

Divided We Stand: RSVP and the Limits of Coherence in AI Safety

Flyxion

September 18, 2025

Abstract

This article situates the Relativistic Scalar Vector Plenum (RSVP) framework within current debates on artificial intelligence safety. It contrasts Eliezer Yudkowsky and Nate Soares’s book *If Anyone Builds It, Everyone Dies* [Yudkowsky and Soares, 2025] with critiques by Scott Alexander [Alexander, 2025], Clara Collier [Collier, 2025], and Robin Hanson [Hanson, 2024]. The analysis highlights RSVP’s distinctive contribution: resilience through entropy-respecting diversity. Where the Machine Intelligence Research Institute (MIRI) emphasizes singular clarity and diamond-hard priors, RSVP insists that survival depends on recursive adaptability, re-bubbling, and dispersion across scalar, vector, and entropic fields (Φ, v, S) . Building on foundational thermodynamic and quantum theories [Jacobson, 1995, Verlinde, 2011, Barandes, 2022], RSVP offers a unified architecture for physics, cognition, and governance, providing a counterpoint to MIRI’s over-coherent worldview.

Contents

1	Introduction	2
2	The Core Framework	3
3	Beyond the Core Framework	3
3.1	Topological Domain	3
3.2	Quantum Domain	4
3.3	Ethical Domain	4
3.4	Socio-Political Domain	4
3.5	Integrative Closure	4
4	Comparative Frameworks	5
4.1	Incrementalism vs. Radical Clarity	5
4.2	Parables and Archetypes	5
4.3	The Problem of Non-Update (Collier’s Critique)	5
4.4	Implications for Governance	5

4.5	Cultural Drift and the Perils of Monoculture (Hanson’s Critique)	6
4.6	Material Privilege and Strategic Posturing	6
4.7	Alternative Entropy-Injection Strategies	6
5	Normative Implications	7
6	Conclusion	7
	Appendix A: Variational Structure	9
	Appendix B: Entropic Stability Criterion	9
	Appendix C: Topological Gluing and MERGE	9
	Appendix D: Entropy-Injection Operators	10

1 Introduction

The publication of Yudkowsky and Soares’s *If Anyone Builds It, Everyone Dies* [Yudkowsky and Soares, 2025] has intensified debate about the fate of artificial intelligence and humanity. The Machine Intelligence Research Institute (MIRI) advances a position of near-total certainty: superintelligence entails extinction, and only a global ban on AI development can avert it. Critics have responded from multiple directions. Scott Alexander [Alexander, 2025] questions the practicality of this radical clarity, describing MIRI’s stance as rhetorically powerful but politically inert. Clara Collier [Collier, 2025] presses further, arguing that the book represents a regression: the intelligence explosion hypothesis is barely defended, and the MIRI worldview has failed to update since the early blogosphere era. Robin Hanson [Hanson, 2024] adds a cultural critique, warning that over-coherence produces fragility by eroding diversity. Taken together, these responses illuminate the political, epistemic, and cultural costs of excessive coherence.

Against this backdrop, the Relativistic Scalar Vector Plenum (RSVP) framework proposes a different strategy. RSVP models reality as a tripartite field (Φ, v, S) of scalar capacity, vector flow, and entropy. It extends beyond physics into cognition, semantics, and governance, formalizing the principle that coherence without entropy is collapse. RSVP builds on prior theoretical advances: Jacobson’s derivation of Einstein’s equations from thermodynamics [Jacobson, 1995], Verlinde’s entropic account of gravity [Verlinde, 2011], and Barandes’s unistochastic reformulation of quantum theory [Barandes, 2022]. Each of these contributions reinterprets fundamental laws as emergent from deeper statistical or probabilistic substrates. RSVP synthesizes these insights into a universal architecture, one that encodes adaptability through entropy-respecting mechanisms.

This article places MIRI’s arguments and their critiques in dialogue with RSVP. Where Yudkowsky and Soares call for singular clarity and diamond-hard priors, RSVP insists that survival depends on recursive adaptability, re-bubbling, and dispersion. The central thesis is that divided we stand: diversity and entropy are not liabilities but the conditions of resilience across physics, cognition, and society.

2 The Core Framework

RSVP models reality in terms of scalar semantic density Φ , vector flow v , and entropy S . Its action functional is defined as

$$F[\Phi, v, S] = \int_{\Omega} \left(\frac{\kappa_{\Phi}}{2} |\nabla \Phi|^2 + \frac{\kappa_v}{2} |\nabla \times v|^2 + \frac{\kappa_S}{2} |\nabla S|^2 - \lambda \Phi S \right) dx.$$

These fields jointly describe cosmological entropic redshift, neural and cognitive dynamics, and semantic infrastructures. In cosmology, Φ represents matter density, v velocity flows, and S entropic expansion [Jacobson, 1995]. In cognition, Φ encodes semantic coherence, v attentional flows, and S unresolved ambiguity [Friston, 2010]. In semantic and social infrastructures, Φ corresponds to stored knowledge or institutional memory, v to the circulation of information and influence, and S to the dissipation of coordination.

The functional F encodes a balance between gradients, curls, and divergences, ensuring that scalar, vector, and entropic contributions remain coupled rather than collapsing into a single dominant term. The positive quadratic terms penalize sharp fluctuations in each field, while the interaction term $-\lambda \Phi S$ enforces reciprocity between scalar concentration and entropy production. Locally, this produces a smoothing dynamic analogous to diffusion, but globally it permits the emergence of coherent structures that persist under perturbation.

In this way, RSVP provides a minimal but generalizable field structure. Its strength lies in capturing systems where order and disorder must coexist—gravity constrained by entropy, cognition stabilized by ambiguity, and governance structured by informational flux. The universality of RSVP depends on extending this local partial differential equation structure into broader domains: topological, quantum, ethical, and socio-political.

3 Beyond the Core Framework

RSVP generalizes through four principal extensions, each addressing a distinct domain while preserving the tripartite structure.

3.1 Topological Domain

Local patches of cosmological cells, neural states, or semantic contexts require sheaf-theoretic gluing. RSVP encodes this in the MERGE operator:

$$\text{Emerge} = \int_{\Omega} \left(\lambda_{\Phi} \|\Phi_1 - \Phi_2\|^2 + \lambda_v \|v_1 - v_2\|^2 + \lambda_S \text{KL}(S_1 \| S_2) \right) dx.$$

This MERGE operator functions as an entropic homotopy colimit, minimizing inconsistency across overlapping covers [Lurie, 2009]. Conceptually, it ensures that local field behaviors—whether in physics, cognition, or semantics—can be coherently extended to global structures without imposing artificial uniformity.

3.2 Quantum Domain

Transitions between RSVP states are modeled as unistochastic maps, embedding unitary quantum evolution as a zero-entropy subset within the RSVP action functional [Barandes, 2022]. This perspective treats ordinary quantum mechanics not as a fundamental law but as a constrained case where entropy fluxes vanish. RSVP thus provides a broader framework in which probabilistic transitions, coherence loss, and entropy production are natural rather than anomalous.

3.3 Ethical Domain

RSVP defines an entropy-production functional along worldlines:

$$\Sigma[\Gamma] = \int_{\Gamma} \max(0, \dot{S}) dt,$$

formalizing irreversible harm as positive entropy growth along trajectories, aligning with thermodynamic ethics [Verlinde, 2011]. This formulation distinguishes between reversible fluctuations, which may be ethically neutral, and irreversible increases in disorder, which encode lasting cost. Ethical evaluation therefore becomes inseparable from entropy management, grounding moral reasoning in the same substrate that governs physics.

3.4 Socio-Political Domain

Governance is reframed as an entropic budgeting problem. The Landauer inequality applies directly:

$$\Delta S(\Omega) \geq (\ln 2) N_{\text{erased}},$$

treating institutional erasure (e.g., data deletion, strategy suppression) as an entropic cost [Landauer, 1961]. Decision-making and policy formation consume informational resources and produce entropy, making stability contingent on careful allocation of finite entropic budgets. Societies that attempt to enforce coherence beyond their entropic means risk collapse, whereas those that budget dispersion effectively can sustain adaptive resilience.

3.5 Integrative Closure

All four domains—topological, quantum, ethical, and socio-political—are projections of (Φ, v, S) , demonstrating RSVP’s methodological parsimony: a single substrate suffices to generate multiple disciplinary projections. Where MIRI envisions safety through static clarity and diamond-hard priors, RSVP emphasizes adaptive coherence achieved through entropy-aware balancing. Survival is thus reinterpreted not as the elimination of uncertainty, but as the cultivation of systems robust enough to integrate and redirect it.

4 Comparative Frameworks

4.1 Incrementalism vs. Radical Clarity

Scott Alexander [Alexander, 2025] situates *If Anyone Builds It, Everyone Dies* within the divide between moderate incrementalists and MIRI’s radical clarity. Incrementalists advocate for transparency, institutional accountability, and incremental safety research, acknowledging a non-trivial probability of catastrophic failure. Yudkowsky and Soares, by contrast, assert that superintelligence entails extinction with near-certainty, necessitating an immediate global ban on AI development [Yudkowsky and Soares, 2025]. Alexander likens MIRI’s strategy to “wearing a hard hat against an asteroid impact,” rhetorically compelling but practically inert. From the RSVP perspective, Alexander’s critique highlights the limitations of singular interventions. A hard ban, like a hard hat, is a single-attractor strategy that fails to integrate entropic feedbacks. RSVP prioritizes recursive adaptability: multiple partial interventions that inject entropy, diversify responses, and avoid brittle over-commitments.

4.2 Parables and Archetypes

The narrative force of *If Anyone Builds It* lies in its parables: gods playing evolutionary games or misaligned chatbots pursuing perverse goals [Yudkowsky and Soares, 2025]. These echo the tragic machine Mima in Harry Martinson’s *Aniara* [Martinson, 1956], which collapses under semantic overload, and the Tower of Babel [The Hebrew Bible], where over-coherence of language is disrupted by dispersion. RSVP interprets these not as mere stories but as field dynamics: runaway attractors (over-concentrated Φ , over-aligned v , suppressed S) collapse without entropy-injection mechanisms.

4.3 The Problem of Non-Update (Collier’s Critique)

Clara Collier [Collier, 2025] identifies epistemic rigidity as the book’s deepest flaw. The intelligence explosion hypothesis, the load-bearing premise of MIRI’s argument, receives only two sentences, despite radical shifts in AI since 2008—from hand-coded architectures to empirically tractable deep learning systems governed by scaling laws. Yet, Yudkowsky and Soares repeat early blogosphere-era arguments almost verbatim, failing to engage with new evidence [Collier, 2025]. Collier calls this a “regression,” noting that MIRI shadowboxes outdated opponents rather than contemporary debates. From the RSVP vantage, this reflects epistemic over-coherence. RSVP encodes update mechanisms: the MERGE operator integrates local evidence into global coherence, and the entropy-production functional $\Sigma[\Gamma]$ formalizes responsiveness to novel conditions [Verlinde, 2011]. Where MIRI’s priors have the “strength of diamond,” RSVP treats adaptability as a principle of resilience.

4.4 Implications for Governance

Both Alexander and Collier converge on governance. Alexander doubts the feasibility of global GPU bans [Alexander, 2025]; Collier questions the empirical assumptions motivating them [Collier, 2025]. RSVP reframes governance as an entropic budget allocation problem:

stability arises from distributing flows across multiple channels, preserving resilience through diversity [Landauer, 1961]. Just as Babel scattered language to prevent collapse [The Hebrew Bible], RSVP prescribes dispersion of semantic, technological, and institutional pathways to avert brittle failure.

4.5 Cultural Drift and the Perils of Monoculture (Hanson’s Critique)

Robin Hanson [Hanson, 2024] warns that cultural drift erodes resilience, leaving societies brittle through monoculture. In evolutionary dynamics, variation across lineages preserves adaptability; in markets, heterogeneous strategies sustain liquidity; in cognition, divergent perspectives prevent epistemic lock-in [Hanson, 2024]. Excessive alignment collapses adaptive capacity, producing fragile systems. RSVP formalizes this as entropic resilience: perturbations in S preserve adaptive capacity against runaway attractors. The Tower of Babel [The Hebrew Bible] serves as a mythological precedent: a unified semantic field Φ enabled a collective vector v toward heaven, but divine intervention scattered Φ into multiple languages, injecting entropy S and restoring resilience. Modern entropy hacks—student-designed ciphers, Arabic-script reform, retro hardware revival—extend this principle, fragmenting corpora and introducing vector drag to prevent over-coherence [Hanson, 2024].

4.6 Material Privilege and Strategic Posturing

A central difficulty with *If Anyone Builds It, Everyone Dies* is not merely the content of its arguments but the position from which they are made. Both Yudkowsky and Soares operate within institutions that guarantee stable income, prestige, and the freedom to publish radical claims without fear of material reprisal. This insulation allows them to frame “radical clarity” as a form of intellectual heroism while remaining materially secure. For those outside this privileged enclave, the risks of technological disruption and policy misfire are borne directly, not abstractly. From the RSVP perspective, this asymmetry exemplifies a failure of entropic budgeting: the burdens of coherence are exported to others, while MIRI maintains internal stability. A more honest engagement would acknowledge the uneven distribution of epistemic and material costs, and the necessity of distributing adaptive capacity across institutions rather than concentrating moral authority in one.

4.7 Alternative Entropy-Injection Strategies

If MIRI prescribes only bans, RSVP proposes dispersion. One strategy is linguistic diversification: replacing English-language centralization with staged adoption of alternative scripts (e.g., Arabic) that resist monocultural capture. Pilot experiments with transliterating English into Arabic script demonstrate the feasibility of increased phonetic transparency while reducing reliance on Latin-script corpora. Another is educational individuation: issuing each student a personalized cipher for their textbooks, cultivated across years of study. What begins as cryptographic opacity becomes, with practice, simply another typeface. By fragmenting linguistic inputs, such policies prevent runaway coherence in training data and foster resilience through pluralism. Similarly, mandating the revival of “obsolete”

hardware—cassettes, CRTs, mechanical typewriters—would slow unidirectional acceleration while uncovering unexploited affordances of earlier technologies. Each of these proposals is intentionally impractical from the standpoint of present technocracy, yet they highlight the principle: survival may require absurd-seeming entropy injections rather than rigid prohibitions. The lesson is not that governments should literally adopt these reforms, but that governance must learn to value pluralization, slowdown, and dispersion as strategies in their own right.

5 Normative Implications

RSVP extends beyond critique by grounding a positive program for survival. First, it re-frames ethics thermodynamically: irreversible harm is modeled as positive entropy growth, making moral reasoning inseparable from physical costs. Second, it redefines governance as an entropic budgeting problem: institutions must distribute informational and material flows across multiple channels, rather than concentrating coherence until systems become brittle. Third, it promotes a strategy of re-bubbling: creating new attractors of meaning, language, and cooperation to counteract collapse.

Here, the earlier critiques converge. The privilege of MIRI’s radical clarity illustrates how concentrated epistemic authority externalizes risk, demanding counter-balances that disperse both power and uncertainty. Entropy-injection proposals—whether through linguistic pluralization, individualized ciphers, or the revival of neglected hardware—serve as metaphors and potential practices for preserving heterogeneity against capture. Their purpose is not to prescribe literal policy but to demonstrate how resilience requires pluralization, slowdown, and dispersion.

Against MIRI’s one-shot framing [Yudkowsky and Soares, 2025], RSVP emphasizes multi-shot adaptation. Its scalar, vector, and entropic fields (Φ, v, S) encode recursive updating, recursive pluralization, and recursive resilience. Normatively, this means survival does not depend on eliminating uncertainty but on continually reconfiguring around it. Diversity is not noise to be suppressed but signal to be preserved: the substrate of adaptability across physics, cognition, and governance.

6 Conclusion

Divided we stand. The debate around AI safety, as crystallized in Yudkowsky and Soares’s *If Anyone Builds It, Everyone Dies* [Yudkowsky and Soares, 2025], reveals the costs of excessive coherence. Scott Alexander’s critique shows the political impracticality of MIRI’s radical clarity: hard bans and absolute prohibitions cannot attract broad coalitions [Alexander, 2025]. Clara Collier’s review identifies the epistemic cost: the refusal to update assumptions in light of new evidence renders the argument brittle and regressive [Collier, 2025]. Robin Hanson’s warning about cultural drift highlights the social cost: over-coherence erodes resilience, leaving civilizations vulnerable to unforeseen shocks [Hanson, 2024].

Taken together, these critiques form a triangulation: political impracticality, epistemic rigidity, and cultural fragility. Each identifies a dimension of collapse that arises when entropy is excluded. RSVP responds by formalizing entropy as the substrate of survival. Its

MERGE operator encodes continual updating [Lurie, 2009]; its entropy-production functional $\Sigma[\Gamma]$ enforces responsiveness [Verlinde, 2011]; and its socio-political framing treats governance as entropic budgeting rather than prohibition [Landauer, 1961].

Beyond these general formalisms, concrete strategies illustrate the value of entropy injection. Script reform, whether through staged adoption of Arabic script or the creation of personalized ciphers for educational materials, fragments linguistic monoculture and builds redundancy into semantic channels. Retro-technological revivals—cassettes, CRT displays, mechanical typewriters—add temporal diversity, slowing technological lock-in and reopening paths of exploration that premature obsolescence foreclosed. These proposals are deliberately heterodox, even satirical, yet they demonstrate the principle that resilience lies not in eliminating uncertainty but in cultivating dispersion. Their formalization appears in Appendix D, where they are cast as entropy-injection operators acting on (Φ, v, S) .

Where MIRI insists that humanity’s only hope is a single shot at perfect alignment [Yudkowsky and Soares, 2025], RSVP offers a different vision: survival depends on diversity, dispersion, and recursive adaptability. Against monoculture, RSVP prescribes re-bubbling; against rigidity, responsiveness; against brittle singular clarity, entropic resilience. Physics, cognition, and society converge on the same principle: coherence without entropy is collapse.

Appendix A: Variational Structure

The RSVP action functional is

$$F[\Phi, v, S] = \int_{\Omega} \left(\frac{\kappa_{\Phi}}{2} |\nabla \Phi|^2 + \frac{\kappa_v}{2} |\nabla \times v|^2 + \frac{\kappa_S}{2} |\nabla S|^2 - \lambda \Phi S \right) dx.$$

Stationarity under variations yields coupled Euler–Lagrange equations:

$$\begin{aligned} -\kappa_{\Phi} \Delta \Phi - \lambda S &= 0, \\ -\kappa_v \nabla \times (\nabla \times v) &= 0, \\ -\kappa_S \Delta S - \lambda \Phi &= 0. \end{aligned}$$

These equations formalize the reciprocal coupling between scalar concentration Φ and entropy S , while vector circulation v is constrained by torsion-free flow.

Appendix B: Entropic Stability Criterion

For a worldline Γ with entropy production functional

$$\Sigma[\Gamma] = \int_{\Gamma} \max(0, \dot{S}) dt,$$

stability requires

$$\frac{d}{dt} (\|\Phi\|^2 + \|v\|^2) \leq \alpha \dot{S},$$

for some $\alpha > 0$. This inequality ensures that growth in scalar or vector intensity is compensated by corresponding entropy dissipation. Violations correspond to runaway attractors (over-coherent priors, cultural lock-in).

Appendix C: Topological Gluing and MERGE

Given a cover $\{U_i\}$ of Ω with local fields (Φ_i, v_i, S_i) , RSVP constructs a global plenum via the MERGE functional:

$$\text{Emerge} = \int_{\Omega} \left(\lambda_{\Phi} \|\Phi_i - \Phi_j\|^2 + \lambda_v \|v_i - v_j\|^2 + \lambda_S \text{KL}(S_i \| S_j) \right) dx.$$

The condition

$$\text{Emerge} \rightarrow 0 \quad \text{as overlaps refine}$$

is equivalent to existence of a homotopy colimit

$$\text{hocolim}_{i \in I} \mathcal{F}_i,$$

where \mathcal{F}_i are local RSVP sheaves. Failure of this limit corresponds to “entropy tears,” structural instabilities where no coherent global extension exists.

Appendix D: Entropy-Injection Operators

RSVP interprets cultural and technological interventions as controlled injections of entropy, designed to prevent over-coherence and monoculture collapse. Let \mathcal{E} denote an entropy-injection operator acting on the tripartite fields (Φ, v, S) . Three classes of \mathcal{E} are illustrative:

Script Diversification Operator

The adoption of alternative scripts (e.g., staged transition to Arabic orthography) functions as

$$\mathcal{E}_{\text{script}}(\Phi) = \Phi \circ \pi,$$

where π is a nontrivial permutation of graphemic encodings. The effect is to redistribute scalar semantic density Φ into new representational bases, increasing entropy S by reducing uniformity across corpora.

Cipher Personalization Operator

Individualized ciphers for students, each defining a bijective mapping $\chi : \Sigma \rightarrow \Sigma$, are modeled as

$$\mathcal{E}_{\text{cipher}}(\Phi_i) = \Phi_i \circ \chi_i, \quad \chi_i \neq \chi_j \text{ for } i \neq j.$$

This introduces vector drag by fragmenting alignment v , ensuring that semantic convergence requires translation across heterogeneous bases. Global Φ remains coherent only through MERGE, preserving adaptability.

Retro-Technology Operator

The mandated revival of legacy media (cassettes, CRT tubes, typewriters) corresponds to

$$\mathcal{E}_{\text{retro}}(v) = v + \eta(t),$$

where $\eta(t)$ is a delay kernel modeling friction introduced by slower, material-bound flows. This operator increases temporal dispersion of v , injecting entropy S by decelerating feedback loops otherwise prone to runaway acceleration.

Interpretation

Together, $\mathcal{E}_{\text{script}}, \mathcal{E}_{\text{cipher}}, \mathcal{E}_{\text{retro}}$ illustrate RSVP's claim that resilience arises not from prohibitions but from entropy-aware interventions. They deliberately fragment Φ , retard v , and diversify S , ensuring that systemic coherence never hardens into brittle monoculture. Such operators are not literal policy prescriptions but mathematical demonstrations of how entropy-injection sustains long-term adaptability.

References

- Scott Alexander. Book review: If Anyone Builds It, Everyone Dies. *Astral Codex Ten*, September 2025. URL <https://astralcodexten.substack.com/p/book-review-if-anyone-builds-it-everyone>. Review of Yudkowsky & Soares, If Anyone Builds It, Everyone Dies.
- Jacob A. Barandes. Unistochastic quantum theory. *Foundations of Physics*, 52(6):94, 2022. doi: 10.1007/s10701-022-00623-4.
- Clara Collier. More was possible: A review of If Anyone Builds It, Everyone Dies. *Asterisk Magazine*, September 2025. URL <https://asteriskmag.com/issues/11/iabied>. Review of Yudkowsky & Soares, If Anyone Builds It, Everyone Dies.
- Karl Friston. The free-energy principle: A unified brain theory? *Nature Reviews Neuroscience*, 11:127–138, 2010. doi: 10.1038/nrn2787.
- Robin Hanson. Beware cultural drift. YouTube, Science, Technology, and the Future, May 2024. URL <https://www.youtube.com/watch?v=Z7C6VC3UxkM>. Talk on the risks of cultural monoculture and societal drift.
- Ted Jacobson. Thermodynamics of spacetime: The einstein equation of state. *Physical Review Letters*, 75:1260–1263, 1995. doi: 10.1103/PhysRevLett.75.1260.
- Rolf Landauer. Irreversibility and heat generation in the computing process. *IBM Journal of Research and Development*, 5(3):183–191, 1961. doi: 10.1147/rd.53.0183.
- Jacob Lurie. *Higher Topos Theory*. Princeton University Press, Princeton, NJ, 2009. doi: 10.1515/9781400830558.
- Harry Martinson. *Aniara: An Epic Science Fiction Poem*. Bonniers, 1956. English translation by Hugh MacDiarmid and Elspeth Harley Schubert, 1963.
- The Hebrew Bible. Genesis 11:1–9, the tower of babel. In *The Old Testament*. Various. Standard reference to the Babel story.
- Erik Verlinde. On the origin of gravity and the laws of newton. *Journal of High Energy Physics*, 2011(4):29, 2011. doi: 10.1007/JHEP04(2011)029.
- Eliezer Yudkowsky and Nate Soares. *If Anyone Builds It, Everyone Dies*. Machine Intelligence Research Institute, Berkeley, CA, 2025. Published September 2025.