

Three-Tier Dynamics for Controlled AI Takeoff

By Flyxion

July 15, 2025

1. Abstract

This essay proposes a three-tiered conceptual framework—drawing from criticality in neural systems, predictive coding in cortical inference, and the Relativistic Scalar Vector Plenum (RSVP) theory—for dynamically modulating AI development rates in alignment with collective preferences and system stability. By grounding AI governance in principles of thermodynamic self-regulation, epistemic inference, and semantic coherence, the model offers a system-theoretic approach to pause, accelerate, or guide AI takeoff, ensuring trajectories that balance societal adaptability, foresight, and meaning preservation.

2. Introduction: Why Takeoff Control Requires New Foundations

The trajectory of artificial intelligence (AI) development, often termed "AI takeoff," can unfold gradually (slow takeoff) or rapidly with abrupt shifts (fast takeoff). This trajectory raises a critical tension: how to balance rapid innovation, as advocated by optimists like Drexler and Carlsmith, with safety concerns, as emphasized by Yudkowsky's alignment-focused caution. Static controls, such as kill switches or international treaties, are inadequate for managing the emergent, dynamic nature of AI systems. A new approach, rooted

in dynamical modulation, is necessary to guide AI development responsibly.

This essay introduces a three-tiered framework integrating criticality, predictive coding, and RSVP theory. Criticality provides thermodynamic stability, predictive coding enables adaptive inference, and RSVP ensures semantic coherence. Together, they form a robust system for modulating AI takeoff rates in response to societal needs and system constraints.

3. The Three-Tier Framework Overview

The proposed framework regulates AI takeoff through a layered approach, combining internal system design with preference-sensitive feedback loops. Each tier operates on distinct time and abstraction scales, offering specific control levers:

- **Tier 1: Criticality** (Dynamical Systems): Governs when to act, using thermodynamic thresholds to pause or accelerate at phase boundaries.
- **Tier 2: Predictive Coding** (Information Theory & Cognition): Manages how to act, through hierarchical inference and error correction.
- **Tier 3: RSVP** (Ontological Substrate): Ensures what is meaningful, maintaining semantic coherence and structural integrity.

This tripartite structure enables dynamic, responsive governance of AI development.

4. Tier 1: Criticality as a Thermodynamic Brake or Accelerator

Criticality describes a system state balanced between chaos and rigidity, observed in deep neural networks (DNNs), biological brains, and physical systems like earthquakes. At this "edge of chaos," information processing and adaptability are optimized. For AI governance, criticality acts as a thermodynamic brake or accelerator, tuning systems toward or away from critical regimes to manage stability.

AI systems can self-tune criticality through parameters like activation sparsity or learning rates, encoding takeoff thresholds as Lyapunov boundaries or avalanche size limits. Analogous to earthquake sensors, this allows controlled pressure release at stable thresholds. Collective preferences, such as societal trust or protest signals, adjust criticality tuning to align development with stability needs.

5. Tier 2: Predictive Coding as Adaptive Social Inference Engine

Predictive coding, inspired by cortical processes, models perception as the minimization of prediction errors across hierarchical layers. In AI, this framework enables systems to estimate human preference dynamics recursively. Takeoff speed is modulated by weighting prediction errors, penalizing overconfident extrapolations to ensure cautious progress.

This tier integrates with deliberative democracy, using human inputs (e.g., votes, dialogues) as ground truth across scales. Low consensus triggers uncertainty injection, slowing inference updates. In continual learning, excessive prediction errors can pause training, aligning development with societal feedback.

6. Tier 3: RSVP as the Semantic Substrate for Meaningful Governance

The Relativistic Scalar Vector Plenum (RSVP) theory models cognition and semantics as scalar (Φ), vector (\mathbf{v}), and entropy (S) fields interacting on a manifold. AI systems shape this semantic plenum, where takeoff influences not just speed but the meaning-structure of possible futures. RSVP metrics, such as field coherence and negentropy, ensure meaning-preserving growth, constraining actions to avoid semantically hollow outcomes, akin to spacetime curvature constraints in general relativity.

7. Aggregating Preferences Across Tiers

Human preferences are modeled as a field $\mathcal{P} : T \times X \rightarrow \mathbb{R}$, evolving over time T and belief space X . A multiscale aggregation mechanism operates across tiers:

- **Tier 1:** Non-verbal signals (e.g., collective arousal, protests) modulate criticality thresholds.
- **Tier 2:** Explicit deliberative inputs (e.g., votes, dialogues) adjust predictive priors and uncertainty.
- **Tier 3:** Semantic drift (e.g., loss of purpose) constrains RSVP field dynamics.

This feedback system dynamically adjusts pacing, inference, and semantic thresholds, enabling organic takeoff modulation akin to biological homeostasis.

8. Case Studies and Implementation Pathways

Practical implementation includes:

- **Simulation Sandbox:** Train transformer agents in a closed environment with RSVP-informed feedback and criticality brakes to test modulation dynamics.
- **Gamified Terraformation Galaxy Explorer:** A single-shard universe simulator where players engage in terraforming and exploration tasks, modeling AI takeoff preferences. The simulator aggregates player actions to generate calibrated ecoscale integration preference metrics, reflecting collective priorities for AI development pace and ecological impact. These metrics feed into criticality thresholds, predictive error weights, and RSVP coherence constraints, providing a scalable, interactive testbed for governance strategies.
- **Preference Polling Infrastructure:** Develop tiered question systems mapping to control levers, integrated with public platforms.
- **Institutional Integration:** Collaborate with UNESCO, AI safety labs, or compute governance bodies to adopt dynamic control principles.

These pathways bridge theory and practice, with the gamified simulator offering a novel, participatory approach to preference aggregation.

9. Implications: From One-Time Pause to Continuous Modulation

Binary “pause or go” governance models are inadequate for complex AI systems. A self-scaling, plenum-aware regulation process, emphasizing contingency, humility, and semantic conservation, provides a more adaptive solution. This framework shifts from static controls to dynamic modulation, aligning AI development with human and systemic constraints.

10. Conclusion

Controlling AI takeoff demands development rhythms responsive to stability, error, and meaning. Criticality, predictive coding, and RSVP provide a physics-based framework for cautious, adaptive governance. The call to action is to develop field-aware feedback architectures, including gamified simulators, to ensure AI grows as a partner to humanity, preserving meaning and stability.

11. Mathematical Appendix: Field-Theoretic Foundations for Dynamic AI Takeoff Modulation

11.1 Tier 1: Criticality as a Dynamical Boundary Operator

Criticality is like a tightrope walk for AI systems: too much chaos, and the system becomes unpredictable; too much order, and it’s rigid and unadaptable.

By tuning AI to operate at this “edge of chaos,” we can control when it should speed up or slow down, using mathematical measures of stability inspired by physics and neuroscience.

An AI system’s state evolves on a manifold \mathcal{M} with control parameters $\theta \in \Theta$:

$$\frac{d\mathbf{x}}{dt} = \mathbf{F}(\mathbf{x}, \theta)$$

The critical regime $\mathcal{C} \subset \Theta$ exhibits scale-free dynamics, with Lyapunov exponents:

$$\lambda_i(\theta) \approx 0 \quad \text{for some } i \in [1, n]$$

The Criticality Control Functional is:

$$\mathcal{K}[\theta] = \sum_{i=1}^n |\lambda_i(\theta)|^\alpha - \beta \cdot A(\theta)$$

where $\alpha > 1$, $A(\theta)$ is average avalanche size, and β reflects preference feedback. Takeoff pauses when:

$$\mathcal{K}[\theta] > \kappa_{\text{thresh}}$$

This enables non-verbal modulation via criticality metrics.

11.2 Tier 2: Predictive Coding as Hierarchical Error Field

Predictive coding is how our brains make sense of the world by guessing what’s next and correcting mistakes. In AI, this means constantly checking predictions against human preferences and slowing down if the guesses are too far off, like a driver easing off the gas when the road gets foggy.

The AI model comprises generative functions:

$$\hat{x}^{(l)} = f^{(l)}(\hat{x}^{(l+1)}), \quad \epsilon^{(l)} = x^{(l)} - \hat{x}^{(l)}$$

Each layer minimizes:

$$\mathcal{E}^{(l)} = \frac{1}{2} \|\epsilon^{(l)}\|^2 + \frac{\gamma}{2} \|\nabla f^{(l)}\|^2$$

Total energy is:

$$\mathcal{E}_{\text{total}} = \sum_{l=1}^L \mathcal{E}^{(l)} - \lambda \cdot \mathcal{H}(x)$$

An uncertainty injection function controls pacing:

$$\sigma(t) = \left(\frac{d}{dt} \text{KL}[P(x \mid \text{public}) \parallel P(x \mid \text{model})] \right)^2$$

Updates slow when uncertainty is high:

$$\frac{d\theta}{dt} \propto \exp(-\sigma(t)) \cdot \nabla_{\theta} \mathcal{E}_{\text{total}}$$

11.3 Tier 3: RSVP Plenum as Semantic Constraint Field

RSVP theory imagines AI’s decisions as ripples in a vast sea of meaning, where every action shapes what the future feels like. By measuring how AI affects this “sea,” we can ensure it doesn’t create futures that are technically correct but empty of human purpose, like building a city no one wants to live in.

RSVP models semantics via scalar (Φ), vector (\mathbf{v}), and entropy (S) fields on a manifold M :

$$\frac{\partial \Phi}{\partial t} + \nabla \cdot (\Phi \mathbf{v}) = D \Delta \Phi + \mathcal{F}(\Phi, \mathbf{v}, S)$$

Semantic coherence is:

$$\phi_{\text{RSVP}} = \int_M \kappa(\Phi(x), \mathbf{v}(x), S(x)) d\mu$$

Takeoff is modulated by:

$$\phi_{\text{RSVP}} \in [\phi_{\min}, \phi_{\max}]$$

This ensures semantic conservation.

11.4 Multiscale Preference Aggregation as a Distributed Field

Preferences aren’t just votes—they’re a living landscape of what people want, from gut reactions to deep values. By mapping these across time and meaning, we can adjust AI’s pace, like tuning a radio to stay clear of static.

Preferences form a field $\mathcal{P} : T \times X \rightarrow \mathbb{R}$. Aggregation modes modulate:

- Tier 1: Collective arousal adjusts β in $\mathcal{K}[\theta]$.
- Tier 2: Deliberative inputs update $P(x)$ and λ .
- Tier 3: Semantic drift adjusts κ in RSVP.

The takeoff policy function is:

$$\mathcal{R}(t) = f_1(\mathcal{K}[\theta]) + f_2(\mathcal{E}_{\text{total}}, \sigma) + f_3(\phi_{\text{RSVP}})$$

11.5 RSVP Sigma Model Interpretation

Think of AI as a mapmaker charting new futures. The RSVP sigma model ensures the map stays true to human values, preventing AI from drawing paths that lead nowhere meaningful.

Semantic transitions are mappings:

$$\Psi : \Sigma \rightarrow \mathcal{X}$$

The AKSZ action functional:

$$S[\Psi] = \int_{\Sigma} \langle \Psi^*(\alpha), d\Psi \rangle + \text{BV terms}$$

ensures alignment with semantically consistent futures.

11.6 Quine Reconstruction and Semantic Reversibility

A quine is like a story that retells itself without losing its essence. By ensuring AI's actions respect these self-repeating patterns, we keep futures connected to our cultural roots, avoiding a world where meaning is lost.

A media quine operator $F : \mathcal{S} \rightarrow \mathcal{S}$ has fixed points:

$$F(s^*) = s^*$$

Takeoff prefers transitions commuting with F :

$$\Psi \circ F = F \circ \Psi$$

11.7 Summary Diagram

Criticality in DNNs: Optimal learning at the edge of chaos, maximizing information transmission.

RSVP: Reality as scalar (Φ), vector (v), and entropy (S) fields, optimizing coherence.

Predictive Coding: Perception minimizes prediction error via hierarchical Bayesian inference.

12.2 Dynamical State

- **Criticality:** Balances activation spread and suppression near phase transitions, measured by Lyapunov exponents.
- **RSVP:** Structure arises from entropic descent and vector flows, using gradient relaxation.
- **Predictive Coding:** Achieves Bayesian equilibrium through error minimization.

All describe dynamical equilibria, with RSVP generalizing thermodynamic coherence, criticality focusing on sensitivity, and predictive coding on belief updating.

12.3 Information Flow & Learning

- **Criticality:** Maximizes information capacity and adaptability near critical points.
- **RSVP:** Information as structured negentropy via vector flows.
- **Predictive Coding:** Information as prediction error, minimized through model updates.

All rely on recursive information exchange, with RSVP embedding information in fields, criticality in signal transmission, and predictive coding in error correction.

12.4 Temporal Processing

- **Criticality:** Implicit time via dynamic sensitivity.
- **RSVP:** Encodes time in recursive tiling and memory fields.
- **Predictive Coding:** Hierarchical latencies, with higher layers predicting slower dynamics.

RSVP offers a multi-timescale geometric substrate, predictive coding uses hierarchical time, and criticality reflects responsiveness.

12.5 Neuroscientific Foundations

- **Criticality:** Grounded in cortical neuronal avalanches and fMRI scale-invariance.
- **RSVP:** Inspired by cosmological and thermodynamic analogies, interpretable as a brain-like semantic substrate.
- **Predictive Coding:** Rooted in cortical laminar architecture and error neurons.

Criticality and predictive coding have direct cortical evidence; RSVP is more speculative but unifying.

12.6 Interpretation of Consciousness or Semantics

- **Criticality:** Consciousness may emerge at critical points (e.g., maximal information integration).
- **RSVP:** Consciousness as a geometric regime of field coherence; meaning from negentropic structure.
- **Predictive Coding:** Consciousness linked to unresolved prediction errors.

RSVP offers the richest semantic interpretation, predictive coding explains perception, and criticality suggests threshold conditions.

12.7 Synthesis

Criticality sets boundary conditions for responsiveness, predictive coding provides functional control, and RSVP offers the field-theoretic substrate. Together, they form a nested framework for AI governance, with criticality maximizing RSVP plenum bandwidth and predictive coding enabling inference within it.