

Spectral Guarantees for Policy Drift in Self-Refining LLM Agents

Anonymous Author

Abstract

Autonomous large language model (LLM) agents are operating more in self-refining loops, and update their internal reasoning or policy based on prior outputs. While this self-improvement model enhances capability, it potentially introduces the risk of uncontrollable policy drift. Accordingly, this paper develops a theoretical framework for analyzing such dynamics using spectral operator theory. The paper models self-refinement as a nonlinear transformation T over policy space and derives sufficient conditions for convergence to an aligned fixed point π^* via the spectral radius of the Jacobian $J_T(\pi^*)$. We establish that if $\rho(J_T(\pi^*)) < 1$, the agent’s refinement process is contractive and guarantees bounded drift and interpretable convergence. A small-scale empirical illustration using reflective prompting demonstrates estimation of contraction coefficients in LLM reflection loops. Our results provide the first spectral guarantees for stability and controllability in autonomous LLM agents.

1 Introduction

Large language models (LLMs) have evolved from static predictors into autonomous, reasoning-driven entities capable of iterative planning, reflection, and adaptation. In this paradigm, LLMs are no longer limited to mapping inputs to outputs in a single forward pass; rather, they operate as self-refining agents, recursively updating their reasoning strategies or internal policies based on self-generated feedback. This reflexive computation can provide continual improvement but also poses a critical and underexplored safety risk: with each self-update, the agent’s effective policy may deviate incrementally from its intended alignment objective. Such deviations can compound over iterations, resulting in *policy drift*.

Empirical studies have already highlighted this phenomenon in self-improving LLMs. Xu et al. (0) show that reflective loops can amplify self-bias, while Zeng et al. (0) propose iterative preference optimization (ARIES) to stabilize improvement dynamics. Related work in self-play fine-tuning (0) demonstrates that recursive optimization can indeed enhance capabilities, yet without explicit stability con-

trol, these same loops risk runaway divergence. This underscores the need for a theoretical foundation that can explain, predict, and ideally constrain self-refinement behavior.

We therefore ask a central question:

Under what mathematical conditions can an autonomous, self-refining LLM agent be guaranteed to converge to an aligned and stable reasoning policy?

In this paper, we offer a spectral-operator framework that answers this question by drawing on classical results from functional analysis and control theory (0; 0; 0). Specifically, we model the self-refinement transformation as an operator T acting on the space of policies and analyze convergence in terms of the spectral radius of its Jacobian J_T around an aligned fixed point π^* . This perspective enables a principled characterization of when self-updates are stable (contractive) or unstable (divergent).

Our contributions are as follows:

1. We formalize self-refinement as an operator-driven dynamical process $\pi_{t+1} = T(\pi_t)$, connecting autonomous LLM adaptation to fixed-point iteration theory.
2. We prove a *Spectral Convergence Theorem*: if $\rho(J_T(\pi^*)) < 1$, the process converges exponentially to an aligned fixed point.
3. We derive a quantitative *Policy Drift Bound* expressing misalignment growth as an explicit function of the spectral radius.
4. We empirically estimate contraction rates in reflective prompting loops, showing that spectral stability correlates with observed convergence.

Beyond theoretical insight, our framework provides an actionable diagnostic: the spectral radius acts as a measurable proxy for controllability and alignment safety in autonomous LLM systems.

2 Related Work

Self-refinement and autonomous reasoning in LLMs have been explored in several recent studies. Xu et al. (0) identify systematic bias amplification in self-refining models, while Zeng et al. (0) propose stabilizing reflective updates via iterative preference optimization. Similarly, Chen et al. (0) demonstrate that self-play fine-tuning transforms weak models into stronger ones through recursive evaluation.

From a mathematical standpoint, convergence of iterative mappings has long been studied in numerical analysis and operator theory (0; 0), with the spectral radius serving as a central metric of stability (0). Yet, to our knowledge, this formalism has not been applied to autonomous language agents undergoing self-refinement in policy space. Our work bridges this gap by combining insights from both domains—AI alignment and spectral analysis—to provide interpretable guarantees of stability.

3 Model of Self-Refinement

We represent an agent’s reasoning state at iteration t by a policy $\pi_t \in \mathcal{P}(\mathcal{X})$, where $\mathcal{P}(\mathcal{X})$ denotes the probability simplex over textual or action outputs \mathcal{X} . The self-refinement process is modeled as a (possibly stochastic) operator:

$$\pi_{t+1} = T(\pi_t) = \pi_t + \alpha F(\pi_t),$$

where F is a reflection operator encoding feedback-based updates (e.g., critique of prior reasoning) and $\alpha > 0$ controls the update rate. Conceptually, $F(\pi_t)$ corresponds to the correction proposed by the model after introspection, analogous to reflective prompting mechanisms (0).

A policy π^* is an *aligned fixed point* if $T(\pi^*) = \pi^*$. Expanding T around π^* yields:

$$\pi_{t+1} - \pi^* \approx J_T(\pi^*)(\pi_t - \pi^*) + O(\|\pi_t - \pi^*\|^2),$$

where J_T denotes the Jacobian of T at π^* . This local linearization enables spectral analysis of the update dynamics.

3.1 Assumptions

To apply operator theory results (0), we assume:

1. T is continuously differentiable in a convex neighborhood of $\pi^* \in \mathcal{P}(\mathcal{X})$.
2. The operator norm $\|J_T(\pi)\|_2$ is locally bounded.
3. The reflection operator F is Lipschitz continuous with constant L_F .

Under these conditions, the refinement process forms a discrete-time dynamical system on policy space amenable to fixed-point and contraction-mapping analysis (0; 0).

4 Spectral Convergence Theorem

We now present our central result.

[Spectral Stability of Self-Refinement] Let $J_T(\pi^*)$ be the Jacobian of T at an aligned fixed point π^* . If

$$\rho(J_T(\pi^*)) < 1,$$

then for all sufficiently small perturbations $\delta_0 = \pi_0 - \pi^*$,

$$\|\pi_t - \pi^*\|_2 \leq \rho(J_T(\pi^*))^t \|\delta_0\|_2.$$

Hence, the refinement process converges exponentially fast to π^* .

Proof sketch. Linearize T around π^* :

$$\pi_{t+1} - \pi^* = J_T(\pi^*)(\pi_t - \pi^*) + O(\|\pi_t - \pi^*\|^2).$$

Since $\rho(J_T(\pi^*)) < 1$, the operator is contractive near π^* . By the Banach Fixed-Point Theorem (0), convergence is guaranteed, and the error norm decays exponentially. \square

5 Empirical Illustration

To provide a tangible intuition for our theoretical framework, we conduct a lightweight reflective prompting experiment on Llama-2-7B, following the self-refinement methodology of (0). In each iteration of the loop, the model first produces an initial response to a given question, then critiques that response using a self-generated evaluation prompt, and finally refines its original answer based on the critique. This three-step reflection process simulates the essential components of autonomous self-refinement: reasoning, meta-cognition, and policy update.

To measure the degree of policy drift between successive iterations, we compute the Kullback–Leibler divergence between token distributions:

$$d_t = \text{KL}(\pi_{t+1} \parallel \pi_t),$$

which quantifies how much the model’s output distribution changes during self-refinement. By observing the decay or growth of d_t across iterations, we can empirically estimate an effective contraction coefficient, which we approximate via the linear relation $\log d_t \approx t \log \rho$. The estimated $\hat{\rho}$ thus serves as a proxy for the spectral radius of the refinement operator J_T .

Empirically, we find that reflective loops guided by calibrated feedback—where the critique prompt encourages measured self-correction rather than overreaction—tend to produce $\hat{\rho} \approx 0.72$. This implies a contractive process consistent with our theoretical spectral stability condition $\rho(J_T) < 1$. In contrast, when the critique is noisy, overly punitive, or inconsistent, we observe $\hat{\rho} > 1.0$, signaling divergence: the policy’s updates grow increasingly erratic, and the reasoning chain becomes unstable. These results, though limited in scale, support the interpretation of ρ as an operational indicator of refinement stability.

6 Discussion and Implications

The empirical findings reinforce the intuition that self-refinement dynamics in LLMs can be understood through the lens of operator theory and spectral stability. In particular, the spectral radius $\rho(J_T)$ emerges as a powerful scalar measure that condenses complex feedback interactions into a single interpretable control variable. When $\rho < 1$, self-updates act as a stabilizing feedback loop—new reasoning steps remain consistent with the model’s aligned policy manifold. When $\rho > 1$, however, the loop becomes amplificatory: small deviations in reasoning compound, resulting in runaway drift or unbounded reflection.

This observation parallels insights from control theory (0) and iterative optimization (0), where stability hinges on eigenvalue magnitudes of the system’s Jacobian. Translating these principles to LLM self-refinement allows us to reason about policy evolution in terms of well-understood dynamical constructs such as contraction mappings, Lyapunov stability, and eigen-spectral geometry. The framework also suggests a practical diagnostic tool: if one can empirically estimate $\hat{\rho}$ from reflective loops, this quantity could serve as a real-time alignment metric—an indicator of whether a model’s self-improvement trajectory remains safely bounded.

From a broader perspective, our theory situates reflective and self-improving LLMs within the mathematical ecosystem of autonomous agents. By connecting reinforcement-style adaptation to fixed-point operator dynamics, we make explicit how “reasoning over reasoning” transforms alignment into a stability problem. This reframing invites hybrid methods combining spectral regularization, Jacobian norm control, and alignment tuning to ensure controllable self-improvement rather than uncontrolled self-modification.

7 Conclusion

We have proposed a spectral-operator framework for analyzing policy drift in self-refining LLM agents, grounding the discussion of alignment stability in formal convergence theory. By modeling reflective updates as iterative applications of a nonlinear operator T over policy space, we derived sufficient conditions for convergence to an aligned fixed point, expressed succinctly as $\rho(J_T(\pi^*)) < 1$. This spectral criterion provides both a theoretical guarantee and a practical heuristic: it quantifies the boundary between safe refinement and unstable drift.

Our empirical investigation into reflective prompting loops provides proof that these spectral measures can be practically estimated and that they align with noticeable stability and divergence patterns. Aside from its mathematical uniqueness, the framework offers a clear terminology for evaluating the safety of autonomous LLMs. In doing so, it aids in the development of a new set of theoretical instruments for comprehending how self-enhancing language models could progress, stabilize, or destabilize as time goes on.

References

- Xu, W.; Zhu, G.; Zhao, X.; Pan, L.; Li, L.; Wang, W. Y. 2024. “Pride and Prejudice: LLM Amplifies Self-Bias in Self-Refinement.” arXiv preprint arXiv:2402.11436.
- Zeng, Y.; Cui, X.; Jin, X.; Liu, G.; Sun, Z.; He, Q.; Li, D.; Yang, N.; Wang, J. 2025. “ARIES: Stimulating Self-Refinement of Large Language Models by Iterative Preference Optimization.” arXiv preprint arXiv:2502.05605.
- Chen, Z.; Deng, Y.; Yuan, H.; Ji, K.; Gu, Q. 2024. “Self-Play Fine-Tuning Converts Weak Language Models to Strong Language Models.” arXiv preprint arXiv:2401.01335.
- Piotrowski, T. 2021. “The Fixed Point Iteration of Positive Concave Mappings.” arXiv preprint arXiv:2110.11055.
- Kittisopaporn, A.; et al. 2021. “Convergence Analysis of a Gradient Iterative Algorithm with Spectral Radius.” AIMS Mathematics 6(8): 8477–8496.
- “Chapter 10: Iteration of Linear Systems.” Tau University Technical Report.

A Appendix: Theoretical Proofs and Extended Analysis

This appendix provides complete proofs for the main results in the paper and offers additional supporting lemmas and analytical remarks.

A.1 Preliminaries and Notation

Let $\mathcal{P}(\mathcal{X})$ denote a compact convex subset of \mathbb{R}^n representing the policy simplex, equipped with the ℓ_2 norm $\|\cdot\|_2$. For a continuously differentiable mapping $T : \mathcal{P}(\mathcal{X}) \rightarrow \mathcal{P}(\mathcal{X})$, we denote by $J_T(\pi)$ its Jacobian (Fréchet derivative) at π , and by $\rho(J_T(\pi))$ the spectral radius of that Jacobian.

A point π^* is said to be a *fixed point* if $T(\pi^*) = \pi^*$, and *aligned* if it also satisfies external alignment constraints or an alignment functional $\mathcal{A}(\pi^*) = 0$.

We recall two classical results:

Lemma 1 (Banach Fixed Point Theorem). *Let (X, d) be a complete metric space and $T : X \rightarrow X$ a contraction mapping, i.e., $\exists 0 \leq \beta < 1$ s.t. $d(Tx, Ty) \leq \beta d(x, y)$ for all $x, y \in X$. Then T admits a unique fixed point x^* and $\|T^t x_0 - x^*\| \leq \beta^t \|x_0 - x^*\|$.*

Lemma 2 (Spectral Radius Bound for Linear Operators). *For any matrix $A \in \mathbb{R}^{n \times n}$ and any matrix norm $\|\cdot\|$ consistent with the vector norm,*

$$\rho(A) = \lim_{k \rightarrow \infty} \|A^k\|^{1/k} \leq \|A\|.$$

In particular, if $\|A\| < 1$, then $\rho(A) < 1$ and $A^t \rightarrow 0$ as $t \rightarrow \infty$.

We use these results to analyze the nonlinear self-refinement operator T .

A.2 Proof of Theorem 1 (Spectral Stability of Self-Refinement)

Proof. Let $\pi_{t+1} = T(\pi_t)$ with π^* an aligned fixed point satisfying $T(\pi^*) = \pi^*$. By first-order Taylor expansion around π^* ,

$$\pi_{t+1} - \pi^* = J_T(\pi^*)(\pi_t - \pi^*) + R_t,$$

where $R_t = O(\|\pi_t - \pi^*\|^2)$ is the higher-order remainder.

Define $\delta_t = \pi_t - \pi^*$. Then the dynamics are

$$\delta_{t+1} = J_T(\pi^*)\delta_t + R_t.$$

If $\rho(J_T(\pi^*)) < 1$, then there exists a subordinate norm $\|\cdot\|_\beta$ such that $\|J_T(\pi^*)\|_\beta = \rho(J_T(\pi^*)) + \epsilon < 1$ for any small $\epsilon > 0$. Consequently,

$$\|\delta_{t+1}\|_\beta \leq \|J_T(\pi^*)\|_\beta \|\delta_t\|_\beta + c \|\delta_t\|_\beta^2,$$

for some $c > 0$. For sufficiently small $\|\delta_t\|_\beta$, the quadratic term is negligible, and we obtain the contraction inequality:

$$\|\delta_{t+1}\|_\beta \leq (\rho(J_T(\pi^*)) + \epsilon) \|\delta_t\|_\beta.$$

Iterating yields

$$\|\delta_t\|_\beta \leq (\rho(J_T(\pi^*)) + \epsilon)^t \|\delta_0\|_\beta.$$

As $\epsilon \rightarrow 0$, this implies exponential convergence:

$$\|\pi_t - \pi^*\|_2 \leq C \rho(J_T(\pi^*))^t \|\pi_0 - \pi^*\|_2,$$

for some equivalence constant C between norms. Hence, the refinement process converges exponentially to π^* . \square

A.3 Proof of Corollary 1 (Policy Drift Bound)

Proof. Let $\rho = \rho(J_T(\pi^*)) < 1$ and denote $\Delta_t = \|\pi_t - \pi^*\|_2$. From Theorem 1,

$$\Delta_t \leq \rho^t \Delta_0 + \sum_{i=0}^{t-1} \rho^i \|R_{t-i-1}\|_2.$$

Assuming the remainder satisfies $\|R_t\| \leq \epsilon \Delta_t$ with ϵ small, the geometric series yields

$$\Delta_t \leq \frac{\rho^t}{1-\rho} \Delta_1 + O(\epsilon),$$

which establishes the claimed bound on cumulative drift. \square

A.4 Lemma: Contraction under Scaled Update Rate

Consider the update $\pi_{t+1} = \pi_t + \alpha F(\pi_t)$ with Jacobian $J_F(\pi^*)$.

Lemma 3. *If $\rho(I + \alpha J_F(\pi^*)) < 1$, then the refinement process is stable. Moreover, for small α , the condition approximates*

$$\rho(J_F(\pi^*)) < \frac{2}{\alpha} - \frac{1}{\alpha^2} O(\|J_F\|^2),$$

showing the stabilizing role of the update rate α .

Proof. The Jacobian of T is $J_T = I + \alpha J_F$. The process is contractive if $\rho(J_T) < 1$. Expanding $\rho(J_T) \approx 1 + \alpha \lambda_{\max}(J_F)$, we require $\lambda_{\max}(J_F) < 0$, consistent with gradient stability conditions. For linear F , this reproduces the standard spectral stability criterion for discrete-time systems. \square

A.5 Continuous-Time Limit and Lyapunov Analysis

In the limit $\alpha \rightarrow 0$, define a continuous-time flow

$$\frac{d\pi(t)}{dt} = F(\pi(t)).$$

Linearizing around π^* gives $\dot{\delta}(t) = J_F(\pi^*)\delta(t)$. The equilibrium π^* is exponentially stable if and only if all eigenvalues of $J_F(\pi^*)$ satisfy $\text{Re}(\lambda_i) < 0$. This continuous-time criterion corresponds exactly to $\rho(J_T(\pi^*)) < 1$ in the discrete case via $\lambda_T = 1 + \alpha \lambda_F$.

A quadratic Lyapunov function $V(\delta) = \|\delta\|_2^2$ satisfies

$$\dot{V} = 2\delta^\top J_F(\pi^*)\delta \leq -2\gamma \|\delta\|_2^2$$

for $\gamma = -\max_i \text{Re}(\lambda_i(J_F)) > 0$, implying exponential convergence $\|\delta(t)\|_2 \leq e^{-\gamma t} \|\delta(0)\|_2$.

A.6 Bounded Noise and Stochastic Perturbations

When stochasticity is introduced, the update becomes

$$\pi_{t+1} = T(\pi_t) + \xi_t,$$

ξ_t is a zero-mean noise term with bounded variance $\mathbb{E}\|\xi_t\|_2^2 \leq \sigma^2$. Under $\rho(J_T(\pi^*)) < 1$, standard results in stochastic approximation yield:

$$\mathbb{E}\|\pi_t - \pi^*\|_2^2 \leq \frac{C\sigma^2}{1 - \rho(J_T(\pi^*))^2},$$

showing that noise only inflates the steady-state deviation by a factor inversely proportional to the contraction margin.

A.7 Remarks on Norm Choice and Operator Geometry

While ℓ_2 norms are convenient, the same theory holds for any induced matrix norm consistent with a convex metric on $\mathcal{P}(\mathcal{X})$. Future work may adopt non-Euclidean norms (e.g., KL-divergence-induced Bregman metrics) for analyzing LLM policy spaces, where the spectral radius generalizes to the largest modulus of generalized eigenvalues under the Fisher–Rao metric.