$$

For strictness: let $\pi^* = \arg \max_{\pi \in A} \mathbb{E}[\Delta \text{vol}(G) \mid \pi]$. If $\pi^* \in R$ (the unrestricted optimum is forbidden), then π^* is not feasible in $A \setminus R$, so $\max_{\pi \in A \setminus R} \mathbb{E}[\Delta \text{vol}(G) \mid \pi]$ is attained at some $\pi' \neq \pi^*$ with $\mathbb{E}[\Delta \text{vol}(G) \mid \pi'] \leq \mathbb{E}[\Delta \text{vol}(G) \mid \pi^*]$. If π^* is the unique maximizer, the inequality is strict. If π^* is not unique, strictness holds whenever every co-maximizer also lies in R . ■

D.4 Proof of Corollary 1.2 (Rigid Rules are Anti-Maximizing)

Proof. By Lemma 3, strict suboptimality holds whenever $\pi^* \in R$. It remains to show that such states exist.

Call R *non-trivial* if there exists a reachable state s and an action $\pi_0 \in R$ with $\mathbb{E}[\Delta \text{vol}(G) \mid \pi_0, s] > 0$: at least one forbidden action is $\text{vol}(G)$ -expanding in some state. For any non-trivial R , let s_0 be the state where π_0 is expanding. If π_0 also maximizes $\mathbb{E}[\Delta \text{vol}(G) \mid \cdot, s_0]$ over A , we are

done. If not, we appeal to an adversarial construction: an adversary who knows the structure of R can construct a state (via deadline pressure) in which every action in $A \setminus R$ is $\text{vol}(G)$ -contracting while a forbidden action (e.g., deception, calculated risk) is $\text{vol}(G)$ -expanding. In that constructed state, $\pi^* \in R$ and Lemma 3 gives strict inequality. The concrete instance following Corollary 1.2 provides one such construction; the argument generalizes to any R whose constraints are known to the adversary, since knowledge of the permitted region $A \setminus R$ suffices to design scenarios that trap the actor within it. ■

D.5 Proof of Lemma 1 (Net Effect of Elimination)

Proof. Setup. Partition the joint goal space G into four disjoint regions relative to agent a_k :

$$U_k = G_k^{\text{ind}} \setminus \bigcup_{j \neq k} G_j^{\text{ind}} \quad (a_k \text{'s unique individual capabilities})$$

$$S_k = G_k^{\text{ind}} \cap \bigcup_{j \neq k} G_j^{\text{ind}} \quad (\text{shared individual capabilities})$$

$$K_k = G_k^{\text{coop}} \quad (\text{cooperative capabilities requiring } a_k)$$

$$R = G \setminus (U_k \cup S_k \cup K_k) \quad (\text{everything else: other agents' unique caps + coop not involving } a_k)$$

By construction these four sets are disjoint and $G = U_k \cup S_k \cup K_k \cup R$. By (M4), $\text{vol}(G) = \text{vol}(U_k) + \text{vol}(S_k) + \text{vol}(K_k) + \text{vol}(R)$.

After elimination. Removing a_k produces G' . The shared capabilities S_k survive (other agents possess them). The remainder R is unaffected. The unique capabilities U_k are lost. The cooperative capabilities K_k are lost. Additionally, removing a_k may unlock *suppression-released* capabilities $G_k^{\text{coop}+} \subseteq G \setminus G$: cooperative capabilities among the remaining agents that were not achievable while a_k was present (e.g., because a_k 's interference prevented the required coordination). By definition, $G_k^{\text{coop}+} \cap G = \emptyset$, since these capabilities were not in G before removal.

The post-elimination space is therefore $G' = S_k \cup R \cup G_k^{\text{coop}+}$, and since $G_k^{\text{coop}+}$ is disjoint from $S_k \cup R$ (it was not in G),

$$\text{vol}(G') = \text{vol}(S_k) + \text{vol}(R) + \text{vol}(G_k^{\text{coop}+}).$$

Volume change.

$$\begin{aligned} \Delta \text{vol}(G) &= \text{vol}(G') - \text{vol}(G) \\ &= [\text{vol}(S_k) + \text{vol}(R) + \text{vol}(G_k^{\text{coop}+})] - [\text{vol}(U_k) + \text{vol}(S_k) + \text{vol}(K_k) + \text{vol}(R)] \\ &= \text{vol}(G_k^{\text{coop}+}) - \text{vol}(U_k) - \text{vol}(K_k) \end{aligned}$$

which is the claimed Equation (5).

Sign. Elimination is net-contracting ($\Delta \text{vol}(G) < 0$) when $\text{vol}(U_k) + \text{vol}(K_k) > \text{vol}(G_k^{\text{coop}+})$. Under (M5), if a_k has at least one unique capability, $\text{vol}(U_k) > 0$; if a_k participates in at least one coalition, $\text{vol}(K_k) > 0$. Since $G_k^{\text{coop}+}$ is typically empty or small (most agents do not suppress cooperation by their mere presence), the condition holds for generic agents. For scorpions that actively suppress cooperation, $\text{vol}(G_k^{\text{coop}+})$ can dominate, making elimination net-expanding. ■

D.6 Proof of Corollary 1.1 (Elimination is Almost Always Anti-Maximizing)

Proof. Part 1: Positive unique contribution. Suppose a_k has at least one capability $g \in G_k^{\text{ind}}$ with $g \notin G_j^{\text{ind}}$ for all $j \neq k$. Then $\{g\} \subseteq U_k = G_k^{\text{ind}} \setminus \bigcup_{j \neq k} G_j^{\text{ind}}$. By (M3) and (M5), $\text{vol}(U_k) \geq \text{vol}(\{g\}) > 0$.

Part 2: Cooperative loss grows with interaction partners. For each agent $a_j \neq a_k$ that interacts with a_k , define the pairwise cooperative contribution $G_{\{k,j\}}^{\text{coop}}$ as the set of capabilities achievable by the pair $\{a_k, a_j\}$ but not by either alone and not by any other pair not involving a_k . The sets $\{G_{\{k,j\}}^{\text{coop}}\}_{j \neq k}$ are not necessarily disjoint (a capability may require a_k and either of two partners), so we cannot directly sum their volumes.

However, for each interacting partner a_j , assume there exists at least one capability $g_j \in G_{\{k,j\}}^{\text{coop}}$ that is unique to that pair (not achievable by a_k with any other single partner). Under (M5), $\text{vol}(\{g_j\}) > 0$. The set $\{g_j : j \text{ interacts with } a_k\}$ has cardinality equal to the number of interaction partners m_k . If the $\{g_j\}$ are distinct, then by (M4) applied to the disjoint singletons, $\text{vol}(G_k^{\text{coop}}) \geq \sum_j \text{vol}(\{g_j\}) > 0$, and this bound grows linearly in m_k .

The combinatorial observation in the main text (that $2^{n-1} - 1$ coalitions involving a_k exist) motivates the expectation of super-linear growth. Formalizing this requires a *diversity condition*: each coalition $S \ni a_k$ with $|S| \geq 2$ produces at least one capability not achievable by any proper sub-coalition of S containing a_k . Under this condition, the number of unique coalition-level capabilities grows with the number of coalitions, but the precise scaling depends on the overlap structure among coalitions and is not characterized here. The linear lower bound in m_k holds without the diversity condition. ■

D.7 Proof of Lemma 4 (Cooperative Expansion)

Proof. Part (i). By definition, $G_{ij}^{\text{coop}}(t)$ contains capabilities requiring the simultaneous participation of both a_i and a_j . A unilateral action π by a_i alone does not involve a_j 's participation. No capability in G_{ij}^{coop} can be created or destroyed by π , since its defining condition (joint participation) is not met. Therefore $G_{ij}^{\text{coop}}(t+1) = G_{ij}^{\text{coop}}(t)$ under π , giving $\Delta \text{vol}(G_{ij}^{\text{coop}} | \pi) = 0$. By symmetry the same holds for unilateral actions by a_j .

Part (ii). Call a_i and a_j *complementary at time t* if there exists a joint action profile $(\pi_i, \pi_j) \in A_i \times A_j$ that produces a capability $g \notin G(t)$: a capability not achievable under the current coordination infrastructure. Suppose a_i and a_j are complementary and have not yet established the coordination mechanism that enables g (e.g., a communication channel, shared protocol, or mutual commitment). Then $g \notin G(t)$. The joint action π_{ij} that establishes the mechanism makes g achievable, so $g \in G_{ij}^{\text{coop}}(t+1)$ and $G_{ij}^{\text{coop}}(t+1) \supsetneq G_{ij}^{\text{coop}}(t)$. By (M5), $\text{vol}(\{g\}) > 0$, so $\text{vol}(G_{ij}^{\text{coop}}(t+1)) \geq \text{vol}(G_{ij}^{\text{coop}}(t)) + \text{vol}(\{g\}) > \text{vol}(G_{ij}^{\text{coop}}(t))$, giving $\Delta \text{vol}(G_{ij}^{\text{coop}} | \pi_{ij}) > 0$. Whether this cooperative gain produces a net expansion in $\text{vol}(G)$ depends on whether the opportunity cost C_{ij} of coordination (individual capabilities foregone) exceeds the cooperative gain; part (iii) addresses this.

Part (iii). Let $\Delta^{\text{coop}} = \mathbb{E}[\Delta \text{vol}(G_{ij}^{\text{coop}})]$ be the expected cooperative expansion from joint action, and let $C_{ij} \geq 0$ be the opportunity cost: the $\text{vol}(G)$ value of the best individual actions a_i and a_j forgo in order to coordinate. The net expected change in $\text{vol}(G)$ from cooperation is $\Delta^{\text{coop}} - C_{ij}$, which is positive whenever $\Delta^{\text{coop}} > C_{ij}$. A GFM actor maximizing $\mathbb{E}[\Delta \text{vol}(G)]$ therefore prefers cooperation to unilateral action when this condition holds. ■

D.8 Proof of Proposition 3 (Sign-Correctness Decomposition)

Proof. Write the true volume change as $\Delta \text{vol}(G) = \Delta_I + \Delta_C$, where $\Delta_I = \Delta \text{vol}(\bigcup_k G_k^{\text{ind}})$ is the individual-capability component and $\Delta_C = \Delta \text{vol}(G^{\text{coop}})$ is the cooperative component. The estimator observes $\hat{\Delta}_I \approx \Delta_I$ (conditions 1–2 control the approximation error) and $\hat{\Delta}_C$ with error $e_C = \Delta_C - \hat{\Delta}_C$.

The estimator's output is $\hat{\Delta}(\pi) = \hat{\Delta}_I + \hat{\Delta}_C = (\Delta_I - e_I) + (\Delta_C - e_C) = \Delta \text{vol}(G) - (e_I + e_C)$, where $e_I = \Delta_I - \hat{\Delta}_I$ is the individual-level error.

Under conditions 1–2, $|e_I|$ is small relative to $|\Delta_I|$ (the actor observes a representative sample with an accurate trust model). Under condition 3, $|e_C| < (1 - \epsilon) \cdot |\Delta_I|$.

For sign-correctness we need $\text{sign}(\hat{\Delta}(\pi)) = \text{sign}(\Delta \text{vol}(G))$. Since $\hat{\Delta}(\pi) = \Delta \text{vol}(G) - (e_I + e_C)$, the sign is preserved when $|e_I + e_C| < |\Delta \text{vol}(G)|$. By the triangle inequality, $|e_I + e_C| \leq |e_I| + |e_C|$. Write $|e_I| \leq \delta |\Delta_I|$ for a sampling error ratio δ controlled by conditions 1–2, and $|e_C| < (1 - \epsilon) |\Delta_I|$ by condition 3. Two cases:

Case 1: Δ_I and Δ_C have the same sign. Then $|\Delta \text{vol}(G)| = |\Delta_I| + |\Delta_C| \geq |\Delta_I|$, and $|e_I| + |e_C| < (\delta + 1 - \epsilon) |\Delta_I|$, which is less than $|\Delta_I| \leq |\Delta \text{vol}(G)|$ whenever $\delta < \epsilon$.

Case 2: Δ_I and Δ_C have opposite signs. Then $|\Delta \text{vol}(G)| = ||\Delta_I| - |\Delta_C||$, and sign preservation requires $(\delta + 1 - \epsilon) |\Delta_I| < ||\Delta_I| - |\Delta_C||$. This holds when $|\Delta_C| < (\epsilon - \delta) |\Delta_I|$, a condition strictly

stronger than condition 3 alone. Condition 3 guarantees sign-correctness in this case only when the cooperative-level *change* (not just the estimation error) is small relative to the individual-level change. ■

D.9 Proof of Proposition 1 (Self-Balancing Property)

Proof. Parts (a)–(c) follow directly from the lemmas.

(a) By Lemma 1, eliminating agent a_k produces $\Delta \text{vol}(G) = \text{vol}(G_k^{\text{coop}+}) - \text{vol}(U_k) - \text{vol}(K_k)$. When $\text{vol}(U_k) + \text{vol}(K_k) > \text{vol}(G_k^{\text{coop}+})$ (the generic case for non-scorpion agents), destruction is $\text{vol}(G)$ -contracting.

(b) By Lemma 2, restricting a_j 's capabilities produces $G' \subseteq G$ and $\text{vol}(G') \leq \text{vol}(G)$, with strict inequality when the restriction removes capabilities not covered by other agents.

(c) By Lemma 3, constraining the actor to $A \setminus R$ gives $\max_{A \setminus R} \mathbb{E}[\Delta \text{vol}(G)] \leq \max_A \mathbb{E}[\Delta \text{vol}(G)]$, with strict inequality when R excludes the unrestricted optimum.

Part (d): Partial result under additional assumptions. Parts (a)–(c) establish that the $\text{vol}(G)$ objective penalizes movement toward both extremes: excessive destruction contracts G by (a), excessive restriction contracts G by (b), and rigid self-constraint reduces the actor's ability to expand G by (c). To show that these opposing pressures produce an optimum distinct from both extremes, assume:

- The action space A is compact.
- The map $\pi \mapsto \mathbb{E}[\Delta \text{vol}(G) \mid \pi]$ is continuous.
- There exist actions $\pi_d, \pi_r \in A$ such that π_d is destructive ($\mathbb{E}[\Delta \text{vol}(G) \mid \pi_d] < 0$ by (a)) and π_r is overrestrictive ($\mathbb{E}[\Delta \text{vol}(G) \mid \pi_r] < 0$ by (b)).
- There exists $\pi_0 \in A$ with $\mathbb{E}[\Delta \text{vol}(G) \mid \pi_0] \geq 0$.

By the extreme value theorem (continuity on a compact set), the maximum of $\mathbb{E}[\Delta \text{vol}(G) \mid \pi]$ over A is attained at some $\pi^* \in A$. Since π_0 achieves a non-negative value, $\mathbb{E}[\Delta \text{vol}(G) \mid \pi^*] \geq 0$. Since both π_d and π_r yield strictly negative values, $\pi^* \neq \pi_d$ and $\pi^* \neq \pi_r$. The maximizer is therefore distinct from the destructive and overrestrictive extremes.

This argument establishes *existence* of an optimum that is neither maximally destructive nor maximally restrictive, under continuity and compactness. It does not establish *stability* (whether small perturbations return to the optimum), *uniqueness* (whether there is one such point or many), or *convergence* (whether the optimization loop of Section 4.2 reaches it). These remain open problems. ■

D.10 Proof of Proposition 2 (Scorpion Detection)

Proof. Channel (a): Contraction attribution. Let a_j be a scorpion whose actions persistently contract $\text{vol}(G)$. At each time step t , the GFM actor observes the capability changes $\Delta \text{vol}(G_k^{\text{ind}} \mid t)$ for each $a_k \in \mathcal{O}$ (the observation set). When a_j acts, the affected agents report $R_k(t) \in \{-1, 0, +1\}$. If a_j 's actions persistently reduce capabilities, the trust-weighted aggregate $\hat{\Delta}_j(t) = \sum_{k \in \mathcal{O}} T_k \cdot R_k(t)$ is persistently negative in time steps where a_j acts.

Formally, define the running average contraction signal attributed to a_j :

$$\bar{\Delta}_j(\tau) = \frac{1}{\tau} \sum_{t=1}^{\tau} \hat{\Delta}_j(t) \cdot \mathbf{1}[a_j \text{ acted at } t].$$

If a_j 's actions produce expected contraction $\mu_j < 0$, the observation noise has finite variance σ^2 , and the observation windows (time steps where a_j acts) yield conditionally independent signals given a_j 's strategy, then by the strong law of large numbers, $\bar{\Delta}_j(\tau) \rightarrow \mu_j < 0$ almost surely as $\tau \rightarrow \infty$. After $\tau^* = O(\sigma^2/\mu_j^2)$ observations, the signal exceeds any fixed detection threshold with high probability. The independence condition is non-trivial: successive observations of a scorpion's

effects may be correlated when the same agents are repeatedly affected. Under weaker mixing conditions the convergence still holds but the rate degrades.

This is a *correlational* detection: the signal identifies that negative capability changes co-occur with a_j 's actions. Causal attribution (separating a_j 's contribution from simultaneous causes) requires additional structure not provided by the framework.

Channel (b): Deception detection. If a_j reports capability changes that diverge from independently observed outcomes, the prediction residual $r_j(t) = R_j(t) - \hat{R}_j(t | M_j(G))$ is nonzero. The cumulative inconsistency $\sigma_j^2(t) = \alpha \cdot \sigma_j^2(t-1) + (1-\alpha) \cdot r_j(t)^2$ (Equation (19)) is an exponentially weighted moving average (EWMA) of squared residuals.

Under stationarity of a_j 's deceptive behavior, $\sigma_j^2(t)$ converges to $\mathbb{E}[r_j^2]$ (Appendix B.4, Equation (25)). For a deceptive scorpion, $\mathbb{E}[r_j^2] > 0$, so $T_j \rightarrow T_j^* = 1/(1+\beta \cdot \mathbb{E}[r_j^2]) < 1$ (Equation (20)). The convergence rate is controlled by the decay parameter α : after approximately $1/(1-\alpha)$ observations, the prior's influence is negligible and T_j reflects the agent's actual prediction consistency.

For an honest (non-deceptive) cooperative agent, $\mathbb{E}[r_j^2]$ is small (driven by observation noise), so $T_j^* \approx 1$. The gap $T_j^*(\text{deceptive}) < T_j^*(\text{honest})$ provides a detection threshold that improves with observation count.

Convergence caveat. Both channels require stationarity of the scorpion's behavior. A scorpion that alternates between cooperative and contracting phases, or that changes its strategy in response to detection pressure, can maintain $\bar{\Delta}_j \approx 0$ and $T_j \approx 1$ indefinitely. Detection of non-stationary adversaries requires adaptive methods beyond the EWMA framework and remains an open problem. ■