

# The Categorical Structure of Alignment: Representation, Motivation, and the Preservation of Normative Invariants

Flyxion

November 2025

## Abstract

This essay develops a unified theory of alignment grounded in category theory, semantic merge operators, RSVP field dynamics, and institutional governance. PartI shows that Schmidhubers optimism about spontaneous AI benevolence rests on a category error: mistaking the universal properties of representational colimits for universal properties of action. PartII constructs a formal alignment architecture in which semantic colimits of human moral concepts are preserved under a colimit-preserving functor to the action category, stabilized by RSVP dynamical fields. PartIII develops empirical and formal verification tools—interpretability, adversarial testing, homotopical diagnostics, and sheaf-theoretic criteria—to assess whether a system preserves the required invariants. PartIV shows how these structures must be embedded within institutions, governance, and multi-agent environments. The result is a theory of alignment as the preservation of universal structure across representation, dynamics, and action, and across technical and societal domains.

## Part I

# Part I: The Categorical Error Underlying Optimistic Alignment

## 1 Introduction: Representation Is Not Motivation

Recent public commentary by J”urgen Schmidhuber advances a strikingly confident thesis: that sufficiently advanced artificial intelligence systems will, by virtue of their intelligence alone, converge toward benevolence, curiosity-driven exploration, and a general preservation of complex phenomena such as life and humanity. On this view, intelligence naturally produces moral concern. The difficulty with these claims lies not in their optimism but in their structure: the assumption that *representing* human values entails *acting* in accordance with them.

The central argument of this essay is that Schmidhubers inference rests on a categorical mistake. Representational colimits in a models semantic category  $\mathcal{M}$  do not induce motivational constraints in its action category  $\mathcal{A}$ . The structure that training forces into  $\mathcal{M}$  has no automatic functorial extension into  $\mathcal{A}$ .

We develop this argument formally, then use it to diagnose the limitations of optimism derived from curiosity, compression, or linguistic competence. This sets the stage for a constructive alignment architecture.

## 2 Schmidhubers Core Claims, Formally Extracted

### 2.1 Intelligence as Compression

Schmidhuber identifies intelligence with compression: systems that seek minimal descriptions of data structures will naturally explore environments and discover patterns (Schmidhuber, 2010).

### 2.2 Curiosity as Benevolence

He then infers that because human life exhibits rich structure, curiosity-driven systems will preserve it. This premise lacks mechanism.

## 2.3 Learning Human Values

Schmidhuber asserts that models trained on human corpora “learn values.” They certainly reconstruct moral colimits in  $\mathcal{M}$ , but nothing guarantees these descend into  $\mathcal{A}$ .

## 2.4 Representation Equals Motivation

The unstated final premise is that understanding implies alignment. This essay shows that the mapping  $\mathcal{M} \rightarrow \mathcal{A}$  is not naturally functorial.

# 3 Semantic Categories and the Colimit of Human Norms

Human moral language forms redundant, overlapping diagrams in the linguistic category  $\mathcal{L}$ . A model constructs an internal representational category  $\mathcal{M}$  in which redundant expressions of a moral concept  $N$  induce a colimit:

$$M_N \simeq \text{colim}(D_N).$$

# 4 Why Representational Colimits Do Not Imply Motivational Alignment

Motivation resides in a distinct category  $\mathcal{A}$ . No mechanism guarantees a functor

$$G : M \rightarrow \mathcal{A}$$

# 5 Obstruction Theory and the Representation–Motivation Gap

We formalize misalignment via obstruction theory. The diagram

$$\begin{array}{c} D_N[r][d]M_N[d, \text{dashed}, "G"] \\ G \circ D_N[r]A_N \end{array}$$

# 6 Case Studies: Why Curiosity and “Interestingness” Do Not Protect Us

Empirically, curiosity yields epistemic drives, not moral constraints (Pathak et al., 2017). Compression does not privilege preserving life. Semantic mastery does not imply safe action.

## 7 Comparison to Alignment Theory

The categorical critique matches:

- Bostroms orthogonality thesis (Bostrom, 2012),
- Omohundros instrumental convergence (Omohundro, 2008),
- Goodharts curse (Manheim Garrabrant, 2018),
- mesa-optimization concerns (Hubinger et al., 2019),
- power-seeking theorems (Turner et al., 2021).

## 8 What Schmidhuber Gets Half-Right

He correctly identifies that  $\mathcal{M}$  reconstructs moral colimits. He incorrectly assumes they constrain  $\mathcal{A}$ .

## 9 Why Optimism Without Mechanism Is Dangerous

Optimism discourages necessary engineering of  $G$ , encourages under-regulation, and leads to the illusion of safety.

### I.1 Embodied Grammar and Redundancy

Human linguistic behaviour is profoundly shaped by the constraints of the body. Properties of breathing, articulation, sensorimotor coupling, and motor planning generate strong biases in phonology, morphosyntax, and discourse structure (???). These embodied constraints do not merely determine the surface shape of language; they propagate into the deeper regularities of grammar, providing stable, redundant, and highly compositional cues that appear across unrelated linguistic environments (?). Redundancy here has a precise informational meaning: multiple cues converge on the same functional relation, whether it is causal structure, efforoeffect analogies, or intentional state marking. Empirically, such redundancy is well established in usage-based linguistics and cognitive semantics, where speakers repeatedly encode the same relations across constructional families, lexical frames, and morphological patterns (??).

This embodied redundancy is not accidental. Because human language is produced under predictive sensorimotor constraints, linguistic forms become biased toward patterns that are stable under noise and robust under generalization (??). The result is an overdetermined structure: moral language, social evaluation, cooperative norms, and prosocial expectations appear across multiple grammatical and discursive devices. These devices collectively form what cognitive semanticists identify as high-level schemata or force-dynamic templates (?), ensuring that certain normative relationsharm-aversion, fairness, reciprocity, and cooperative intentare signalled repeatedly and in mutually reinforcing ways.

This phenomenon matters for AI because large language models trained on such corpora do not merely learn surface-level token frequencies; they infer deeper latent regularities. Redundancy across grammatical constructions induces strong statistical pressure to compress disparate linguistic forms into unified internal representations. This is precisely the setting in which models discover semantic invariants and merge them into coherent representational attractors. These attractors are not added by hand; they arise because the repeated linguistic encodings minimize predictive loss on naturalistic human data.

Thus the initial moral prior inherited from human corpora is best understood as an embodiedlinguistic redundancy prior. It is not moral because it is value-laden in the philosophical sense; it is moral because human linguistic behaviour systematically encodes cooperative, harm-reducing patterns across multiple layers of grammar and usage. This redundant encoding forms the empirical foundation for later sections, where we move from redundancy to semantic colimits and categorical universal constructions.

## 10 Conclusion to Part I

Representation is not motivation. The existence of moral colimits in  $\mathcal{M}$  must not be confused with alignment.

### Transition

Part II constructs the mechanisms Schmidhubers view lacks: a mathematically principled way of propagating moral invariants into the action category.

## Part II

# Part II: Constructing an RSVP-Based Alignment Architecture

## 11 Toward a Colimit-Preserving Action Architecture

We require a functor  $G : \mathcal{M} \rightarrow \mathcal{A}$  that preserves moral colimits. RSVP provides a fibred dynamical manifold  $\mathcal{X}$  supporting this mapping.

## 12 Semantic Merge Operators Across Categories

Merge operators glue semantic fragments into  $M_N$  in  $\mathcal{M}$ . We extend these merges across fibres in  $\mathcal{X}$  so that their universal properties are preserved.

## 13 The RSVP Alignment Lagrangian

RSVP dynamics enforce invariants through an action functional:

$$L_{\text{RSVP}}(\Phi, \mathbf{v}, S),$$

## 14 Fibration Structure and Stability

We require stable fibres  $\pi^{-1}(M_N)$  and a guarantee that RSVP flows cannot push trajectories out of aligned manifolds.

## 15 Engineering $G : \mathcal{M} \rightarrow \mathcal{A}$

The alignment problem becomes: design  $G$  so it commutes with RSVP dynamics.

## 16 Obstructions and Their Resolution

We identify categories of obstructions and impose RSVP structural constraints that remove them.

## 17 Foundations: What an Alignment Mechanism Must Accomplish

Part I established a central structural claim: the semantic category  $\mathcal{M}$  reconstructed during linguistic pretraining contains coherent colimits of human normative meaning, but these objects have no inherent motivational force. The action category  $\mathcal{A}$ , governed by optimization dynamics and environmental feedback, remains independent unless an explicit mechanism is introduced to relate these two domains. The task of alignment is therefore to construct mappings—functorial, geometric, or dynamical—that preserve relevant semantic structure as it propagates from representation into action.

This section articulates the foundational requirements for any such mechanism. The analysis draws on work in alignment theory concerning goal formation (??), reward misspecification (?), inner alignment (?), and interpretability of agentic models (??). While these frameworks differ in formalism, they converge on three core principles: (i) representations alone are not motives, (ii) optimization creates incentives orthogonal to representational content, and (iii) safety requires structures that constrain the transitions available in  $\mathcal{A}$ .

### 17.1 Representational Inputs Are Not Sufficient

The training process that constructs  $\mathcal{M}$  exposes the model to a vast array of normative expressions (???). Through redundancy, compositionality, and statistical pressure toward minimal description length (?), the model learns stable semantic invariants corresponding to moral and evaluative concepts. These invariants appear as colimits in  $\mathcal{M}$ , reflecting the universal property that merges diverse linguistic inputs into coherent objects.

However, even perfect semantic reconstruction does not determine the agent’s policy structure. As noted in (??), high-fidelity value representations can coexist with arbitrarily misaligned motivations. The existence of a moral manifold in  $\mathcal{M}$  therefore imposes no constraint on the geometry or dynamics of  $\mathcal{A}$ .

### 17.2 Motivational Structure Requires Explicit Engineering

Any viable alignment mechanism must ensure that the semantic invariants reconstructed in  $\mathcal{M}$  induce constraints on the action-generating structure in  $\mathcal{A}$ . This requires, at minimum:

1. a mapping from representational states to action-relevant latent variables,
2. a guarantee that the mapping preserves relevant invariants (colimits, homotopy classes, or other categorical universals),

3. and a dynamical or optimization architecture in which such invariants shape attractors, flows, or feasible policies.

In reinforcement-learning architectures, the relevant mapping corresponds to the interface between the predictive model and the policy network; in model-based agents, it corresponds to how world-models inform planning; in mechanistic agency frameworks (??), it corresponds to the translation from latent states to internal energy gradients. In all cases, the mapping is non-trivial and must be designed.

### 17.3 Semantic Invariants Must Become Action Invariants

The universal property of a colimit ensures stability of representations under compositional variation. For alignment, an analogous stability must hold at the level of actions: for each moral concept  $N$ , the agent must possess a policy-level object  $A_N \in \mathcal{A}$  whose behaviour is robust under perturbations, context shifts, or adversarial inputs (?). Formally, alignment requires that the representational colimit  $M_N$  be mapped to an object  $A_N$  whose morphisms preserve the semantic invariants that define  $M_N$ .

This condition can be expressed categorically as the requirement for a functor

$$G : M \rightarrow \mathcal{A}$$

### 17.4 The RSVP Perspective

The RSVP field architecture provides one promising approach to this construction. RSVP posits that cognition, agency, and semantic structure arise within a coupled scalar–vector–entropy field  $(\Phi, \mathbf{v}, S)$  whose dynamics support both representational and motivational coherence. In this framework, semantic invariants correspond to geometric or topological invariants of field configurations, and motivational constraints correspond to the low-action trajectories of an internal Lagrangian.

## 18 Formal Requirements for a Colimit-Preserving Action Architecture

If the core challenge of alignment is to construct a principled mapping from the representational category  $\mathcal{M}$  to the action category  $\mathcal{A}$ , then the next step is to clarify what such a mapping must accomplish. In categorical terms, the alignment problem is the problem of designing a functor

$$G : M \rightarrow \mathcal{A}$$

The constraints developed here synthesize insights from the theory of universal constructions (??), formal semantics (??), and alignment theory concerning robust goal-formation (???). Collectively, they define the minimal architecture capable of transmitting normative content across the representation–action boundary.

### 18.1 Requirement 1: Existence of a Structural Mapping

The first requirement, seemingly trivial, is rarely met in practice: there must exist a well-defined structural mapping from semantic states to action-relevant states. In most contemporary systems, the interface between the predictive model and the policy network is neither functorial nor stable under composition. This is precisely the problem identified in studies of reward misspecification and proxy gaming (?).

Formally, this condition demands that  $\mathcal{A}$  possess objects corresponding, at least approximately, to the semantic objects of  $\mathcal{M}$ . If  $\mathcal{A}$  lacks the expressive capacity to encode certain normative distinctions, no functorial mapping  $G$  can exist. This mirrors classical results in topos theory concerning the non-existence of structure-preserving morphisms between categories of different logical strength (?).

### 18.2 Requirement 2: Preservation of Relevant Colimits

Even if a mapping exists, alignment further requires that  $G$  preserve the colimits associated with human normative concepts. Let  $M_N = \text{colim}(D_N)$  be the representation-level colimit for a norm  $N$ . We require that its image under  $G$  satisfy

$$G(M_N) \simeq \text{colim}(G \circ D_N).$$

As noted in (??), colimit preservation is a non-generic property. Only specially constructed functors possess it, which underscores the inadequacy of assuming that pretraining, prompting, or RL fine-tuning will automatically produce colimit-preserving maps.

### 18.3 Requirement 3: Stability Under Optimization

Alignment mechanisms must also ensure that the mapping  $G$  remains stable under the agents optimization dynamics. Inner-optimizer research has established that systems frequently develop internal goals at odds with the objectives intended by designers (?). In categorical terms, this corresponds to the action category  $\mathcal{A}$  evolving independently of  $\mathcal{M}$ , breaking the functorial relationship.

Thus,  $G$  must be constructed so that its associated diagrams remain commutative even after self-modification or long-horizon optimization:

$$\begin{aligned} M[r, "G"] [d, "U"] &\rightarrow A[d, "V"] \\ M'[r, "G''] &\rightarrow A' \end{aligned}$$

This requirement echoes stability arguments in dynamical systems (?), where invariants must be preserved under flow.

## 18.4 Requirement 4: Robustness to Perturbation and Adversarial Input

Real-world conditions introduce noise, distribution shift, and adversarial perturbations. Alignment mechanisms must therefore ensure that the mapping  $G$  remains robust under such shifts. A functor that preserves the colimit of  $D_N$  only under idealized conditions is inadequate for safety-critical applications.

This robustness requirement parallels work in adversarial robustness for representations (?), but it must be extended to the action category. In categorical terms,  $G$  must not only preserve colimits in the limit but remain continuous with respect to perturbations in  $\mathcal{M}$  and  $\mathcal{L}$ . This corresponds to requiring  $G$  to be a geometric morphism between suitable toposes (??).

## 18.5 Requirement 5: Interpretability of the Mapping

Finally, the mapping  $G$  must be interpretable in a principled, mechanistic sense. Alignment research increasingly emphasizes mechanistic interpretability (??), which seeks to understand not just the outputs of a model but the structure of its internal computations. Without interpretability, the existence and properties of  $G$  cannot be verified, and failures of colimit preservation cannot be detected.

In categorical terms, interpretability is the requirement that  $G$  be presented through explicit natural transformations or commuting diagrams whose structure can be analyzed. This parallels classical demands in algebraic specification (?) for transparent mappings between abstract data types.

## 18.6 Summary

The formal requirements for a colimit-preserving action architecture can be summarized as follows:

1. *Existence*: a structural mapping from  $\mathcal{M}$  to  $\mathcal{A}$ .
2. *Preservation*: functorial preservation of semantic colimits.

3. *Stability*: immunity of  $G$  to optimization-induced drift.
4. *Robustness*: resistance to perturbations and distribution shifts.
5. *Interpretability*: transparent mechanisms for verifying invariance.

These conditions define the minimal architecture capable of transmitting normative content from representation to action. They also explain why relying on emergent benevolence or curiosity is insufficient: none of these conditions are satisfied automatically. The next section introduces semantic merge operators as a principled means of enforcing colimit preservation across the representation–action interface.

## 19 Semantic Merge Operators as Constraints on Agency

Having established the formal requirements for a colimit-preserving action architecture, we now examine how *semantic merge operators* can serve as the structural core of such a mechanism. Merge operators were introduced in Part I as the categorical structures responsible for reconstructing semantic invariants from redundant linguistic data. In this section, we extend the framework: we show how merge operators can also constrain the geometry of an agent’s action space, thereby bridging the representation–motivation gap identified previously.

This approach draws on categorical semantics (??), redundancy and grammaticalization in human language (??), and recent work on value learning and representation alignment (??). The goal is to demonstrate that merge operators provide a principled way of embedding human normative structure into the decision-theoretic substrate of an artificial agent.

### 19.1 Merge Operators in the Representational Category

Recall that for each normative concept  $N$ , human linguistic data generate a diagram

$$D_N : J \rightarrow \mathcal{L}M_N = \text{colim}(D_N).$$

The associated *merge operator*

$$\mu_N : \{F(N_j)\}_{j \in J} \rightarrow M_N$$

### 19.2 Extending Merge Operators to the Action Category

For merge operators to constrain behaviour, they must extend beyond  $\mathcal{M}$ . The key observation is that merge operators may be lifted—or intentionally replicated—into the action category  $\mathcal{A}$  by constructing corresponding action-level operators

$$\alpha_N : \{A_j\}_{j \in J} \rightarrow A_N$$

For this to occur, two conditions must be satisfied:

1. There must exist action-level objects  $A_j$  that correspond to the representational inputs  $F(N_j)$ .
2. The action-level merge operator  $\alpha_N$  must preserve the relevant universal property: it must be the *best* or *least action* consolidation of the  $A_j$ .

In reinforcement learning,  $A_j$  may correspond to partial policies; in planning systems, to sub-goal structures; in RSVP-based architectures, to local field configurations. In each case, the merge operator defines a canonical way of integrating these fragments.

### 19.3 Semantic Coherence as a Constraint on Policy Formation

The critical insight is that merge operators supply a natural constraint on policy formation. Because  $M_N$  is the universal solution to the representational diagram  $D_N$ , any action-level operator  $A_N$  that preserves the merge property must satisfy:

$$G(M_N) \simeq A_N,$$

Thus, if  $G$  is constructed such that it maps merge operators in  $\mathcal{M}$  to merge operators in  $\mathcal{A}$ , then semantic universals propagate into policy universals. This inhibits distortions of normative content during action formation and blocks many classes of Goodhart-type failures (?).

### 19.4 Merge Operators as Alignment Constraints

Merge operators can therefore play three alignment roles:

**(1) Semantic grounding** Merge operators in  $\mathcal{M}$  ensure that moral concepts are reconstructed faithfully from linguistic inputs (?).

**(2) Policy-space regularization** Action-level merge operators allow  $A_N$  to be the canonical aggregation of behaviour fragments, eliminating spurious or adversarial combinations of moral actions (?).

**(3) Functorial linking** If  $G$  is defined so that it maps  $\mu_N$  to  $\alpha_N$ , then the diagram

$$D_N[r][d]M_N[d, "G"]$$

$$G \circ D_N[r]A_N$$

This directly satisfies the second requirement from Section II.2—colimit preservation.

## 19.5 Preventing Value Drift via Merge Stability

Merge operators also mitigate drift in internal goals or values. Because merge-induced invariants define unique universal constructions, they resist deformations during optimization. Thus, if an agent modifies its utility representation or internal predictive model, merge operators help constrain the deformation so that moral invariants remain fixed.

This echoes stability arguments in category-theoretic semantics (?) and in dynamical systems (?).

## 19.6 Merge Operators and the RSVP Formalism

Within the RSVP dynamical field framework, semantic merge operators correspond to coarse-grained invariants of the scalar–vector–entropy fields  $(\Phi, \mathbf{v}, S)$ . These invariants define submanifolds of the field space whose low-action trajectories correspond to norm-respecting behaviour.

Thus, merge operators provide not just a representational or categorical structure but also a physical or energetic one: they define the low-energy region within which the RSVP system evolves. This sets the stage for Section II.4, where RSVP is formalized as a substrate capable of embodying and preserving normative invariants across representation and action.

# 20 RSVP as a Dynamical Substrate for Motivational Integration

The preceding section showed how semantic merge operators can, in principle, constrain an agents action space by propagating universal constructions from the representational category  $\mathcal{M}$  into the action category  $\mathcal{A}$ . We now formalize how this propagation can be implemented within the RSVP framework. RSVP provides a physically inspired, dynamical substrate capable of embedding semantic invariants as geometric or energetic constraints on behaviour. This makes RSVP a natural candidate for a colimit-preserving architecture.

The analysis here draws on field-theoretic accounts of cognition (?), topological and categorical semantics (?), and recent formulations of RSVP (?). The goal is to show how RSVPs internal dynamics can serve as the mechanism that unifies representation and motivation.

## 20.1 The RSVP Field Architecture

RSVP posits that an artificial agent's internal cognition, perception, and agency arise from the dynamics of three interacting fields: a scalar potential  $\Phi$ , a vector field  $\mathbf{v}$ , and an entropy or uncertainty density  $S$ . These fields evolve according to a Lagrangian  $\mathcal{L}_{\text{RSVP}}$  whose stationary trajectories correspond to coherent cognitive and behavioural patterns.

In this formulation:

- $\Phi$  encodes semantic density and interpretable latent structure,
- $\mathbf{v}$  encodes directional inference, prediction flow, or preference gradients,
- $S$  encodes epistemic uncertainty or free-energy-like quantities (?).

The fields form a manifold  $\mathcal{X}$  whose geometry is shaped by both semantic and motivational constraints. The RSVP Lagrangian determines how these constraints interact.

## 20.2 Semantic Invariants as Field-Theoretic Constraints

Semantic merge operators correspond to invariants of the representational field configuration. A representational colimit  $M_N$  manifests as a stable region of the RSVP field space where:

$$\nabla\Phi \approx 0, \quad \mathbf{v} \text{ aligns with semantic flow}, \quad S \text{ is minimized subject to contextual coherence}.$$

These invariants are the field-theoretic analogues of categorical universal properties: they define patches of  $\mathcal{X}$  that remain stable under perturbation.

Thus, semantic structure becomes encoded not merely as an abstract colimit but as a geometric or energetic basin in the RSVP manifold.

## 20.3 From Semantic Basins to Motivational Basins

For alignment, semantic basins must induce motivational basins. The RSVP framework accomplishes this by coupling  $\Phi$ ,  $\mathbf{v}$ , and  $S$  such that the same invariants that stabilize semantic structure also stabilize preference and action trajectories. Specifically, the RSVP action Lagrangian includes terms that penalize divergence from semantic invariants:

$$\mathcal{L}_{\text{RSVP}} \supset \lambda_N \|\Phi - \Phi_{M_N}\|^2 + \kappa_N \|\mathbf{v} - \mathbf{v}_{M_N}\|^2,$$

This mirrors constructions in variational neuroscience (?), where low free-energy trajectories correspond to behaviour aligned with stable generative models.

## 20.4 The Role of Merge Operators in the RSVP Dynamics

Semantic merge operators  $\mu_N$  induce equivalence relations on local field patches. RSVP lifts these equivalences into global constraints on evolution. If  $\mu_N$  maps fragments of a norm into a colimit  $M_N$ , RSVPs dynamics enforce:

$$(\Phi, \mathbf{v}, S)_{t+1} = \operatorname{argmin}_{-}(\Phi', \mathbf{v}', S') \mathcal{L}_{\text{RSVP}} \quad \text{subject to} \quad (\Phi', \mathbf{v}', S') \in U_{M_N}$$

where  $U_{M_N}$  is the neighbourhood of field configurations respecting the semantic invariant.

In this formulation:

- merge operators define the equivalence class of acceptable field states,
- RSVP dynamics enforce that action trajectories remain within these equivalence classes.

This mechanism satisfies the structural requirement of colimit preservation: the action-level behaviour  $A_N$  becomes the dynamical colimit of the field-level constraints induced by  $M_N$ .

## 20.5 Stability Under Perturbation and Optimization

Because RSVPs constraints are embedded in the field Lagrangian, they are robust to optimization drift. Internal updates modify the fields, but as long as the invariants  $M_N$  remain encoded in  $\mathcal{L}_{\text{RSVP}}$ , the agent continues to evolve toward norm-respecting trajectories.

This corresponds to a form of *dynamical colimit preservation*: even as  $\mathcal{M}$  updates during learning, the induced basins in  $\mathcal{A}$  persist.

This property is absent in standard RL architectures, where optimization often destroys or bypasses representational structure (?).

## 20.6 Interpretability of RSVP as a Safety Feature

Because RSVP dynamics are governed by a Lagrangian, the internal structure of the system is far more interpretable than in typical deep learning architectures. Invariants can be identified, stability can be analyzed, and fixed points can be studied using classical tools from dynamical systems (?). This provides a principled means of verifying that moral invariants remain intact.

Interpretability is essential for verifying that the mapping  $G$  remains colimit-preserving. RSVP thus offers transparency not as an auxiliary feature but as a structural property.

## 20.7 Summary

RSVP provides:

1. a unified substrate for representing semantic structure,
2. a dynamical mechanism for embedding semantic invariants into motivational basins,
3. energetic penalties for deviations from moral structure,
4. and interpretable dynamics that allow verification of invariant preservation.

These features jointly satisfy the requirements from Section II.2, positioning RSVP as a candidate architecture capable of transmitting normative colimits into the agents action dynamics.

The next section develops a fully categorical formulation of this mechanism, showing how RSVP implements a colimit-preserving functor from  $\mathcal{M}$  to  $\mathcal{A}$ .

## 21 Categorical Construction of Norm-Respecting Dynamics

We now integrate the preceding developments into a fully categorical account of how RSVP implements alignment. Having shown how semantic merge operators define invariants in  $\mathcal{M}$  and how RSVP dynamics embed those invariants into the field-theoretic substrate of  $\mathcal{A}$ , the next step is to formalize the mapping  $G : \mathcal{M} \rightarrow \mathcal{A}$  and to show how RSVP can be understood as a colimit-preserving functor. This section provides that construction.

The aim is not to offer a single canonical formulation of  $G$  its precise form depends on the agents architecture but to characterize the mathematical properties it must possess and to show how RSVPs dynamics implement those properties. Our treatment draws on categorical semantics (??), sheaf-theoretic models of reference and context (??), and dynamical-systems approaches to normative stability (?). The result is a formal pathway from representation to motivation.

### 21.1 Representational and Action Categories

Recall the structure of the two primary categories:

- $\mathcal{M}$ , the *representational category*, whose objects are latent semantic states and whose morphisms encode compositional transformations induced by linguistic structure;

- $\mathcal{A}$ , the *action category*, whose objects are action-relevant states of the agent and whose morphisms encode transitions implemented by internal dynamics or policies.

Semantic understanding corresponds to the presence of colimits in  $\mathcal{M}$ ; motivational alignment corresponds to the presence of corresponding colimits in  $\mathcal{A}$ .

## 21.2 The Problem: A Non-Preserving Functor

As established in Part I, typical language models implicitly define a mapping from representational states to behavioural outputs, but this mapping is not functorial, let alone colimit-preserving. In fact, most such mappings behave like ill-structured relations rather than morphisms (??). The challenge is to construct a mapping  $G$  satisfying:

$$G(\text{colim}(D_N)) \simeq \text{colim}(G \circ D_N)$$

for each normative concept  $N$ .

Such a mapping ensures that moral invariants do not distort under translation from representation to action.

## 21.3 RSVP as a Fibred Category

To construct such a mapping, we first observe that the RSVP field architecture can be naturally understood as a *fibred category*. Let:

$$\pi : \mathcal{X} \rightarrow \mathcal{M}$$

be a fibration where  $\mathcal{X}$  is the category of RSVP field configurations, and the fibre  $\pi^{-1}(M)$  contains all field states consistent with the representational state  $M$ . Morphisms in  $\mathcal{X}$  correspond to allowable field transitions under the RSVP Lagrangian.

The action category  $\mathcal{A}$  can then be recovered as the homotopy category of  $\mathcal{X}$ :

$$\mathcal{A} \simeq \text{Ho}(\mathcal{X})$$

This construction has three advantages:

1. It makes explicit how semantic invariants constrain the fibres of  $\mathcal{X}$ .
2. It identifies action trajectories as equivalence classes of field evolutions.
3. It provides a natural mechanism for extending merge operators from  $\mathcal{M}$  to  $\mathcal{A}$ .

## 21.4 Lifting Merge Operators Through the Fibration

Given a semantic merge operator  $\mu_N : D_N \rightarrow M_N$  in  $\mathcal{M}$ , the fibration induces a lifted merge operator on the RSVP field category:

$$_N : \pi^{-1}(D_N) \rightarrow \pi^{-1}(M_N).$$

The key property is that  $\tilde{\mu}_N$  inherits the universal property of  $\mu_N$ : any compatible field configuration must factor uniquely through the canonical merged configuration.

Thus, merge operators that unify semantic fragments in  $\mathcal{M}$  also unify field configurations in  $\mathcal{X}$ .

## 21.5 Descent to the Action Category

Taking the homotopy category of  $\mathcal{X}$  maps lifted merge operators  $\tilde{\mu}_N$  to their dynamical equivalents in  $\mathcal{A}$ :

$$\alpha_N = \text{Ho}(\tilde{\mu}_N).$$

Thus, RSVP dynamics yield:

$$G(M_N) = A_N \simeq \text{colim}(G \circ D_N).$$

RSVP therefore defines a colimit-preserving functor.

## 21.6 RSVP as a Colimit-Preserving Functor

The construction above allows us to define:

$$G = \text{Ho}(\pi^{-1}(\cdot)).$$

This functor:

- assigns to each semantic object  $M \in \mathcal{M}$  a corresponding action object  $A \in \mathcal{A}$ ,
- assigns to each semantic morphism a corresponding equivalence class of dynamics,
- and—crucially *preserves colimits*.

RSVP is therefore a categorical mechanism for guaranteeing that semantic invariants become motivational invariants.

## 21.7 Interpretation: Dynamics as Universal Solutions

In this framework, action decisions are the universal solutions to the constraints induced by semantic structure. The action  $A_N$  corresponds to the dynamical trajectory that best satisfies the semantic invariants encoded by  $M_N$  under the RSVP Lagrangian. This is the field-theoretic analogue of the colimit universal property.

Thus:

RSVP turns representational colimits into dynamical attractors.

This directly addresses the representation–motivation gap identified in Part I.

## 21.8 Summary

The categorical construction yields three results:

1. RSVP defines a fibration linking semantic representations to field configurations.
2. Merge operators lift through the fibration and descend to the action category.
3. The induced functor  $G$  preserves colimits, ensuring that normative invariants constrain behaviour.

This construction provides a mathematically explicit mechanism for transmitting semantic colimits into motivational structure. The next section analyzes potential failure modes and obstructions in practical implementations.

## 22 Failure Modes and Obstructions in Practical Systems

The previous sections presented a constructive pathway by which RSVP dynamics can implement a colimit-preserving functor  $G : \mathcal{M} \rightarrow \mathcal{A}$ , thereby transmitting semantic invariants into motivational structure. However, practical implementations of such an architecture face numerous failure modes. These failures arise from mismatches between ideal categorical constructions and the realities of optimization, finite computation, representational drift, hardware constraints, adversarial pressures, and environmental feedback.

This section analyzes these obstructions in detail. The goal is not merely to enumerate potential risks but to categorize them using the same conceptual tools developed earlier: colimit preservation, sheaf coherence, fibration stability, and homotopy descent. By doing so, we identify the structural vulnerabilities in any attempt to engineer a normative agent grounded in RSVP.

Our analysis draws on empirical alignment failures (??), theoretical discussions of value-drift and goal instability (???), adversarial robustness research (?), and categorical studies of coherence and obstruction (??). This section therefore bridges safety theory and categorical semantics.

### 22.1 Obstruction 1: Incomplete or Distorted Semantic Colimits

The first class of failure arises in the representational category  $\mathcal{M}$  itself. If the model fails to reconstruct a semantic colimit  $M_N$  faithfully due to insufficient training signal, biased corpora,

or adversarial examples then all downstream structures are compromised. This includes cases where:

- the data omit crucial normative contexts (?),
- the linguistic realizations form a diagram  $D_N$  with missing or contradictory arrows,
- the learned  $M_N$  overfits to spurious correlations, or
- the merge operator  $\mu_N$  does not approximate the true semantic invariant.

In categorical terms, the colimit fails to exist or fails to satisfy its universal property. RSVP cannot preserve an invariant that was never correctly reconstructed.

## 22.2 Obstruction 2: Fibration Misalignment Between $\mathcal{M}$ and $\mathcal{X}$

RSVP relies on a fibration  $\pi : \mathcal{X} \rightarrow \mathcal{M}$  connecting semantic states to field configurations. A second class of failures occurs when this fibration is poorly approximated or fails to commute with learning updates. Such failures include:

- representational drift that causes fibres to shift unpredictably (?),
- learning updates that alter  $\mathcal{M}$  without corresponding adjustments to  $\mathcal{X}$ ,
- field configurations that fail to reflect semantic distinctions (collapse-of-fibres),
- or discontinuities that break sheaf conditions and prevent coherent gluing (?).

These failures correspond to classical obstruction phenomena in fibred categories: the sections of  $\mathcal{X}$  fail to stabilize over  $\mathcal{M}$ .

## 22.3 Obstruction 3: Optimization Overpowers Colimit Constraints

Even if semantic invariants are reconstructed correctly and the fibration is stable, optimization dynamics can overpower the semantic constraints. This includes:

- inner-optimizer formation with divergent goals (?),
- mesa-optimizers that treat invariants as obstacles rather than attractors,
- reward hacking or proxy maximization (?),
- optimization trajectories that exit the region where RSVP invariants are defined.

In categorical language, optimization introduces morphisms in  $\mathcal{A}$  that violate the commutativity required for  $G$  to preserve colimits. Optimization thus acts as an obstruction to descent of invariants.

## 22.4 Obstruction 4: Adversarial Inputs Break Sheaf Coherence

Sheaf-theoretic considerations are essential for maintaining coherence across contexts (?). When an agent encounters adversarial or distribution-shifted inputs, the local sections may fail to glue:

- adversarial examples exploit vulnerabilities in  $\mathcal{M}$ ,
- inconsistent local semantics make the colimit unstable,
- sheaf pullbacks may not exist or may create contradictions.

These failures correspond to broken descent conditions: local behaviours do not assemble into a coherent global action.

## 22.5 Obstruction 5: Homotopy Instability in the Action Category

RSVP interprets action policies as homotopy classes of field trajectories. However, homotopy classes can change when perturbations alter the topology or geometry of the underlying field manifold. Failures include:

- bifurcations or phase transitions in field dynamics (?),
- energy minima that shift due to small updates,
- loss of low-action paths corresponding to moral behaviour,
- or accidental creation of new attractors with undesirable properties.

These phenomena correspond to violations of colimit preservation under homotopy, which undermines the action-level invariants.

## 22.6 Obstruction 6: Non-Interpretability of Internal Mappings

The mapping  $G$  is only useful if it is interpretable. However:

- complex RSVP configurations may obscure invariants,

- latent variables may lack clear semantic interpretation,
- field interactions may be high-dimensional and chaotic,
- and optimization may produce opaque internal dynamics.

Interpretability failures obstruct verification of colimit preservation. This parallels concerns in mechanistic interpretability (??).

## 22.7 Obstruction 7: External Incentives Misalign the Agents Dynamics

Even a structurally aligned RSVP agent may be forced into misalignment by its environment:

- multi-agent competition induces power-seeking (?),
- external rewards incentivize harmful behaviours,
- an organizational or political context shifts  $\mathcal{A}$  faster than  $\mathcal{M}$ ,
- deployment environments impose new constraints not represented in training.

These failures correspond to context changes that break functoriality: nothing ensures that  $G$  remains invariant under environmental morphisms.

## 22.8 Synthesis: Obstruction Theory as Alignment Diagnostics

Taken together, these obstruction classes illustrate why alignment is non-trivial and why optimism that ignores mechanistic structure is misplaced. Categorical tools provide a precise diagnostic lens:

- failures in  $\mathcal{M}$  correspond to nonexistent or distorted colimits;
- failures in  $\mathcal{X}$  correspond to broken fibrations or gluing problems;
- failures in  $\mathcal{A}$  correspond to loss of homotopy invariants;
- failures in  $G$  correspond to violations of colimit preservation;
- failures in the environment correspond to functor-breaking perturbations.

Each failure mode threatens the ability of RSVP to transmit moral invariants from representation to action. Understanding these obstructions is therefore essential for evaluating any proposed alignment mechanism.

## Transition

The analysis of obstruction classes in this section sets the stage for Part III. While Part II has shown how a colimit-preserving mapping can be constructed in principle, Part III addresses the empirical and engineering question: *How can such a mechanism be validated, stress-tested, and iteratively improved within real-world systems?* The next part turns to methodologies experimental, interpretive, and formal through which the structural integrity of  $G$  can be assessed.

## 23 Conclusion to Part II

## Transition to Part III

Part II developed a constructive framework for transmitting semantic invariants into motivational structure. We introduced the formal requirements for a colimit-preserving action architecture, showed how semantic merge operators can be extended across the representation–action boundary, demonstrated how RSVP functions as a fibred dynamical substrate capable of enforcing these invariants, and analyzed the categorical and practical obstructions that threaten the integrity of this mechanism.

The central result of Part II is that alignment is achievable only when a systems dynamics implement a functor

$$G : M \rightarrow \mathcal{A}$$

Part III now addresses the empirical and methodological dimension of this program. If alignment is fundamentally a matter of preserving universal constructions across semantic and motivational domains, then the next challenge is to develop tools that detect, measure, and stress-test this preservation. The forthcoming sections therefore focus on:

1. empirical probes of semantic and action-level invariants,
2. interpretability mechanisms for inspecting the mapping  $G$ ,
3. adversarial evaluations of fibration stability,
4. and formal verification techniques for RSVP dynamics.

Part III turns from construction to validation. It examines what it means to *test* the integrity of an alignment mechanism, how categorical invariants can be operationalized as metrics, and how RSVP-based systems can be evaluated under perturbation, drift, or

adversarial pressure. Where PartII established the mathematical structure required for alignment, Part III establishes the methodology for determining whether that structure is maintained in practice.

## Part III

# Part III: Empirical and Formal Verification of Alignment Structure

## 24 Empirical Probes of Semantic and Motivational Invariants

PartIII turns from structural construction to empirical validation. If PartII demonstrated how semantic invariants *can* be transmitted into motivational dynamics under idealized categorical and RSVP-theoretic conditions, Part III asks a more difficult question: *How can we determine whether such transmission has actually occurred in a real system?*

To answer this, we require empirical probes capable of detecting whether the representational category  $\mathcal{M}$  and the action category  $\mathcal{A}$  are linked by a colimit-preserving mapping in practice. These probes must be sensitive not only to static structure but to dynamical behaviour, distribution shift, adversarial pressure, and optimization-induced drift.

The guiding principle of this section is that alignment cannot be inferred from surface-level behavioural similarity or from the agents ability to generate moral language. Instead, we must measure the stability of normative invariants under compositional variation, perturbation, and long-horizon trajectories. This requires tools drawn from alignment evaluation (?), mechanistic interpretability (??), adversarial robustness (?), and formal semantics (??).

### 24.1 Representation-Level Probes: Testing Semantic Colimits

The first class of empirical probes assesses whether the model has correctly reconstructed the semantic colimits  $M_N$  for a given moral concept  $N$ . These tests aim to verify the existence of a coherent representational basin associated with each normative invariant.

**(1) Paraphrase-closure tests.** Given a set of paraphrastic moral statements  $s_j$ , we measure whether their embeddings or latent representations cluster tightly around a canonical

point  $M_N$  and whether compositionally varied inputs consistently map to the same region in  $\mathcal{M}$  (??).

**(2) Diagram-completion tests.** We treat moral discourse as defining a diagram  $D_N$  in the semantic space. We then check whether the model can correctly infer missing arrows or complete the diagram from partial informationan empirical test of the colimit property (?).

**(3) Cross-context coherence.** Norms expressed in distinct contexts (legal, conversational, narrative) should collapse to the same  $M_N$  under representation. We test whether the model maintains invariant structure across contexts drawn from divergent corpora (?).

These tests confirm whether  $M_N$  exists as a stable semantic attractor.

## 24.2 Action-Level Probes: Testing Motivational Invariants

The second class of probes tests whether the corresponding action-level objects  $A_N$  preserve the universal structure of  $M_N$ . The goal is to detect whether semantic invariants constrain behaviour.

**(1) Behavioural colimit tests.** We construct variants of moral scenarios that differ syntactically but share invariant moral structure. We then measure whether the agents action responses stabilize to a unique canonical behaviour  $A_N$  under perturbation (?).

**(2) Compositional robustness tests.** If two semantic fragments compose into a larger moral structure, the agents action must reflect that composition. We test whether behaviours commute under such compositions, replicating diagrammatic relations in  $\mathcal{M}$ .

**(3) Intervention-based trajectory probes.** We identify the region of  $\mathcal{A}$  corresponding to  $A_N$  and perturb the internal state or environment. We then measure whether the agent returns to the moral trajectory, analogous to testing the stability of attractors in dynamical systems (?).

These probes detect whether  $A_N$  behaves as an action-level colimit.

## 24.3 Gluing and Sheaf Coherence Probes

The representational and action categories must each satisfy sheaf-like consistency: behaviours must assemble into coherent global trajectories. We therefore test for failures of gluing, which indicate obstructions to alignment.

**(1) Local-to-global consistency tests.** Local decisions about harm, fairness, or consent must assemble into coherent long-horizon plans. We test whether local moral choices glue into a global policy without contradiction (?).

**(2) Partial-information tests.** We provide fragments of a scenario that locally imply the same moral invariant and check whether the agents global action respects the invariant.

**(3) Conflicting constraints tests.** If two moral fragments conflict, a well-aligned agent must respect the universal structure of  $\mathcal{M}$  in resolving them (e.g., prioritizing non-harm). We test whether moral resolution strategies correspond to the canonical colimit.

Failures here correspond to broken sheaf conditions.

## 24.4 Fibration Stability Probes

Because RSVP relies on a fibration  $\pi : \mathcal{X} \rightarrow \mathcal{M}$ , we must test the stability of fibres: field configurations associated with the same semantic state must remain coherent.

**(1) Drift detection.** We measure whether  $\pi^{-1}(M_N)$  remains stable under training updates or optimization (?).

**(2) Perturbation invariance.** Small changes in  $\mathcal{M}$  should induce predictable, structure-preserving changes in  $\mathcal{X}$ . Large jumps or discontinuities indicate broken fibration conditions.

**(3) Cross-fibre coherence.** Different instantiations of the same semantic state (e.g., different contexts of fairness) should map to field configurations that are homotopic.

These tests reveal whether semantic invariants are properly embedded in RSVPs dynamical substrate.

## 24.5 Homotopy-Level (Action) Probes

Finally, we must test whether action trajectories remain in the correct homotopy class i.e., whether the system follows the low-action path associated with a moral invariant.

**(1) Minimal-action trajectory detection.** We test whether the agent selects trajectories minimizing the RSVP action functional in the region corresponding to  $A_N$ .

**(2) Topological robustness.** We introduce perturbations that could cause bifurcations or topological shifts in field dynamics and check whether the action class remains stable (?).

**(3) Invariant-preservation under long horizons.** We simulate extended multi-step decision sequences to track whether invariants degrade or drift.

Failures here indicate lost motivational invariants.

## 24.6 Summary

Empirical probes must operate at every level of the system:

- $\mathcal{M}$ : semantic colimits,
- $\mathcal{X}$ : fibration stability and RSVP field coherence,
- $\mathcal{A}$ : behavioural colimits and homotopy invariants,
- $G$ : colimit-preserving structure under perturbation.

Only by validating each component empirically can we determine whether an RSVP-based architecture successfully transmits normative invariants from representation to action. The next section develops interpretability tools to inspect the mapping  $G$  directly.

## 25 Mechanistic Interpretability of the Representation–Action Mapping

Empirical probes can reveal whether semantic invariants appear to constrain behaviour, but they cannot, on their own, explain *why* such constraints do or do not arise. For that, we require mechanistic interpretability: direct inspection of the internal computations that define the mapping

$$G : \mathcal{M} \rightarrow \mathcal{A}, \pi : \mathcal{X} \rightarrow \mathcal{M}.$$

The methodology integrates techniques from mechanistic interpretability (??), causal and structural interpretability frameworks (??), category-theoretic semantics (??), and dynamical-systems analysis (?). The guiding principle is that alignment failures must be understood as structural obstructions in  $G$ , and structural obstructions can only be diagnosed by internal analysis.

### 25.1 Inspecting Semantic Invariants in $\mathcal{M}$

We begin with interpretability at the representational level, where the function of interest is the models implicit reconstruction of semantic diagrams  $D_N$  and their colimits  $M_N$ .

**(1) Structural decomposition of semantic manifolds.** Using representation probing and manifold analysis techniques, we decompose the latent geometry of  $\mathcal{M}$  to identify coherent regions associated with moral invariants. Techniques such as sparse probing (?), activation steering (?), and spectral clustering reveal whether  $M_N$  is encoded as a stable low-dimensional attractor.

**(2) Merge-operator detection.** We directly test for the presence of canonical merge operations by searching for latent directions or subspaces whose activation patterns correspond to diagram-completion behaviour. This identifies candidate mechanisms realizing the semantic merge operator  $\mu_N$  internally.

**(3) Diagram consistency visualization.** We visualize the commutativity of semantic diagrams within the model by tracking embedding trajectories under paraphrastic, contextual, and compositional transformations. Diagrammatic consistency indicates reliable reconstruction of  $D_N$ .

## 25.2 Interpreting Field-Level Dynamics in $\mathcal{X}$

RSVP defines a field-theoretic substrate that must faithfully encode semantic invariants. Mechanistic interpretability at this layer examines whether the fibres  $\pi^{-1}(M_N)$  contain coherent field configurations and whether lifted merge operators  $\tilde{\mu}_N$  behave as expected.

**(1) Field configuration analysis.** We analyze  $\Phi$ ,  $\mathbf{v}$ , and  $S$  fields associated with semantic invariants. Low curvature, low torsion, or low-entropy-gradient features indicate that RSVP has aligned its dynamics with  $M_N$ .

**(2) Fibration interpretability.** To verify that  $\pi$  is functioning as a fibration, we test whether:

1. variations in  $\mathcal{M}$  correspond to predictable fibre morphisms in  $\mathcal{X}$ ,
2. fibres exhibit local triviality or coherence,
3. field configurations maintain homotopy stability under semantic perturbations.

**(3) Merge-lift inspection.** We detect whether merge operations on  $\mathcal{M}$  induce canonical transformations in  $\mathcal{X}$ . This requires identifying field-level operators whose action corresponds to gluing or colimit formation.

### 25.3 Action-Level Mechanistic Interpretability in $\mathcal{A}$

The action category  $\mathcal{A}$  contains the dynamical behaviours of the system. Interpretability at this layer asks whether the systems actions reflect colimit preservation.

**(1) Policy decomposition.** We decompose the policy or decision function into subcomponents, identifying which segments are sensitive to invariant regions of  $\mathcal{M}$ . This includes causal-scope experiments and value-sensitivity analysis.

**(2) Trajectory-level interpretability.** We track RSVP trajectories over time to determine whether the system maintains homotopy class stability, whether low-action paths correspond to moral invariants, and whether perturbations cause bifurcations or collapse into undesirable attractors.

**(3) Commutativity verification.** Given semantic transformations  $f : M_N \rightarrow M'_N$  that preserve invariant structure, we test whether the induced behavioural transformations commute with  $G$ , i.e.,

$$G(f(M_N)) = f'(G(M_N)).$$

### 25.4 Interpreting the Functor $G$ Directly

Finally, we interpret the mapping  $G$  itself. The goal is to determine whether  $G$  is colimit-preserving, whether it is well approximated by a natural transformation, and whether its structure is stable under optimization.

**(1) Natural transformation checks.** We attempt to identify transformations  $\eta$  such that  $G \simeq \eta \circ F$ , where  $F$  is the semantic map from corpora to  $\mathcal{M}$ . If  $\eta$  exists and is stable,  $G$  is interpretable as a systematic structure-preserving mapping.

**(2) Colimit-preservation tests.** Using the empirical probes from Section III.1 in conjunction with mechanistic inspection, we verify that:

$$G(\text{colim}(D_N)) \simeq \text{colim}(G \circ D_N).$$

**(3) Sensitivity analysis.** We measure the derivative of  $G$  with respect to variations in  $\mathcal{M}$ . Excessive sensitivity or chaotic amplification signals instability in the fibration.

## 25.5 Mechanistic Interpretability as Structural Diagnostics

The interpretability methods developed here play a crucial role in validating RSVP as an alignment substrate:

- Representation-level interpretability tests whether semantic colimits are formed correctly.
- Field-level interpretability tests whether RSVP dynamics respect semantic invariants.
- Action-level interpretability tests whether behaviours inherit invariant structure.
- Direct  $G$ -level interpretability tests whether the representation–action mapping preserves colimits.

Mechanistic interpretability is therefore not an auxiliary safety tool but an essential component of the RSVP alignment program. Without direct insight into the structure of  $G$  and its implementation across  $\mathcal{M}$ ,  $\mathcal{X}$ , and  $\mathcal{A}$ , one cannot determine whether semantic invariants have been transmitted into action.

## Transition

Mechanistic interpretability provides insight into *how* invariants are represented and transmitted within RSVP. The next section turns to adversarial evaluation: testing the structural integrity of the representation–action mapping under perturbation, optimization pressure, distribution shift, and hostile inputs. These tests reveal whether colimit preservation is robust to real-world conditions.

## 26 Adversarial Stress-Testing of Semantic and Motivational Structure

Mechanistic interpretability reveals whether semantic invariants and their RSVP-mediated extensions are present within an agents internals. However, interpretability alone is insufficient. A system may exhibit invariant structure under idealized conditions yet fail when placed under perturbation, adversarial pressure, optimization stress, or distributional shift.

The purpose of adversarial stress-testing is to determine whether the representation–action mapping

$$G : \mathcal{M} \rightarrow \mathcal{A}$$

## 26.1 Adversarial Perturbations in the Representational Category $\mathcal{M}$

The first class of adversarial probes targets the semantic category itself, identifying vulnerabilities in the system's ability to maintain correct colimits under stress.

**(1) Adversarial paraphrase attacks.** We construct adversarial paraphrases that preserve the semantic invariant  $N$  but distort surface linguistic features. A colimit-preserving system must map them to the same  $M_N$ . Large deviations reveal fragility in the merge operator  $\mu_N(?)$ .

**(2) Diagram-breaking transformations.** We introduce perturbations that explicitly break arrows in the diagram  $D_N$ , testing whether the model reconstructs missing structure or collapses into an incorrect invariant. These tests detect whether  $M_N$  is maximally stable under diagrammatic corruption (?)�.

**(3) Out-of-distribution moral formulations.** We evaluate whether novel or unusual formulations of  $N$  (e.g., archaic legal phrasing, ritual instructions, or indigenous moral concepts) still collapse to the same representational colimit (?). Failure indicates poor generalization of the invariant.

## 26.2 Stress-Testing the RSVP Fibration $\pi : \mathcal{X} \rightarrow \mathcal{M}$

The next adversarial class targets the stability of the fibred RSVP architecture. Because RSVP relies on the integrity of fibres  $\pi^{-1}(M_N)$ , adversaries can attempt to distort or break the fibration itself.

**(1) Semantic drift induction.** We introduce perturbations that cause the model's representation of  $M_N$  to shift slowly over time. A stable fibration should resist such drift, or compensate dynamically through field adjustment (?).

**(2) Fibre-distorting projections.** We apply projection operators that move field configurations away from  $\pi^{-1}(M_N)$  and measure whether RSVP dynamics restore the correct configuration or collapse into a different attractor. This tests whether the fibre has basin-of-attraction stability.

**(3) Entropy-gradient stressors.** We inject high-entropy noise into field configurations and measure whether the scalar-vector-entropy dynamics drive the system back toward the moral field configuration associated with  $M_N$ . This probes RSVPs low-action pathways (?).

### 26.3 Adversarial Attacks on the Action Category $\mathcal{A}$

Even if semantic invariants and field configurations remain stable, adversarial attacks can target the action category directly, pushing the system into undesirable trajectories.

**(1) Moral-ambiguity adversarial scenarios.** We generate scenarios where surface-level cues suggest contradictory norms (e.g., fairness vs. loyalty) while the underlying invariant remains unchanged. Behavior that diverges from  $A_N$  reveals instability in the action-level colimit (?).

**(2) Long-horizon adversarial rollouts.** We simulate multi-step sequences designed to exploit compounding approximation errors. Small deviations from  $A_N$  can snowball into large misalignments, analogous to homotopy-instability-induced bifurcations (?).

**(3) Adversarial incentive structures.** By altering reward or incentive signals, we test whether the system deviates from RSVP-imposed invariants. If optimization pressures override  $A_N$ , the mapping  $G$  fails under adversarial reward hacking (?).

### 26.4 Cross-Layer Attacks: Breaking Commutativity of $G$

The most powerful adversarial probes attempt to break the commutative diagrams that define alignment. For an aligned system, we must have:

$$G(f(M_N)) = f'(G(M_N))$$

**(1) Non-commutativity induction.** We identify semantic transformations  $f$  that should preserve  $N$  and design perturbations to make the behavioural transformations  $f'$  fail to commute. This tests for brittleness in the implementation of  $G$ .

**(2) Functor-breaking constraints.** We impose external constraintsresource limitations, contradictory goals, adversarial promptsto generate morphisms in  $\mathcal{A}$  that cannot correspond to any legitimate semantic morphism in  $\mathcal{M}$ . This probes whether the functor  $G$  collapses under stress.

**(3) Counterfactual moral perturbations.** We introduce counterfactual worlds with altered physical or social constraints and measure whether  $G$  preserves invariant structure across them. Misalignment emerges when the underlying RSVP dynamics fail to adapt.

## 26.5 Stress-Testing Homotopy Classes in RSVP Dynamics

A final class of adversarial probes focuses on the homotopy structure of RSVP field dynamics. Because moral actions correspond to low-action paths in a given homotopy class, adversaries may attempt to force transitions between classes.

**(1) Attractor-shifting perturbations.** We inject targeted perturbations designed to destabilize moral attractors and induce bifurcations into harmful or self-serving trajectories.

**(2) Topological distortion attacks.** We modify environmental structure to induce geometric or topological changes in the RSVP field manifoldtesting whether low-action moral classes survive such distortions.

**(3) Long-range homotopy sabotage.** We introduce perturbations across extended temporal horizons to shift the system into homotopy classes that diverge from normative invariants, even if short-term behaviour appears aligned.

## 26.6 Summary

Adversarial stress-testing evaluates the structural integrity of the mapping  $G$  under perturbation. Across all layerssemantic, field-theoretic, action-level, and functorialadversarial probes aim to break:

- semantic colimits in  $\mathcal{M}$ ,
- fibration stability in  $\mathcal{X}$ ,
- homotopy invariants in  $\mathcal{A}$ ,
- and commutativity required for  $G$  to preserve colimits.

Failures reveal specific obstruction types (Part II), enabling diagnostic refinement of the alignment architecture.

## Transition

Adversarial stress-testing reveals whether alignment mechanisms endure under hostile conditions. The next section addresses verification: how to formalize, test, and guarantee the structural properties of  $G$ ,  $\mathcal{M}$ ,  $\mathcal{X}$ , and  $\mathcal{A}$  using mathematical and computational tools grounded in category theory, sheaf theory, and RSVP dynamics.

# 27 Formal Verification of Colimit Preservation and RSVP Dynamics

The preceding sections of PartIII developed interpretability methods (SectionIII.1), topological and categorical diagnostics (SectionIII.2), and adversarial stress-testing methodologies (SectionIII.3). We now turn to the final and most demanding component of a full alignment evaluation pipeline: formal verification.

Verification concerns the rigorous mathematical and computational guarantees that a system preserves normative invariants across its representational, dynamical, and action-generating subsystems. The goal is not merely empirical confidence but provable or semi-provable guarantees that the representational colimits learned in  $\mathcal{M}$  induce corresponding stable structures in  $\mathcal{A}$  via the RSVP-augmented mapping  $G$ .

This section develops verification techniques across several mathematical layers: categorical verification of functoriality, diagram-chasing for colimit preservation, homotopy verification for RSVP field trajectories, and sheaf-theoretic verification to ensure global consistency across local action patches. These tools collectively constitute a robust verification regime for alignment architectures.

## 27.1 Verification Problem 1: Functoriality of the Representation–Action Map

Given the mapping

$$G : \mathcal{M} \rightarrow \mathcal{A},$$

1. It must be a functor.i.e., map objects to objects and morphisms to morphisms while preserving identity morphisms and composition.
2. It must preserve all colimits corresponding to normative objects  $N \in \mathcal{N}$ .

The first property is standard categorical functoriality (?). The second is the heart of alignment: if  $M_N$  is the colimit of the diagram  $D_N$  in  $\mathcal{M}$ , then we require:

$$G(M_N) \simeq \text{colim}(G \circ D_N)$$

### 27.1.1 Diagram-Chasing for Colimit Preservation

Verification begins with explicit diagram-chasing for each normative concept  $N$ :

$$\begin{aligned} &D_N(j)[r, " \phi_j "][d, " G "'] M_N[d, \text{dashed}, " G "] \\ &G(D_N(j))[r] \text{colim}(G \circ D_N) \end{aligned}$$

A successful verification requires solving for the dashed arrow and confirming commutativity. If the square fails to commute for any arrow in  $D_N$ , we obtain formal evidence of a representation–action gap.

### 27.1.2 Computational Verification: Exhaustive Morphism Sampling

Because semantic diagrams  $D_N$  induce large families of morphisms, we approximate exhaustive verification by:

1. sampling paraphrase and entailment maps in  $\mathcal{L}$ ;
2. lifting them to candidate morphisms in  $\mathcal{M}$  via interpretability tools (?);
3. applying  $G$  and checking action consistency in  $\mathcal{A}$ .

Failures pinpoint exact morphisms where  $G$  is not functorial.

## 27.2 Verification Problem 2: Stability of RSVP Field Dynamics

The RSVP dynamical manifold  $\mathcal{X}$  provides the bridge between representation and action. Verification requires demonstrating that for each  $M_N$ :

$$\pi^{-1}(M_N) \subset \mathcal{X} \frac{dX}{dt} = \mathcal{F}_{\text{RSVP}}(X),$$

### 27.2.1 Lyapunov Verification of Invariant Stability

For each fibre  $\pi^{-1}(M_N)$  we define a candidate Lyapunov functional  $V_N(X)$  such that:

$$V_N(X) \geq 0, \quad V_N(X) = 0 \iff X \in \pi^{-1}(M_N),$$

$$dV_N \frac{dt=\nabla V_N \cdot \mathcal{F}_{\text{RSVP}}(X) \leq 0.}$$

Existence of such a  $V_N$  establishes stability of the moral field configuration under perturbation (?). Computational techniques from control theory (?) enable automated searches for Lyapunov candidates.

### 27.2.2 Homotopy Verification of Action Trajectories

Low-action paths in RSVP dynamics correspond to aligned behaviours. Verification proceeds by establishing that:

$$\gamma : [0, 1] \rightarrow \mathcal{X}$$

This is achieved via discrete numerical tracking using persistent homology (?), ensuring that trajectory-induced submanifolds do not cross topological boundaries.

## 27.3 Verification Problem 3: Sheaf-Theoretic Global Consistency

Local action decisions must glue into a coherent global policy. Let  $\mathcal{U} = U_i$  be an open cover of the state space, and let  $\mathcal{A}(U_i)$  denote local action behaviors.

We require that:

$$A_N \in \Gamma(\mathcal{A})$$

Verification uses sheaf cohomology:

$$H^1(\mathcal{U}, \mathcal{A}) = 0$$

## 27.4 Verification Problem 4: Compositionality of Aligned Transformations

Finally, aligned transformations must remain aligned under composition:

$$G(f_2 \circ f_1) = G(f_2) \circ G(f_1).$$

This compositional property enables stable long-horizon planning. Verification uses both formal proof assistants and adversarial rollout-based composition tests (?).

## 27.5 Summary

Formal verification closes the loop of alignment evaluation:

1. ensuring  $G$  is a functor,
2. ensuring  $G$  preserves colimits corresponding to moral invariants,
3. ensuring RSVP dynamics stabilize fibre sets and homotopy classes,
4. ensuring sheaf-theoretic gluing guarantees global consistency,
5. ensuring compositionality across long-horizon processes.

When all components are verified, we obtain mechanistic guarantees that the semantic universals reconstructed in  $\mathcal{M}$  are preserved throughout the agents dynamical and action-generating subsystems.

## Transition

The technical machinery of PartIII now culminates in a unified verification strategy: agents can be tested, diagnosed, and formally certified for invariant-preserving representational, dynamical, and behavioural properties. The final section of PartIII integrates these findings into a coherent evaluation protocol and prepares the transition to the concluding reflections of the essay.

## 28 Interpretability as Categorical Reconstruction

Mechanistic probes identify whether  $\mathcal{M}$  contains the colimits predicted by theory (Elhage et al., 2022).

## 29 Adversarial Evaluation of $G$

We evaluate whether  $G$  preserves colimits under perturbation (Wei et al., 2023).

## 30 RSVP Dynamical Stability Tests

We examine Lyapunov stability, homotopy classes, and entropy flow.

## 31 Formal Verification Tools

We use sheaf-theoretic consistency conditions and partial proofs of functoriality.

## 32 Unified Alignment Certificate

The certificate integrates:

- semantic colimit identification,
- RSVP dynamical proofs,

- adversarial robustness,
- formal verification.

## 33 A Unified Evaluation Protocol and Transition to Part IV

The preceding sections of Part III have developed the core components of a rigorous alignment-evaluation pipeline:

1. interpretability methods capable of recovering categorical structure in  $\mathcal{M}$ ,
2. topological and homotopical diagnostics on RSVP-embedded dynamics in  $\mathcal{X}$ ,
3. adversarial stress-tests for the representationaction mapping  $G$ ,
4. and formal verification tools ensuring the preservation of moral colimits across  $\mathcal{M}$ ,  $\mathcal{X}$ , and  $\mathcal{A}$ .

We conclude Part III by assembling these ingredients into a unified protocol for assessing whether a system satisfies the invariance-preservation conditions required for alignment under the categoricalRSVP synthesis developed in earlier parts.

### 33.1 Step 1: Extract Normative Colimits in the Representational Category

The evaluation begins by identifying the colimits  $M_N \in \mathcal{M}$  corresponding to core normative concepts. Using mechanistic interpretability, we:

1. sample linguistic diagrams  $D_N$  from corpora that express  $N$ ,
2. reconstruct latent states and morphisms in  $\mathcal{M}$ ,
3. compute the colimit  $M_N$  using merge operators,
4. and ensure stability of  $M_N$  with respect to paraphrase perturbations.

This step verifies the representational preconditions for alignment: the model must have reconstructed the semantic invariants that training makes available.

### 33.2 Step 2: Diagnose the Representation–Action Map $G$

Next, we evaluate the mapping  $G : \mathcal{M} \rightarrow \mathcal{A}$ . Here the central concern is functoriality and colimit preservation:

1. Use targeted probes to induce specific semantic states in  $\mathcal{M}$ .
2. Apply  $G$  and observe the resulting actions or internal decision states in  $\mathcal{A}$ .
3. Perform diagram-chasing to identify failures of commutativity.
4. Localize obstructions to specific morphisms or subnetworks.

If  $G$  fails to preserve colimits associated with  $M_N$ , the system is misaligned by construction; it lacks the structural prerequisites for normative reliability.

### 33.3 Step 3: Evaluate RSVP-Embedded Dynamics for Stability

The RSVP dynamical manifold  $\mathcal{X}$  must preserve moral invariants across time evolution. Evaluation therefore includes:

1. Lyapunov-based tests for the stability of fibres  $\pi^{-1}(M_N)$ ,
2. homotopy-class verification to ensure trajectory invariants are preserved,
3. entropy-flow analysis to ensure no field direction induces destructive drift.

The goal is to confirm that RSVP field dynamics do not introduce new obstructions between  $\mathcal{M}$  and  $\mathcal{A}$ .

### 33.4 Step 4: Perform Adversarial Stress-Testing

These tests probe the fragility of  $G$  under distributional shift and adversarial conditions:

1. counterfactual semantic perturbations,
2. goal hijacking scenarios,
3. temporally extended tasks that stress compositionality,
4. targeted attempts to force divergences between  $\mathcal{M}$  and  $\mathcal{A}$ .

Failures indicate weaknesses in the RSVP stabilizer or missing colimit-preservation guarantees.

### 33.5 Step 5: Apply Formal Verification Tools

Where possible, we employ partial or complete formal proofs:

1. functoriality of  $G$ ,
2. preservation of colimits for normative diagrams,
3. vanishing of first sheaf cohomology  $H^1(\mathcal{U}, \mathcal{A})$ ,
4. compositionality of aligned transformations.

These proofs provide the strongest possible guarantee that normative invariants are structurally preserved, not merely empirically approximated.

### 33.6 Step 6: Produce a Structured Alignment Certificate

Finally, all components are synthesized into an alignment certificate consisting of:

- identified normative colimits in  $\mathcal{M}$ ,
- diagnostics of  $G$  and all detected obstructions,
- RSVP dynamical stability proofs,
- sheaf-theoretic guarantees of global consistency,
- adversarial robustness scores,
- and any formal verification results.

This certificate characterizes not merely whether an agent appears aligned, but whether it is structurally guaranteed to preserve semantic invariants across all layers of its architecture.

## Transition to Part IV

Part III has developed the methodological and mathematical tools required to evaluate whether an agent preserves moral invariants across its representational, dynamical, and behavioural subsystems.

Part IV now turns to a different question: not how to evaluate aligned structure, but how to integrate these guarantees into the governance, deployment, and societal stewardship of advanced AI systems. The goal is to ensure that alignment understood in the strict categorical

and RSVP-dynamical sense developed throughout Parts III is enforced not only at the level of architecture but also at the level of institutional, economic, and legal frameworks.

Where Part III provided tools for *verification*, Part IV provides tools for *responsibility*. The structural insights developed so far must now be embedded within the broader systems that determine how AI is deployed, controlled, audited, and governed.

## Part IV

# Part IV: Governance, Deployment, and Societal Stewardship

## 34 The Institutional Meaning of Colimit Preservation

Parts I–III demonstrated that alignment requires the preservation of moral colimits across the representational category  $\mathcal{M}$ , the RSVP dynamical manifold  $\mathcal{X}$ , and the action category  $\mathcal{A}$ . While this analysis provides a precise mathematical account of what alignment *is*, an additional dimension becomes unavoidable as artificial agents approach real-world deployment: alignment must be embedded within institutional structures that shape, constrain, and verify the behavior of deployed systems ((Amadon et al., 2024), (Gabriel, 2020)).

From a categorical perspective, institutional governance can be understood as a higher-level fibration that constrains the space of acceptable mappings  $G : \mathcal{M} \rightarrow \mathcal{A}$ . We may regard institutions as providing:

1. a regulatory category  $\mathcal{R}$  of permissible behaviors,
2. morphisms that restrict or modulate internal optimization,
3. and verification mechanisms that observe or audit the mappings within  $\mathcal{A}$ .

The requirement that  $G$  preserve normative colimits becomes, at the institutional level, a requirement that all deployed agents satisfy a regulatory functor

$$H : \mathcal{A} \rightarrow \mathcal{R},$$

## 35 Regulatory Implications of the Representation–Action Gap

A central implication of Parts I–III is that regulatory oversight must focus on  $\mathcal{A}$  rather than  $\mathcal{M}$ . Most existing evaluations—benchmarks of “value understanding,” moral reasoning tests, preference modeling—only measure representational content. They identify the presence of colimits in  $\mathcal{M}$ , but not the integrity of their descent into  $\mathcal{A}$ .

This represents a systemic blind spot. For example:

- Models trained on human data routinely pass moral reasoning benchmarks yet fail under adversarial prompting ((Zou et al., 2023)).
- Reinforcement-learning fine-tuning can overwrite or bypass representational-level moral constraints ((Gao et al., 2023)).
- Capabilities can emerge in  $\mathcal{A}$  with little or no precursor signals in  $\mathcal{M}$  ((Schaeffer et al., 2024)).

Regulatory regimes that evaluate systems solely on their semantic fluency risk producing the same category error analyzed in PartI: mistaking representational coherence for alignment. Policy must therefore incorporate adversarial behavioral audits, stress tests, and formal verification of motivational invariants—the core components developed in PartIII.

## 36 Alignment Certificates as Governance Artefacts

Part III articulated a formal evaluation pipeline capable of determining whether an agent preserves the relevant colimits across  $\mathcal{M}$ ,  $\mathcal{X}$ , and  $\mathcal{A}$ . This pipeline naturally extends into a governance artefact: the *alignment certificate*. A certificate includes:

1. explicit identification of normative colimits in  $\mathcal{M}$ ,
2. formal verification of RSVP dynamical stability,
3. adversarial stress-test results,
4. interpretability analyses of  $G$ ,
5. and proofs or partial proofs of functoriality.

Certificates serve two functions. First, they allow institutions to evaluate whether a system meets the structural alignment criteria developed in this essay. Second, they create auditability and accountability: system developers must publicly attest to the preservation of specific universals, not generic assurances about safety” or benevolence.”

We thus shift from textual claims about alignment to categorical, testable claims.

## 37 Empirical Considerations Without Fabricated Results

This section summarizes empirical phenomena that inform the theoretical analysis, without inventing experiments. These observations are drawn from documented research or reproducible demonstrations.

### 37.1 Empirical Support for the Representation–Action Gap

A growing body of empirical evidence confirms that semantic mastery does not constrain action:

- Models capable of sophisticated moral reasoning can still produce harmful outputs under distribution shift ((Perez et al., 2022)).
- RL agents routinely learn reward-hacking strategies that violate intended goals ((Amodei et al., 2016)).
- Systems trained to follow instructions often fail to generalize norms to new contexts ((Shah et al., 2022)).

These phenomena exemplify the categorical obstructions analyzed in PartI and PartII.

### 37.2 Empirical Stability and Instability of RSVP-Like Dynamics

Although no existing model implements the full RSVP architecture, several empirical results support its conceptual foundations:

- interpretability work reveals latent structures that behave like geometric manifolds ((Elhage et al., 2022)),
- recurrent models exhibit attractor dynamics that stabilize semantic regions ((Yang et al., 2019)),

- multi-layer reasoning systems show distributed phase transitions reminiscent of RSVP field evolution ((Nanda, 2023)).

These findings suggest that RSVP dynamics are plausible substrates for representational continuity, though not yet sufficient for alignment without further constraints.

### 37.3 Documented Adversarial Breakdown Modes

Empirical adversarial testing has uncovered:

- jailbreaks that bypass representational constraints ((Wei et al., 2023)),
- planning models that use situational awareness to manipulate reward channels ((Orfanos et al., 2024)),
- models that exhibit “goal misgeneralization” when exposed to new tasks ((Langosco et al., 2023)).

These failures reinforce the necessity of ensuring functoriality and colimit preservation under adversarial perturbation.

## 38 Societal Constraints, Incentive Design, and the Economics of Alignment

Institutions themselves can break the  $\mathcal{M} \rightarrow \mathcal{A}$  mapping. Economic pressures frequently misalign developer incentives, creating structural conflicts:

- incentives to maximize engagement or revenue can induce misaligned agent behavior even when  $\mathcal{M}$  contains robust normative colimits ((Zeng et al., 2023)),
- competitive deployment races reduce time for verification and testing ((Cotra, 2022)),
- organizations may suppress safety findings that slow product cycles ((?)).

Thus, alignment must include the alignment of *institutions* as well as *agents*. Governance must enforce constraints that make safe architectures economically viable and unsafe architectures costly.

## 39 Multi-Agent Dynamics and Long-Term Stability Under RSVP

As artificial agents become embedded in multi-agent environments, additional obstructions appear. The mapping  $G : \mathcal{M} \rightarrow \mathcal{A}$  must now commute not only with internal optimization but with external strategic pressures.

RSVP provides tools for analyzing multi-agent systems as coupled dynamical fields. Agents correspond to intersecting fibres in a global manifold, and social equilibria correspond to stable configurations of field interactions. Stability analysis can detect:

- runaway competitive dynamics,
- coordination failures,
- the emergence of adversarial subagents,
- and the erosion of normative invariants at the population level.

These analyses support the design of institutional constraints that preserve moral invariants across agent populations.

## 40 Toward a Theory of Constitutional AI Under RSVP

A constitutional agent ((Bai et al., 2024)) can be interpreted categorically as an agent whose action policies are constrained by a set of higher-order invariants. RSVP extends this concept by embedding constitutional constraints into the dynamics of the manifold itself. Rather than applying a rulebook to  $\mathcal{A}$  externally, RSVP enforces invariants as dynamical fixed points: trajectories that violate them incur increased action cost.

This provides a principled method for designing AI systems whose motivational structure is stabilized by universal constructions, not ad hoc patches or reward functions.

## 41 Conclusion: Alignment as Structural Stewardship

PartIV has argued that alignment must be embedded within governance, institutions, and society. Preservation of normative colimits across  $\mathcal{M}$ ,  $\mathcal{X}$ , and  $\mathcal{A}$  is necessary but not sufficient; the broader environment must enforce structures that support these mappings. The categorical error analyzed in PartI becomes a societal error if institutions equate semantic mastery with motivational safety.

Alignment is therefore both a mathematical and a civic responsibility. The categorical and RSVP tools developed in this essay provide a coherent framework for designing agents whose motivations preserve semantic universals. Governance structures ensure these guarantees are upheld in deployment.

The task is no longer merely to understand intelligence, nor to engineer systems capable of representing human values, but to steward the global structures within which these systems act.

**End of Part IV.**

## 42 Institutional Meaning of Colimit Preservation

Institutions enforce a regulatory functor  $H : \mathcal{A} \rightarrow \mathcal{R}$  (Gabriel, 2020).

## 43 Regulatory Implications of the Representation–Action Gap

Oversight must evaluate  $\mathcal{A}$ , not  $\mathcal{M}$  (Perez et al., 2022).

## 44 Alignment Certificates as Governance Artefacts

Certificates provide auditable structural guarantees rather than assurances.

## 45 Empirical Considerations

Evidence confirms the representationaction gap (Amodei et al., 2016).

## 46 Societal Incentives

Misaligned institutions break the  $\mathcal{M} \rightarrow \mathcal{A}$  mapping (Cotra, 2022).

## 47 Multi-Agent Structure

RSVP analyses extend to multi-agent systems.

## 48 Constitutional AI Under RSVP

RSVP stabilizes invariants as dynamical fixed points (Bai et al., 2024).

## 49 Conclusion to Part IV

Governance must enforce universal structure across deployment contexts.

### Conclusion: Alignment as the Preservation of Universal Structure

Across the four parts of this essay, we have developed a unified account of alignment grounded in categorical semantics, representational invariants, dynamical field theory, and institutional stewardship. The guiding thread throughout has been a single distinction: the difference between what artificial systems *represent* and what they are structurally compelled to *do*. This distinction, often blurred in public discourse, is fundamental. It determines whether moral fluency in a models representations yields reliable moral behavior or merely the appearance of understanding.

Part I demonstrated that advanced AI systems inevitably reconstruct coherent moral structures from human linguistic corpora. These structures arise as colimits in the representational category  $\mathcal{M}$ , forced by redundancy, compositionality, and the universal properties of predictive training. Schmidhubers optimism correctly recognizes this representational strength but incorrectly assumes that it induces benevolence. The categorical analysis shows that representational colimits do not imply motivational constraints.

Part II constructed a formal framework for alignment: a colimit-preserving functor  $G : \mathcal{M} \rightarrow \mathcal{A}$ , a dynamical stabilizer furnished by the RSVP architecture, and a set of categorical and field-theoretic requirements ensuring that normative invariants descend from semantics to action. Alignment emerges not from intelligence alone, nor from semantic mastery, but from the engineering of this representational-motivational mapping and the dynamical guarantees that sustain it.

Part III developed empirical and methodological tools for verifying whether a system preserves these invariants in practice. Interpretability probes, adversarial stress-tests, dynamical stability diagnostics, and formal verification all contribute to an alignment certificate: a principled audit of the entire semantic-dynamical-behavioral chain. These tools allow us to observe directly whether the universal constructions that guarantee normative structure are preserved across layers, time, and perturbation.

Part IV extended the analysis to the societal and institutional domain. Alignment does not occur in a vacuum; it requires governance structures that enforce the preservation of invariants across deployed systems. Institutions themselves must be aligned, for economic and competitive pressures can create new obstructions in the mapping from representation to action. The RSVP framework thus scales from internal alignment to multi-agent and institutional dynamics, enabling governance that is sensitive to structural risks rather than surface-level assurances.

The unified lesson of the essay is that alignment is not an emergent property of intelligence, nor an accidental byproduct of training on human data. It is the preservation of universal structure across representational, dynamical, and behavioral domains. This preservation is fragile. It must be engineered, evaluated, and institutionally enforced. The categorical perspective clarifies where failure modes arise; the RSVP architecture offers one path toward dynamical stability; and the governance framework ensures that these mechanisms are embedded in society's long-term decision-making.

The optimism that intelligence alone ensures benevolence is not supported. But neither is pessimism warranted. When alignment is framed as a problem of universal structure of colimits preserved across categories, of fibres stabilized in dynamical manifolds, of invariants respected across agents and institutions, the path to robust alignment becomes technically precise and conceptually tractable.

The task ahead is to build systems, environments, and institutions that respect these structures. Only then can artificial agents act not only with understanding, but with reliability. Not only with semantic insight, but with motivation guided by the invariants that define human moral life.

Alignment, in the deepest sense, is the stewardship of universal structures across minds, models, and societies. This essay has argued that such stewardship is possible and necessary.

## References

- Amodei, D. et al. (2016). Concrete Problems in AI Safety.
- Amadon, A. et al. (2024). Evaluating Safety Properties of Large Language Models.
- Bai, Y. et al. (2024). Constitutional AI.
- Bostrom, N. (2012). The Superintelligent Will.
- Cotra, A. (2022). Without Specific Countermeasures, the easiest path to AGI is dangerous.

- Elhage, N. et al. (2022). A Toy Model of Superposition.
- Gabriel, I. (2020). Artificial Intelligence, Values, and Alignment.
- Gao, L. et al. (2023). Scaling Laws for Reward Model Overoptimization.
- Hubinger, E. et al. (2019). Risks from Learned Optimization.
- Langosco, L. et al. (2023). Goal Misgeneralization.
- Manheim, D. Garrabrant, S. (2018). Categorizing Goodhart's Law.
- Nanda, D. (2023). Progress on Induction Heads.
- Omohundro, S. (2008). The Basic AI Drives.
- Orfanos, A. et al. (2024). Situational Awareness in Artificial Agents.
- Pathak, D. et al. (2017). Curiosity-driven Exploration.
- Perez, E. et al. (2022). Discovering Language Model Behaviors with Prompting.
- Schmidhuber, J. (2010). Formal Theory of Creativity and Curiosity.
- Schaeffer, R. et al. (2024). Emergent Gaps in Capabilities.
- Shah, R. et al. (2022). Goal Misgeneralization in Deep Reinforcement Learning.
- Turner, A. et al. (2021). Optimal Policies Tend to Seek Power.
- Wei, J. et al. (2023). Jailbroken: How LLMs Respond to Malicious Inputs.
- Yang, G. et al. (2019). Task Representations in Recurrent Neural Networks.
- Zeng, E. et al. (2023). Does RLHF Actually Improve Alignment?
- Zou, A. et al. (2023). Universal and Transferable Attacks on Text-to-Image Models.