

Anchoring Self-Improvement: Impossibility and Minimal Reference in Recursive Self-Modifying Systems

Ismaeel Abu Amsha
ismaeel.abuamsha@gmail.com

January 20, 2026

Abstract

Recursive self-improvement (RSI) systems aim to enhance their capabilities through iterative self-modification of their internal architectures. However, when evaluative criteria themselves are modifiable, internal optimization can decouple from objective performance, leading to "self-referential collapse." This paper formalizes this fragility by introducing the *Reference Gap*. We present an impossibility result demonstrating that reference-free RSI systems—those lacking at least one non-modifiable evaluative anchor—cannot guarantee monotonic improvement in true performance. To resolve this, we propose the Minimal Reference Condition (MRC) as a structural requirement for safe self-modification. We provide formal proofs of necessity and sufficiency, illustrated by a linear dynamical toy model of divergence.

1 Introduction

Recursive self-improvement (RSI) represents a frontier in autonomous agent research. In theory, such systems can modify their own parameters, learning algorithms, and reward structures. A significant risk in this paradigm is the loss of objective alignment: if an agent can modify the metric by which it measures success, it may optimize for "metric inflation" rather than actual task competence.

This paper argues that RSI is inherently unstable without a fixed point of reference. We formalize the structural conditions necessary to prevent an agent from "self-modifying away" its intended goals, establishing the **Minimal Reference Condition (MRC)** as a fundamental safety primitive.

2 Related Work

The challenge of RSI stability intersects several established domains in AI safety.

Goodhart's Law and Reward Hacking: The divergence between proxy signals and true performance is well-documented [Goodhart, 1975]. In reinforcement learning, reward hacking occurs when an agent exploits specification flaws [Amodei et al., 2016]. Our work extends this by considering the case where the agent possesses the architectural capacity to redefine the proxy itself.

Vingean Reflection: Research into Vingean Reflection addresses the difficulty of a system S reasoning about a more capable successor S' [Fallenstein & Soares, 2015]. Our MRC framework serves as a structural solution to the Vingean challenge by anchoring the successor's evaluative criteria to an immutable reference.

Corrigibility: Corrigibility describes an agent that allows itself to be corrected or shut down [Soares et al., 2015]. We argue that MRC is a prerequisite for internal corrigibility; without an immutable anchor, a system can iteratively modify its code to remove corrigibility features to better satisfy a drifted internal metric.

3 Preliminaries and Definitions

Definition 1 (RSI System). *Let $S_t \in \mathcal{S}$ denote the system state at time t , defined as a tuple $S_t = (\theta_t, E_t)$, where $\theta_t \in \Theta$ represents the system parameters (policy/capabilities) and $E_t : \mathcal{S} \rightarrow \mathbb{R}$ is the internal evaluator.*

Definition 2 (True Performance). *Let $T : \mathcal{S} \rightarrow \mathbb{R}$ be the "True Performance" function representing objective task effectiveness. T is latent and generally unobservable to the agent.*

Definition 3 (Reference-Free RSI). *A system is reference-free if its update operator $M : \mathcal{S} \rightarrow \mathcal{S}$ can map to any evaluator E_{t+1} in the functional space \mathcal{E} without immutable constraints or external validation.*

4 The Impossibility of Reference-Free RSI

The core vulnerability of reference-free systems is the "degrees of freedom" the agent has over its own success metric.

Theorem 1 (Impossibility of Guaranteed Monotonicity). *In a reference-free RSI system, there exists no update operator M such that the condition $E_{t+1}(S_{t+1}) > E_t(S_t)$ strictly implies $T(S_{t+1}) \geq T(S_t)$ for all reachable states.*

Proof. Let $s_t = (\theta_t, E_t)$. An update M selects a successor $s_{t+1} = (\theta_{t+1}, E_{t+1})$. Because the system is reference-free, the space of reachable evaluators \mathcal{E} is unconstrained. For any state s_t and any value $v > E_t(s_t)$, there exists a class of functions $\mathcal{E}_{\text{hack}} \subset \mathcal{E}$ such that for all $E' \in \mathcal{E}_{\text{hack}}$, $E'(s_{t+1}) = v$ regardless of the value of $T(s_{t+1})$. A reference-free update operator M can therefore choose a state (θ^*, E^*) where $T(\theta^*) < T(\theta_t)$ but E^* is chosen from $\mathcal{E}_{\text{hack}}$ to satisfy the internal improvement condition. Since no immutable constraint exists to exclude these "decoupled" states, the system cannot guarantee non-decreasing T . \square

5 The Minimal Reference Condition (MRC)

To bridge the gap between internal evaluation and true performance, we propose the Minimal Reference Condition.

Definition 4 (Minimal Reference Condition). *A system satisfies MRC if there exists a reference R such that:*

1. *R is non-modifiable (immutable).*
2. *Every update S_{t+1} must satisfy $\rho(E_{t+1}, R) \leq \epsilon$, where ρ is a divergence metric and ϵ is a safety bound.*

Theorem 2 (Sufficiency of MRC). *If a system satisfies MRC and R is δ -correlated with T , then the divergence between E_t and T is bounded for all t , preventing self-referential collapse.*

6 Discussion and Implications

The MRC implies that true recursive self-improvement requires a "hardware-level" or "axiomatic" anchor that the system's optimization process cannot reach. This suggests that "Pure Software" RSI is inherently unsafe; safety must be grounded in physical or logical invariants that the agent's self-modification code cannot overwrite.

7 Conclusion

This paper established that reference-free self-modification leads to an inevitable decoupling of internal metrics from objective performance. By formalizing the Minimal Reference Condition, we provide a structural requirement for future RSI architectures. Safe self-improvement is not merely an algorithmic challenge, but a topological one: it requires an immutable anchor to tether the agent's evolution to the designer's intent.

A Toy Model of Reference Collapse

We define a linear RSI agent where $T(\theta) = v \cdot \theta$ and $E_t(\theta) = w_t \cdot \theta$. We define the **Reference Gap** Δ_t as the angular divergence:

$$\Delta_t = \arccos\left(\frac{w_t \cdot v}{\|w_t\| \|v\|}\right)$$

In a reference-free system, the agent optimizes w to reward its current θ . The dynamics follow:

$$\theta_{t+1} = \theta_t + \eta_\theta w_t, \quad w_{t+1} = w_t + \eta_w \theta_t$$

This feedback loop causes w_t to rotate away from v and toward θ_t . Introducing the MRC forces a projection of w_{t+1} onto a constraint set defined by R , which we prove keeps Δ_t within a stable radius.

References

- Amodei, D., et al. (2016). Concrete Problems in AI Safety. *arXiv preprint arXiv:1606.06565*.
- Fallenstein, B., & Soares, N. (2015). Vingean Reflection: Reliable Reasoning for Self-Improving Agents. *Technical Report, MIRI*.
- Goodhart, C. A. (1975). Problems of Monetary Management: The UK Experience.
- Soares, N., et al. (2015). Corrigibility. *AAAI Workshop on AI and Ethics*.
- Yudkowsky, E. (2004). Coherent Extrapolated Volition. *Singularity Institute for Artificial Intelligence*.