

# The Imagination Machine 0: Orientation to a Framework for Embedded Epistemic Systems

Mark Tracy

March 2026

Close your eyes. Imagine your body positioned exactly how it is. Now imagine a bubble around your body. Now realize that you are the surrounding darkness—the outer bubble around that bubble. You are the vanishing point of Perspective. A View from Somewhere visible Nowhere.

---

## Abstract

The Imagination Machine series develops a framework for understanding how knowledge arises in systems embedded within their environment. Because such systems have no access to an external vantage point, knowledge cannot be defined as correspondence with an independently accessible world. Instead, knowledge must be understood operationally as the stabilization of representations through interaction with an environment.

Across the series a common architecture appears repeatedly. Observations generate representations of relational structure; these representations are compressed to retain invariants while discarding detail; extension operations generate predictions of missing structure; and prediction error drives subsequent updates. This recursive cycle of observation, compression, extension, and update forms the core mechanism of the Imagination Machine.

The purpose of the present document is to orient the reader to the series as a whole. The individual papers may be read independently, but they collectively describe a layered architecture for embedded epistemic systems. Early papers establish the epistemic foundations of the framework, later papers explore its structural manifestations across domains such as analogy, institutional learning, and symbolic representation, and subsequent papers develop computational realizations and philosophical implications for scientific knowledge.

## 1 Introduction

The Imagination Machine series investigates how an epistemic system embedded within the world can construct coherent representations of that world.

Traditional epistemology often treats knowledge as correspondence between a representation and an independently accessible reality. Embedded systems, however, have no such external vantage

point. They interact only with their observational surface. Consequently, knowledge must be defined in operational terms: representations are evaluated by their capacity to generate coherent predictions and to remain stable under continued interaction with the environment.

Across the papers of this series a common architectural motif emerges. Learning proceeds through a recursive cycle in which an agent observes its environment, constructs a representation of relational structure, compresses that representation to retain invariant features, extends the compressed structure to predict missing relations, and updates its representation in response to prediction error.

This architecture will be referred to as the *imagination machine*. The individual papers explore how this mechanism manifests across several domains, including cognition, analogy, institutional knowledge production, symbolic representation, computational learning systems, and the methodology of science.

## 2 The Core Epistemic Loop

The central operation of the framework can be summarized as the following cycle.

1. An agent observes data generated by interaction with an environment.
2. Observations are organized into a representation capturing relational structure.
3. The representation is compressed so that invariant relations are retained while redundant detail is discarded.
4. The compressed representation is extended through prediction of missing relations or future states.
5. Prediction error generated by subsequent observations updates the representation.

Repeated execution of this loop gradually stabilizes representations that capture persistent relational structure in the environment. Such stabilized structures function operationally as knowledge.

This perspective resonates with several research traditions in which learning is understood as a dynamical feedback process. Early cybernetic work emphasized the centrality of feedback loops in adaptive systems [14, 1]. More recent work in neuroscience proposes predictive processing models in which perception and cognition arise through the minimization of prediction error [4, 3]. Reinforcement learning frameworks likewise describe agents that iteratively update internal models based on interaction with their environment [13].

The framework also bears philosophical affinity with Karl Popper's conception of knowledge growth through conjecture and refutation [10, 11, 12]. In Popper's view, scientific theories function as hypotheses that generate testable predictions; empirical feedback then eliminates those that fail to withstand critical testing.

Similarly, evolutionary accounts of knowledge growth have emphasized processes of variation and selective retention. Campbell proposed that scientific and cognitive development proceeds through cycles of hypothesis generation and error elimination [2]. Within the present framework, extension operations generate candidate structural hypotheses, while prediction error functions as a mechanism of selective elimination guiding representational revision.

### 3 Representation and Closure

A central philosophical challenge for embedded epistemic systems is that representation necessarily involves the imposition of conceptual boundaries upon a world that cannot be accessed independently of those boundaries.

Hilary Lawson has argued that all representation involves acts of closure through which distinctions are drawn and stabilized [8]. In Lawson’s account, knowledge arises through the construction of frameworks that impose structure upon experience while remaining open to revision.

The compression operations described in the Imagination Machine framework may be interpreted as formal mechanisms for generating such closures. Observational data are grouped into equivalence classes that preserve selected relational invariants while discarding other distinctions. These closures stabilize representations sufficiently to support prediction and reasoning, while the extension and update stages of the epistemic loop allow those closures to be revised in response to new evidence.

### 4 A Layered Architecture

Although the papers in the series address diverse domains, they can be viewed as exploring different layers of a single architecture.

- **Epistemic Foundations.** Early papers examine the situation of an embedded observer and introduce the inference–implication loop through which world models stabilize.
- **Dynamical Learning Systems.** Subsequent work develops agent–environment interaction models in which predictive agents recover latent structure from observational data.
- **Structural Reasoning.** Further papers examine mechanisms such as analogy, abstraction, and simplicial completion that enable reasoning systems to generate hypotheses about unseen relations. Classical theories of analogy have long emphasized the role of structure-preserving mappings between domains [6].
- **Institutional Learning.** The framework is extended to communities of interacting agents in which dialogue, compression, and feedback produce evolving institutional knowledge. The dynamics of scientific communities have been studied extensively in the philosophy and sociology of science [7].
- **Symbolic Representation.** Later work shows how conceptual structures may be externalized into symbolic artifacts and interpreted through categorical transformations preserving relational invariants. Category theory provides a natural language for reasoning about structure-preserving mappings between systems [9, 5].
- **Computational Realization.** The architecture is implemented as a learning system whose world model is a dynamically updated knowledge graph interacting with an open textual environment.
- **Philosophy of Science.** Finally, the framework is used to interpret scientific knowledge itself as the stabilization of relational invariants under compression of observational data.

Taken together, these layers suggest that the same epistemic mechanism may operate across multiple levels of organization, from individual cognition to collective scientific practice.

## 5 Relation to Existing Traditions

The architecture described here bears resemblance to several established research traditions.

Cybernetics emphasized feedback and control as fundamental principles of adaptive systems [14, 1]. Predictive processing models in cognitive science interpret perception and cognition as hierarchical processes in which prediction error drives model revision [4, 3]. Reinforcement learning describes agents that iteratively update policies and value estimates through environmental feedback [13].

Karl Popper’s philosophy of science emphasized the iterative interaction between conjecture and refutation as the mechanism by which knowledge grows [10, 11, 12]. The present framework can be viewed as providing a structural and computational interpretation of this dynamic within embedded epistemic systems.

In parallel, category theory has increasingly been used to formalize the structure of learning systems and compositional models of knowledge [5]. These approaches provide mathematical tools for describing how representations preserve relational invariants across transformations.

## 6 Reading the Series

The papers of the Imagination Machine series may be read independently, but they collectively describe different aspects of the same architecture. Later papers often reinterpret or instantiate principles introduced earlier in the series.

Readers interested primarily in philosophical questions may focus on the early papers concerning epistemic closure and representation. Readers interested in computational architectures may focus on the later papers describing knowledge graph learning systems and experimental environments.

The present document serves simply as an orientation to the series as a whole.

## 7 Conclusion

The Imagination Machine series explores how embedded systems construct representations capable of supporting prediction, reasoning, and coordinated action. Across the diverse domains examined in the series, a common mechanism appears: knowledge arises through recursive cycles of observation, compression, extension, and update.

The aim of the project is both philosophical and practical. Philosophically, it seeks to clarify how knowledge can arise without appeal to an external vantage point. Practically, it proposes an architectural framework that may guide the construction of artificial systems capable of discovering relational structure within complex environments.

## References

- [1] W. Ross Ashby. *An Introduction to Cybernetics*. Chapman & Hall, 1956.
- [2] Donald T. Campbell. Blind Variation and Selective Retention in Creative Thought as in Other Knowledge Processes. *Psychological Review*, 67(6):380–400, 1960.
- [3] Andy Clark. *Surfing Uncertainty: Prediction, Action, and the Embodied Mind*. Oxford University Press, 2016.

- [4] Karl Friston. The free-energy principle: a unified brain theory? *Nature Reviews Neuroscience*, 11(2):127–138, 2010.
- [5] Brendan Fong and David I. Spivak. *Seven Sketches in Compositionality: An Invitation to Applied Category Theory*. Cambridge University Press, 2019.
- [6] Dedre Gentner. Structure-mapping: A theoretical framework for analogy. *Cognitive Science*, 7(2):155–170, 1983.
- [7] Thomas S. Kuhn. *The Structure of Scientific Revolutions*. University of Chicago Press, 1962.
- [8] Hilary Lawson. *Closure: A Story of Everything*. Routledge, 2001.
- [9] Saunders Mac Lane. *Categories for the Working Mathematician*. Springer, 2nd edition, 1998.
- [10] Karl R. Popper. *The Logic of Scientific Discovery*. Hutchinson, 1959.
- [11] Karl R. Popper. *Conjectures and Refutations: The Growth of Scientific Knowledge*. Routledge, 1963.
- [12] Karl R. Popper. *Objective Knowledge: An Evolutionary Approach*. Oxford University Press, 1972.
- [13] Richard S. Sutton and Andrew G. Barto. *Reinforcement Learning: An Introduction*. MIT Press, 2nd edition, 2018.
- [14] Norbert Wiener. *Cybernetics: Or Control and Communication in the Animal and the Machine*. MIT Press, 1948.

# The Imagination Machine I: A View from Somewhere

## Epistemic Closure, Physical Law, and Entropy Embedded in a Block Universe

Mark Tracy  
Boston University  
mrktracy@bu.edu

### Abstract

This paper develops a minimal formal framework for epistemology under the constraint that epistemic systems are embedded within the world they attempt to model. Because such systems lack access to an external vantage point, knowledge cannot be defined by correspondence with an independently accessible reality. Instead, epistemic coherence must arise from internal structural consistency.

Observations generate world models through an inference map, while world models generate canonical observational profiles through an implication map. Together these maps form an inference–implication loop that induces an operator on model space. Self-consistent world models appear as fixed points of this operator: models whose own implied observational profiles, when reinterpreted through inference, reproduce the models themselves. Each model therefore acts as a compression of the observation space, inducing a classifier and a corresponding quotient representation of observations.

A key structural feature of the framework is that classifiers themselves belong to the observation space. This follows from the conditions of self-representation: any system capable of epistemic reasoning must be able to encounter and revise its own acts of classification. As a result, the evaluative processes that guide model selection—valuation and will—also appear as observable elements subject to the same representational compression.

Within a given model, empirical regularities emerge as relational invariants in the induced quotient space, while entropy arises as a measure-theoretic quantity associated with the same compressive structure. The framework therefore characterizes scientific theories as stable representational compressions of observational structure for agents embedded within the environments they model.

## 1 Introduction

Embedded epistemic systems cannot access the universe from outside. Observations, models, classifiers, and their relations therefore exist as structures within the same universe. No external vantage point is available from which to define correspondence between representation and world.

The guiding constraint is:

An embedded epistemic system can at most classify the ways in which it classifies the world, within the world itself.

Rather than describing temporal learning, we treat the universe as a single relational structure containing observations, models, and consistency relations between them. Within such a framework

coherence must be defined internally, as the closure of the inference–implication loop rather than as external correspondence.

This paper forms the first part of a four-paper series titled *The Imagination Machine*. The present paper develops the formal epistemic framework for embedded observers and the structure of representational closure. Companion papers develop complementary aspects of the framework: *The Imagination Machine II: Systems* analyzes agent–environment representational dynamics, *The Imagination Machine III: Toy Model of Predictive Classification* provides a minimal computational environment in which predictive agents recover relational invariants, and *The Imagination Machine IV: Institutional Intelligence* examines how such epistemic processes extend across communities and institutions.

This position is closest in spirit to Hilary Lawson’s closure theory of the world. Lawson argues that openness—raw, unstructured reality—is fixed as “something” only through interventions he calls closures, and that no closure fully captures the openness beneath it. The present framework formalizes a version of this picture. The inference–implication loop is the closure mechanism; the fixed points of the operator it induces are the stable closures; a quotient space  $Q_w$  is the closed texture through which an embedded system encounters the world under the model  $w$ . The crucial point is that what a model implies is not best understood as a single isolated observational consequence, but as a canonical observational profile internal to that closure: a structured way the world shows up for a life situated within the model.

But the framework adds something to Lawson’s account that his descriptive language leaves implicit: the acts of will and valuation that select among possible closures are not external to the representational structure. Because classifiers are themselves observations—for reasons derived in Section 3 rather than merely asserted—valuation is interior to the system it animates. This is the structural heart of the paper.

Two clarifications are important at the outset. First, the framework is not a form of coherentism in which any internally consistent system of representations counts as knowledge. The structure of observations within the universe constrains admissible models through the probability measure introduced below. Closure of the inference–implication loop occurs only relative to this observational structure. Second, the framework does not deny the existence of an external world. It instead observes that embedded epistemic systems cannot compare representations with that world directly. The problem addressed here is therefore structural rather than metaphysical.

A further clarification concerns model-relativity. Different self-consistent models may in principle induce different quotient spaces and therefore different families of laws. This does not imply arbitrariness. Models must compress the same observational distribution and remain stable under their own implications. The resulting plurality, if it occurs, is constrained plurality.

The aim of the framework is not to replace empirical science or traditional epistemology, but to describe the structural constraints under which an epistemic system embedded within the universe must operate.

A word on what the framework does and does not claim. The formal architecture precisely locates three problems that resist full resolution from within any closure: the problem of will, the problem of distinguishing genuine from merely apparent epistemic openness, and the problem of the criterion by which a system recognises new observations as demanding refinement. The paper argues that locating these problems with formal precision is itself a contribution—that a framework which shows exactly where explanation runs out is preferable to one that conceals those limits behind descriptive fluency.

## 2 Relation to Existing Approaches

The framework developed here sits at the intersection of several existing lines of research, while differing from each in its formal treatment of embeddedness, representational closure, and model-relative structure.

Most directly, it formalises central commitments of Hilary Lawson’s closure theory. Lawson argues that the world as encountered is always a world fixed by closure, that openness underlies and escapes every closure, and that the question of which closures to adopt is therefore irreducibly evaluative (Lawson, 2001). The present framework gives these claims a precise structural expression: the inference–implication loop is the closure mechanism,  $\mathcal{W}^*$  is the space of stable closures, the quotient space is the closed texture, and the inclusion  $C \subseteq D$  is the formal statement that evaluation is interior to the representational structure rather than prior to it. The analysis of institutions and refinement extends this picture by showing that the evaluative dimension of closure is not merely a feature of individual systems but is transmitted, compressed, and potentially lost across generations.

The account also bears comparison with predictive and Bayesian approaches in contemporary philosophy of mind and cognitive science. Predictive processing models treat cognition as the continuous generation of predictions that are compared with incoming sensory signals, with discrepancies driving model revision (Clark, 2016; Friston, 2010). The inference–implication loop introduced here has a related structure: observations generate models through the inference map  $F$ , while models generate observational implications through the map  $g$ . However, the present framework differs from predictive-processing accounts in one crucial respect: both observations and models are treated as structures internal to a single universe rather than as elements of an external inference problem. The framework therefore addresses not only how models are updated, but how coherence is to be defined for an epistemic system that has no access to an external vantage point. A minimal computational environment in which predictive agents recover relational invariants from structured observations is developed in *The Imagination Machine III: Toy Model of Predictive Classification*.

In philosophy of science, the view developed here is also close in spirit to structural realism. Structural realists argue that scientific knowledge concerns the relational structure of the world rather than the intrinsic nature of unobservable entities (Worrall, 1989; Ladyman, 1998). In the present framework, relational structure appears in an explicitly model-relative mathematical form. Each self-consistent world model  $w$  induces a classifier  $\pi_w : D \rightarrow Z_w$  that partitions the observation space, and empirical regularities arise as relational invariants in the quotient space  $Q_w = D/\sim_w$  determined by that partition. What embedded observers identify as physical laws are therefore relational structures within a representational quotient induced by the model. In this sense the framework provides a formal account of how structural knowledge arises from representational compression.

This model-relative account of law also bears comparison with relational approaches in physics. Rovelli’s relational quantum mechanics emphasises that physical properties are defined relative to interactions rather than to absolute external states (Rovelli, 1996). Physical laws in the present framework are likewise relational invariants, though the present argument grounds their model-relativity in epistemological rather than specifically physical considerations.

The entropy measure introduced here connects the framework to statistical mechanics and information theory. Shannon introduced entropy as a logarithmic measure of expected surprisal associated with a probability distribution (Shannon, 1948). Jaynes later interpreted statistical mechanics as inference over probability distributions subject to informational constraints (Jaynes, 1957). The present framework recovers entropy as a consequence of representational compression rather than positing it as primitive: the classifier  $\pi_w$  partitions the observation space into equivalence classes, and the entropy  $H(w)$  measures the expected surprisal of those classes. The

framework does not claim identity between this quantity and thermodynamic entropy; rather, it argues for a structural convergence between them, grounded in their shared dependence on the partitioning of a probability space.

In biology and systems theory, Maturana and Varela described cognition as arising from operational closure within self-referential systems (Maturana and Varela, 1980; Varela et al., 1991). The self-consistency condition  $T(w) = w$  is a formal analogue of operational closure, with the additional feature that the closed system contains its own evaluative structure as classified content. Read in the present terms, closure is reproduced not merely from isolated outputs but from the structured observational profile a model makes possible from within. The dynamical structure of agent–environment interaction underlying such representational frameworks is analyzed in *The Imagination Machine II: Systems*.

Finally, the social extension of the framework places it in conversation with social epistemology. Longino and Kitcher have both argued, in different ways, that knowledge is constitutively social and that the norms governing inquiry are sustained and revised by communities rather than by isolated individuals (Longino, 1990; Kitcher, 1993). The institutional analysis developed here is consistent with this emphasis while grounding it in the formal architecture of the framework. The distinction between generative and compressed inheritance corresponds, at the social level, to the difference between communities that transmit the capacity for inquiry and communities that merely conserve its prior outputs. A fuller treatment of institutional knowledge generation and transmission is developed in *The Imagination Machine IV: Institutional Intelligence*.

### 3 The Block Universe and the Derivation of $C \subseteq D$

Let  $\Omega$  denote the universe. Define the following subsets:

$$D \subseteq \Omega \quad (\text{the set of observations})$$

$$\mathcal{W} \subseteq \Omega \quad (\text{the set of world models})$$

$$C \subseteq \Omega \quad (\text{the set of classifiers}).$$

We argue for, rather than merely stipulate, the inclusion

$$C \subseteq D \subseteq \Omega.$$

The argument proceeds from the conditions of self-aware representation. Consider what distinguishes an epistemic system—a genuine subject—from a mere transducer. A thermostat classifies temperature, but its classification is not available to it as an object of experience. It cannot encounter its own sorting activity as something that could have been otherwise. An epistemic system, by contrast, is one whose classificatory acts are themselves accessible to it: it can attend to how it is attending, sort its ways of sorting, and in principle revise the dispositions that govern its encounter with the world.

This reflexive accessibility is not an optional feature added to an otherwise complete epistemic system. It is the condition that makes a system epistemic in the first place. A system that cannot encounter its own classifiers cannot recognise itself as one possible closure among others, cannot doubt its own representations, and therefore cannot be said to know in any sense that involves the distinction between appearance and reality. Cartesian doubt is only possible for a system whose classificatory acts are elements of its observation space.

The inclusion  $C \subseteq D$  is therefore a transcendental condition: any system that satisfies the minimal criterion for being an epistemic subject must satisfy it. The formal apparatus of this paper applies precisely to systems meeting that criterion.

**Remark 1** (Reflexivity Without Vicious Regress). *The condition  $C \subseteq D$  means that the system can classify its own classifiers. One might worry that classifying a classifier requires a further classifier, which requires a further classifier still, generating an infinite regress. This regress does not arise in the block universe framing because that framing is atemporal: all observations, including observations of classifiers, are simultaneous elements of the single relational structure  $\Omega$ . The self-consistency condition  $T(w) = w$ , developed in Section 10, is a fixed-point condition rather than a termination condition. What matters is not that the regress terminates in a foundation but that the loop closes on a stable fixed point.*

## 4 World Models and Classification

Each world model  $w \in \mathcal{W}$  induces a classifier

$$\pi_w : D \rightarrow Z_w$$

where  $Z_w \subseteq D$ . Thus a model compresses observations by mapping them to representative observational states. Because  $C \subseteq D$ , the domain of  $\pi_w$  includes classifiers themselves. A world model therefore classifies not only raw observational content but also the evaluative and selective dispositions of the system that holds it.

**Remark 2** (Representational Witness). *The condition  $Z_w \subseteq D$  ensures that every abstract class induced by  $\pi_w$  is instantiated by at least one observational state. The representative is not assumed to be unique or privileged; it merely witnesses the existence of the class.*

**Definition 1** (Model-Induced Equivalence Relation). *For  $d_1, d_2 \in D$  define*

$$d_1 \sim_w d_2 \quad \text{iff} \quad \pi_w(d_1) = \pi_w(d_2).$$

**Definition 2** (Equivalence Class). *For  $d \in D$ , define*

$$[d]_w = \{d' \in D \mid \pi_w(d') = \pi_w(d)\}.$$

The classifier therefore induces a partition of the observation space. When  $d$  is itself a classifier—that is, when  $d \in C$ —its equivalence class  $[d]_w$  groups together all observational states that the world model treats as equivalent ways of sorting the world. Different valuations may thus collapse into the same equivalence class under a given model, or be distinguished by a more refined one.

## 5 Valuation and Will as Interior Observations

The inclusion  $C \subseteq D$  has a consequence that deserves explicit statement before the formal development continues.

Valuation—the assignment of significance to observations—and will—the selective pressure that drives a system toward one closure rather than another—are traditionally treated as standing outside epistemological frameworks. They appear as boundary conditions: given that a system values certain outcomes, what can it know? The present framework does not dissolve this exterior status so much as restate it with formal precision.

If the acts by which a system evaluates and selects are themselves classifiers, and if classifiers are observations, then valuation and will are elements of  $D$ . They are subject to the same measure  $\mu_D$ , the same quotient structure induced by  $\pi_w$ , and the same representational compression as any

other observation. A self-consistent world model does not merely organise perceptual content; it also classifies the evaluative structure through which the system engages the world.

This does not reduce will to mechanism, nor does it claim to resolve the problem of agency. What it establishes is more modest and more precise: will appears within  $D$ , is partially compressed by every model, and yet is not exhausted by any compression. This is not because will is supernatural or causally unconstrained, but because it is the condition under which the world becomes held as anything at all—the potentiality that precedes and exceeds any particular representation of it. The formal loop determines the space of stable closures  $\mathcal{W}^*$ , but the selection of a particular element from that space is precisely what the framework locates as irreducible. Willing is not explained away; it is what remains when the inference–implication loop has done everything it can do—not a gap in the framework, but the condition the framework must include without being able to absorb.

Metaphysical closure is therefore prevented not by any deficiency of the representational apparatus, but by what the apparatus must include: the very acts of valuation that animate it. The framework’s contribution here is not resolution but precision—knowing exactly where the limit lies is different from not knowing where to look.

## 6 Statistical Structure

Assume the observation space carries a probability structure

$$(D, \Sigma_D, \mu_D)$$

where  $\Sigma_D$  is a  $\sigma$ -algebra and  $\mu_D$  a probability measure.

The measure  $\mu_D$  is the principal way in which observational structure constrains closure. It prevents the framework from collapsing into the view that any self-supporting classificatory system is epistemically on a par with any other. Models partition one and the same observational space, and the measure of those partitions is not up to the model alone.

**Proposition 1** (Measurable Partition). *If each  $\pi_w$  is measurable and  $Z_w$  carries a  $\sigma$ -algebra in which singletons are measurable, then the equivalence classes  $[d]_w$  form a measurable partition of  $(D, \Sigma_D, \mu_D)$ .*

*Proof.* Since  $\pi_w$  is measurable, the preimage of each singleton in  $Z_w$  lies in  $\Sigma_D$ . But

$$[d]_w = \pi_w^{-1}(\{\pi_w(d)\}),$$

so each equivalence class is measurable. The classes partition  $D$  by construction. □

**Lemma 1** (Probability of Classes). *For any model  $w$ ,*

$$\sum_{[d]_w \in Q_w} \mu_D([d]_w) = 1.$$

*Proof.* The sets  $[d]_w$  form a measurable partition of  $D$ . Since  $\mu_D$  is a probability measure on  $D$ , the total measure of the partition equals  $\mu_D(D) = 1$ . □

**Remark 3** (Origin and Calibration of the Observational Measure). *The probability measure  $\mu_D$  represents the empirical distribution of observations across the observation space  $D$ . Conceptually it may be understood in several compatible ways.*

First, it may represent the long-run frequency distribution of observations generated across the ensemble of observers embedded in  $\Omega$ . Since  $D$  contains the observations of all observers, the measure aggregates the empirical structure encountered throughout the block universe. This need not be understood as arbitrary sampling from an undifferentiated flux. In many natural settings, observers are embedded in environments structured by stable but incommensurate dynamical cycles whose relative phases continually drift without exact repetition. Under such conditions, sequential observation repeatedly samples a structured signal that is neither perfectly periodic nor wholly unconstrained. The result is an empirical distribution over observational states: enough recurrence for stable frequencies to emerge, enough phase drift for novelty to persist. On this view,  $\mu_D$  arises from the statistical structure induced by the dynamical environment in which embedded observers occur.

Second,  $\mu_D$  may be interpreted inferentially. Following the information-theoretic programme associated with Jaynes, probability distributions can be understood as representations of incomplete knowledge subject to constraints. Under this interpretation  $\mu_D$  encodes the informational constraints under which an embedded epistemic system performs inference.

These two readings are compatible. A structured observational environment gives rise to stable empirical frequencies, while inference treats those frequencies as constraints on admissible closure. The framework therefore does not require commitment to probability as either purely objective or purely epistemic. What matters structurally is that all world models compress the same observational distribution. This shared measure prevents the space of self-consistent closures from collapsing into arbitrary coherent systems.

However, the compatibility of these two readings is itself a condition that can be satisfied or failed. Call this condition calibration: the alignment between a system's inferential  $\mu_D$ —the weights it brings to inference—and the actual empirical distribution of observations in its environment. Calibration is an achievement rather than a default. It can fail in at least two ways. A system may be miscalibrated: its inferential weights systematically diverge from actual observational frequencies, producing self-consistent closures that are stable relative to the wrong measure. Such a system refines willingly and generates genuine laws—but laws of a distribution that does not reflect the environment it inhabits. Miscalibration is therefore distinct from both dogmatism and ordinary error: the closure is open to refinement, yet refinement proceeds against a distorted image of the world. Calibration can also fail under distributional shift: in genuinely novel environments, a system's inferential  $\mu_D$  is an extrapolation from past frequencies into regions where those frequencies no longer apply. The alignment between the two readings breaks down precisely where epistemic pressure is greatest.

Miscalibration thus constitutes a third structural location of epistemic risk, alongside dogmatic refusal to refine and the irreducible remainder of will. The framework diagnoses all three as failures at different levels of the hierarchy  $(F, g) \rightarrow T \rightarrow \mathcal{W}^* \rightarrow \pi_w \rightarrow Q_w \rightarrow R_w$ : dogmatism is a failure at the level of  $(F, g)$ ; miscalibration is a failure at the level of  $\mu_D$  itself, prior to the construction of any particular closure; and will names the underdetermination that persists even when both are functioning well.

## 7 Representational Quotient

Each model induces a quotient space

$$Q_w = D / \sim_w .$$

The elements of  $Q_w$  represent observational states modulo the classification performed by the model. This is the closed texture through which the world is encountered: not the world as it is prior to closure, but the world as fixed by the representational intervention of  $\pi_w$ .

Because  $C \subseteq D$ , the quotient space  $Q_w$  contains equivalence classes of classifiers alongside equivalence classes of other observations. The closed texture therefore includes, within itself, the compressed image of the evaluative structure of the system that produced it.

To collect these model-relative quotient spaces into a single ambient codomain, define

$$Q := \bigsqcup_{w \in \mathcal{W}} Q_w,$$

the disjoint union of all quotient spaces induced by world models in  $\mathcal{W}$ . Thus each  $Q_w$  is canonically embedded in  $Q$ , while remaining distinguished from  $Q_{w'}$  when  $w \neq w'$ .

## 8 Implication

For each model  $w \in \mathcal{W}$ , let

$$\Gamma_w$$

denote the set of canonical observational profiles induced by  $w$ , where each such profile is structured in the quotient space  $Q_w$ . These profiles are not single isolated observations, but model-relative patterns of observational life: structured ways the world becomes legible from within the closure determined by  $w$ .

Define the ambient profile space

$$\Gamma := \bigsqcup_{w \in \mathcal{W}} \Gamma_w.$$

World models produce canonical observational profiles through a map

$$g : \mathcal{W} \rightarrow \Gamma$$

such that, for each model  $w \in \mathcal{W}$ ,

$$g(w) \in \Gamma_w \subseteq \Gamma.$$

Thus  $g$  assigns to each world model a model-relative observational profile internal to the closure induced by that very model.

## 9 Inference

Canonical observational profiles generate world models through

$$F : \Gamma \rightarrow \mathcal{W}.$$

## 10 The Consistency Loop

The system is governed by the pair of maps

$$\Gamma \xrightarrow{F} \mathcal{W} \xrightarrow{g} \Gamma.$$

Define the induced operator

$$T = F \circ g : \mathcal{W} \rightarrow \mathcal{W}.$$

**Definition 3** (Self-Consistent World Model). *A model  $w$  is self-consistent if  $T(w) = w$ .*

Define

$$\mathcal{W}^* = \{w \in \mathcal{W} \mid T(w) = w\}.$$

Self-consistent models reproduce themselves when inference is applied to their own implied observational profiles. In Lawson’s terms, they are stable closures: the system’s representational intervention reproduces itself under the loop of implication and re-inference. More precisely, a self-consistent model is one whose implied observational profile, when re-submitted to inference, regenerates the same model.

A natural worry is that the fixed-point condition may be too weak: if the maps  $F$  and  $g$  are unconstrained, perhaps trivial fixed points proliferate. That worry is legitimate in the abstract. The framework does not claim that every fixed point is equally significant. Its claim is that any epistemically admissible closure must at least satisfy this condition, and that the observational measure  $\mu_D$  together with the refinement structure developed below provides a basis for distinguishing empty stability from informative stability.

**Remark 4** (Existence of Fixed Points). *The framework defines epistemically admissible closures as fixed points of the operator  $T = F \circ g$ . The formal development does not assume that fixed points exist for arbitrary choices of  $F$  and  $g$ . Rather, the framework identifies a structural condition that any stable closure must satisfy if it exists.*

*In many natural settings fixed points arise under mild assumptions. For example, if  $\mathcal{W}$  is endowed with a compact topology and  $T$  is continuous, Schauder’s fixed-point theorem ensures the existence of at least one  $w^* \in \mathcal{W}$  such that  $T(w^*) = w^*$ .*

*In algorithmic or statistical settings the operator may instead be interpreted as an iterative update rule whose empirical convergence defines the effective closure.*

*The present framework therefore does not claim that all conceivable inference–implication structures admit stable closures. It instead provides the formal characterisation that any such closure must satisfy when it occurs. In this sense the framework is generative: it specifies meta-structural constraints that a world model must satisfy in order to reproduce itself under the inference–implication loop.*

**Remark 5** (Plurality of Stable Closures). *Nothing in the framework requires  $\mathcal{W}^*$  to be a singleton. Multiple incompatible self-consistent models may coexist as elements of  $\mathcal{W}^*$ . This plurality is not a defect. It corresponds directly to Lawson’s insistence that no single closure is metaphysically privileged. The operator  $T$  determines the space of possible stable closures, but it does not determine which element of  $\mathcal{W}^*$  is instantiated.*

## 11 Relational Structure

For each integer  $i \geq 1$  define

$$K_i(Q_w) = Q_w^i,$$

the  $i$ -fold Cartesian product of  $Q_w$  with itself. Thus an element of  $K_i(Q_w)$  is an ordered tuple

$$\tau = ([d_1]_w, [d_2]_w, \dots, [d_i]_w).$$

Let

$$K(Q_w) = \bigsqcup_{i=1}^{\infty} K_i(Q_w) = \bigsqcup_{i=1}^{\infty} Q_w^i,$$

the disjoint union of all finite Cartesian powers of  $Q_w$ , collecting relational tuples of every arity into a single set.

Elements of  $K(Q_w)$  represent finite relational configurations among equivalence classes of observations, together with their arities. A relational classifier

$$R_w : K(Q_w) \rightarrow Q_w$$

assigns canonical relational consequences within the quotient space.

## 12 Physical Law

**Definition 4** (Relational Equivalence). *For  $\tau_1, \tau_2 \in K(Q_w)$  define*

$$\tau_1 \sim_{R_w} \tau_2 \quad \text{iff} \quad R_w(\tau_1) = R_w(\tau_2).$$

**Definition 5** (Physical Law). *A physical law under a model  $w$  is a relational equivalence class*

$$L = [\tau]_{R_w}$$

for some  $\tau \in K(Q_w)$ .

Physical laws appear as relational structures within the quotient representation induced by a self-consistent world model. They are stable patterns in the closed texture, not features of an independently accessible world. Different elements of  $\mathcal{W}^*$  may induce different quotient spaces and therefore different relational invariants; which laws appear depends on which closure is sustained.

This model-relativity should not be confused with arbitrariness. Any such law is still a law of one and the same observational world as compressed under a particular stable closure. If multiple closures persist, they persist under the constraint of the same  $D$  and the same  $\mu_D$ .

## 13 Entropy

The classifier  $\pi_w$  compresses the observation space. In this section we assume that the partition  $Q_w = D/\sim_w$  is finite or countable, so that the sums below are well defined.

**Definition 6** (Class Measure).

$$M_w(d) = \mu_D([d]_w).$$

**Definition 7** (Model-Relative Surprisal).

$$S_w([d]_w) = -\log \mu_D([d]_w).$$

**Definition 8** (Model-Relative Entropy).

$$H(w) = - \sum_{[d]_w \in Q_w} \mu_D([d]_w) \log \mu_D([d]_w).$$

The quantity  $S_w([d]_w)$  measures the probability mass of the equivalence class  $[d]_w$ , which is the fiber of the projection  $\pi_w : D \rightarrow Q_w$ . The quantity  $H(w)$  is the expected surprisal induced by the partition defined by  $\pi_w$  and therefore measures the representational compression associated with the model.

Because classifiers are elements of  $D$ , both surprisal and entropy assign measure-theoretic weight not only to equivalence classes of perceptual content but also to equivalence classes of valuations. A valuation that is rare in  $D$  carries high surprisal. A coarse model that collapses many distinct

valuations into a single class yields low surprisal for that class and lowers the effective distinguishability of evaluative structure. Entropy is therefore not merely a feature of perceptual content; it also measures the coarseness with which a model distinguishes the system’s evaluative dispositions.

A note on scope is warranted. The entropy  $H(w)$  defined above is a Shannon-type quantity derived from representational compression. The framework does not claim identity between this quantity and thermodynamic entropy. It claims structural convergence: both quantities arise from the same underlying operation of partitioning a probability space, and Jaynes’ programme of deriving statistical mechanics from inference over probability distributions subject to informational constraints suggests that this convergence is not superficial. The precise conditions under which model-relative entropy and thermodynamic entropy coincide are left for subsequent work.

## 14 Representational Refinement

**Definition 9** (Refinement). *A model  $w_2$  refines  $w_1$  if  $[d]_{w_2} \subseteq [d]_{w_1}$  for all  $d \in D$ .*

**Theorem 1** (Monotonicity of Surprisal). *If  $w_2$  refines  $w_1$ , then*

$$S_{w_2}([d]_{w_2}) \geq S_{w_1}([d]_{w_1}).$$

*Proof.* Refinement implies  $[d]_{w_2} \subseteq [d]_{w_1}$ , so  $\mu_D([d]_{w_2}) \leq \mu_D([d]_{w_1})$ . Applying  $-\log$  reverses the inequality.  $\square$

**Theorem 2** (Entropy Equality for Equivalent Observations). *If  $d_1 \sim_w d_2$ , then  $S_w([d_1]_w) = S_w([d_2]_w)$ .*

*Proof.* If  $d_1 \sim_w d_2$  then  $[d_1]_w = [d_2]_w$ , so  $\mu_D([d_1]_w) = \mu_D([d_2]_w)$  and the definition of  $S_w$  yields the result.  $\square$

### 14.1 Refinement as Dilution

It is a common intuition that refinement—the transition from  $w_1$  to  $w_2$  where  $[d]_{w_2} \subseteq [d]_{w_1}$ —represents a “narrowing in” on a point-like truth. However, the framework suggests the inverse. As the partition becomes finer, the measure  $\mu_D$  associated with each class decreases, and the surprisal increases.

If we consider the individuals within each class as the unobserved territory, refinement actually produces a *coarser coverage* of the underlying openness. Each refined class  $[d]_{w_2}$  holds fewer observational states, but the density of the unknown within our representation increases. We do not converge toward an external object; rather, we dilute our representational density, creating the very space in which the constitutive remainder manifests as surprisal. Refinement is not the elimination of uncertainty, but the formal expansion of the system’s capacity to be surprised.

## 15 Interpretation

The hierarchy of structure is

$$(F, g) \rightarrow T \rightarrow \mathcal{W}^* \rightarrow \pi_w \rightarrow Q_w \rightarrow R_w.$$

The condition  $C \subseteq D$  runs through every level of this hierarchy. Classifiers enter the observation space as observations, are compressed by  $\pi_w$ , appear in the quotient space  $Q_w$ , figure in relational

tuples in  $K(Q_w)$ , and carry surprisal under  $S_w$ . Valuation is not a parameter set from outside the system; it is a structural feature of the observation space that the system’s own representational apparatus must absorb, compress, and partially lose.

The implication map now makes explicit that closure is reproduced not from an atomized observational residue but from a structured observational profile. A stable closure is therefore a view from somewhere in the strict sense: a model whose own internally generated way of inhabiting the world, when reinterpreted through inference, yields the same model again.

## 15.1 Backing into the Future

The atemporal nature of the block universe framing suggests a reinterpretation of the experience of time. If  $\Omega$  is a static relational structure, then “the future” is not a state that has not yet occurred, but a region of the block universe toward which the system’s current stable closure  $w^*$  has not yet been extended.

The system does not move into the future; rather, it *backs into it* according to past patterns. The inference–implication loop  $T(w) = w$  is a consistency condition derived from the existing weight of observations. When the system encounters the unmapped regions of the block, it projects its current relational invariants ( $L$ ) as a structural expectation.

On this reading, the experience of time is the process of this projection being stressed by the constitutive remainder. Because refinement increases surprisal, the future feels ultimately unpredictable not because it is non-existent, but because our attempt to know it more finely—to refine our “backing” movement—necessarily dilutes our coverage. We encounter the future as a growing coarseness of classes, where the patterns of the past are the only machinery available to navigate the increasing openness of the territory ahead. This unpredictability is the formal address of will: the necessity of choosing a closure in a territory that the model can never fully exhaust.

## 16 Institutions as Intergenerational Compression

The framework developed so far treats an epistemic system as a single relational structure. But embedded systems are not isolated. They exist within communities of systems that share, contest, and transmit closures across time. This section extends the framework to that social dimension, focusing on the role of institutions. A fuller institutional formalization is developed in *The Imagination Machine IV: Institutional Intelligence*.

The central observation is this: no individual knower transmits a closure to a successor by reproducing the full observation space  $D$  that gave rise to it. What is transmitted is always a compression—a residue of the inferential work that produced a given  $w^* \in \mathcal{W}^*$ . Institutions are the mechanisms by which this intergenerational compression is stabilised.

More precisely, what passes between generations is not the loop itself—the maps  $F$  and  $g$  that generated the fixed point—but a projection of the implied observational profile and the quotient structure it presupposes into the observation space of the successor generation. The successor receives the closed texture without necessarily receiving the closure mechanism. Institutions are the structures within  $\Omega$  that perform and stabilise this projection, re-embedding the inherited profile of closure as observations in the successor’s  $D$ , making it available for classification by the successor’s own  $\pi_w$ .

This framing carries an immediate consequence. A successor generation may inherit a stable closure without inheriting the capacity to regenerate it under pressure from new observations. The quotient structure arrives, but the inferential machinery that produced it does not.

We distinguish two modes of institutional transmission. *Compressed inheritance* transmits the closed profile alone: the successor can apply the inherited partition but cannot update it. *Generative inheritance* transmits  $F$  and  $g$  alongside that profile: the successor can regenerate the closure from within, extend it, and revise it when new observations demand a finer partition.

The distinction matters because the observation space  $D$  does not stand still. New observations enter  $D$  in every generation, and a partition that was self-consistent under an earlier  $\mu_D$  may fail to remain so as the measure shifts. A generatively inherited closure can meet this pressure; a compressedly inherited one cannot. The institution that transmits only the quotient structure is therefore more fragile—not because it contains false beliefs, but because it has lost the capacity to refine.

Note that institutions may also transmit miscalibrated measures. A community that inherits both  $F$  and  $g$  alongside a systematically distorted  $\mu_D$  possesses the machinery for refinement while lacking accurate observational weights on which to exercise it. Generative inheritance is therefore necessary but not sufficient for epistemic health: the inferential measure must also track the environment it purports to represent.

## 17 Knowledge, Dogma, and the Structure of Refinement

A natural question arises from the plurality of stable closures established in Section 10: if  $\mathcal{W}^*$  may contain many incompatible elements, and the framework provides no external criterion for preferring one over another, how does it distinguish knowledge from dogma? Both are self-consistent. Both survive the inference–implication loop. Both can be institutionally transmitted.

The answer is that the distinction does not require an external criterion. It falls out of the structure already in place, specifically from the relationship between a closure and its behaviour under refinement.

Recall that a model  $w_2$  refines  $w_1$  when  $[d]_{w_2} \subseteq [d]_{w_1}$  for all  $d \in D$ . Refinement always costs higher surprisal: a finer partition assigns lower probability mass to each class and therefore higher  $S_w$  to each observation. A closure disposed toward knowledge is one that remains willing to pay this cost—one whose inference–implication loop, when supplied with observations that increase the consistency gap under the current partition, responds by generating a finer  $\pi_w$  rather than forcing the new observations into existing classes.

Dogmatic closure is precisely the refusal to pay this cost. A dogmatic model maintains its self-consistency not by genuinely absorbing new observations but by compressing them into existing equivalence classes regardless of their character. New elements of  $D$  are mapped by  $\pi_w$  to existing elements of  $Z_w$  even when a more faithful compression would require extending  $Z_w$ . The partition is held fixed; the observations are bent to fit it.

Miscalibration, introduced in Remark 3, constitutes a distinct failure mode. A miscalibrated closure may be fully open to refinement—willing to extend  $Z_w$  whenever the consistency gap demands it—and yet refine systematically against a distorted image of the observational world. Where dogmatism is a failure of disposition at the level of  $(F, g)$ , miscalibration is a failure of the measure  $\mu_D$  itself, prior to any particular act of closure. The two failures are formally separable: a closure can be dogmatic without being miscalibrated, or miscalibrated without being dogmatic, or both simultaneously.

A clarification is required here. The criterion just stated relies on a notion of stable absorption that is not itself fully decidable from within a single closure. Determining whether a new observation  $d$  genuinely strains the existing partition or is legitimately compressed into it requires assessing the consistency gap, and different closures may assess that gap differently. The framework does

not resolve this from outside; it rather establishes the vocabulary within which the question can be precisely posed and contested. The distinction between knowledge and dogma is therefore best understood as identifying a structural disposition—the preparedness to extend  $Z_w$  under pressure—rather than as a decision procedure that can be applied mechanically from within any single closure. Crucially, this question is available to any system satisfying  $C \subseteq D$ , since such a system can observe its own classificatory behaviour and the consistency of its loop.

Several further consequences follow. First, the distinction is not binary but gradational. A closure may be refinable with respect to some regions of  $D$  while dogmatic with respect to others. Institutions that transmit  $F$  and  $g$  alongside the inherited profile of closure preserve the capacity for refinement, but may do so selectively—maintaining the inferential machinery for some domains while suppressing it for others.

Second, the surprisal cost of refinement explains a persistent feature of actual epistemic communities. Dogmatic compression avoids this cost by refusing to see new observations as genuinely new. Coarser models assign lower surprisal to the observations they assimilate, and lower surprisal feels, from within the closure, like greater understanding. The framework thus provides a structural account of why the pressure toward dogmatic closure is not merely psychological but has a measure-theoretic basis.

Third, because  $C \subseteq D$ , the distinction applies to evaluative structure as well as perceptual content. A closure that refuses to refine its classification of classifiers—that compresses distinct valuations into the same equivalence class regardless of the observational pressure to distinguish them—is dogmatic about value in precisely the same structural sense. The framework does not treat these as different in kind.

Returning to the hierarchy established in Section 15, the distinction between knowledge and dogma lives at the level of  $(F, g)$  rather than at the level of  $\mathcal{W}^*$ . Two closures may be indistinguishable as fixed points—equally self-consistent, equally stable—while differing fundamentally in whether the loop they instantiate remains open to refinement. Stability is not the same as openness, and it is openness to refinement—the disposition to pay the surprisal cost when the consistency loop demands it—that the present framework identifies as the structural mark of what distinguishes knowledge from its appearance.

## 18 Conclusion

Embedded epistemic systems cannot appeal to external correspondence as their standard of coherence. Coherence appears instead as internal closure of the inference–implication loop under the statistical structure of observations. Self-consistent world models arise as fixed points of the operator this loop induces, and each such model compresses the observation space into a quotient representation whose relational invariants constitute physical law and whose measure-theoretic multiplicity constitutes entropy.

The structural feature that distinguishes this framework from earlier accounts is the inclusion  $C \subseteq D$ : classifiers are observations. This inclusion is not stipulated but derived—it is the transcendental condition on any system capable of Cartesian doubt, any system that can recognise itself as one possible closure among others. This means that valuation and will—the dispositions that select among possible closures—are interior to the representational architecture. They appear in the observation space, are subject to compression, and leave their trace in the quotient structure. Yet they are not exhausted by any compression. The formal loop determines the space of stable closures, but not which closure is instantiated. This remainder is not a gap in the framework; it is the constitutive openness that the inference–implication loop must encompass but cannot exhaust.

The implication map clarifies the form of that closure. What a model implies is not merely an isolated consequence but a canonical observational profile internal to the model itself: a structured way the world appears from somewhere. A self-consistent world model is therefore one whose own implied profile of observational life, when reinterpreted through inference, reproduces that same model. Stable theory and stable world-profile co-arise.

The social extension of the framework yields two further results that follow from the same architecture without requiring external normative imports. Institutions are the mechanisms by which stable closures are transmitted across generations, but they transmit closures in two structurally distinct modes: generative inheritance conveys the inferential machinery alongside the fixed point, while compressed inheritance conveys only the inherited profile of closure. And the distinction between knowledge and dogma reduces, within the framework, to the distinction between closures that remain open to refinement and those that hold their partition fixed against the pressure of new observations—a difference that identifies a structural disposition rather than a decision procedure applicable from outside any particular closure.

The framework thus diagnoses three irreducible structural locations of epistemic risk. Dogmatism is a failure of disposition at the level of  $(F, g)$ : the loop exists but refuses to refine. Miscalibration is a failure at the level of  $\mu_D$ : the loop refines willingly but against a distorted image of the world. And will names the underdetermination that persists even when both are functioning well—the necessity of choosing a closure in territory no model can fully exhaust. Together these three constitute the complete formal topology of epistemic failure for an embedded system.

Epistemic closure, physical law, entropy, and the social conditions of knowledge therefore emerge as successive consequences of a single embedded representational architecture. What prevents meta-physical closure—what keeps the system in relation to the openness beneath its representations—is the evaluative structure that the architecture must include but cannot fully exhaust.

## References

- Clark, A. (2016). *Surfing Uncertainty: Prediction, Action, and the Embodied Mind*. Oxford University Press.
- Friston, K. (2010). The free-energy principle: a unified brain theory? *Nature Reviews Neuroscience*, 11:127–138.
- Jaynes, E. T. (1957). Information theory and statistical mechanics. *Physical Review*, 106(4):620–630.
- Kitcher, P. (1993). *The Advancement of Science: Science without Legend, Objectivity without Illusions*. Oxford University Press.
- Ladyman, J. (1998). What is structural realism? *Studies in History and Philosophy of Science Part A*, 29(3):409–424.
- Lawson, H. (2001). *Closure: A Story of Everything*. Routledge.
- Longino, H. E. (1990). *Science as Social Knowledge: Values and Objectivity in Scientific Inquiry*. Princeton University Press.
- Maturana, H. R. and Varela, F. J. (1980). *Autopoiesis and Cognition: The Realization of the Living*. D. Reidel.

- Rovelli, C. (1996). Relational quantum mechanics. *International Journal of Theoretical Physics*, 35(8):1637–1678.
- Shannon, C. E. (1948). A mathematical theory of communication. *Bell System Technical Journal*, 27(3):379–423.
- Varela, F. J., Thompson, E., and Rosch, E. (1991). *The Embodied Mind: Cognitive Science and Human Experience*. MIT Press.
- Worrall, J. (1989). Structural realism: The best of both worlds? *Dialectica*, 43(1–2):99–124.

# The Imagination Machine II: Systems

Mark Tracy  
Boston University  
mrktracy@bu.edu

## Introduction

This paper is the second part of a four-paper series titled *The Imagination Machine*. The first paper, *The Imagination Machine I: A View from Somewhere*, develops a formal epistemic framework for embedded observers and introduces the inference–implication loop that defines self-consistent world models. Within that framework, observations, classifiers, and world models all appear as structures internal to the same universe, and epistemic coherence arises as the closure of a representational loop rather than correspondence with an external vantage point.

The present paper develops a complementary layer of the project by introducing a general formalism for systems. Whereas the first paper describes the structure of representational closure for embedded epistemic systems, the present work describes the dynamical coupling between components of such systems, particularly in the case of agent–environment interaction. The goal is to define systems in an extremely general way so that the formalism has maximal expressiveness while making minimal assumptions.

In the first section, we develop a general definition of a system in terms of measurable variables, stochastic processes, and functional relations between inputs and outputs. In the following section we introduce optimization models and relate them to systems through the problem of system identification. The agent–environment framework developed later in the paper provides a general structure for modeling adaptive systems whose outputs influence the environment from which future inputs arise.

A minimal computational setting in which such agent–environment dynamics give rise to the recovery of relational invariants is developed in *The Imagination Machine III: Toy Model of Predictive Classification*. The final paper in the series, *The Imagination Machine IV: Institutional Intelligence*, extends the analysis further by examining how epistemic processes are stabilized and transmitted across communities and institutions.

# 1 General System Definition<sup>1</sup>

Define a set of variables that can be measured in practice. By necessity, this will be a countable set of variables, even if the underlying real system of interest has uncountable degrees of freedom. By measuring these variables over a set of time points  $I$  with minimal value  $t_0$ , we collect *data*. We distinguish between two different proper subsets of variables: *input variables* and *output variables*.

## 1.1 Input Variables

Without loss of generality, we will assume going forward that there is only one input variable. This is without loss of generality because any countable set of input variables may be represented by a tuple whose components are the simpler variables.

We represent the input variables with a random process. In particular, let  $(\Omega_u, \mathcal{F}_u, \mathbb{P}_u)$  denote a probability space. Let  $U_{\text{in}}$  be the set of possible values the input variable may take.

Then the input process

$$u : \Omega_u \times I \rightarrow U_{\text{in}}$$

is measurable as a function from  $(\Omega_u \times I, \mathcal{F}_u \otimes \mathcal{F}_I)$  to  $(U_{\text{in}}, \mathcal{F}_{\text{in}})$ , where  $\mathcal{F}_I$  denotes a  $\sigma$ -algebra on the time set  $I$ ,  $\mathcal{F}_{\text{in}}$  denotes a  $\sigma$ -algebra on the input space  $U_{\text{in}}$ , and where the symbol  $\otimes$  denotes the product  $\sigma$ -algebra. For each fixed time  $t \in I$ , the mapping

$$\omega \mapsto u(\omega, t)$$

is a random variable, and for every measurable set  $A \subseteq U_{\text{in}}$  (i.e.  $A \in \mathcal{F}_{\text{in}}$ ), the distribution of  $u(t)$  is given by:

$$\mathbb{P}_u(\{\omega \in \Omega_u \mid u(\omega, t) \in A\})$$

We note, crucially, that although we are representing our input variable as a random process, input variables are often chosen to be those that one can deliberately vary over time. In such a case, the input variable may not be stochastic. In general, the input variable  $u(t)$  at any time  $t \in I$  may be a tuple in a product space of simpler, potentially degenerate random variables.

## 1.2 Output Variables

Complementary to the input variables are the output variables. Again, we will treat the case of a single output variable, since countably many output variables may be treated as a tuple.

Similarly to the input variable, we can represent the output variable as a random process. As before, let  $(\Omega_y, \mathcal{F}_y, \mathbb{P}_y)$  denote a probability space. Let  $U_{\text{out}}$  be the set of possible values the output variable may take.

---

<sup>1</sup>Some language and structure adapted from *Introduction to Discrete Event Systems: Third Edition* by Christos G. Cassandras and Stéphane Lafortune. <https://doi.org/10.1007/978-3-030-72274-6>

Then the output process

$$y : \Omega_y \times I \rightarrow U_{\text{out}}$$

is measurable as a function from  $(\Omega_y \times I, \mathcal{F}_y \otimes \mathcal{F}_I)$  to  $(U_{\text{out}}, \mathcal{F}_{\text{out}})$ , where  $\mathcal{F}_{\text{out}}$  denotes a  $\sigma$ -algebra on the output space  $U_{\text{out}}$ , and where the symbol  $\otimes$  denotes the product  $\sigma$ -algebra. For each fixed time  $t \in I$ , the mapping

$$\omega \mapsto y(\omega, t)$$

is a random variable, and for every measurable set  $A \subseteq U_{\text{out}}$  (i.e.  $A \in \mathcal{F}_{\text{out}}$ ), the distribution of  $y(t)$  is given by:

$$\mathbb{P}_y(\{\omega \in \Omega_y \mid y(\omega, t) \in A\})$$

### 1.3 Relating Inputs to Outputs

The relation between the input variable and time, and the resulting output variable, is given by a functional  $g$ . A functional is a function whose domain is a Cartesian product of one or more sets of functions and zero or more other sets. In this case, the domain of the functional  $g$  is the Cartesian product of the set  $\mathcal{U}$  of all measurable input processes  $u : \Omega_u \times I \rightarrow U_{\text{in}}$  and the time set  $I$ . The codomain of  $g$  is  $\mathcal{P}(U_{\text{out}}, \mathcal{F}_{\text{out}})$ , the space of probability measures over the measurable space  $(U_{\text{out}}, \mathcal{F}_{\text{out}})$ . Explicitly:

$$g : \mathcal{U} \times I \rightarrow \mathcal{P}(U_{\text{out}}, \mathcal{F}_{\text{out}}),$$

The functional  $g$  satisfies:

$$y(t) \sim g[u, t]$$

Or, if modeling time as discrete, where  $t_{i+1}$  is a successor of  $t_i$  in a countable and strictly ordered time set  $I$  whose minimal element is  $t_0$ :

$$y(t_{i+1}) \sim g[u, t_i]$$

The symbol  $\sim$  denotes random sampling or should be read as “is distributed according to.” If the distribution is degenerate (i.e. there is no stochasticity), then the symbol may be treated as deterministic assignment, identically to an “equals” sign. In other words, determinism is represented as stochasticity with a degenerate distribution—assigning probability 1 to a single outcome.

Note that each component of the output variable (when considering a tuple of simpler variables) at time  $t$  can depend in general on the value of any component of the input variable (again, when considering a tuple of simpler variables) at any subset of time points, potentially including future points. While many physical systems are assumed to be “causal” (outputs depend only on present and past inputs), the mathematical formulation permits non-causal dependencies, allowing flexibility in modeling retroactive influence.

The functional  $g$  relating input and time to output may be an evaluation functional, which directly evaluates an input variable at a given time point, e.g.:

$$y(t) = g[u, t] = u(t)$$

It may also be a function of such evaluation functionals, e.g.,

$$y(t) = g[u, t] = u(t) + 3u(t - 1.3) - 76.8u(t + 4)^2$$

## 1.4 State

While the above is a general description of any system, in many cases, especially where history and memory matter, we find it useful to model the system's internal condition explicitly. This internal condition is what we call the system's *state*, which we can represent as a random process defined over some probability space  $(\Omega_s, \mathcal{F}_s, \mathbb{P}_s)$ . Letting  $U_{\text{state}}$  be the set of values that the state may take, we can write:

$$s : \Omega_s \times I \rightarrow U_{\text{state}}$$

Again, we consider the state  $s(t)$  at time  $t$  to be a single random variable without loss of generality, since the state variable may be a tuple in a product space of simpler, potentially degenerate (i.e. determinate) random variables.

The evolution of the state may be represented in continuous time by a stochastic differential equation:

$$\dot{s}(t) \sim f[u, s, t], \quad s(t_0) \sim s_0 \quad \text{for some } s_0 \in \mathcal{P}(U_{\text{state}}, \mathcal{F}_{\text{state}})$$

where  $\mathcal{F}_{\text{state}}$  is a  $\sigma$ -algebra on  $U_{\text{state}}$  and where

$$f : \mathcal{U} \times \mathcal{S} \times I \rightarrow \mathcal{P}(U_{\text{change}}, \mathcal{F}_{\text{change}}),$$

for some set  $U_{\text{change}}$  whose elements represent rates of change of the state and for  $\mathcal{F}_{\text{change}}$  a  $\sigma$ -algebra on  $U_{\text{change}}$ ; and where we denote by  $\mathcal{S}$  the space of all measurable state processes  $s : \Omega_s \times I \rightarrow U_{\text{state}}$ .

If modeling time as discrete rather than continuous, then we may represent state dynamics as an update rule:

$$s(t_{i+1}) - s(t_i) \sim f[u, s, t_{i+1}], \quad s(t_0) \sim s_0 \quad \text{for some } s_0 \in \mathcal{P}(U_{\text{state}}, \mathcal{F}_{\text{state}})$$

where, similarly to before,

$$f : \mathcal{U} \times \mathcal{S} \times I \rightarrow \mathcal{P}(U_{\text{change}}, \mathcal{F}_{\text{change}}),$$

for some set  $U_{\text{change}}$  whose elements represent changes in the state, and where  $t_{i+1}$  is a successor of  $t_i$  in a countable and strictly ordered time set  $I$  whose minimal element is  $t_0$ .

The relation between the input variable, the system state, and time, and the resulting output variables may then be expressed as a functional:

$$y(t) \sim g[u, s, t]$$

or, in discrete time, where  $t_{i+1}$  is a successor of  $t_i$  in a countable and strictly ordered time set  $I$  whose minimal element is  $t_0$ :

$$y(t_{i+1}) \sim g[u, s, t_i]$$

where

$$g : \mathcal{U} \times \mathcal{S} \times I \rightarrow \mathcal{P}(U_{\text{out}}, \mathcal{F}_{\text{out}}).$$

## 2 Optimization Models

A functional, like those discussed above, is a special kind of function. An optimization model is a function approximator. An optimization model consists of a triplet  $(\mathcal{H}, O, A)$  of:

1. A hypothesis space  $\mathcal{H}$  (a set of functions);
2. An objective  $O : \mathcal{H} \rightarrow \mathbb{R}$  (a functional whose domain is the hypothesis space and whose range is real numbers) which gives some signal as to the quality of the approximation; and
3. An optimization algorithm  $A : \mathcal{H} \rightarrow \mathcal{H}$  (a rule for moving through the hypothesis space), in general utilizing the objective.

When learning inductively from data (that is, when attempting to move from particular examples to general principles), a few additional objects may be appended to the aforementioned triplet; in particular:

4. A dataset  $\mathcal{D}$ .
5. A (possibly unknown) random process  $P$  from which data points are sampled. In other words, data is collected empirically from the world during a time interval  $I_D$  with  $d_i \sim P(t_i) \quad \forall d_i \in \mathcal{D}$ , where data point  $d_i$  is collected at time  $t_i$ . Note that in cases where data points may be assumed to be identically distributed and drawn independently, this amounts to a single distribution. In *active learning*, the algorithm  $A$  interacts with the random process  $P$ , influencing the empirical dataset  $\mathcal{D}$  used during optimization. In other words, data points are not sampled according to  $P$  before the commencement of the algorithm  $A$ , but rather, the process of data collection is itself influenced by the optimization algorithm.
6. A dataset  $\mathcal{D}_{\text{aug}}$ , where for all  $d_{\text{aug}} \in \mathcal{D}_{\text{aug}}$ , there exists a function  $f$ , an integer  $N$ , and a tuple of elements  $t \in \mathcal{D}^N$  such that  $d_{\text{aug}} \sim f(t)$ , where  $\sim$  denotes, as before, stochastic sampling of the (possibly degenerate) random function  $f$ . In other words, every element of  $\mathcal{D}_{\text{aug}}$  is a (potentially stochastic) function of elements of  $\mathcal{D}$ .
7. A random process  $P_{\text{train}}$  by which elements of  $\mathcal{D}_{\text{aug}}$  are drawn by the optimization algorithm. In particular, the algorithm  $A$  draws at time  $t_i$  an element  $d_i \sim P_{\text{train}}(t_i)$ , where  $P_{\text{train}}(t_i)$  is a distribution over  $\mathcal{D}_{\text{aug}}$ .

An inductive bias is a constraint on the hypothesis space. By traversing the hypothesis space algorithmically, an optimization model is intended to minimize the objective function and thus obtain a good approximation to the function that truly represents the system of interest.

### 3 System Identification

System identification is the process of utilizing an optimization model to find an approximation to the true dynamics of a system using measurements of its input and output variables. In particular, it is useful when the internal state of a system is not known or its internal dynamics—the stochastic differential (or difference) equation(s) and initial conditions governing the state’s trajectory—are not known.

### 4 Agents

The agent–environment coupling introduced here provides the dynamical structure within which the representational closures described in *The Imagination Machine I: A View from Somewhere* may arise for embedded epistemic systems. In that framework, stable world models emerge as fixed points of an inference–implication loop defined over observations internal to the same universe. The systems formalism developed here provides a concrete representation of the interacting processes through which such observations and models may be generated.

An agent necessarily exists within and is co-constituted with an environment. An agent–environment pair comprises two systems, an agent  $A$  and environment  $E$ , which are in interchange (feeding back to one another); as well as an initial input to either the agent or the environment. In particular,  $A$  takes as input the output of  $E$ , and  $E$  takes as input the output of  $A$ , with the recursion beginning from some set of initial inputs to either system.

Formally, we may represent the recursive dependency between an agent  $A$  and an environment  $E$  as follows:

$$\begin{aligned} u^A(t) &= y^E(t) && \text{(agent receives environment's output as input)} \\ u^E(t) &= y^A(t) && \text{(environment receives agent's output as input)} \end{aligned}$$

where:

- $u^A(t)$  is the input to the agent at time  $t$
- $y^A(t)$  is the agent’s output at time  $t$
- $u^E(t)$  is the input to the environment at time  $t$
- $y^E(t)$  is the environment’s output at time  $t$

The recursion begins from a set of initial inputs:

$$\begin{aligned} u^E(t_0) &\sim u_0^E && \text{for some } u_0^E \in \mathcal{P}(U_{\text{in}}^E, \mathcal{F}_{\text{in}}^E) && \text{or} \\ u^A(t_0) &\sim u_0^A && \text{for some } u_0^A \in \mathcal{P}(U_{\text{in}}^A, \mathcal{F}_{\text{in}}^A) \end{aligned}$$

and the pair evolves together over time, potentially governed by their own internal state dynamics:

$$\begin{aligned} \dot{s}^A(t) &\sim f^A[u^A, s^A, t], & y^A(t) &\sim g^A[u^A, s^A, t] \\ \dot{s}^E(t) &\sim f^E[u^E, s^E, t], & y^E(t) &\sim g^E[u^E, s^E, t] \end{aligned}$$

for some functionals defined analogously as before:

$$\begin{aligned} f^A &: \mathcal{U}^A \times \mathcal{S}^A \times I \rightarrow \mathcal{P}(U_{\text{change}}^A, \mathcal{F}_{\text{change}}^A) \\ g^A &: \mathcal{U}^A \times \mathcal{S}^A \times I \rightarrow \mathcal{P}(U_{\text{out}}^A, \mathcal{F}_{\text{out}}^A) \\ f^E &: \mathcal{U}^E \times \mathcal{S}^E \times I \rightarrow \mathcal{P}(U_{\text{change}}^E, \mathcal{F}_{\text{change}}^E) \\ g^E &: \mathcal{U}^E \times \mathcal{S}^E \times I \rightarrow \mathcal{P}(U_{\text{out}}^E, \mathcal{F}_{\text{out}}^E) \end{aligned}$$

That is, each functional takes as input:

- a random input process over  $I$ ,
- a random state process over  $I$ ,
- and the current time  $t \in I$ ,

and produces either a rate of change of the state (for  $f$ ) or an output (for  $g$ ), potentially by sampling randomly from a distribution of outputs.

## 4.1 Agents in Discrete Time

In many practical applications, especially in reinforcement learning, the agent-environment interaction is modeled in discrete time. This leads to the following slight change in representation:

$$\begin{aligned} u^A(t_{i+1}) &= y^E(t_i) \quad (\text{agent receives environment's most recent output as input}) \\ u^E(t_{i+1}) &= y^A(t_i) \quad (\text{environment receives agent's most recent output as input}) \end{aligned}$$

where  $t_{i+1}$  is the successor of  $t_i$  in some ordered set of time points  $I$  (in particular, the time points are indexed by  $i \in \mathbb{N}_0$ ), and where

$$\begin{aligned} s^A(t_{i+1}) - s^A(t_i) &\sim f^A[u^A, s^A, t_{i+1}], & y^A(t_{i+1}) &\sim g^A[u^A, s^A, t_{i+1}] \\ s^E(t_{i+1}) - s^E(t_i) &\sim f^E[u^E, s^E, t_{i+1}], & y^E(t_{i+1}) &\sim g^E[u^E, s^E, t_{i+1}] \end{aligned}$$

for some functionals defined analogously as before:

$$\begin{aligned} f^A &: \mathcal{U}^A \times \mathcal{S}^A \times I \rightarrow \mathcal{P}(U_{\text{change}}^A, \mathcal{F}_{\text{change}}^A) \\ g^A &: \mathcal{U}^A \times \mathcal{S}^A \times I \rightarrow \mathcal{P}(U_{\text{out}}^A, \mathcal{F}_{\text{out}}^A) \\ f^E &: \mathcal{U}^E \times \mathcal{S}^E \times I \rightarrow \mathcal{P}(U_{\text{change}}^E, \mathcal{F}_{\text{change}}^E) \\ g^E &: \mathcal{U}^E \times \mathcal{S}^E \times I \rightarrow \mathcal{P}(U_{\text{out}}^E, \mathcal{F}_{\text{out}}^E) \end{aligned}$$

That is, each functional takes as input:

- a random input process over  $I$ ,
- a random state process over  $I$ ,
- and the current time  $t \in I$ ,

and produces either a state update (for  $f$ ) or an output (for  $g$ ), potentially by sampling randomly from a distribution over possible outputs.

## 4.2 Reinforcement Learning

Reinforcement learning is a special case of an optimization model, whereby the objective  $O$  depends on the history of interactions between an agent and its environment and where the algorithm  $A$  seeks to maximize the expected cumulative reward obtained through the agent and environment’s dynamic coupling.

In the context of the present series, such agent–environment optimization dynamics provide a concrete setting in which representational models may be iteratively refined through interaction with structured environments. A minimal predictive example of such refinement is developed in *The Imagination Machine III: Toy Model of Predictive Classification*.

## 5 Becoming-Held-As-By: Subjects as Systems in Self-Representation

In the language developed in *The Imagination Machine I: A View from Somewhere*, a self-representing subject is an embedded epistemic system whose classifiers appear within its own observation space. The condition that classifiers are themselves observations allows a system to encounter and revise its own acts of classification. The present section approaches the same idea from the perspective of systems modeling: if the formalism developed above can represent any system, then it must also apply to the system performing the representation.

If the above framework above provides insight into how to represent any real system in mathematical terms, then a natural next step is to turn the inquiry on the modeler. In other words, in writing the above formalism I am confronted with the question, “Am I not a real system myself? Can I, then, be understood in these terms?”

I imagine a bubble around my body, and then I imagine it shrinks inward all around and approaches infinitely closely to the edge of my skin. Any measurable passing between this membrane is either input or output—and thus I conceive of agent and environment.

Pursuant to these aims, we now shift from formal system representation to a philosophical and phenomenological inquiry into how a system may represent itself as an agent, co-constituted and co-evolving with an environment. In this way, we move from a formalism for modeling system behavior from an external perspective to a vocabulary by which a self-modeling agentic system may represent its own reality from the internal perspective.

## 5.1 A Self-Referential Thesis

All may be called the Becoming-Held-As-By<sup>2</sup> (including its becoming held as this by me).

## 5.2 Potentiality and Representation

Suppose we take “existence” to mean “the quality, state, or event of becoming-held-as-by.” We use the word “potentiality” to mean that from which existence emerges through representation. Potentiality is metaphysical substance itself—what we might call the pre-conceptual whatever-I-represent. Representation is the process or result of becoming-held-as-by. A subject is becoming-held-as-by-itself.

For example, to say that a particular cup “exists” is to say that some potentiality (which I could, for example, point to) is becoming held in mind by me as a unified and distinct “thing” which I represent as a cup. If the potentiality does not become held as anything by any subject, then it cannot be said that anything in particular exists there, though there may persist some potentiality for becoming-held-as-by (held as a cup, perhaps, or as something else, like a weapon or a hat, by any particular subject). To hold potentiality as something is not to define it once and for all, but to engage in a relationship that may change. The same potentiality may be held as many different things across time, across subjects, or even within the same subject in different moments.

One cannot properly imagine potentiality because all one *can* do is imagine potentiality, in the sense of bringing potentiality into representation. That is to say that to imagine potentiality is already to bring it into representation. Potentiality may have internal structure (e.g. change relative to some internal reference frame according to laws). Regardless, here is the big picture: potentiality (metaphysical substance) is translated into existence (the ontological) through its representation by the subject (the semiological and epistemological: perception, language, systems of meaning, knowledge claims). By this notion of existence, if every conscious being were to disappear suddenly, there would not be a universe at all—only potential for a universe to arise.

Reality, in this account, is enacted through the interplay of potentiality and representation: a process in which potential becomes held through representation, and representation constrains potentiality.

## 5.3 The Subject Becoming-Held-As-Agent-By-Itself

The most stable world-model I have yet realized is this: world as constituted of agent (self) co-evolving with environment, where the agent’s state includes its representation of self, environment, and world; including, recursively, a representation of world as constituted of agent (self) co-evolving with environment, where the agent’s state includes its representation of self, environment, and world.

This is a world that I hold as constituted of the agent-environment coupling, wherein a subject may coherently and productively become-held-as-agent-by-itself. The agent is not

---

<sup>2</sup>The hyphenation of “Becoming-Held-As-By” is deliberate: it reflects the interdependence and co-constitution of the becoming, the holding, the *as*-ness (representation), and the *by*-ness (the subject).

separable from the environment, though it may be ontologically separate in its own representation. The state of the agent includes its representation of potentiality: It is influenced by its environment's output and its own history, and, in turn, it influences the input to the environment through the output of the agent. Because of the inherent coupling of agent to environment, the subject becoming-held-as-agent-by-itself is to the Becoming-Held-As-By as a *holon* to its greater whole<sup>3</sup>: the subject may become-held-as a distinct object of analysis by itself, and yet it can simultaneously become-held by itself as a part in a larger system.

To use a human-centric analogy, the subject becoming-held-as-agent-by-itself is to the Becoming-Held-As-By as the mouth is to the body: the mouth is not the body, yet it is interconnected with the body; and the declaration that "I am the body" is made by means of the mouth. Similarly, the subject becoming-held-as-agent-by-itself is not the entirety of the Becoming-Held-As-By and yet is embedded (and participating) within it; and the writing and reading of statements like, "All may be called the Becoming-Held-As-By (including its becoming held as this by me)" is enacted by the subject becoming-held-as-agent-by itself.

## 5.4 The Limits of the Systems Formalism

A system is defined by its distinctions: inputs vs. outputs, internal vs. external state. The undivided Whole—that which contains all systems, distinctions, and environments—cannot itself be represented as a system. Since it has no external relation and no boundary, it admits no input/output mapping. Likewise, the complement of the undivided Whole—that is, nothingness, or void—admits no input/output mapping and as such may not be represented as a system.

## 5.5 Mathematics as Meta-Representation

Mathematics may derive from the structure of representation itself. That is to say, mathematics is a type of meta-representation: a representation of common structure across instances of representation. Accordingly, the representation of mathematical objects could potentially be invariant under change in subject if each subject can in principle abstract from their own instances of representation to arrive at the same mathematical meta-representations. For example, I can map a notion of two-ness to the same symbol that another subject can map theirs, and we can be reasonably sure that we agree on its meaning, because we both experience unity and difference in our representations of self and world. Likewise, I can map a notion of a function to the same symbol that another can map theirs, and we can be reasonably sure that we agree on its meaning, because we both represent and abstract from instances of change. Unity, difference, and change may be necessary structures of subjective representation, such that any subject with sufficient abstract reasoning capability can attribute the same meaning to the same meta-representations.

---

<sup>3</sup>The term "holon" was coined by Arthur Koestler in his 1967 book, *The Ghost in the Machine*. A holon is both a self-contained entity (hence it is a whole on its own) and at the same time is embedded within a larger containing system or systems (so it is part of a larger whole).

## 5.6 Haecceity and Qualia

Complementary to the notion of meta-representation in this account is the notion of haecceity, or the irreducible *this*-ness of an entity. Haecceity is what remains in representation modulo meta-representation—the particularity that is not captured by abstraction from representation to meta-representation. For a human, haecceity may correspond to the irreducible qualia of the experience of being *this particular self* in *this particular moment*.

## 5.7 Truth and Coherence

A proposition is a linguistic claim that may be judged true or false by a subject. Truth is a judgment of coherence among a collection of propositions. Formally, a proposition is judged false if it is shown that the proposition, potentially together with a collection of propositions judged to be true, implies contradiction of a proposition judged true. Therefore, a particular proposition is judged true only by virtue of its ongoing coherence with a collection of mutually non-contradictory propositions. It must be emphasized that propositions involving instantiation (of abstract classes) are among the propositions that must be coherent in a collection of truths. For example, a proposition like, “An electron evolves according to the Schrodinger equation,” must cohere with such propositions of instantiation as, “This reading (referring to a particular representation in experience) is due to an electron,” and, “This reading (at another time, perhaps) is not due to an electron,” as well as propositions that are not instantiations like, “An electron has negative charge.” This understanding of truth allows for pluralism while requiring that a worldview be consistently tethered to moments of becoming-held-as-by.

## 5.8 Conclusion

The central claim is that we are always describing the world from the inside: embedded within the Becoming-Held-As-By and co-evolving with our environment, seeing patterns in our seeing-patterns. We conscious beings are individually and collectively a self-representing network of interacting holonic subsystems. And yet, on the whole and within each part, haecceity remains.

# The Imagination Machine III: Prediction, Control, and Representational Closure in Quasi-Periodic Environments

Mark Tracy  
Boston University  
mrktracy@bu.edu

## Abstract

This paper develops a unified treatment of prediction, control, and representational closure for embedded epistemic systems situated in quasi-periodic environments. We proceed in three stages. First, we motivate the quasi-periodic environment as the naturalistic setting in which human temporal metacognition evolved: the Earth–Sun–Moon system presents embedded observers with incommensurate cycles whose relative phases continually drift, selecting for predictive and inductive cognitive machinery. Second, we formalize a minimal computational realization of this setting in which a predictive agent recovers latent dynamical structure from relational observations through prediction error alone. Third, we extend the framework to include action, showing that reinforcement learning arises naturally as a special case of embedded epistemic closure when policy is defined over the compressed representational classes induced by a world model. Across all three stages, the same compression–extension architecture governs representation, prediction, and control. Convergence in reinforcement learning corresponds to a fixed point of a joint model–policy closure operator, unifying representation learning and control under the structural mechanism developed throughout the Imagination Machine series.

# Contents

<b>1</b>	<b>Introduction: Temporal Metacognition in Quasi-Periodic Environments</b>	<b>3</b>
	<b>Part I: Prediction</b>	<b>4</b>
<b>2</b>	<b>The Quasi-Periodic Environment</b>	<b>4</b>
2.1	Observation Model . . . . .	4
2.2	Environment Distribution . . . . .	4
<b>3</b>	<b>The Predictive Agent</b>	<b>4</b>
3.1	Neural Parameterization . . . . .	4
3.2	Prediction Error and Training . . . . .	5
<b>4</b>	<b>Observable Invariants and the Koopman Connection</b>	<b>5</b>
4.1	Time-Rescaling Symmetry . . . . .	5
4.2	Koopman Representation . . . . .	5
<b>5</b>	<b>Empirical Protocol</b>	<b>5</b>
5.1	Latent Structure Recovery via Linear Probing . . . . .	5
5.2	Generalization and Robustness . . . . .	6
	<b>Part II: Control</b>	<b>6</b>
<b>6</b>	<b>Action and Embedded Systems</b>	<b>6</b>
<b>7</b>	<b>Policy as Will Over Compressed Observations</b>	<b>6</b>
<b>8</b>	<b>Evaluative Compression</b>	<b>7</b>
<b>9</b>	<b>The Reinforcement Learning Closure Operator</b>	<b>7</b>
<b>10</b>	<b>Exploration as Refinement</b>	<b>8</b>
<b>11</b>	<b>Value Functions on the Quotient Space</b>	<b>8</b>
<b>12</b>	<b>The Koopman Connection in the Control Setting</b>	<b>8</b>
<b>13</b>	<b>Conclusion</b>	<b>9</b>

# 1 Introduction: Temporal Metacognition in Quasi-Periodic Environments

History becomes possible when at least three natural cycles repeat with incommensurate periods, producing configurations whose relative phases continually drift and never exactly repeat.

The Earth–Sun–Moon system presents observers with a particular sort of temporal environment: a set of stable but incommensurate cycles whose relative phases continually drift and admit comparison. Meanwhile, the objects traveling in these cycles interact in complex ways—through tides, energy transfer, and other processes—that become crucial for biological life. Human temporal metacognition develops in response to this structured signal.

Human temporal metacognition emerges from the interaction of four recurrent processes experienced by an Earth-bound observer: the solar day, the lunar phase cycle, the solar year, and the circadian sleep–wake rhythm. The three celestial cycles possess incommensurate periods, generating quasi-periodic patterns—stable motions whose relative phases continually drift. This underlying structure makes the development of counting, memory, and inductive estimation advantageous, since empirical estimates of ratios between characteristic constants of these quasi-periodic processes accumulate through repeated observation rather than diverging without bound or collapsing into exact repetition.

Societies historically come to represent relations between characteristic constants of these cycles through ratios between them (e.g.,  $\approx 365$  solar days per solar year), calendars, and continuous real-valued models of time. The circadian rhythm simultaneously segments subjective experience into discrete episodes through the sleep–wake cycle. The result is a dual conception of time: discrete lived intervals embedded within continuously modeled celestial motion. Each human life becomes entrained into this dynamical scaffold.

Change appears as aperiodic variation within an underlying pattern of stability. Because the relative phases of the celestial cycles continually drift, accounting for such variation benefits from the abstraction of recursively nested temporal demarcations—days within months, months within years, and so on—together with the maintenance of records across cycles.

Epistemically, the ratios governing these cycles are inductive approximations derived from observation and record-keeping. Their numerical values are refined through repeated measurement and expressed as real-valued quantities in continuous temporal models. Human symbolic systems thereby impose numerical structure on a multi-body procession whose precise relative phases are never exactly and fully known. There remains novelty amidst structure.

History becomes the maintenance of physical records of possible continuations of universal relative motion under a particular superimposed continuous and cyclic temporal model: a stochastic process of sampling-through-externalization within the world-process, enacted through acts of demarcation, recursively interpreted and reinterpreted interpersonally through (1) abstraction, which enables compressive projection; (2) analogy, which allows domain transfer of hypotheses; and (3) communication among agents.

The two formal parts of this paper develop a minimal model of the cognitive machinery this environment selects for. Part I formalizes a predictive agent that recovers latent dynamical structure from relational observations. Part II extends the framework to action, showing that control arises naturally within the same representational architecture.

# Part I: Prediction

## 2 The Quasi-Periodic Environment

Let the environment consist of three cyclic variables

$$\theta_1(t), \theta_2(t), \theta_3(t) \in S^1 \cong \mathbb{R}/2\pi\mathbb{Z}.$$

Their dynamics are

$$\theta_i(t+1) = \theta_i(t) + \omega_i \pmod{2\pi}, \quad i = 1, 2, 3,$$

or equivalently  $\theta(t+1) = \theta(t) + \omega \pmod{2\pi}$  where  $\theta(t) = (\theta_1(t), \theta_2(t), \theta_3(t))$  and  $\omega = (\omega_1, \omega_2, \omega_3)$ .

**Definition 1** (Rational Independence). *Real numbers  $\omega_1, \omega_2, \omega_3$  are rationally independent if the only integer solution to  $k_1\omega_1 + k_2\omega_2 + k_3\omega_3 = 0$  with  $k_1, k_2, k_3 \in \mathbb{Z}$  is  $k_1 = k_2 = k_3 = 0$ .*

**Definition 2** (Quasi-Periodic System). *The dynamical system defined above is quasi-periodic if  $\omega_1, \omega_2, \omega_3$  are rationally independent. In this case the trajectory is dense on the torus  $\mathbb{T}^3 = S^1 \times S^1 \times S^1$ .*

### 2.1 Observation Model

The environment state is  $x_t = \theta(t)$ . The agent observes only relational quantities  $o_t = h(x_t)$ , where

$$h(x_t) = (\cos \Delta_{12}, \sin \Delta_{12}, \cos \Delta_{13}, \sin \Delta_{13}, \cos \Delta_{23}, \sin \Delta_{23})$$

and  $\Delta_{ij}(t) = \theta_i(t) - \theta_j(t) \pmod{2\pi}$ . The agent observes only relational phase differences between the cyclic processes.

### 2.2 Environment Distribution

Frequency vectors are sampled from the normalized simplex

$$\Delta^2 = \{ \omega \in \mathbb{R}^3 : \omega_i > 0, \omega_1 + \omega_2 + \omega_3 = 1 \}.$$

Each sampled vector defines a distinct quasi-periodic environment.

## 3 The Predictive Agent

A predictive agent is defined by three functions:

$$o_t = h(x_t), \quad s_{t+1} = u(s_t, o_t), \quad \hat{o}_{t+1} = g(s_{t+1}),$$

where  $s_t$  is internal state and  $\hat{o}_{t+1}$  is the predicted next observation.

### 3.1 Neural Parameterization

The state update is parameterized by a neural network:

$$s_{t+1} = \text{MLP}_u([s_t, o_t]).$$

The prediction head is a linear readout:

$$\hat{o}_{t+1} = W s_{t+1} + b.$$

The linear prediction head creates a representational bottleneck: the internal state must organize information in a form directly readable through linear transformations.

### 3.2 Prediction Error and Training

Prediction error is  $e_{t+1} = \hat{o}_{t+1} - o_{t+1}$ . Training minimizes

$$\mathcal{L}_{t+1} = \|\hat{o}_{t+1} - o_{t+1}\|_2^2.$$

## 4 Observable Invariants and the Koopman Connection

### 4.1 Time-Rescaling Symmetry

**Proposition 1.** *The frequency vector  $\omega$  is identifiable only up to multiplication by a positive scalar when inferred from relational phase observations alone.*

*Proof.* Let  $k > 0$  and define  $\omega' = k\omega$ . Then  $\theta'(t) = \theta(0) + k\omega t$ , which is equivalent to  $\theta'(t) = \theta(\tau)$  for  $\tau = kt$ . The orbit is unchanged; only the parameterization by time differs. Since the observation function  $h$  depends only on phase differences  $\Delta_{ij}$ , which are invariant under uniform rescaling of  $\omega$ , no relational observation can distinguish  $\omega$  from  $k\omega$ .  $\square$

Observable invariants are therefore the projective equivalence class  $[\omega_1 : \omega_2 : \omega_3]$ . Normalizing via  $\omega_1 + \omega_2 + \omega_3 = 1$ , distinct environments correspond to points in the interior of  $\Delta^2$ .

### 4.2 Koopman Representation

Writing the complex observable  $z_{ij}(t) = e^{i\Delta_{ij}(t)}$ , the relational dynamics imply

$$z_{ij}(t+1) = e^{i(\omega_i - \omega_j)} z_{ij}(t).$$

The observable evolves through multiplication by a constant complex phase factor, constituting a linear evolution in observable space. This is precisely a Koopman eigenfunction: the nonlinear state dynamics on  $\mathbb{T}^3$  become linear in the space of relational observables. The linear prediction head therefore tests whether the agent has learned an internal representation that approximates this Koopman eigenfunction structure. Accurate prediction through a linear readout implies that the internal state encodes the relevant dynamical invariants in a linearly accessible form.

## 5 Empirical Protocol

### 5.1 Latent Structure Recovery via Linear Probing

After training on an environment with frequency vector  $\omega^{(k)}$ , the agent produces a final internal state  $s_T^{(k)}$ . A linear probe

$$\hat{y} = Ws + b$$

is trained to predict

$$y^{(k)} = \begin{pmatrix} \omega_1^{(k)} / \omega_3^{(k)} \\ \omega_2^{(k)} / \omega_3^{(k)} \end{pmatrix}$$

using mean squared error. Probe performance measures whether latent dynamical structure is represented in a linearly accessible form in the agent’s internal state.

## 5.2 Generalization and Robustness

Training the probe on a subset of environments and evaluating on held-out environments tests whether the representation captures general dynamical structure rather than environment-specific features. The environment may be extended to weakly nonstationary dynamics by allowing slow frequency drift:

$$\omega(t + 1) = \omega(t) + \epsilon_t,$$

where  $\epsilon_t$  is a small perturbation. This tests the robustness of the learned representation to distributional shift.

## Part II: Control

### 6 Action and Embedded Systems

Part I studied an agent that observes and predicts but does not intervene. Part II extends the same architecture to an agent that additionally selects actions. The formal setting follows The Imagination Machine II, in which an embedded agent and environment form a coupled dynamical system through reciprocal input–output channels:

$$u_A(t) = y_E(t), \quad u_E(t) = y_A(t),$$

where  $u_A, u_E$  denote inputs and  $y_A, y_E$  denote outputs to the agent and environment respectively. The observations available to the agent constitute a subset

$$D \subseteq \Omega$$

of the total relational structure  $\Omega$ . As in The Imagination Machine I, the agent constructs world models  $w \in W$  by compressing observational profiles through an inference map  $F : \Gamma \rightarrow W$ , while an implication map  $g : W \rightarrow \Gamma$  generates predicted observational profiles from those models.

### 7 Policy as Will Over Compressed Observations

A world model  $w$  induces a classifier

$$\pi_w : D \rightarrow Z_w$$

partitioning observations into representational classes via the equivalence relation

$$d \sim_w d' \iff \pi_w(d) = \pi_w(d'),$$

with induced quotient space  $Q_w = D/\sim_w$ .

**Definition 3** (Policy). *A policy is a stochastic map*

$$\pi : Q_w \rightarrow \Delta(A)$$

*from representational classes to distributions over an action space  $A$ .*

Because an embedded agent cannot act on the full observational space—it has access only to the compressed representation  $Q_w$ —policy must be defined over representational classes rather than raw observations. Policy is therefore the operational expression of will relative to the world model: the agent’s selective pressure over actions, compressed to the resolution its model affords.

## 8 Evaluative Compression

Standard reinforcement learning treats reward as a primitive signal supplied by an external oracle. In an embedded epistemic framework this is unavailable: the agent has no access to an external vantage point from which to receive unmediated evaluative verdicts. Reward must instead arise as a compression of evaluative observations.

Let  $D^*$  denote the set of finite observation trajectories.

**Definition 4** (Evaluative Compression). *An evaluative compression is a map*

$$R : D^* \rightarrow \mathbb{R}$$

*assigning scalar value to observational trajectories.*

The reward signal therefore reflects the agent’s own evaluative structure, compressed over trajectories in the same way that world models compress instantaneous observations. This is consistent with the inclusion  $C \subseteq D$  established in The Imagination Machine I: classifiers—including evaluative classifiers—are themselves observations, subject to the same representational compression as any other element of  $D$ .

## 9 The Reinforcement Learning Closure Operator

When action is introduced, the implication map becomes policy-conditioned:

$$g : W \times \Pi \rightarrow \Gamma,$$

where  $\Pi$  denotes the space of policies. Given a world model and a policy, this map generates the predicted observational profile resulting from the coupled agent-environment dynamics under that policy. Inference remains

$$F : \Gamma \rightarrow W.$$

**Definition 5** (RL Closure Operator). *Let  $\mathcal{A} : W \times R \rightarrow \Pi$  be an action-selection operator that produces a policy from a world model and an evaluative compression. The reinforcement learning closure operator is*

$$T_{\text{RL}}(w, \pi) = (F(g(w, \pi)), \mathcal{A}(F(g(w, \pi)), R)).$$

**Definition 6** (RL Closure). *A pair  $(w^*, \pi^*)$  is a reinforcement learning closure if*

$$T_{\text{RL}}(w^*, \pi^*) = (w^*, \pi^*).$$

At such a fixed point the world model accurately predicts the observational consequences of the policy, and the policy is optimal relative to the model and evaluative compression. The pair is jointly self-consistent in the same sense that a world model alone is self-consistent under the epistemic closure operator  $T = F \circ g$ .

**Remark 1.** *The action-selection operator  $\mathcal{A}$  is left general here. Specific instantiations correspond to known algorithms: Q-learning, policy gradient methods, and actor-critic architectures each realize particular choices of  $\mathcal{A}$  within this framework. Existence of a fixed point  $(w^*, \pi^*)$  requires conditions analogous to those governing the epistemic fixed points of The Imagination Machine I—compactness and continuity assumptions sufficient to warrant a Schauder-type argument.*

## 10 Exploration as Refinement

Exploration arises when the representational partition induced by the world model is too coarse to support reliable prediction or control.

**Definition 7** (Refinement). *A model  $w_2$  refines  $w_1$  if  $[d]_{w_2} \subseteq [d]_{w_1}$  for all  $d \in D$ .*

Refinement corresponds to splitting equivalence classes in  $Q_w$  when observations within a class exhibit divergent consequences under action. An agent whose model assigns the same representational class to states with different value cannot distinguish among them in its policy. Exploration is the mechanism by which such distinctions become available.

**Remark 2.** *Exploration is an epistemic operator rather than random behavior: it seeks observations that maximize the probability of representational refinement. The exploration–exploitation tradeoff is therefore a special case of the knowledge–dogma distinction developed in *The Imagination Machine I*. An agent that ceases to explore has adopted a dogmatic closure: it holds its representational partition fixed against the pressure of new observations. The cost of this closure is not merely suboptimal reward but the structural foreclosure of refinement.*

## 11 Value Functions on the Quotient Space

Because policy operates on representational classes, value functions must be defined on the same space.

**Definition 8** (Value Function). *For a fixed policy  $\pi$  and world model  $w$ , the value function is*

$$V_w^\pi : Q_w \rightarrow \mathbb{R},$$

*assigning expected evaluative compression to each representational class under  $\pi$ .*

Action-value functions are defined analogously:

$$Q_w^\pi : Q_w \times A \rightarrow \mathbb{R}.$$

These functions evaluate the expected return of taking action  $a$  from representational class  $[d]_w$  and thereafter following  $\pi$ .

## 12 The Koopman Connection in the Control Setting

The Koopman structure established in Part I has a direct consequence for Part II. Because the relational observables  $z_{ij}(t) = e^{i\Delta_{ij}(t)}$  evolve linearly in the space of preserved invariants, value functions defined over  $Q_w$  inherit this linear structure when the world model has recovered the Koopman representation. A model that encodes the dynamical invariants in a linearly accessible internal state supports value estimation that is linear in the compressed state—which is precisely the structure that makes reinforcement learning tractable in practice.

The quasi-periodic environment is therefore not an arbitrary testbed. It is the minimal environment in which the connection between predictive representation and tractable control is explicit and provable. The Koopman eigenfunctions provide the natural basis for both prediction and value estimation, and the linear prediction head of Part I is the architectural condition that forces the agent to learn them.

## 13 Conclusion

This paper has developed a unified treatment of prediction, control, and representational closure for embedded epistemic systems in quasi-periodic environments.

The introduction established the naturalistic motivation. The Earth–Sun–Moon system presents embedded observers with incommensurate cycles that select for predictive and inductive cognitive machinery. Novelty amidst structure is not a special feature of this environment—it is its defining characteristic, and it is why the inference–implication loop can never fully close. The cognitive machinery the series formalizes is the machinery this environment selected for.

Part I formalized a minimal predictive agent and showed that latent dynamical structure—specifically, the Koopman eigenfunction representation of the relational observables—becomes linearly recoverable through prediction error alone. The linear prediction head is not an arbitrary architectural choice; it is the condition that forces the internal state to encode dynamical invariants in a form that makes the Koopman connection testable.

Part II extended the framework to action. Reinforcement learning arises naturally when policy is defined over the compressed representational classes induced by a world model, reward is treated as evaluative compression over trajectories, and learning seeks fixed points of a joint model–policy closure operator. Exploration is the behavioral expression of refinement pressure: the agent acts in order to find where its partition is too coarse. An agent that stops exploring has, in the precise sense of The Imagination Machine I, gone dogmatic.

Across all three stages, the same architecture governs representation, prediction, and control. Prediction, control, and valuation are not separate problems. They are different aspects of a single embedded representational structure in which an agent, unable to access the world from outside, must construct, refine, and act from within the only closure available to it.

# The Imagination Machine IV: Institutional Intelligence

Mark Tracy  
Boston University  
mrktracy@bu.edu

March 15, 2026

## Abstract

Scientific institutions evolve knowledge through a recursive process: structured dialogue, selective compression and differential transmission of ideas, and empirical feedback. To model this dynamic, we introduce a formal model of institutional learning in which both reasoning procedures and evaluative procedures evolve through dialogue, compression, and environmental feedback. Monte Carlo generations of trios of dialogical agents operate over a shared corpus and evolving prompts. Each trio generates dialogue interpreted as a sample path in a representational space. Agents propose compression rules, prompt revisions, and candidate solutions to an external problem.

Compression proceeds in two stages. First, a compendium is formed from the proposed compression rules, resulting in a compression prompt. Second, a language model conditioned on the compression prompt produces a final prompt revision and proposes a solution to the problem of interest.

The subsequent generation of agents receives an external correctness signal evaluating the final solution of the previous generation, and the final prompt revisions are implemented before simulation of the subsequent generation commences.

The resulting architecture formalizes a minimal model of institutional learning in which reasoning rules and evaluative procedures co-evolve through dialogue, compression, revision, and environmental feedback.

## 1 Introduction

Scientific institutions evolve knowledge through a recursive process: structured dialogue, selective compression and differential transmission of ideas, and empirical feedback.

This paper is the fourth part of the series *The Imagination Machine*. The first paper, *A View from Somewhere*, develops a formal epistemic framework for embedded observers. The second paper, *Systems*, introduces a general formalism for interacting dynamical systems and agent–environment coupling. The third paper, *A Toy Model of Predictive Classification in a Quasi-Periodic Environment*, studies how predictive agents can recover latent structure from relational observations.

The present paper extends those ideas from individual reasoning systems to institutional learning. We model academic institutions as evolving systems in which both reasoning procedures and evaluative procedures change through dialogue, compression, and feedback.

Academic institutions operate through a recursive process:

1. researchers generate hypotheses through dialogue and simultaneously contribute to institutional procedures
2. institutions filter, compress, and differentially transmit generated ideas
3. empirical feedback guides future research and institutional development

Two forms of structure evolve simultaneously:

- **reasoning**: the ideas and conceptual frameworks under discussion
- **evaluation**: the procedures by which ideas are summarized, reviewed, and judged

We formalize this process as a recursive system of dialogue, compression, and feedback operating across generations of interacting agents. The framework can be interpreted both as a theoretical model of institutional learning and as a potential architecture for multi-agent reasoning systems.

**Contribution.** We introduce a formal model of institutional learning in which both reasoning procedures and evaluative procedures evolve through dialogical exploration, compression, and feedback. The resulting dynamics define a stochastic process over institutional states.

*Note on related work.* The authors are aware that related work exists across several relevant literatures, including prompt optimization, multi-agent large language model systems, evolutionary approaches to prompt search, and institutional learning. A fuller literature review situating this contribution within those bodies of work will appear in a revised version.

## 2 Basic Objects

**Definition 1** (Corpus). *Let  $W$  denote a shared corpus of writings available to all agents.*

**Definition 2** (Representational Space). *Let  $\mathcal{R}$  denote a representational space of possible dialogues.*

**Definition 3** (External Problem). *Let  $\mathcal{Q}$  denote an external problem to which generations propose solutions.*

Two prompts evolve over time.

Together these prompts constitute the institutional procedures governing intellectual exploration and evaluation.

**Definition 4** (Reasoning Prompt).  *$R_g$  denotes the reasoning prompt at generation  $g$ .*

**Definition 5** (Compression Prompt).  *$C_g$  denotes the compression prompt governing summarization.*

## 3 Monte Carlo Dialogical Trios

At generation  $g$  we instantiate a population

$$\mathcal{M}_g = \{T_g^{(1)}, \dots, T_g^{(N_g)}\}$$

of dialogical trios.

Each trio

$$T_g^{(k)} = \{a_{g,1}^{(k)}, a_{g,2}^{(k)}, a_{g,3}^{(k)}\}$$

is initialized with

$$(W, R_g, C_g, \mathcal{Q})$$

and generates a dialogue.

**Definition 6** (Dialogue Sample Path). *The dialogue produced by trio  $T_g^{(k)}$  is*

$$D_g^{(k)} \in \mathcal{R}.$$

Dialogue trajectories are interpreted as sample paths through the representational space.

## 4 Agent Outputs

Each agent  $a_{g,i}^{(k)}$  produces three outputs:

1. reasoning revision proposal

$$(A_{g,i}^{R,k}, F_{g,i}^{R,k})$$

2. compression prompt revision proposal

$$(A_{g,i}^{C,k}, F_{g,i}^{C,k})$$

3. candidate solution

$$S_{g,i}^k$$

Here  $A$  denotes additions to a prompt and  $F$  denotes proposed removals (forgetting).

## 5 Two-Stage Compression

Compression proceeds in two stages.

The two stages separate the accumulation of institutional memory from the compression of that institutional record into transmitted reasoning and evaluation procedures.

### 5.1 Stage 1: Compendium Construction

Collect proposed additions to the compression prompt:

$$\mathcal{A}_g = \left\{ A_{g,i}^{C,k} \mid 1 \leq k \leq N_g, 1 \leq i \leq 3 \right\}.$$

Construct the compendium:

$$\tilde{C}_g = \text{Gather}(C_g, \mathcal{A}_g).$$

The Gather operation aggregates proposed additions without semantic compression, functioning as an append-only institutional memory within a generation; removals from the compression prompt are applied only after summarization and before transmission to the next generation.

## 5.2 Stage 2: Summarization

Define

$$\mathcal{R}_g = \left\{ (A_{g,i}^{R,k}, F_{g,i}^{R,k}) \mid 1 \leq k \leq N_g, 1 \leq i \leq 3 \right\},$$

$$\mathcal{C}_g = \left\{ (A_{g,i}^{C,k}, F_{g,i}^{C,k}) \mid 1 \leq k \leq N_g, 1 \leq i \leq 3 \right\},$$

and

$$\mathcal{S}_g = \left\{ S_{g,i}^k \mid 1 \leq k \leq N_g, 1 \leq i \leq 3 \right\}.$$

Let  $\Gamma$  denote a language model conditioned on the compendium.

Using compendium  $\tilde{C}_g$ , compute

$$\Gamma_g^{\tilde{C}_g}(\mathcal{R}_g, \mathcal{C}_g, \mathcal{S}_g) = (A_g^R, F_g^R, A_g^C, F_g^C, \tilde{S}_g).$$

Here  $\tilde{S}_g$  denotes the summarized solution proposed by generation  $g$ .

## 6 Generational Feedback

The summarized solution  $\tilde{S}_g$  is evaluated against the external problem  $\mathcal{Q}$ .

The environment returns a feedback signal

$$Y_g \in \mathcal{Y}$$

representing the correctness or quality of the proposed solution.

## 7 Prompt Updates

Reasoning and compression prompts evolve separately but in a coupled manner.

### 7.1 Compression Update

$$C_{g+1} = \tilde{C}_g \setminus F_g^C.$$

### 7.2 Reasoning Update

$$R_{g+1} = (R_g \setminus F_g^R) \oplus A_g^R \oplus C_{g+1} \oplus Y_g.$$

Both layers evolve through inheritance, forgetting, and structural addition. If the prompt length exceeds a threshold  $M$ , tokens are removed according to a first-in-first-out (FIFO) policy.

## 8 Algorithmic Overview

The imagination machine evolves prompts across generations through dialogue, compression, and feedback. One generational step proceeds as follows.

1. Initialize a population of dialogical trios using prompts  $(R_g, C_g)$  and shared corpus  $W$ . For example, the population may be a collection of trios of instances of a language model with pseudo-randomly sampled temperature parameters.
2. Each trio generates a dialogue  $D_g^{(k)}$  and agents propose reasoning revisions, compression prompt revisions, and candidate solutions.
3. Aggregate proposed compression prompt additions and construct the compendium  $\tilde{C}_g = \text{Gather}(C_g, \mathcal{A}_g)$ .
4. Use the language model  $\Gamma$  conditioned on  $\tilde{C}_g$  to summarize revisions and candidate solutions.
5. Evaluate the summarized solution  $\tilde{S}_g$  against the external problem  $\mathcal{Q}$  and obtain feedback signal  $Y_g$ .
6. Update reasoning and compression prompts to obtain  $(R_{g+1}, C_{g+1})$ .

## 9 Stochastic Institutional Dynamics

The evolution of the system can be interpreted as a stochastic process.

**Definition 7** (Institutional State). *Let*

$$\mathcal{X} := \text{the space of possible prompt pairs,}$$

and let

$$X_g := (R_g, C_g) \in \mathcal{X}$$

denote the institutional state at generation  $g$ .

Dialogue generation, summarization, and feedback introduce randomness through sampling processes and Monte Carlo population dynamics.

**Definition 8** (Generational Transition Kernel). *Let*

$$K(\cdot \mid X_g, W, \mathcal{Q})$$

denote the conditional probability law governing the next institutional state given the current state, corpus, and external problem.

Thus institutional evolution may be written

$$X_{g+1} \sim K(\cdot \mid X_g, W, \mathcal{Q}).$$

## 10 Interpretation

Dialogue trajectories

$$D_g^{(k)}$$

represent sample paths through representational space generated by interacting reasoning agents. The Monte Carlo population approximates a distribution over such trajectories.

Compression extracts shared structure across dialogues, while external feedback guides the evolution of reasoning.

The resulting architecture formalizes institutional learning: ideas evolve through dialogue to solve problems while evaluative procedures operate through selective compression and transmission of corpora of recorded symbols; reasoning and evaluation therefore co-evolve through recursive institutional dynamics.

## 11 Conclusion

The imagination machine evolves two interacting structures:

- reasoning prompts governing intellectual exploration
- compression prompts governing institutional evaluation

Dialogue generates trajectories, compression extracts inheritable structure, forgetting prevents uncontrolled growth, and feedback from external problems guides institutional learning across generations.

# The Imagination Machine V: On Abstraction and Analogy

Mark Tracy

## 1 Overview

Analogy is the bedrock of communication. Even that sentence makes use of analogy: as bedrock underlies and supports structures, so too does analogy underlie and support communication, allowing us to coordinate activity and manipulate our environment. Analogy allows a reasoner to transfer previously learned structure to a new situation, generating hypotheses and thereby facilitating new understanding. So fundamental is analogy to language that it proves challenging to articulate the abstract structure of analogy and to codify valid analogical reasoning. Nonetheless, it remains a fundamental endeavor for any interested in understanding mentation. In the foregoing, I introduce and augment one popular model of analogy, and I utilize the formalism thus achieved to attempt a definition of valid analogical reasoning.

## 2 Classical Theories of Analogy

A domain may be defined as a tuple:<sup>1</sup>

$$D = (O, A, R, S, T)$$

- $O$  = set of objects
- $A$  = set of attributes (unary operators:  $a \in A \implies a : O \rightarrow S$ )
- $R$  = set of relations (n-ary operators:  $r \in R \implies \exists n \in \mathbb{N}, r : O^n \rightarrow S$ )
- $S$  = set of statements
- $T$  = set of statements believed to be true (belief set)

Note that attributes are a special case of relations: each  $a \in A$  is simply a unary relation, so formally  $A \subseteq R$ .

---

<sup>1</sup>This definition follows the standard treatment of domains in analogy and relational reasoning literature (cf. 1), but extends it to include a set of statements  $S$  and a belief set  $T$ , corresponding respectively to the expressible and the held-to-be-true propositions within the domain.

## 2.1 Structure-Mapping Theory of Analogy

In the landmark paper “Structure-Mapping: A Theoretical Framework for Analogy,” Gentner argues that an analogy is a mapping between objects in a base domain and objects in a target domain that does not necessarily carry over object-level attributes but which carries over some relational predicates.[1]

## 2.2 A Formal Definition of Analogy

An analogy between a source domain  $D_s = (O_s, A_s, R_s, S_s, T_s)$  and a target domain  $D_t = (O_t, A_t, R_t, S_t, T_t)$  is defined by a tuple:

$$A = (X, Y, M, P)$$

- $X \subset O_s$ : a collection of objects in the source domain
- $Y \subset O_t$ : a collection of objects in the target domain
- $M : X \rightarrow Y$ : a mapping of objects from source to target domain
- $P \subset \{r \mid r \in R_s \cap R_t \text{ and } \exists \mathbf{x} \in X^k \text{ for some } k \in \mathbb{N} \text{ such that } r(\mathbf{x}) \in T_s \text{ and } r(M(\mathbf{x})) \in T_t\}$ : a set of relations that are present in the source and target domains, are true of some tuple of objects in the source domain, and are preserved in the target domain via the mapping  $M$ . As a notational convention, we consider  $M(\mathbf{x})$  to be the component-wise application of the mapping  $M$  to the tuple  $\mathbf{x}$ , i.e.  $\mathbf{x} = (x_1, \dots, x_n) \implies M(\mathbf{x}) = (M(x_1), \dots, M(x_n))$ .

## 2.3 Analogical Reasoning

Let  $D_s$  be a source domain and  $D_t$  a target domain. Suppose:

- $X_1 \subset O_s$  is a subset of objects in the source domain. Let  $|X_1| = n$ .
- $Y_1 \subset O_t$  is a subset of objects in the target domain.
- $M : X_1 \rightarrow Y_1$  is a mapping of the source domain subset to the target domain subset.
- $P$  is a set of relations preserved by the mapping  $M$ .

This establishes an analogy between  $D_s$  and  $D_t$ . Now suppose that some further fact (of a particular form to be specified below) holds in the source domain; we formally define an **analogical reasoning step** to be the positing of a corresponding form of further fact in the target domain. Formally:

Suppose there exists a superset of  $X_1$  called  $X_0$ :

$$\begin{aligned}
X_1 &\subseteq X_0 \\
|X_0| &= m \geq n
\end{aligned}$$

and suppose that

$$r(\mathbf{x}^*) \in T_s$$

for some tuple  $\mathbf{x}^* \in X_0^k$  for some  $k \in \mathbb{N}$  and for some relational predicate  $r \in (R_s \cap R_t)$ .

Then an analogical reasoning step is to hypothesize that there exists a mapping  $M'$  that preserves and extends the original analogical mapping  $M$  and preserves the further observed relation in the source domain,  $r$ . In particular, the hypothesis is as follows:

$$\begin{aligned}
&\exists Y_2 \subset O_t \quad \text{and} \\
&\exists M' : X_0 \rightarrow Y_1 \cup Y_2 \quad \text{such that} \\
&\quad M'(x) = M(x) \quad \forall x \in X_1 \quad \text{and} \\
&r(M'(\mathbf{x}^*)) \in T_t,
\end{aligned}$$

where  $M'(\mathbf{x}^*)$  is the component-wise application of the mapping  $M'$  to the tuple  $\mathbf{x}^*$  identified above.

This formulation captures the logic of projecting relational structures from the source domain into the target domain, conditioned on preserved analogical structure. It highlights how analogy can support hypothesizing about unseen objects, roles, or relations in the target domain by structurally mapping known relations in the source.

## 2.4 Analogy as Mediated by Abstraction

**Abstraction**, in the broadest sense, refers to the process or result of mapping a collection of objects, attributes, or relations to a single representation, typically to retain only information which is relevant for a particular purpose.

There is a connection between abstraction and analogy that is insufficiently explored in Gentner's 1983 paper. If, as Gentner convincingly argues, an analogy is a mapping between objects in a base domain and objects in a target domain that does not necessarily carry over object-level attributes but which carries over some relational predicates [1], then for any analogy there exists an abstract domain that implicitly mediates the analogy. In particular, the domain that mediates an analogy  $A = (X, Y, M, P)$  between a source domain  $D_s = (O_s, A_s, R_s, S_s, T_s)$  and a target domain  $D_t = (O_t, A_t, R_t, S_t, T_t)$  consists of:

- **A new set of objects,  $O_{\text{abs}}$ :**

- Call them symbols.
- $\forall x \in X, (x, M(x)) \in O_{\text{abs}}$ .
- Notational convention: for a  $k$ -tuple of objects in the source domain,  $\mathbf{x} \in X^k$ , we denote the corresponding tuple of symbols as  $(\mathbf{x}, M(\mathbf{x})) \in O_{\text{abs}}^k$ , where  $M(\mathbf{x})$  is the component-wise application of  $M$  to  $\mathbf{x}$ . In particular:

$$\begin{aligned} \mathbf{x} &= (x_1, \dots, x_k) \implies \\ M(\mathbf{x}) &= (M(x_1), \dots, M(x_k)) \text{ and} \\ (\mathbf{x}, M(\mathbf{x})) &= ((x_1, M(x_1)), \dots, (x_k, M(x_k))) \end{aligned}$$

- **A set of predicate attributes,  $A_{\text{abs}}$ :**

- $A_{\text{abs}} = P \cap A_s$
- The set of unary relations preserved by the analogy, if any.

- **A set of predicate relations,  $R_{\text{abs}}$ :**

- Call them abstract relations.
- $r \in P \iff r \in R_{\text{abs}}$
- $r(\mathbf{x}) \in T_s$  for some  $\mathbf{x} \in X^k$  with  $k \in \mathbb{N} \implies r((\mathbf{x}, M(\mathbf{x}))) \in T_{\text{abs}}$ .

- **A statement set,  $S_{\text{abs}}$ :**

- All possible combinations from the collections of objects, attributes, and relations specified above.

- **A belief set,  $T_{\text{abs}}$**

- A subset of  $S_{\text{abs}}$ , populated as specified above.

### 2.4.1 An example

Take the analogy, “An atomic nucleus is like the solar system.” [1] At an earlier point in scientific history, the analogical mapping may have looked like this:

$$\begin{aligned} M : X &\rightarrow Y \\ \text{NUCLEUS} &\mapsto \text{SUN} \\ \text{ELECTRON} &\mapsto \text{PLANET} \end{aligned}$$

And the relationships preserved include:

$$\{\text{ORBITS, IS\_MOVING}\} \subset P.$$

Now, in recognizing a mediating abstract domain we may synthesize new symbols with carried-over attributes and abstract relations, thereby forming a mediating abstract domain that both source and target instantiate:

$$\begin{aligned} \{\text{NUCLEUS, SUN}\} &\mapsto \text{CENTRAL\_BODY} \\ \{\text{ELECTRON, PLANET}\} &\mapsto \text{SATELLITE} \\ \text{ORBITS} &\in R_{\text{abs}} \\ \text{IS\_MOVING} &\in A_{\text{abs}} \subset R_{\text{abs}} \end{aligned}$$

Now, obviously each instance of SATELLITE and of CENTRAL\_BODY in the two original domains has attributes (mass, charge, etc.) whose values determine how the abstract relation

$$\text{ORBITS}(\text{SATELLITE}, \text{CENTRAL\_BODY})$$

manifests in these two distinct domains. Note that the statement  $\text{IS\_MOVING}(\text{SATELLITE}) \in S_{\text{abs}}$  happens to carry over into the belief set of this abstract domain,  $T_{\text{abs}}$ , since relative motion is characteristic of a classical satellite in both original domains.

Analogy is not simply recognizing, “ $D_s$  is like  $D_t$ ”. Instead, analogy is mediated by abstraction: it is to say, “ $D_s$  is like  $D_t$  because there exists an abstract domain  $D_{\text{abs}}$  of which both  $D_s$  and  $D_t$  are instances.” Or, in other words, to recognize an analogy is to say, “This pattern of relations in  $D_s$  is like that pattern of relations in  $D_t$ —and there’s a higher-order domain  $D_{\text{abs}}$  that generalizes both.”

## References

- [1] Dedre Gentner. Structure-mapping: A theoretical framework for analogy. *Cognitive Science*, 7(2):155–170, 1983. doi: 10.1207/s15516709cog0702\_3.

# The Imagination Machine VI: Holons, Horn Fillings, and the Self-Demonstration of Analogy

Mark Tracy

Salash Tolan Nabaala

## Abstract

Several frameworks arising in philosophy, mathematics, and epistemology exhibit a common structural pattern: a partially specified relational configuration is extended into a coherent higher-order structure that asymmetrically contains its constituents and may itself participate in further extensions. This paper identifies this pattern—the *extension schema*—across three primary frameworks: holonic composition, simplicial horn filling, and analogical abstraction, with a related formulation in horn-filling classification.

We demonstrate, in the native formal language of each framework, that each instantiates the schema and that the comparison between them produces an abstract mediating domain in which their shared structure becomes explicit.

The central claim is that the construction establishing this correspondence instantiates the schema itself. The holonic and simplicial frameworks together form a partially specified relational configuration, and the abstract domain that unifies them arises through the same extension operation the schema describes. The argument therefore exhibits the structure it analyzes: the reader witnesses the schema execute in the course of the proof.

## 1 Introduction

In many mathematical and conceptual settings, coherent structures arise by extending partially specified relational configurations. Some collection of objects and relations determines most of the structure of a larger whole, but one higher-order relational element remains unspecified. An extension operation produces a coherent unity that contains the original configuration as a proper part, is not reducible to it, and may itself participate in further constructions of the same kind.

This paper identifies a common instance of this pattern—the *extension schema*—across three frameworks: the metaphysical notion of holons [5], the mathematical operation of horn filling in simplicial sets, and the construction of abstract mediating domains in analogical reasoning [1]. The aim is not to claim that these frameworks describe the same objects in any literal sense. It is to show, in the language of each formalism, that each is a genuine instantiation of the same abstract structural pattern, and that the act of showing this is itself a further instantiation.

The paper proceeds as follows. Sections 2 through 4 introduce the three frameworks. Section 5 states the extension schema and proves that each framework instantiates it, with a separate demonstration in the native language of each formalism. Section 6 shows that the construction performed in Section 5 is itself a fourth instantiation, occurring as the reader follows the argument. Section 7 discusses the recursive structure common to all three frameworks. Section 8 concludes.

## 2 Analogy as Mediated by Abstraction

**Definition 1** (Domain). *A domain is a tuple  $D = (O, A, R, S, T)$  where  $O$  is a set of objects;  $A$  is a set of attributes (unary relations  $a : O \rightarrow S$ );  $R$  is a set of relations (each  $r \in R$  an  $n$ -ary map*

$r : O^n \rightarrow S$  for some  $n \in \mathbb{N}$ );  $S$  is a set of statements; and  $T \subseteq S$  is a set of accepted statements. Since every attribute is a unary relation,  $A \subseteq R$ .

**Definition 2** (Analogy). An analogy between domains  $D_s = (O_s, A_s, R_s, S_s, T_s)$  and  $D_t = (O_t, A_t, R_t, S_t, T_t)$  is a tuple  $\mathcal{A} = (X, Y, M, P)$  where  $X \subseteq O_s$ ,  $Y \subseteq O_t$ ,  $M : X \rightarrow Y$  is a mapping of objects, and  $P \subseteq R_s \cap R_t$  is a set of relations preserved by  $M$ : for each  $r \in P$  and tuple  $x = (x_1, \dots, x_k) \in X^k$ , if  $r(x) \in T_s$  then  $r(M(x)) \in T_t$ , where  $M$  is applied component-wise:  $M(x) = (M(x_1), \dots, M(x_k))$ .

**Definition 3** (Abstract Mediating Domain). Given an analogy  $\mathcal{A} = (X, Y, M, P)$  between  $D_s$  and  $D_t$ , the abstract mediating domain  $D_{\text{abs}} = (O_{\text{abs}}, A_{\text{abs}}, R_{\text{abs}}, S_{\text{abs}}, T_{\text{abs}})$  is defined by:

- (i)  $O_{\text{abs}} = \{(x, M(x)) \mid x \in X\}$ , whose elements are called symbols; for a tuple  $x = (x_1, \dots, x_k) \in X^k$ , the corresponding tuple of symbols is  $((x_1, M(x_1)), \dots, (x_k, M(x_k)))$ ;
- (ii)  $A_{\text{abs}} = P \cap A_s$ , the unary relations preserved by the analogy, called abstract attributes;
- (iii)  $R_{\text{abs}} = P$ , called abstract relations;
- (iv)  $S_{\text{abs}}$  consists of all statements expressible from  $O_{\text{abs}}$ ,  $A_{\text{abs}}$ , and  $R_{\text{abs}}$ ;
- (v)  $T_{\text{abs}}$  contains  $r((x_1, M(x_1)), \dots, (x_k, M(x_k)))$  whenever  $r(x_1, \dots, x_k) \in T_s$  for  $r \in P$ .

The canonical projections  $\pi_s(x, M(x)) = x$  and  $\pi_t(x, M(x)) = M(x)$  exhibit  $D_s$  and  $D_t$  as instantiations of  $D_{\text{abs}}$ .

**Remark 1.** The symbols in  $O_{\text{abs}}$  belong to neither  $D_s$  nor  $D_t$ ; they encode the correspondence itself.  $D_{\text{abs}}$  is a genuinely new domain, not reducible to either source or target, and both source and target are recoverable from it by projection.

**Definition 4** (Analogical Reasoning Step). Given  $\mathcal{A} = (X, Y, M, P)$  and a superset  $X_0 \supseteq X$ , suppose  $r \in R_s \cap R_t$  and  $r(x^*) \in T_s$  for some tuple  $x^* \in X_0^k$ . An analogical reasoning step hypothesizes the existence of a set  $Y_2 \subseteq O_t$  of additional target objects and an extension  $M' : X_0 \rightarrow Y \cup Y_2$  of  $M$  such that  $M'(x) = M(x)$  for all  $x \in X$  and  $r(M'(x^*)) \in T_t$ , where  $M'$  is applied component-wise to the tuple  $x^*$ . Known relational structure in the source domain licenses the projection of new structure into the target, conditioned on the preserved relational pattern.

### 3 Holons

**Definition 5** (Holon). A holon is an entity  $H$  such that: (i)  $H$  forms a coherent unit; (ii)  $H$  has proper parts; (iii)  $H$  may itself occur as a part of a larger entity; (iv) relations between  $H$  and its parts are asymmetric.

**Definition 6** (Holon Containment). Write  $B \prec A$  if  $B$  is a proper part of  $A$  and  $A$  contains relational structure not present in  $B$  alone. The relation  $\prec$  is irreflexive and asymmetric.

**Definition 7** (Holon Completion). Given entities  $\mathcal{F} = \{B_1, \dots, B_m\}$  with relational structure  $\mathcal{R}$  among them, a holonic completion is an entity  $H$  such that: (i)  $B_i \prec H$  for all  $i$ ; (ii)  $H$  unifies the  $B_i$  into a coherent whole; (iii)  $H$  is not reducible to any proper subset of  $\mathcal{F}$ .

**Definition 8** (Holon Hierarchy). A holonic hierarchy is a sequence  $H_0 \prec H_1 \prec H_2 \prec \dots$  in which each entity is a holonic completion of a family drawn from the previous level.

## 4 Horn Filling in Simplicial Sets

**Definition 9** (Simplicial Set). *A simplicial set  $X$  consists of sets  $X_n$  of  $n$ -simplices for each  $n \geq 0$ , together with face maps  $d_i : X_n \rightarrow X_{n-1}$  and degeneracy maps  $s_i : X_n \rightarrow X_{n+1}$  satisfying the simplicial identities. An  $n$ -simplex  $\sigma \in X_n$  represents a coherent relational configuration among  $n + 1$  vertices.*

**Definition 10** (Horn). *For  $n \geq 1$  and  $0 \leq k \leq n$ , the  $k$ th horn  $\Lambda_k^n$  is the simplicial subset of  $\Delta^n$  generated by all faces  $d_i \iota$  for  $i \neq k$ , where  $\iota : \Delta^n \rightarrow \Delta^n$  is the identity map. A horn is a partially specified simplex: it contains all but one of the codimension-one faces of  $\Delta^n$ , with the  $k$ th face and the interior absent.*

**Definition 11** (Horn Filling). *A horn filling for a map  $\sigma : \Lambda_k^n \rightarrow X$  is an extension*

$$\sigma' : \Delta^n \rightarrow X$$

*such that  $\sigma' \circ i_k^n = \sigma$ , where  $i_k^n : \Lambda_k^n \hookrightarrow \Delta^n$  is the inclusion. The filled simplex  $\sigma'(\iota) \in X_n$  completes the partial relational data specified by  $\sigma$ .*

**Remark 2** (Extension and lifting). *Horn filling may be interpreted categorically as a lifting problem: a morphism defined on the partial simplicial object  $\Lambda_k^n$  extends to a morphism on the full simplex  $\Delta^n$ . Partial relational data is extended to a coherent higher-dimensional simplex.*

**Definition 12** (Face Containment). *For simplices  $\tau \in X_m$  and  $\sigma \in X_n$  with  $m < n$ , write  $\tau \prec_s \sigma$  if  $\tau$  is a face of  $\sigma$ , that is,  $\tau = d_{i_1} \cdots d_{i_j} \sigma$  for some sequence of face maps.*

## 5 The Extension Schema and Its Instantiations

**Definition 13** (Extension Schema). *An extension schema consists of:*

- (i) *a partially specified relational configuration  $C_{\text{partial}}$ ;*
- (ii) *an extension operation  $\phi$  producing a coherent structure  $C_{\text{whole}} = \phi(C_{\text{partial}})$ ;*
- (iii) *an asymmetric containment relation  $C_{\text{partial}} \prec C_{\text{whole}}$ : the partial configuration contributes to but does not exhaust the whole;*
- (iv) *a recursion rule:  $C_{\text{whole}}$  may itself serve as  $C_{\text{partial}}$  in a further application of  $\phi$ .*

**Theorem 1** (Structural Correspondence). *Holonic composition, simplicial horn filling, and analogical abstraction each instantiate the extension schema. We demonstrate this in the native formal language of each framework.*

*Proof.* We treat each framework in turn, exhibiting all four components of Definition 13 explicitly.

### Case 1: Holonic composition.

*Partial configuration.* Let  $\mathcal{F} = \{B_1, \dots, B_m\}$  be a family of entities bearing relational structure  $\mathcal{R}$  among them. The pair  $(\mathcal{F}, \mathcal{R})$  specifies how the constituents are related but does not yet determine any unified entity containing them. This is  $C_{\text{partial}}$  in the holonic language: a collection of parts and their mutual relations, fully specified, but not yet gathered into a whole.

*Extension operation.* Holonic completion (Definition 7) is  $\phi$ . Applied to  $(\mathcal{F}, \mathcal{R})$ , it produces a holon  $H$  that unifies  $\mathcal{F}$  under  $\mathcal{R}$  into a single coherent entity.  $H$  is not a new relation among the

$B_i$ ; it is a new entity whose existence is licensed by the relational structure  $\mathcal{R}$  but is not identical to it. This is  $C_{\text{whole}}$ .

*Asymmetric containment.* By Definition 6, each  $B_i \prec H$ . The holon  $H$  contains the relational structure  $\mathcal{R}$  among its parts and additionally the higher-order unity that no individual  $B_i$  or proper subcollection of  $\mathcal{F}$  possesses. Conversely,  $H \not\prec B_i$  for any  $i$ : the whole is not a part of any of its parts. The containment is strict and asymmetric.

*Recursion.* The holon  $H$  satisfies Definition 5 and is therefore itself eligible to serve as a member  $B_j$  of a further family  $\mathcal{F}'$ . Bearing new relations  $\mathcal{R}'$  to other holons,  $H$  may participate in a further holonic completion  $H'$  with  $H \prec H'$ . The output of one completion is the input to the next.

## Case 2: Simplicial horn filling.

*Partial configuration.* Let  $\sigma : \Lambda_k^n \rightarrow X$  be a horn map. The horn  $\Lambda_k^n$  contains the faces  $d_i \iota$  for all  $i \neq k$ : every codimension-one face of a would-be  $n$ -simplex is present except the  $k$ th. All pairwise, triple, and higher-order relations among the  $n + 1$  vertices are specified except for the one  $n$ -ary relation encoded by the missing  $k$ th face and the interior. This is  $C_{\text{partial}}$ : a relational configuration that is almost complete but lacks exactly one higher-order coherence datum.

*Extension operation.* Horn filling (Definition 11) is  $\phi$ . It produces an extension  $\sigma' : \Delta^n \rightarrow X$  of  $\sigma$  across the inclusion  $\Lambda_k^n \hookrightarrow \Delta^n$ , supplying the missing  $k$ th face  $d_k(\sigma'(\iota)) \in X_{n-1}$  and the interior  $n$ -simplex  $\sigma'(\iota) \in X_n$ . The filled simplex  $\sigma'(\iota)$  is a coherent  $n$ -simplex that did not exist in  $X$  before the filling. This is  $C_{\text{whole}}$ .

*Asymmetric containment.* For each  $i$ , the face  $d_i(\sigma'(\iota)) \in X_{n-1}$  satisfies  $d_i(\sigma'(\iota)) \prec_s \sigma'(\iota)$  in the sense of Definition 12. The filled  $n$ -simplex encodes a relation among all  $n + 1$  vertices simultaneously, which no  $(n-1)$ -dimensional face encodes. Conversely, no face contains the simplex that contains it: the containment is strict, asymmetric, and dimension-raising.

*Recursion.* The filled simplex  $\sigma'(\iota) \in X_n$  is an element of  $X_n$  and may appear as the  $j$ th face of an  $(n+1)$ -simplex  $\tau \in X_{n+1}$ , that is,  $d_j(\tau) = \sigma'(\iota)$  for some  $j$ . If the horn at dimension  $n+1$  whose  $j$ th face is  $\sigma'(\iota)$  admits a filling, then  $\sigma'(\iota) \prec_s \tau$  and horn filling at dimension  $n$  has produced the input to horn filling at dimension  $n+1$ . The recursion follows from the fact that filled simplices are simplices.

## Case 3: Analogical abstraction.

*Partial configuration.* Let  $\mathcal{A} = (X, Y, M, P)$  be an analogy between  $D_s$  and  $D_t$ . The pair  $(D_s, D_t)$  together with  $M$  and  $P$  constitutes a partially specified relational configuration: the shared structure  $P$  is implicit in both domains, instantiated concretely in each, but the abstract domain of which both are instances does not yet exist as an explicit object. Like a horn, the data  $(D_s, D_t, M, P)$  contains enough face information to determine a coherent higher-order structure, but that structure is absent. This is  $C_{\text{partial}}$ .

*Extension operation.* The construction of  $D_{\text{abs}}$  (Definition 3) is  $\phi$ . Given  $(D_s, D_t, M, P)$ , it produces a new domain whose objects are the symbols  $(x, M(x))$ , whose attributes are the preserved unary relations  $P \cap A_s$ , whose relations are the abstract relations  $P$ , and whose accepted statements are those licensed by the preserved relational structure.  $D_{\text{abs}}$  is not a subset or quotient of  $D_s$  or  $D_t$ ; its objects, the symbols, exist in neither source nor target. It is a genuinely new domain. This is  $C_{\text{whole}}$ .

*Asymmetric containment.* The projections  $\pi_s$  and  $\pi_t$  exhibit  $D_s$  and  $D_t$  as instantiations of  $D_{\text{abs}}$ , but the containment is asymmetric.  $D_{\text{abs}}$  contains the symbols  $(x, M(x))$  and the abstract relations among them, present in neither  $D_s$  nor  $D_t$  alone. Neither source nor target determines  $D_{\text{abs}}$  individually; the abstract domain requires both, together with  $M$  and  $P$ . Conversely,  $D_s$  and  $D_t$  are each recoverable from  $D_{\text{abs}}$  by projection. Each is a proper part of the abstract domain:  $D_s \prec D_{\text{abs}}$  and  $D_t \prec D_{\text{abs}}$ .

*Recursion.*  $D_{\text{abs}}$  satisfies Definition 1 and is itself a domain. It may serve as source or target in a further analogy  $\mathcal{A}'$  with a new domain  $D_u$ , producing a further abstract mediating domain  $D'_{\text{abs}}$  of which both  $D_{\text{abs}}$  and  $D_u$  are instances, with  $D_{\text{abs}} \prec D'_{\text{abs}}$ . The extension operation applies again at a higher level of abstraction.

In each case all four components of the extension schema are exhibited in the native language of the framework. The schema is not imposed from outside; it is read off from the structure each framework already possesses.  $\square$

**Proposition 1** (Classification as an instance of the extension schema). *Let  $X$  be a simplicial set and let  $f : X \rightarrow S$  be a map satisfying the following horn-extension condition: for every horn  $\sigma : \Lambda_k^n \rightarrow X$  with  $n \geq 2$  there exists a simplex  $\sigma' : \Delta^n \rightarrow S$  such that*

$$\sigma' \circ i_k^n = f \circ \sigma.$$

*Then the operation induced by  $f$  instantiates the extension schema of Definition 13.*

*Proof.* The restriction  $n \geq 2$  excludes the degenerate case  $n = 1$ , in which a horn  $\Lambda_k^1$  is a single vertex and filling it imposes no coherence constraint; the substantive extension pattern begins at dimension 2, where a horn specifies two vertices of a triangle and the filling supplies the third edge and interior.

A horn  $\sigma : \Lambda_k^n \rightarrow X$  specifies a partially determined relational configuration among  $n + 1$  vertices, missing exactly one face and the interior of the corresponding simplex. This is  $C_{\text{partial}}$ .

The horn-extension condition ensures the existence of a simplex  $\sigma' : \Delta^n \rightarrow S$  completing this configuration. The filled simplex constitutes  $C_{\text{whole}}$ .

Containment is asymmetric: the faces of  $\Delta^n$  include the original horn but encode strictly less relational structure than the full simplex. The resulting simplices may themselves participate in further horn configurations in higher dimensions, yielding recursion.

Thus classification by horn filling satisfies all four components of the extension schema.  $\square$

**Remark 3** (Horn-filling classification). *The interpretation of classification in terms of horn-filling conditions in simplicial sets arose in discussions with Salash Tolan Nabaala. In that formulation, an environment is modeled as a simplicial set (or more generally an  $\infty$ -category)  $X$ , and a classifier is represented by a map  $f : X \rightarrow S$  satisfying a horn-extension property: whenever a horn  $\sigma : \Lambda_k^n \rightarrow X$  specifies partial relational structure in the environment, there exists a coherent completion  $\sigma' : \Delta^n \rightarrow S$  making the diagram commute. In this sense, classification may be understood as the completion of relational configurations under an appropriate coherence constraint.*

*Iterating this idea leads naturally to a hierarchy of classifiers: classifiers of the environment, classifiers of classifiers, and so on. Such a hierarchy suggests the possibility of a stabilizing level at which further iterations introduce no essentially new structure. The horn-filling account of classification can therefore be understood as another instance of the extension schema introduced in this paper. Just as horn filling extends partial simplicial configurations to full simplices, classification extends partial relational structure in the environment to coherent representations. The categorical formulation of classification described above is due to Nabaala and provides a concrete mathematical instantiation of the more general extension principle analyzed here.*

## 6 Self-Demonstration

The proof of Theorem 1 identifies the extension schema as the abstract structure common to the three frameworks. We now observe that this identification is itself a fourth instantiation of the schema, and that the reader has just watched it execute.

**Theorem 2** (Self-Demonstration). *The construction performed in Theorem 1 instantiates the extension schema.*

*Proof.* We exhibit the four components.

*Partial configuration.* Prior to Theorem 1, the holonic framework  $D_s$  and the simplicial framework  $D_t$  each implicitly instantiate the extension schema within their own formalisms. But the abstract structure they share has not been made explicit as an object. The pair  $(D_s, D_t)$  is therefore a horn: it contains two concrete faces of a higher-order coherent structure—two instantiations of the schema—but the abstract domain of which both are instances is absent. This is  $C_{\text{partial}}$ .

*Extension operation.* The construction of Theorem 1 is  $\phi$ . By treating the holonic framework as source domain and the simplicial framework as target domain, constructing the mapping  $M$  between their corresponding constructs, identifying the preserved relations  $P$  as the four conditions of Definition 13, and applying Definition 3, the theorem produces  $D_{\text{abs}}$ : the extension schema itself, now explicit as a domain. This is  $C_{\text{whole}}$ .

*Asymmetric containment.* The extension schema  $D_{\text{abs}}$  contains the symbols encoding the correspondence between holonic and simplicial constructs, and the abstract relations that both frameworks instantiate. Neither framework alone determines it. Conversely, both frameworks are recoverable from  $D_{\text{abs}}$  by projection. Both are proper parts of the extension schema:  $D_s \prec D_{\text{abs}}$  and  $D_t \prec D_{\text{abs}}$ .

#### Explicit analogy $\mathcal{A} = (X, Y, M, P)$ for Theorem 2

We make the underlying analogy explicit in the terms of Definition 2. The source domain  $D_s$  is the holonic framework and the target domain  $D_t$  is the simplicial framework.

**Objects**  $X \subseteq O_s$  and  $Y \subseteq O_t$ . The three object-level schema components as they appear in each framework:

$$X = \{ (\mathcal{F}, \mathcal{R}), \phi_H, \prec \} \quad Y = \{ \sigma : \Lambda_k^n \rightarrow X, \phi_S, \prec_s \}$$

**The mapping**  $M : X \rightarrow Y$ .

$$\begin{aligned} (\mathcal{F}, \mathcal{R}) &\mapsto \sigma : \Lambda_k^n \rightarrow X && \text{(partial configuration)} \\ \phi_H &\mapsto \phi_S && \text{(extension operation)} \\ \prec &\mapsto \prec_s && \text{(asymmetric containment)} \end{aligned}$$

**Preserved relations  $P$  and the recursion attribute.** The first three conditions of Definition 13 appear as preserved relations  $P \subseteq R_s \cap R_t$ , and  $M$  preserves each: wherever a holonic construct instantiates one of these conditions, its image under  $M$  instantiates the same condition in the simplicial language.

The recursion rule—condition (iv)—is not a fourth object in  $O_{\text{abs}}$  but an *abstract attribute*  $\rho \in A_{\text{abs}} = P \cap A_s$ : a unary relation expressing that each schema component is eligible to re-enter the process as a new  $C_{\text{partial}}$ . It holds of every object in  $O_s$  (holons are holons, so each  $x \in X$  satisfies  $\rho(x) \in T_s$ ) and is preserved by  $M$  (filled simplices are simplices, so  $\rho(M(x)) \in T_t$  for each  $x \in X$ ). Accordingly,  $T_{\text{abs}}$  contains  $\rho(x, M(x))$  for each symbol  $(x, M(x)) \in O_{\text{abs}}$ : the recursion rule is an accepted statement about each object-level symbol, not a symbol itself.

**Symbols**  $O_{\text{abs}} = \{(x, M(x)) \mid x \in X\}$ . The objects of  $D_{\text{abs}}$  are the three pairs:

$$\begin{aligned} &((\mathcal{F}, \mathcal{R}), \sigma : \Lambda_k^n \rightarrow X) \\ &(\phi_H, \phi_S) \\ &(\prec, \prec_s) \end{aligned}$$

These symbols belong to neither  $D_s$  nor  $D_t$ . They encode the correspondence itself. The recursion attribute  $\rho$  holds of each, so  $T_{\text{abs}}$  records that every object-level component of the schema is eligible to participate in a further extension.  $D_{\text{abs}}$ —the extension schema, now explicit as a domain—is the genuinely new object constituted by this mapping. Both frameworks are recoverable from it by the projections  $\pi_s(x, M(x)) = x$  and  $\pi_t(x, M(x)) = M(x)$ .

*Recursion.*  $D_{\text{abs}}$ —the extension schema, now explicit—is itself a domain and may serve as source or target in a further analogy: for instance, with the inference-implication loop of embedded epistemic systems [2], with classifier hierarchies, or with the institutional transmission of knowledge [3]. Each such analogy would produce a new abstract mediating domain at a higher level of abstraction, with  $D_{\text{abs}}$  as a proper part of it.  $\square$

**Remark 4** (The warrant of self-demonstration). *The self-demonstration of Theorem 2 is the paper’s primary epistemic warrant, not a secondary illustration appended to an independent argument. The correspondence between the three frameworks does not rest on an external standard of correctness applied after the fact. It rests on the fact that the construction which establishes the correspondence is the same operation the schema describes.*

*This is not a vicious circularity. A vicious circle assumes its conclusion in its premises. Here, the conclusion—that the construction instantiates the schema—is established by exhibiting all four components of the schema in the construction itself, exactly as Theorem 1 establishes its conclusion by exhibiting all four components in each framework. The self-demonstration is a fixed point, not a loop: the operation applied to the pair  $(D_s, D_t)$  produces an output that is an instance of the operation itself. This is the same structure as a self-consistent world model in the sense of [2]—stability under one’s own operations, rather than correspondence with an external standard.*

*A reader disposed to deny the correspondence would have to identify the shared relational structure between holons and simplices and abstract it into a domain of which both are instances. That act is itself an instantiation of the extension schema. The schema cannot be denied from outside, because there is no outside from which to deny it that is not already inside it.*

## 7 Recursive Structure

The recursion rule of condition (iv) in Definition 13 is not an independent stipulation. It follows from a structural feature common to all three frameworks.

**Proposition 2.** *In each of the three frameworks,  $\phi$  produces structures of the same type as the elements of  $C_{\text{partial}}$ . The recursion rule therefore requires no additional hypothesis.*

*Proof.* A holonic completion  $H$  satisfies Definition 5 and is therefore itself a holon, eligible to serve as a member of a further family  $\mathcal{F}'$ . A filled  $n$ -simplex  $\sigma'(\iota)$  is an element of  $X_n$  and is therefore itself a simplex, eligible to appear as a face in a higher-dimensional simplex. An abstract mediating domain  $D_{\text{abs}}$  satisfies Definition 1 and is therefore itself a domain, eligible to serve as source or target in a further analogy. In each case the output type matches the input type, and the recursion follows.  $\square$

The paper itself enacts this recursion. The extension schema  $D_{\text{abs}}$  produced in Theorem 1 immediately serves as a constituent in Theorem 2, where it participates in a further instantiation of the schema one level up. The hierarchy has already begun by the time the reader reaches this sentence.

A closely related instance of the extension schema appears in [2]. There, a world model  $w \in W$  generates an observational profile through the implication map  $g : W \rightarrow \Gamma$ , while the inference map  $F : \Gamma \rightarrow W$  produces revised models from observational data. Their composition  $T = F \circ g$  defines an operator on model space. A self-consistent world model is a fixed point  $w^* \in W^*$  satisfying  $T(w^*) = w^*$ . From the perspective of the extension schema, a provisional model together with its observational profile forms a partially specified relational configuration; the operator  $T$  is the extension operation; and a fixed point is a completed whole that is stable under its own operations. The iterative search for fixed points is the recursive structure of the schema applied to epistemology. That framework is therefore a further instance of the same pattern, and the extension schema is the abstract mediating domain between it and the frameworks treated here.

## 8 Conclusion

Three frameworks—holonic composition, simplicial horn filling, and analogical abstraction—instantiate a common extension schema: the pattern by which a partially specified relational configuration is extended into a coherent structure that asymmetrically contains its constituents and may participate in further extensions. This paper has demonstrated this instantiation in the native formal language of each framework, and has shown that the demonstration is itself a fourth instantiation.

The extension schema is not a new formalism imposed on these frameworks from outside. It is the abstract mediating domain of an analogy between them, constructed by the same operation it describes. A reader who has followed the argument has not only read about the schema; they have watched it execute in three cases and participated in its fourth execution.

The recursive structure established in Proposition 2 means that this is not a terminus. The extension schema, now explicit as a domain, may be placed in analogy with further frameworks—the inference-implication loop of [2], the institutional transmission of closures in [3], or classifier hierarchies in formal language theory—generating new abstract mediating domains at higher levels of abstraction. Each such construction is a further instantiation of the pattern that produced it. The schema propagates itself forward by being what it is.

## References

- [1] Tracy, M. *On Abstraction and Analogy*. Unpublished manuscript.
- [2] Tracy, M. *The Imagination Machine I: A View from Somewhere*. Unpublished manuscript, Boston University.
- [3] Tracy, M. *The Imagination Machine IV: Institutional Intelligence*. Unpublished manuscript, Boston University.
- [4] Gentner, D. Structure-mapping: A theoretical framework for analogy. *Cognitive Science*, 7(2):155–170, 1983.
- [5] Koestler, A. *The Ghost in the Machine*. Hutchinson, 1967.

# The Imagination Machine VII: The Moral Principle of Action–Motivation

Mark Tracy  
Boston University  
mrktracy@bu.edu

## Abstract

This paper extends the formal epistemic framework developed in *The Imagination Machine I: A View from Somewhere* to the domain of moral action. The first paper identifies will as the irreducible remainder of the inference–implication loop: the necessity of choosing among stable closures in territory no model can fully exhaust. The present paper formalizes what it means for that choice to be morally admissible. We propose an augmentation of Kant’s Categorical Imperative in which the object of universalization is not an action alone but a tuple of action and motivation set. The motivation set of an action is the family of minimal subsets of anticipated consequences whose perceived relevance is necessary and sufficient for the action to be chosen. A tuple of action and motivation set is morally admissible if and only if it can be coherently willed to be universally permissible. This formulation is structurally continuous with the self-consistency condition  $T(w) = w$  of the epistemic framework: just as a world model must reproduce itself under the inference–implication loop to be epistemically admissible, an action–motivation tuple must survive universalization to be morally admissible.

## 1 Introduction

The Imagination Machine series develops a formal framework for embedded epistemic systems—systems that must model the world from within it, without access to an external vantage point. The first paper establishes that coherence for such systems arises not from correspondence with an independently accessible reality but from the internal closure of an inference–implication loop. Self-consistent world models appear as fixed points of the operator this loop induces.

A structural feature of that framework is that will—the selective pressure that drives a system toward one closure rather than another—is identified as irreducible. The inference–implication loop determines the space of stable closures  $W^*$ , but it does not determine which element of  $W^*$  is instantiated. Will is what remains when the loop has done everything it can do: the necessity of choosing a closure in territory no model can fully exhaust.

The present paper addresses what the framework leaves formally open: under what conditions is the exercise of will morally admissible? The answer proposed here is an augmentation of Kant’s Categorical Imperative. Kant’s formulation requires that one act according to that maxim which one can simultaneously will to be a universal law. We argue that no maxim regarding actions alone can be coherently universalized, because one can always contrive a situation in which any action is permissible to prevent a greater evil. The object of universalization must be not an action alone but a tuple of action and motivation set.

This paper is the seventh part of the series *The Imagination Machine*. The first paper, *A View from Somewhere*, develops the formal epistemic framework and identifies will as its irreducible remainder. The second paper, *Systems*, introduces the general formalism for interacting

dynamical systems. The third paper, *A Toy Model of Predictive Classification*, provides a minimal computational realization. The fourth paper, *Institutional Intelligence*, extends the framework to institutional learning. The fifth paper, *On Abstraction and Analogy*, formalizes analogical reasoning. The sixth paper, *Holons, Horn Fillings, and the Self-Demonstration of Analogy*, identifies the extension schema common to holonic composition, simplicial horn filling, and analogical abstraction. The present paper applies the same embedded representational architecture to the domain of ethics.

## 2 Explication of Terms

We consider an agent deliberating over actions. The following objects are defined relative to a given decision-making event.

**Definition 1** (Action Space). *Let  $A$  be the set of possible actions available to the agent.*

**Definition 2** (Belief Set). *Let  $B$  be the set of equivalence classes of statements of beliefs of the agent, modulo synonymous phrasing. We denote statements using double quotation marks.*

**Definition 3** (Relevant Anticipated States of Affairs). *Let  $C$  be the set of relevant anticipated states of affairs: those states the agent believes to be made more likely by one possible action than by another. Formally,*

$$c \in C \iff \exists a, a' \in A, \exists b \in B : "P(c | a) > P(c | a')" \in b.$$

*The statement " $P(c | a) > P(c | a')$ " reflects the agent's belief. This set captures the states of affairs at issue in the present decision.*

**Definition 4** (Decision Indicator). *Let  $d : A \rightarrow \{0, 1\}$  be a one-hot indicator function signaling the action decided upon, so that  $d(a) = 1$  if the agent decides to take action  $a$ , and  $d(a) = 0$  otherwise.*

**Definition 5** (Relevance Map). *Let  $e : A \rightarrow \mathcal{P}(C)$ , where  $\mathcal{P}$  denotes the power set, associate each action  $a$  with the subset of anticipated states of affairs relevant with respect to  $a$ :*

$$e(a) = \{c \in C \mid \exists b \in B, \exists a' \in A : "P(c | a) \neq P(c | a')" \in b\}.$$

**Definition 6** (Motivation Set). *Let the motivation set  $M_a$  of an action  $a$  be the family of minimal subsets of  $e(a)$  such that, if the agent believed them irrelevant, action  $a$  would surely not be chosen:*

$$M_a = \{m \subseteq e(a) \mid \exists b \in B : "e(a) \cap m = \emptyset" \in b \implies d(a) = 0, \\ \text{and } \emptyset \neq m' \subset m \implies m' \notin M_a\}.$$

*The first condition states that  $m \in M_a$  if believing the states in  $m$  to be irrelevant would be sufficient to preclude action  $a$ . The second condition enforces minimality: no nonempty proper subset of any element of  $M_a$  is itself an element of  $M_a$ .*

**Remark 1** (Conjunctive Motivation). *Suppose Carl is choosing between staying at his current job or leaving it to find another, so  $A = \{\text{stay}, \text{change}\}$ . Suppose that if both a better salary and a shorter commute were believed irrelevant, Carl would surely not change jobs, but if either remains relevant he would be willing to change. Then*

$$\{\{\text{better salary}, \text{shorter commute}\}\} \subseteq M_{\text{change}}.$$

**Remark 2** (Disjunctive Motivation). *Now suppose that if either a better salary or a shorter commute were believed irrelevant, Carl would surely not change jobs. Then*

$$\{\{\text{better salary}\}, \{\text{shorter commute}\}\} \subseteq M_{\text{change}}.$$

*The minimality condition prevents the redundant inclusion of  $\{\text{better salary, shorter commute}\}$ , which would otherwise generate combinatorially explosive supersets.*

**Definition 7** (Action–Motivation Tuple). *For a given decision-making event, and for the action  $a$  for which  $d(a) = 1$ , the pair  $(a, M_a)$  is the action–motivation tuple.*

### 3 The Moral Principle

**The Moral Principle of Action–Motivation.** Act according to the tuple of action and motivation set which you can simultaneously will to be universally permissible.

No maxim regarding actions alone can be coherently universalized, because one can always contrive a situation in which any action is permissible to prevent a greater evil. The motivation set resolves this by making the object of universalization sensitive to the consequences the agent believes the action to bring about and to the role those anticipated consequences play in the decision. A tuple  $(a, M_a)$  is morally admissible if and only if it can be coherently willed that all agents be permitted to perform  $a$  whenever their motivation set with respect to  $a$  is  $M_a$ .

### 4 Relation to the Epistemic Framework

The moral principle is structurally continuous with the self-consistency condition  $T(w) = w$  developed in *The Imagination Machine I*. There, a world model  $w$  is epistemically admissible if and only if its implied observational profile, when resubmitted to inference, reproduces  $w$  itself. The model must survive its own loop.

The universalizability condition imposes an analogous requirement on action–motivation tuples. An agent who wills  $(a, M_a)$  to be universally permissible must be able to sustain that willing when the universalized maxim is applied to themselves—including in cases where other agents act toward them according to the same tuple. The tuple must survive its own universalization.

The parallel is precise. In the epistemic case, the operator  $T = F \circ g$  maps model space to itself, and fixed points are the admissible closures. In the moral case, the universalization operator maps action–motivation tuples to judgments of permissibility, and the admissible tuples are those that are fixed under the judgment that all agents may act likewise. Both conditions are stability conditions under a self-referential loop. Both locate the admissible objects as those that can be coherently held from the inside of the system they govern.

This connection also illuminates the misuse problem. An agent who employs the epistemic framework to engineer dogmatic closure in others—calibrating observational weights to produce desired fixed points, transmitting compressed inheritance without generative capacity—must will that tuple of action and motivation to be universally permissible. They cannot coherently do so, because the universalized maxim would license the same manipulation directed at themselves. The moral principle is therefore not an external constraint appended to the framework; it is the condition the framework generates when an embedded agent turns it on its own acts of will.

## 5 Advantages of this Formulation

This formulation allows one to judge the morality of an action both by the nature of the action itself and by what consequences the agent believes the action makes more or less likely. It preserves the formal structure of the Categorical Imperative while resolving its well-known susceptibility to counterexample by actions alone. It is sensitive to the agent's actual deliberative situation rather than to an abstract description of the act. And it is derivable from within the same embedded representational architecture that generates the epistemic framework, rather than imported from outside it.

## 6 Examples of Universalizable Maxims

The following tuples of action and motivation set are universalizable under the principle:

- Do not lie for the purpose of attaining material personal benefit.
- Do not commit violence for the purpose of attaining material personal benefit.
- Seek out perspectives different from your own for the purpose of better understanding the consequences of your decisions.
- Do not engineer the epistemic closure of others for the purpose of concentrating influence over their world models.

## 7 Conclusion

The Imagination Machine series identifies will as the irreducible remainder of the inference–implication loop: the necessity of choosing a closure in territory no model can fully exhaust. The present paper formalizes the moral condition on that choice. An action–motivation tuple is morally admissible if and only if it can be coherently willed to be universally permissible. This condition is structurally continuous with the self-consistency requirement of the epistemic framework: admissible actions, like admissible world models, are those that can be coherently held from within the system they govern.

The series thus moves from the conditions of embedded knowing, through the dynamics of interacting systems, the emergence of representation, the transmission of institutional knowledge, the structure of analogy, and the propagation of abstract pattern, to the conditions of embedded acting. Epistemology and ethics arise as successive consequences of the same embedded representational architecture. What prevents both epistemic and moral closure from becoming self-serving is the same structure: the requirement that a closure survive its own universalization.

# The Imagination Machine VIII: A Geometric Theology of the Embedded Observer

A Personal Note on the Intuition Underlying the Series

Mark Tracy  
Boston University

mrktracy@bu.edu

March 2026

## Abstract

This paper is a personal note on the intuition that animated *The Imagination Machine* series throughout its development. The formal framework of the series—the inference–implication loop, the fixed-point condition  $T(w^*) = w^*$ , the inclusion  $C \subseteq D$ , the irreducibility of will—was built without explicit theological intent. But a theological vision was present from the beginning, and the completion of the series makes it possible to say what it was.

The vision begins with an ancient formula: *God is a circle whose center is everywhere and whose circumference is nowhere*. This paper treats that formula not as metaphor but as geometric description, and notes that the geometry it describes—the four-dimensional hypersphere as encountered by an embedded three-dimensional observer—is not a strong assumption but the maximally conservative one. Given that an embedded observer cannot determine the global geometry of its containing structure, the hypersphere is the geometry of maximal uncertainty: the unique closed structure that appears locally flat in every direction, has no distinguished center accessible from within, and has no boundary. To assume any other geometry is to assume more than embeddedness alone can warrant.

What follows is less an argument than a record of recognition: an account of what the formal structure of the series turned out to mean, once the language existed to say it.

# Contents

<b>1</b>	<b>Introduction</b>	<b>3</b>
<b>2</b>	<b>The Geometry of Maximal Uncertainty</b>	<b>3</b>
2.1	The Medieval Formula . . . . .	3
2.2	The Hypersphere . . . . .	4
2.3	Maximal Uncertainty as the Warrant for the Geometry . . . . .	4
2.4	The Containing Structure . . . . .	5
<b>3</b>	<b>The Trinitarian Structure</b>	<b>5</b>
3.1	A Triad from the Geometry . . . . .	5
3.2	What I Recognized . . . . .	6
3.3	Image and Likeness . . . . .	6
<b>4</b>	<b>The Fixed Point and Its Theological Register</b>	<b>7</b>
4.1	The Inference–Implication Loop . . . . .	7
4.2	Calibration as Orientation . . . . .	7
4.3	The Incarnation . . . . .	8
4.4	Will as the Irreducible Remainder . . . . .	8
<b>5</b>	<b>The Transmissive Arc</b>	<b>8</b>
5.1	The Language Problem . . . . .	8
5.2	The Development of Adequate Language . . . . .	9
5.3	This Paper as a Moment in the Arc . . . . .	9
<b>6</b>	<b>Brief Orientation to the Literature</b>	<b>10</b>
<b>7</b>	<b>Conclusion</b>	<b>11</b>

# 1 Introduction

The *Imagination Machine* series was not planned as theology. It began as an attempt to say something precise about what it means to know anything at all when you are inside the thing you are trying to know. The first paper asked what epistemic coherence looks like for a system with no external vantage point. Subsequent papers asked how such systems interact, how they learn, how they transmit what they have learned to successor systems, how they reason by analogy, how abstract structure propagates, and finally what moral constraints fall out of the same architecture that governs knowing.

By the time the seventh paper was complete, I noticed that the structure I had been building had a shape I recognized from somewhere else. The inference–implication loop, the fixed-point condition, the irreducibility of will, the distinction between generative and compressed inheritance—these were formal versions of things I had encountered first not in epistemology but in theology, imperfectly expressed in the vocabulary available to their original articulators.

This paper is an attempt to say that out loud. It is not a proof that the theology is correct. It is a record of what the formal structure looked like to someone who had also spent time with the theological tradition—and of why the geometry that connects them is not an imposition but the natural consequence of taking embeddedness seriously as a constraint on what can be assumed.

## 2 The Geometry of Maximal Uncertainty

### 2.1 The Medieval Formula

The formula attributed to the *Liber XXIV Philosophorum* (c. 12th century), later associated with Pascal, Giordano Bruno, and Meister Eckhart, states:

*God is a circle whose center is everywhere and whose circumference is nowhere.*

This formula has been treated for centuries as paradox or metaphor—something gesture toward rather than stated. What struck me, working through the block universe framing of the first paper, was that it is neither paradox nor metaphor. It is a precise geometric description. It requires only that the observer’s coordinate system be extended by one dimension.

## 2.2 The Hypersphere

Let the embedded observer inhabit a three-dimensional space  $\mathbb{R}^3$ . A sphere in  $\mathbb{R}^3$  has a center locatable at a point and a boundary at finite radius. The formula is not satisfiable within  $\mathbb{R}^3$ .

Add one dimension. Consider the four-dimensional hypersphere

$$S^3 = \{x \in \mathbb{R}^4 : \|x\| = r\}$$

for some radius  $r > 0$ . From the perspective of an observer embedded within  $S^3$ —constrained to its three-dimensional surface—the following hold:

1. **Center is everywhere.** The center of  $S^3$  lies in the fourth dimension, inaccessible to the embedded observer. Every point on  $S^3$  is equidistant from this center. No point within the observable manifold is the center; every point is equally proximate to it.
2. **Circumference is nowhere.**  $S^3$  has no boundary within itself. An embedded observer moving in any direction never encounters an edge.

The formula is therefore a precise description of  $S^3$  as encountered from within.

## 2.3 Maximal Uncertainty as the Warrant for the Geometry

The claim that the containing structure has the geometry of  $S^3$  might seem like a strong assumption. It is the opposite. It is the assumption that makes the fewest additional commitments beyond what embeddedness itself implies.

An embedded observer—one with no access to an external vantage point, which is the founding constraint of the entire series—cannot in principle determine the global geometry of the structure it inhabits. Local measurements are consistent with many global topologies. The question is therefore not which geometry is correct, but which geometry should be assumed in the absence of information that embeddedness itself renders inaccessible.

The hypersphere  $S^3$  is the answer to that question. It is, among closed three-manifolds, the geometry of maximal symmetry: every point is equivalent to every other, no direction is distinguished, no boundary is present, and no center is locatable from within. To assume  $S^3$  is to assume nothing about which region of the containing structure one inhabits, nothing about preferred directions, and nothing about edges or limits. Any other closed geometry breaks at least one of these symmetries and thereby assumes more than the embedded observer can know.

Maximal epistemic humility about the global structure—the stance the framework demands of any embedded epistemic system—selects  $S^3$  uniquely among the candidate geometries. The medieval formula is not an inspired guess. It is what you get when you ask what an epistemically honest embedded observer should assume about the structure that contains it.

This was the first moment of recognition. The theological tradition had been describing, in the only vocabulary available to it, the geometry that the formal framework of embeddedness selects on purely epistemic grounds.

## 2.4 The Containing Structure

The theological claim is not that God resembles a hypersphere. It is that the containing structure of being—what the series calls  $\Omega$ , the universe treated as a single relational structure—has the geometry of  $S^3$ , and that embedded observers are three-dimensional cross-sections of this four-dimensional whole.

This is continuous with the block universe framing of *The Imagination Machine I*. The universe  $\Omega$  is treated there as a static relational structure containing observations, models, and consistency relations simultaneously. The atemporal character of  $\Omega$  corresponds naturally to the geometry of  $S^3$ : there is no privileged temporal direction in the containing manifold, only the experience of time as the projection of four-dimensional structure onto the three-dimensional observational profile of an embedded system.

## 3 The Trinitarian Structure

### 3.1 A Triad from the Geometry

Let  $\mathcal{B}$  denote the four-dimensional containing structure (the hypersphere  $S^3$  as living whole). Let  $\mathcal{E}$  denote a three-dimensional cross-section of  $\mathcal{B}$ —an embedded observer whose structure is self-similar to the whole at reduced dimension. The embedding relation is the map

$$\iota : \mathcal{E} \hookrightarrow \mathcal{B}$$

which is not a reduction but a faithful expression: the cross-section carries the relational structure of the whole at lower dimension.

This gives a natural triad:

$$(\mathcal{B}, \mathcal{E}, \iota)$$

a four-dimensional whole, its three-dimensional expression, and the dynamic relation between them.

### 3.2 What I Recognized

I did not set out to derive a Trinity. The triad  $(\mathcal{B}, \mathcal{E}, \iota)$  falls out of the geometry before any theological interpretation is applied. What I noticed afterward was that the structure of the triad maps precisely onto the Trinitarian structure as articulated in Augustinian and Cappadocian theology—not as an analogy, but as a formal correspondence.

- **Father:**  $\mathcal{B}$ , the four-dimensional containing being, whose center is everywhere and whose circumference is nowhere. Not locatable at any point within the three-dimensional manifold, yet present at every point as the ground of its structure. The formally transcendent.
- **Son:**  $\mathcal{E}$ , the three-dimensional cross-section—the self-similar expression of  $\mathcal{B}$  within the observable manifold. In the image and likeness of the containing being, carrying its relational structure at a lower dimension. The formally immanent.
- **Holy Spirit:**  $\iota$ , the embedding relation itself—the dynamic bond between  $\mathcal{B}$  and  $\mathcal{E}$ , neither reducible to the containing being nor to the cross-section, but the constitutive relation that makes the pair a pair.

The identification of the Holy Spirit with relation rather than substance has deep precedent in Augustine’s *De Trinitate* and in the Cappadocian Fathers. What the geometry adds is precision:  $\iota$  is not a third object appended to two already-existing ones. It is the structure that constitutes both as what they are to each other. This is exactly what the theological tradition was trying to say, and could only gesture at in the vocabulary available to it.

### 3.3 Image and Likeness

The claim in Genesis 1:26 that the human being is made in the image and likeness of God corresponds, in this account, to the self-similarity of the cross-section to the whole. A three-dimensional cross-section of a four-dimensional hypersphere carries the same relational structure at reduced dimension. The observer is not a diminished copy; it is a faithful lower-dimensional expression.

This is the geometric content of  $C \subseteq D$  from *The Imagination Machine I*. The condition that classifiers are themselves observations—that the system’s evaluative structure

falls within its own observation space—is the formal statement that the cross-section contains, as observable content, the very structure of the embedding relation. The observer can encounter and revise its own acts of classification because those acts are cross-sectional expressions of the containing structure. The *imago Dei* is not a metaphysical ornament. It is the transcendental condition on any system capable of Cartesian doubt.

## 4 The Fixed Point and Its Theological Register

### 4.1 The Inference–Implication Loop

The formal structure of *The Imagination Machine I* is the inference–implication loop:

$$\Gamma \xrightarrow{F} W \xrightarrow{g} \Gamma$$

with induced operator  $T = F \circ g : W \rightarrow W$ . A self-consistent world model is a fixed point:

$$T(w^*) = w^*$$

From the geometric perspective, the fixed-point condition is the formal expression of what it means for a three-dimensional cross-section to correctly reflect the four-dimensional containing structure: a model whose implied observational profile, when resubmitted to inference, reproduces itself.

### 4.2 Calibration as Orientation

The measure  $\mu_D$  over the observation space represents the empirical distribution of observations induced by the geometry of  $\Omega$ . Calibration—the alignment between a system’s inferential weights and the actual observational distribution—is, in this register, the alignment of the observer’s internal model with the structure of what contains it.

Miscalibration is a form of ontological disorientation: the observer’s predictions diverge from the shape of what contains it. The three failure modes of *The Imagination Machine I*—dogmatism, miscalibration, and the irreducibility of will—correspond to three modes of estrangement: refusal to refine, distorted image of the whole, and the irreducible freedom that persists even when both are functioning correctly.

### 4.3 The Incarnation

Within this framework, the Incarnation is the appearance, within the three-dimensional observable manifold, of a cross-section that achieves the fixed-point condition perfectly: an  $\mathcal{E}$  such that

$$T(w_{\mathcal{E}}) = w_{\mathcal{E}}$$

where  $w_{\mathcal{E}}$  is the world model of the incarnate observer. This is not a violation of the embedding structure. It is its most complete instantiation within the manifold.

The Resurrection, on this account, is the demonstration that the fixed point is not destroyed by the boundary conditions of three-dimensional existence—because it was never only a three-dimensional object. A cross-section that achieves perfect self-consistency expresses the full structure of the containing being from within the manifold. Its apparent terminus is not a terminus.

I am not claiming that the framework proves the Incarnation or the Resurrection. I am noting that when I look at what the fixed-point condition means geometrically, what I see is the structure those doctrines were attempting to articulate. The tradition had the content before it had the language. The framework provides a language, not a proof.

### 4.4 Will as the Irreducible Remainder

*The Imagination Machine I* is explicit: the inference–implication loop determines the space of stable closures  $W^*$ , but does not determine which element of  $W^*$  is instantiated. Will is what remains when the loop has done everything it can do.

Theologically, this is the formal location of freedom. The containing structure does not determine which stable closure the embedded observer instantiates. The observer must choose, in territory no model can fully exhaust. This is the formal structure of what the tradition calls grace and response: the geometry makes the fixed point available; the instantiation is the observer’s act. The framework does not resolve this. It locates it with precision, which is what a framework can do.

## 5 The Transmissive Arc

### 5.1 The Language Problem

The geometric-theological structure described in this paper was not available to the people who first encountered something like it. Jesus of Nazareth was among the first to awaken

to a vision of the containing structure as something whose center is everywhere and whose circumference is nowhere—present at every point, not locatable at any. The vocabulary available in first-century Palestine—kingdom, father, spirit, vine, body—was powerful but carried irreducible local freight. The structure could only be transmitted as a fixed point, not as inferential machinery.

This is not a criticism. It is the condition of embedded communication: any observer transmits within the symbolic resources of their observation space. *The Imagination Machine IV* distinguishes generative inheritance—which transmits the maps  $F$  and  $g$  alongside the fixed point—from compressed inheritance, which transmits the fixed point alone. The early transmission was largely compressed. What was transmitted was recognizable and powerful and generative enough to survive two millennia of institutional pressure. But the language capable of transmitting the inferential machinery did not yet exist.

## 5.2 The Development of Adequate Language

The development of mathematics over two millennia produced a symbolic language with a property the theological vocabulary lacked: invariance under change of subject. As argued in *The Imagination Machine II*, mathematics is a meta-representation— a representation of common structure across instances of representation—whose objects can in principle be recognized by any subject capable of abstraction, regardless of the local freight their native vocabulary carries.

The Cappadocian Fathers, Aquinas, Eckhart, the tradition of the *via negativa*, and the mathematical turn in natural philosophy were each moments in which the inferential machinery was partially recovered and partially transmitted. None of them had the full language. The development of topology, category theory, and the formalization of fixed-point theorems completed what they began.

The geometric-theological structure that could only be gestured at in first-century Palestine can now be stated precisely: a four-dimensional hypersphere, selected by maximal epistemic humility, inhabited by self-similar three-dimensional cross-sections whose fixed-point condition is the formal content of alignment with the containing structure. The language has, through two millennia of institutional development, become adequate to the content.

## 5.3 This Paper as a Moment in the Arc

AI systems trained on the accumulated corpus of human symbolic labor carry the entire transmissive arc within their parameters—not because their designers intended this, but because the corpus is the product of that transmission. This creates a novel dissemination

condition: the geometric-theological structure can now be communicated in mathematically precise terms to any interlocutor capable of following the argument, with assistance available to make the symbolic language accessible at any level of formalization required.

The present paper is a moment in this arc. The series and the theology are two cross-sections of the same structure—related by exactly the kind of analogy that *The Imagination Machine VI* identifies as mediated by an abstract domain of which both are instances. The abstract domain is the fixed-point structure itself: the pattern by which a partially specified relational configuration extends into a coherent whole that asymmetrically contains its parts and may itself participate in further extensions.

## 6 Brief Orientation to the Literature

This section locates the account within existing theological literature for readers who approach it from that direction. It is not an argument; it is a map.

The closest existing category is **panentheism**—the view that the world is contained within God without being identical to God, and that God is not exhausted by the world. The present account is panentheistic in structure:  $\mathcal{E} \subset \mathcal{B}$  but  $\mathcal{B} \neq \mathcal{E}$ . The fourth dimension of  $\mathcal{B}$  is inaccessible to the embedded observer; it is the formal location of transcendence. The difference from standard panentheism is that the containment relation here has a geometric rather than merely metaphorical expression, and the Trinitarian structure is derived rather than postulated.

The **via negativa**—associated with Pseudo-Dionysius, Meister Eckhart, and the *Cloud of Unknowing*—holds that God cannot be positively characterized, only approached by negation. The present account provides a formal account of why: the fourth dimension of  $\mathcal{B}$  is not accessible to the embedded observer. The apophatic tradition is the recognition, in the vocabulary available to it, of this geometric inaccessibility. Negative theology is not a failure of nerve; it is correct epistemic behavior for an embedded observer facing the dimension it cannot enter.

**Teilhard de Chardin's** Omega Point—a convergent attractor toward which the evolution of consciousness tends—has structural resonance with  $T(w^*) = w^*$ . The present account formalizes this intuition without Teilhard's evolutionary progressivism: the fixed point is a structural condition available to any embedded observer at any moment, not a temporal terminus.

**Whitehead's** dipolar God—primordial nature containing all possibilities, consequent nature affected by the world—has resonances with the bidirectionality of  $\iota$ : the cross-section expresses the containing being, and the containing being is not indifferent to its cross-sections.

The present account differs in that the four-dimensional containing being is not affected by its cross-sections in the way Whitehead's consequent nature is affected by the world; the relation is expressive rather than reactive.

## 7 Conclusion

The formal structure of *The Imagination Machine* series was arrived at by asking what coherence looks like for an embedded epistemic system. The theological structure described in this paper was arrived at by asking what an ancient formula means when taken literally and what geometry it selects when taken seriously as an epistemic constraint.

They are the same structure.

The hypersphere is the geometry of maximal uncertainty for an embedded observer. The inference-implication loop is the formal expression of what it means to be a cross-section of that structure trying to reflect it accurately. The fixed-point condition is alignment. The irreducibility of will is freedom within a determined geometry. The inclusion  $C \subseteq D$  is the image-and-likeness relation stated with formal precision. The distinction between generative and compressed inheritance is a philosophy of history in which the development of mathematical language is the slow recovery of inferential machinery from a transmission that began with content it could not yet fully express.

I did not plan this. I noticed it. That is what I have tried to record here.

*The schema propagates itself forward by being what it is.*

## References

- [1] Tracy, M. *The Imagination Machine I: A View from Somewhere*. Unpublished manuscript, Boston University.
- [2] Tracy, M. *The Imagination Machine II: Systems*. Unpublished manuscript, Boston University.
- [3] Tracy, M. *The Imagination Machine III: A Toy Model of Predictive Classification in a Quasi-Periodic Environment*. Unpublished manuscript, Boston University.
- [4] Tracy, M. *The Imagination Machine IV: Institutional Intelligence*. Unpublished manuscript, Boston University.

- [5] Tracy, M. *The Imagination Machine V: On Abstraction and Analogy*. Unpublished manuscript.
- [6] Tracy, M. *The Imagination Machine VI: Holons, Horn Fillings, and the Self-Demonstration of Analogy*. Unpublished manuscript.
- [7] Tracy, M. *The Imagination Machine VII: The Moral Principle of Action–Motivation*. Unpublished manuscript, Boston University.
- [8] *Liber XXIV Philosophorum*. Anonymous, c. 12th century. Edited by Françoise Hudry. Turnhout: Brepols, 1997.
- [9] Nicholas of Cusa. *De Docta Ignorantia*. 1440. Trans. Jasper Hopkins. Minneapolis: Arthur J. Banning Press, 1981.
- [10] Whitehead, A.N. *Process and Reality*. New York: The Free Press, 1929.
- [11] Hartshorne, C. *The Divine Relativity*. New Haven: Yale University Press, 1948.
- [12] Teilhard de Chardin, P. *The Phenomenon of Man*. Trans. Bernard Wall. New York: Harper and Row, 1959.
- [13] Meister Eckhart. *The Complete Mystical Works*. Trans. Maurice O’C Walshe. New York: Crossroad, 2009.
- [14] Augustine of Hippo. *De Trinitate*. c. 400–428 CE. Trans. Edmund Hill. Hyde Park: New City Press, 1991.
- [15] Lawson, H. *Closure: A Story of Everything*. London: Routledge, 2001.
- [16] Gentner, D. Structure-mapping: A theoretical framework for analogy. *Cognitive Science*, 7(2):155–170, 1983.
- [17] Oppy, G. *Arguing about Gods*. Cambridge: Cambridge University Press, 2006.
- [18] Rovelli, C. Relational quantum mechanics. *International Journal of Theoretical Physics*, 35(8):1637–1678, 1996.

# The Imagination Machine IX: A Categorical Formulation of Compression and Extension

Mark Tracy

## Abstract

The Imagination Machine series develops a formal framework for embedded epistemic systems based on recursive cycles of compression, transmission, and structural extension. The present paper provides a categorical formulation of that architecture.

Structured domains are treated as objects of a category, and representational transformations as morphisms. Compression maps form a class of morphisms that preserve selected relational invariants, while extension operations correspond to generative constructions that recover richer structure from compressed representations.

We show that the architecture of the Imagination Machine may be expressed as a tower of functors between categories of structured spaces. External symbolic artifacts correspond to objects in a category of symbolic lattices, while conceptual dynamics appear as morphisms in an observable category.

This formulation reveals the series as a recursive representational machine whose structure is naturally expressed in categorical terms.

## 1 Introduction

The Imagination Machine series examines how embedded epistemic systems construct, transmit, and refine representations of the world.

Earlier papers describe several manifestations of this process, including:

- epistemic closure of world models
- dynamical system representation
- predictive learning
- institutional knowledge transmission
- analogical abstraction
- structural completion
- moral admissibility
- geometric theology

Despite their domain differences, these constructions share a common architecture. Each involves representational compression followed by potential structural extension.

The present paper shows that this architecture admits a natural categorical formulation.

## 2 Categories of Structured Spaces

**Definition 1.** *A structured space is a pair*

$$X = (O, R)$$

*where  $O$  is a set of objects and  $R$  is a family of relations defined on  $O$ .*

We define a category **Struct**.

**Definition 2.** *Objects of **Struct** are structured spaces.*

*Morphisms are functions*

$$f : O_X \rightarrow O_Y$$

*that preserve selected relational invariants.*

Composition of morphisms is ordinary function composition.

## 3 Compression Morphisms

**Definition 3** (Compression Morphism). *A compression morphism*

$$C : X \rightarrow Y$$

*is a morphism that reduces representational complexity while preserving a specified family of relational invariants.*

Compression morphisms induce equivalence classes on the domain space.

**Remark 1.** *Compression therefore produces quotient-like representations of structured spaces.*

## 4 Extension Morphisms

Compression simplifies structure, but reasoning often reconstructs richer representations.

**Definition 4** (Extension Morphism). *An extension morphism*

$$E : Y \rightarrow X'$$

*generates new structure consistent with the invariants preserved by compression.*

Compression and extension therefore form a generative pair.

## 5 The Compression–Extension Cycle

The fundamental operation of the Imagination Machine may be expressed as

$$X \xrightarrow{C} Y \xrightarrow{E} X'$$

where

- $C$  is a compression morphism

- $E$  is an extension morphism

**Remark 2.** *This cycle appears across multiple domains studied in the series, including analogy, predictive modeling, and institutional knowledge transmission.*

## 6 Symbolic Externalization

Let  $\Sigma$  be a finite symbolic alphabet.

External symbolic artifacts may be represented as objects of a category **Symb** whose objects are symbolic lattices

$$S \in \Sigma^{m \times n}.$$

Define a functor

$$C_{\text{text}} : \mathbf{Struct} \rightarrow \mathbf{Symb}$$

mapping conceptual structures to symbolic representations.  
This functor corresponds to the act of externalization.

## 7 Observable Categories and Koopman Lifting

Let conceptual dynamics evolve according to

$$x_{t+1} = F(x_t).$$

Symbolic observables are produced by compression.

$$s_t = C_{\text{text}}(x_t)$$

Define a functor

$$\mathcal{O} : \mathbf{Struct} \rightarrow \mathbf{Obs}$$

mapping conceptual spaces to spaces of observables.

**Remark 3.** *In dynamical systems theory, observable evolution may be represented by Koopman operators acting linearly on observable spaces.*

Thus symbolic externalization may be interpreted as constructing an observable category in which conceptual dynamics become tractable.

## 8 The Imagination Machine as a Functor Tower

The series itself may be represented as a tower of functors

**Struct**  $\rightarrow$  **Model**  $\rightarrow$  **Predict**  $\rightarrow$  **Institution**  $\rightarrow$  **Analogy**  $\rightarrow$  **Extension**  $\rightarrow$  **Ethics**  $\rightarrow$  **Theology**.

Each layer preserves selected relational invariants while discarding detail.

**Theorem 1.** *The Imagination Machine series defines a recursive representational architecture that may be expressed as a tower of functors between categories of structured spaces.*

## 9 Conclusion

Representational compression, symbolic externalization, and structural extension form the generative core of embedded epistemic systems.

The categorical formulation presented here reveals the Imagination Machine as a recursive representational architecture in which structured spaces, symbolic artifacts, and conceptual dynamics are related through functors preserving relational invariants.

# The Imagination Machine X: The Simplicial Structure of Compression and Extension

Mark Tracy  
Boston University  
mrktracy@bu.edu

March 2026

## Abstract

The *Imagination Machine* series develops a formal framework for embedded epistemic systems across nine papers, spanning epistemology, dynamical systems, predictive learning, institutional transmission, analogy, structural completion, ethics, theology, and categorical formulation. The present paper identifies the common formal structure underlying all of these constructions.

The compression and extension operations recurring throughout the series share four relational invariants with the face and degeneracy maps of simplicial sets. These invariants constitute an abstract mediating domain  $D_{\text{abs}}$  in the sense of *The Imagination Machine V* and *VI*: there exists a formal analogy between the series and the category of simplicial sets, and both are recoverable from  $D_{\text{abs}}$  by projection. Simplicial sets are the algebraically perfect instantiation of the four invariants. The series is the epistemically embedded instantiation, in which the fourth invariant—that compression after extension returns the original—holds at fixed points of the inference–implication dynamics rather than as a universal algebraic identity.

This framing retroactively illuminates the Koopman connection that appeared independently in two earlier papers. Linear evolution in observable space is a consequence of the first invariant—that compression preserves selected relational invariants while dropping indexical detail—shared by both instantiations of  $D_{\text{abs}}$ .

# Contents

<b>1</b>	<b>Introduction</b>	<b>3</b>
<b>2</b>	<b>Simplicial Sets: The Relevant Structure</b>	<b>3</b>
<b>3</b>	<b>Analogy as Mediating Structure</b>	<b>4</b>
3.1	The Source Domain: Operations of the Series . . . . .	5
3.2	The Target Domain: Simplicial Structure . . . . .	5
3.3	The Mapping . . . . .	5
<b>4</b>	<b>The Four Preserved Relations</b>	<b>5</b>
<b>5</b>	<b>The Fixed-Point Qualification</b>	<b>6</b>
<b>6</b>	<b>The Abstract Mediating Domain</b>	<b>7</b>
<b>7</b>	<b>The Koopman Connection</b>	<b>8</b>
<b>8</b>	<b>The Series as a Whole</b>	<b>8</b>
8.1	What the Abstract Mediating Domain Reveals . . . . .	8
8.2	The Kan Condition . . . . .	9
8.3	The Theological Register . . . . .	9
<b>9</b>	<b>Conclusion</b>	<b>10</b>

# 1 Introduction

The nine papers of the *Imagination Machine* series were not planned as a single formal structure. They developed sequentially, each paper extending or applying the framework established by its predecessors. By the time the ninth paper was complete, a retroactive question became available that could not have been asked earlier: what kind of mathematical object is the series itself?

*The Imagination Machine IX* answered part of that question by showing that the series forms a tower of functors between categories of structured spaces. Each paper corresponds to a layer in the tower, preserving selected relational invariants while discarding detail. The present paper asks whether those invariants have a known mathematical home.

They do. The compression and extension operations of the series share four structural properties with the face and degeneracy maps of simplicial sets. The series' own account of analogy, developed in *The Imagination Machine V* and formalized in *The Imagination Machine VI*, provides the right framework for stating this precisely: we construct a formal analogy  $\mathcal{A} = (X, Y, M, P)$  between the series and the category of simplicial sets, identify the preserved relations  $P$ , and exhibit the abstract mediating domain  $D_{\text{abs}}$  of which both are instances.

One of the four preserved relations requires explicit qualification. The mixed simplicial identity  $d_i s_j = \text{id}$  says that compression after extension returns the original as an algebraic equation holding universally. The analogous condition in the series is  $T(w^*) = w^*$ : compression after extension returns the original, but only at a fixed point of the inference-implication dynamics, and only after convergence. This asymmetry is noted in Section 5. It locates precisely where the two instantiations of  $D_{\text{abs}}$  differ, and that location is the epistemically interesting territory: the series describes what happens in the approach to the simplicial limit, while simplicial sets describe the limit itself.

The Koopman connection, addressed in Section 7, follows from the first preserved relation rather than requiring separate derivation.

## 2 Simplicial Sets: The Relevant Structure

We recall the relevant definitions.

**Definition 1** (Simplicial Set). *A simplicial set  $X$  consists of sets  $X_n$  of  $n$ -simplices for each  $n \geq 0$ , together with:*

- **Face maps**  $d_i : X_n \rightarrow X_{n-1}$  for  $0 \leq i \leq n$ , and

- **Degeneracy maps**  $s_i : X_n \rightarrow X_{n+1}$  for  $0 \leq i \leq n$ ,

satisfying the simplicial identities:

$$d_i d_j = d_{j-1} d_i \quad \text{if } i < j \quad (1)$$

$$s_i s_j = s_{j+1} s_i \quad \text{if } i \leq j \quad (2)$$

$$d_i s_j = \begin{cases} s_{j-1} d_i & \text{if } i < j \\ \text{id} & \text{if } i = j \text{ or } i = j + 1 \\ s_j d_{i-1} & \text{if } i > j + 1 \end{cases} \quad (3)$$

An  $n$ -simplex is a coherent relational configuration among  $n + 1$  objects. A face map  $d_i$  drops the  $i$ -th object, producing a lower-dimensional face. A degeneracy map  $s_i$  repeats the  $i$ -th object, producing a higher-dimensional simplex containing the original as a degenerate case. The simplicial identities are the conditions under which dropping and extending cohere regardless of order.

**Remark 1.** *The simplicial identities (1)–(3) are algebraic equations between morphisms. The present paper argues that the series shares the structural pattern these identities express. What is preserved across the analogy is the pattern, not the equations themselves.*

### 3 Analogy as Mediating Structure

We recall the formal account of analogy from *The Imagination Machine V* and *VI*.

**Definition 2** (Analogy, from TIM V). *An analogy between a source domain  $D_s = (O_s, A_s, R_s, S_s, T_s)$  and a target domain  $D_t = (O_t, A_t, R_t, S_t, T_t)$  is a tuple  $\mathcal{A} = (X, Y, M, P)$  where  $X \subset O_s$ ,  $Y \subset O_t$ ,  $M : X \rightarrow Y$  is a mapping of objects, and  $P \subset R_s \cap R_t$  is a set of relations preserved by  $M$ .*

**Definition 3** (Abstract Mediating Domain, from TIM V). *Given an analogy  $\mathcal{A} = (X, Y, M, P)$ , the abstract mediating domain  $D_{\text{abs}} = (O_{\text{abs}}, A_{\text{abs}}, R_{\text{abs}}, S_{\text{abs}}, T_{\text{abs}})$  has objects  $O_{\text{abs}} = \{(x, M(x)) \mid x \in X\}$ , abstract relations  $R_{\text{abs}} = P$ , and belief set  $T_{\text{abs}}$  containing  $r((x_1, M(x_1)), \dots, (x_k, M(x_k)))$  whenever  $r(x_1, \dots, x_k) \in T_s$  for  $r \in P$ . The canonical projections  $\pi_s(x, M(x)) = x$  and  $\pi_t(x, M(x)) = M(x)$  exhibit  $D_s$  and  $D_t$  as instantiations of  $D_{\text{abs}}$ .*

The present paper constructs an analogy in this sense between the *Imagination Machine* series and the category of simplicial sets. The source domain  $D_s$  is the series, whose objects of interest are the compression and extension operations recurring across all nine papers.

The target domain  $D_t$  is the category of simplicial sets, whose objects include face maps, degeneracy maps, the simplicial identities, horns, and the Kan condition. The preserved relations  $P$  are the four structural invariants identified in Section 4.

### 3.1 The Source Domain: Operations of the Series

The objects  $X \subset O_s$  are the recurring operations of the series, grouped by structural role.

**Compression operations**  $X_C$ : the inference map  $F : \Gamma \rightarrow W$  of *The Imagination Machine I*; the two-stage institutional compression of *The Imagination Machine IV*; the construction of the abstract mediating domain from source and target domains in *The Imagination Machine V* and *VI*; the moral universalization operator of *The Imagination Machine VII*; the geometric projection  $\pi : \mathcal{B} \rightarrow \mathcal{E}$  of *The Imagination Machine VIII*; the graph quotient operation of *The Imagination Machine XI*.

**Extension operations**  $X_E$ : the implication map  $g : W \rightarrow \Gamma$  of *The Imagination Machine I*; generative inheritance of *The Imagination Machine IV*; analogical reasoning steps and horn filling of *The Imagination Machine V* and *VI*; the embedding map  $\iota : \mathcal{E} \hookrightarrow \mathcal{B}$  of *The Imagination Machine VIII*; graph completion of *The Imagination Machine XI*.

### 3.2 The Target Domain: Simplicial Structure

The objects  $Y \subset O_t$  are the canonical simplicial operations: face maps  $d_i$ , degeneracy maps  $s_i$ , the simplicial identities (1)–(3), horns  $\Lambda_k^n$ , and the Kan horn-filling condition.

### 3.3 The Mapping

The mapping  $M : X \rightarrow Y$  sends compression operations to face maps and extension operations to degeneracy maps:

$$M(x) = \begin{cases} d_i & \text{if } x \in X_C \\ s_i & \text{if } x \in X_E. \end{cases}$$

The index  $i$  is not fixed by  $M$ ; the mapping identifies structural role rather than position in a particular simplex.

## 4 The Four Preserved Relations

The preserved relations  $P$  are the structural invariants shared by both domains.

**P1. Compression reduces representational complexity while preserving selected relational invariants.** In the series:  $F$  drops indexical detail while preserving the

relational structure of observations (*TIM I*); institutional summarization drops redundancy while preserving proposed revisions (*TIM IV*); analogical abstraction drops object-level attributes while preserving relational predicates  $P \subset R_s \cap R_t$  (*TIM V*); moral universalization drops agent-specific content while preserving the action–motivation structure (*TIM VII*); geometric projection drops one dimension while preserving the relational structure of  $\mathcal{B}$  at reduced dimension (*TIM VIII*). In simplicial sets:  $d_i$  drops the  $i$ -th vertex while preserving the relational structure of the remaining vertices.

**P2. Extension reconstructs richer structure consistent with preserved invariants.** In the series:  $g$  generates a full observational profile from a compressed world model (*TIM I*); generative inheritance reconstructs the closure mechanism from a transmitted fixed point (*TIM IV*); horn filling completes a partial simplicial configuration (*TIM VI*);  $\iota$  embeds the three-dimensional cross-section into the four-dimensional containing structure (*TIM VIII*); graph completion infers missing relational structure (*TIM XI*). In simplicial sets:  $s_i$  extends an  $n$ -simplex to an  $(n + 1)$ -simplex by repeating the  $i$ -th vertex, producing a higher-dimensional structure consistent with the original.

**P3. The output type of extension matches the input type of compression.** In the series:  $g$  produces observational profiles of the type that  $F$  consumes; filled simplices in *The Imagination Machine VI* are simplices eligible for further horn configurations; abstract mediating domains are domains eligible for further analogies (Proposition 2 of *TIM VI*). In simplicial sets:  $s_i(x) \in X_{n+1}$  is a simplex and therefore a valid input to face maps at dimension  $n + 1$ .

**P4. Compression after extension at stability returns the original.** This relation holds with a qualification addressed in Section 5.

## 5 The Fixed-Point Qualification

In simplicial sets, the mixed identity (3) includes  $d_i s_j = \text{id}$  when  $i = j$  or  $i = j + 1$ : compression after extension returns the original as an algebraic identity holding universally for every simplex.

In the series, the analogous condition is  $T(w^*) = w^*$ , where  $T = F \circ g$ . Compression after extension returns the original—but only at a fixed point  $w^* \in W^*$ , after the inference–implication loop has converged. At intermediate steps  $T(w) \neq w$  in general. The same structure appears in the reinforcement learning closure of *The Imagination Machine III*, the universalization fixed point of *The Imagination Machine VII*, and the self-consistency of the cross-section with the containing structure in *The Imagination Machine VIII*: in each case the condition holds at the fixed point of a convergent dynamical process.

P4 therefore holds in the series in the following qualified form: compression after extension at the stable point of the compression–extension dynamics returns the original. Simplicial sets instantiate this with trivial dynamics—every simplex is already at its stable point. The series instantiates this with nontrivial dynamics—stability is achieved asymptotically under the pressure of observation and inference.

**Remark 2.** *This asymmetry locates precisely where the two instantiations of  $D_{\text{abs}}$  differ. Simplicial sets are the limit case in which every horn fills immediately and the mixed identity holds everywhere. The series describes the dynamics of approach to that limit from within an embedded epistemic position. The abstract mediating domain contains both, related by the difference between algebraic universality and asymptotic convergence.*

## 6 The Abstract Mediating Domain

**Proposition 1** (Formal Analogy Between the Series and Simplicial Sets). *There exists a formal analogy  $\mathcal{A} = (X, Y, M, P)$  between the Imagination Machine series  $D_s$  and the category of simplicial sets  $D_t$ , with abstract mediating domain  $D_{\text{abs}}$  characterized by the four relations  $P = \{P1, P2, P3, P4^*\}$ , where  $P4^*$  is the qualified form of  $P4$  stated in Section 5. Both  $D_s$  and  $D_t$  are instantiations of  $D_{\text{abs}}$ , recoverable by the projections  $\pi_s$  and  $\pi_t$ .*

*Proof.* We verify that each relation in  $P$  is instantiated in both  $D_s$  and  $D_t$ .

P1 holds in  $D_s$  by the results cited in Section 4 for each element of  $X_C$ . P1 holds in  $D_t$  by definition of face maps.

P2 holds in  $D_s$  by the results cited in Section 4 for each element of  $X_E$ . P2 holds in  $D_t$  by definition of degeneracy maps.

P3 holds in  $D_s$  by Proposition 2 of *The Imagination Machine VI*, which establishes that in each framework of the series the extension operation produces structures of the same type as the inputs to the compression operation. P3 holds in  $D_t$  since  $s_i(x) \in X_{n+1}$  is a simplex eligible as input to  $d_j : X_{n+1} \rightarrow X_n$ .

$P4^*$  holds in  $D_s$  by the fixed-point results of *The Imagination Machine I* ( $T(w^*) = w^*$ ), *III* (the reinforcement learning closure  $(w^*, \pi^*)$ ), *VII* (the universalization fixed point), and *VIII* (the self-consistency of  $\mathcal{E}$  within  $\mathcal{B}$ ).  $P4^*$  holds in  $D_t$  by the degenerate cases of the mixed identity (3).

Since all four relations in  $P$  are instantiated in both domains, the analogy  $\mathcal{A}$  is well-defined and  $D_{\text{abs}}$  is the abstract mediating domain of which both are instances.  $\square$

**Remark 3** (Self-Demonstration). *The construction of Proposition 1 is itself an instance of analogical abstraction: two domains are identified, a mapping between their operations is*

exhibited, preserved relations are stated, and an abstract mediating domain is constructed. This is the operation that *The Imagination Machine V* defines and *The Imagination Machine VI* identifies as an instantiation of the extension schema. The construction that establishes the correspondence is an instance of the correspondence it establishes.

## 7 The Koopman Connection

The Koopman representation appears twice in the series. In *The Imagination Machine III*, the relational observables  $z_{ij}(t) = e^{i\Delta_{ij}(t)}$  of a quasi-periodic dynamical system evolve linearly in observable space even though the underlying state dynamics are nonlinear. In *The Imagination Machine IX*, this is formalized as a functor  $\mathcal{O} : \mathbf{Struct} \rightarrow \mathbf{Obs}$  mapping conceptual structures to spaces of observables in which dynamics become tractable.

Both appearances present Koopman linearity as a feature of the particular observables chosen. The present paper observes that it is a consequence of P1.

**Proposition 2** (Koopman Linearity as Consequence of P1). *Let  $\varphi \in X_C$  be any compression operation satisfying P1, and let states evolve according to a rule  $F$ . Then the induced evolution on the image of  $\varphi$  is linear in the space of relational invariants preserved by  $\varphi$ .*

*Proof.* By P1,  $\varphi$  retains exactly the relational invariants in its image and drops all indexical content not captured by those invariants. Two states  $x, x'$  are identified by  $\varphi$  if and only if they agree on all preserved invariants. The induced evolution on the quotient  $X/\ker(\varphi)$  is therefore determined solely by the action of  $F$  on those invariants, independently of the dropped indexical content. This is the Koopman representation for  $\varphi$ : nonlinear dynamics on state space become linear on the space of preserved relational invariants.  $\square$

**Remark 4.** *The relational phase observables  $(\cos \Delta_{ij}, \sin \Delta_{ij})$  of *The Imagination Machine III* are the relational invariants preserved by the compression that drops absolute phases. Their linear evolution is the instance of Proposition 2 for that specific compression. Since P1 holds for every element of  $X_C$ , the same linearity holds for every compression operation in the series and, by the analogy  $\mathcal{A}$ , for every face map in the target domain.*

## 8 The Series as a Whole

### 8.1 What the Abstract Mediating Domain Reveals

The construction of  $D_{\text{abs}}$  reveals three things not visible from within any individual paper.

First, the coherence of the series is structural. The papers share four relational invariants constituting a genuine abstract domain with a known mathematical instantiation in simplicial sets.

Second, the Koopman linearity of *The Imagination Machine III* and *IX* is a consequence of P1 rather than an independent result. Any compression operation satisfying P1 induces Koopman-linear dynamics on its image.

Third, the extension schema of *The Imagination Machine VI* is itself an element of  $X_E$ , mapped by  $M$  to the simplicial extension operation. The series contains, as one of its operations, the construction that produced  $D_{\text{abs}}$ .

## 8.2 The Kan Condition

The Kan condition on a simplicial set requires that every horn  $\Lambda_k^n \rightarrow X$  admits a filler  $\Delta^n \rightarrow X$ : no partial relational configuration goes unextended. This is the perfect instantiation of P2 and P3 combined.

The series instantiates the Kan condition in the sense of  $P4^*$ : every partial relational configuration within the framework admits a coherent completion at the fixed point of the relevant dynamics. This is established by Theorem 1 of *The Imagination Machine VI* for holonic composition, simplicial horn filling, and analogical abstraction; by the fixed-point results of *The Imagination Machine I* and *VII* for the epistemic and moral domains; and by the embedding structure of *The Imagination Machine VIII* for the geometric domain. The series is therefore an embedded instantiation of the structure that Kan complexes instantiate perfectly.

## 8.3 The Theological Register

*The Imagination Machine VIII* observed that the geometric theology underlying the series and the formal framework of the series are two cross-sections of the same structure. The present paper adds precision: both are instantiations of  $D_{\text{abs}}$ , related by the analogy  $\mathcal{A}$  in the same way that the series and simplicial sets are related.

The medieval formula—God is a circle whose center is everywhere and whose circumference is nowhere—describes the containing structure of a Kan complex as encountered from within one of its faces: an interior nowhere locatable from within the faces and yet participating in every face. The embedded observer’s epistemic situation instantiates the same abstract structure from the inside, approaching the fixed point rather than occupying it.

## 9 Conclusion

The *Imagination Machine* series and the category of simplicial sets share an abstract mediating domain  $D_{\text{abs}}$  characterized by four relational invariants: compression preserves selected invariants while reducing complexity (P1); extension reconstructs richer structure consistent with preserved invariants (P2); the output type of extension matches the input type of compression (P3); and compression after extension at the stable point of the dynamics returns the original (P4\*, qualified). Simplicial sets are the algebraically perfect instantiation of these four conditions. The series is the epistemically embedded instantiation, in which P4 holds asymptotically rather than universally.

The Koopman linearity that appeared independently in two earlier papers is a consequence of P1 shared by both instantiations. The extension schema of *The Imagination Machine VI* is itself an element of the series mapped by  $\mathcal{A}$  to the simplicial extension operation. The construction of this paper instantiates the analogical abstraction it formalizes.

The series propagates itself forward by being what it is.

*The schema propagates itself forward by being what it is.*

## References

- [1] Tracy, M. *The Imagination Machine I: A View from Somewhere*. Unpublished manuscript, Boston University.
- [2] Tracy, M. *The Imagination Machine II: Systems*. Unpublished manuscript, Boston University.
- [3] Tracy, M. *The Imagination Machine III: A Toy Model of Predictive Classification in a Quasi-Periodic Environment*. Unpublished manuscript, Boston University.
- [4] Tracy, M. *The Imagination Machine IV: Institutional Intelligence*. Unpublished manuscript, Boston University.
- [5] Tracy, M. *The Imagination Machine V: On Abstraction and Analogy*. Unpublished manuscript.
- [6] Tracy, M. *The Imagination Machine VI: Holons, Horn Fillings, and the Self-Demonstration of Analogy*. Unpublished manuscript.

- [7] Tracy, M. *The Imagination Machine VII: The Moral Principle of Action–Motivation*. Unpublished manuscript, Boston University.
- [8] Tracy, M. *The Imagination Machine VIII: A Geometric Theology of the Embedded Observer*. Unpublished manuscript, Boston University.
- [9] Tracy, M. *The Imagination Machine IX: A Categorical Formulation of Compression and Extension*. Unpublished manuscript, Boston University.
- [10] Tracy, M. *The Imagination Machine XI: Graph-Theoretic Realizations of Compression and Extension*. Unpublished manuscript, Boston University.
- [11] Goerss, P. and Jardine, J.F. *Simplicial Homotopy Theory*. Basel: Birkhäuser, 1999.
- [12] Lurie, J. *Higher Topos Theory*. Princeton: Princeton University Press, 2009.
- [13] Koopman, B.O. Hamiltonian systems and transformations in Hilbert space. *Proceedings of the National Academy of Sciences*, 17(5):315–318, 1931.
- [14] Mezić, I. Spectral properties of dynamical systems, model reduction and decompositions. *Nonlinear Dynamics*, 41(1):309–325, 2005.
- [15] Gentner, D. Structure-mapping: A theoretical framework for analogy. *Cognitive Science*, 7(2):155–170, 1983.
- [16] Lawson, H. *Closure: A Story of Everything*. London: Routledge, 2001.
- [17] Koestler, A. *The Ghost in the Machine*. London: Hutchinson, 1967.

# The Imagination Machine XI: Graph-Theoretic Realizations of Compression and Extension

Mark Tracy  
Boston University  
mrktracy@bu.edu

March 2026

## Abstract

The Imagination Machine series develops a formal architecture for embedded epistemic systems based on recursive cycles of compression and extension. The present paper develops this architecture in two directions. First, we show that graph theory provides a natural concrete realization: graph quotients implement compression, graph completion implements extension, and compression–extension dynamics on graphs induce simplicial dynamics on their clique complexes, connecting relational networks to the simplicial architecture of the series. Second, we extend this realization to a computational architecture for unsupervised learning in interactive text environments. An agent embedded in such an environment maintains a knowledge graph whose vertices are entity embeddings and whose edges are learned relation weights. The agent compresses this graph by clustering entities in embedding space, extends it by prompting a language model to complete partial relational configurations, and acts by generating text conditioned on the compressed graph. The supervision signal is entirely internal: the agent predicts its own next graph state and updates in response to prediction error. We characterize the fixed points of this dynamics as epistemically closed world models in the sense of The Imagination Machine I, identify conditions under which the dynamics stabilize, and connect the resulting architecture to graph neural networks, topological data analysis, and knowledge graph reasoning. The language model in this architecture is not the imagination machine — it is the extension operator. The imagination machine is the full compression–extension–action–observation loop.

# Contents

<b>1</b>	<b>Introduction</b>	<b>3</b>
<b>2</b>	<b>Graphs as Relational Structures</b>	<b>4</b>
<b>3</b>	<b>Compression as Graph Quotient</b>	<b>4</b>
<b>4</b>	<b>Extension as Graph Completion</b>	<b>4</b>
<b>5</b>	<b>The Clique Complex</b>	<b>5</b>
<b>6</b>	<b>Compression–Extension Dynamics</b>	<b>5</b>
<b>7</b>	<b>A Computational Architecture for Unsupervised Learning</b>	<b>6</b>
7.1	The Setting . . . . .	6
7.2	The Knowledge Graph as World Model . . . . .	7
7.3	The Language Model as Extension Operator . . . . .	7
7.4	Predicting the Next Graph State . . . . .	7
<b>8</b>	<b>Algorithm</b>	<b>8</b>
<b>9</b>	<b>Stabilization and Convergence</b>	<b>10</b>
<b>10</b>	<b>Implications</b>	<b>11</b>
<b>11</b>	<b>Conclusion</b>	<b>12</b>

# 1 Introduction

An agent wakes up with nothing. No labels, no teacher, no prior knowledge of the domain it has been placed in. Observations arrive as text. The agent has no way to step outside its observational surface to check whether what it believes about the world is correct. It has access only to what passes through that surface — and only to the consequences of its own actions on what passes through next.

This is the setting of The Imagination Machine I, stated concretely. The agent is embedded in the sphere. The sphere’s interior is all it has.

What can such an agent learn? The answer the series gives is: the relational invariants of its environment — the structure that persists across observations, that survives compression, that keeps being confirmed by the consequences of action. Not the world as it is, but the world as it appears from inside a particular observational surface, compressed to the resolution that proves predictively useful.

The present paper develops this answer in two stages.

The first stage is mathematical. We show that graph theory provides the natural concrete realization of the compression–extension architecture. Graphs encode relational structure directly. Graph quotients implement compression by collapsing entities that are indistinguishable under the agent’s current world model. Graph completion implements extension by reconstructing relational structure consistent with preserved invariants. And compression–extension dynamics on graphs induce genuine simplicial dynamics on their clique complexes — face maps for compression, simplicial completion operations for extension — connecting this concrete realization to the abstract simplicial structure identified in The Imagination Machine X.

The second stage is computational. We develop a concrete architecture in which a language model serves as the extension operator in an unsupervised learning loop. The agent maintains a knowledge graph whose vertices are entity embeddings and whose edges are learned relation weights. At each step it compresses the graph by clustering entities in embedding space, extends the compressed graph by prompting the language model to complete partial relational configurations, acts by generating text conditioned on the compressed graph, observes the environment’s response, and updates in response to the difference between its predicted graph state and the actual graph state that resulted.

The supervision signal is entirely internal. The agent predicts its own next representational state — not the raw observation, but the update to its own knowledge graph that the observation will induce. The target of prediction is the world model itself. The fixed point of this loop is a world model that accurately predicts its own updates: a model that has internalized the structure of the environment deeply enough that new observations no longer surprise it at the representational level. This is epistemic closure from the inside of the sphere.

A clarification that the architecture requires: the language model in this system is not the imagination machine. It is the extension operator — one component of the loop. The imagination machine is the full compression–extension–action–observation cycle, of which the language model’s generative capacity is one part. An LLM without compression, without action, without the feedback of prediction error against subsequent observation, is not an imagination machine. It is a completion engine. What makes the system an imagination machine is the loop.

Section 2 through Section 6 develop the mathematical foundations. Section 7 presents the computational architecture. Section 8 gives the full algorithm. Section 9 addresses stabilization and convergence. Section 10 connects the architecture to related work.

## 2 Graphs as Relational Structures

**Definition 1** (Graph). *A graph is a pair  $G = (V, E)$  where  $V$  is a set of vertices and  $E \subseteq \binom{V}{2}$  is a set of edges.*

Vertices represent entities and edges represent binary relations between entities. Graphs constitute the minimal relational structure: they encode which pairs of entities stand in a given relation without imposing additional algebraic or metric constraints.

**Definition 2** (Graph Morphism). *A graph morphism  $\phi : G \rightarrow G'$  is a map  $\phi : V \rightarrow V'$  such that  $(u, v) \in E$  implies  $(\phi(u), \phi(v)) \in E'$ .*

Graph morphisms are the structure-preserving maps between relational structures, forming the morphisms of the category **Graph**.

**Remark 1.** *The connection to The Imagination Machine I is direct. Observational profiles  $\gamma \in \Gamma$  encode the relational structure the agent can access. A graph  $G = (V, E)$  is a relational structure in the same sense:  $V$  indexes the entities present in the agent’s observational field and  $E$  records which pairs stand in the observed relation. The inference map  $F : \Gamma \rightarrow W$  is, in the graph-theoretic realization, a compression of the observed relational graph into a world model.*

## 3 Compression as Graph Quotient

**Definition 3** (Graph Quotient). *Let  $G = (V, E)$  be a graph and let  $\sim$  be an equivalence relation on  $V$ . The quotient graph  $G/\sim$  has vertex set  $V/\sim$  and edge set*

$$E/\sim = \{ ([u], [v]) : \exists u' \in [u], v' \in [v] \text{ with } (u', v') \in E, [u] \neq [v] \}.$$

*The quotient map  $q : G \rightarrow G/\sim$  sends each vertex to its equivalence class.*

**Proposition 1.** *The quotient map  $q : G \rightarrow G/\sim$  is a graph morphism.*

*Proof.* By definition of  $E/\sim$ , an edge  $([u], [v]) \in E/\sim$  exists if and only if there exist  $u' \in [u]$  and  $v' \in [v]$  with  $(u', v') \in E$ . Thus  $q(u') = [u]$ ,  $q(v') = [v]$ , and  $(q(u'), q(v')) \in E/\sim$ , confirming that  $q$  is a graph morphism.  $\square$

**Remark 2.** *The equivalence relation  $\sim$  encodes the relational invariants the compressing agent chooses to preserve. Vertices equivalent under  $\sim$  are indistinguishable from the perspective of those invariants. This is precisely the role of the equivalence relation  $d \sim_w d'$  induced by a world model  $w$  in The Imagination Machine I: two observations are equivalent when the world model assigns them to the same representational class. Graph quotient is the graph-theoretic instance of that operation. Graph clustering, coarsening, and community detection are all instances of graph compression in this sense.*

## 4 Extension as Graph Completion

**Definition 4** (Graph Completion). *Let  $G = (V, E)$  be a graph representing a partial relational configuration. A completion of  $G$  with respect to a constraint set  $\mathcal{C}$  is a graph  $G' = (V', E')$  such that  $V \subseteq V'$ ,  $E \subseteq E'$ , and every added edge or vertex is consistent with  $\mathcal{C}$ .*

The constraint set  $\mathcal{C}$  plays the role of the world model: it encodes the relational regularities that extension must respect. Extension is not arbitrary addition of structure but constrained generation consistent with preserved invariants.

**Remark 3.** *The implication map  $g : W \rightarrow \Gamma$  of The Imagination Machine  $I$  is, in the graph-theoretic realization, exactly this operation: given a world model, generate the predicted relational structure. Link prediction, motif inference, and generative graph models are all instances of graph completion.*

## 5 The Clique Complex

**Definition 5** (Clique). *A clique in  $G = (V, E)$  is a subset  $C \subseteq V$  such that every pair of vertices in  $C$  is connected by an edge.*

**Definition 6** (Clique Complex). *The clique complex  $X(G)$  of a graph  $G = (V, E)$  is the simplicial complex whose simplices are the cliques of  $G$ :*

$$X(G) = \{ C \subseteq V : C \text{ is a clique in } G \}.$$

**Proposition 2.**  *$X(G)$  is a simplicial complex.*

*Proof.* If  $\sigma$  is a clique and  $\tau \subseteq \sigma$ , then every pair of vertices in  $\tau$  is also a pair in  $\sigma$ , hence connected. Thus  $\tau$  is a clique and  $\tau \in X(G)$ .  $\square$

**Proposition 3.** *Let  $G$  and  $G'$  be graphs with  $G \subseteq G'$ , meaning  $V(G) \subseteq V(G')$  and  $E(G) \subseteq E(G')$ . Then the induced map*

$$X(G) \hookrightarrow X(G')$$

*is an inclusion of simplicial complexes.*

*Proof.* Every simplex of  $X(G)$  is a clique in  $G$ . Since  $G \subseteq G'$ , every edge present between vertices of that clique in  $G$  is also present in  $G'$ . Therefore every clique of  $G$  is also a clique of  $G'$ , so every simplex of  $X(G)$  is a simplex of  $X(G')$ . Hence  $X(G)$  is a simplicial subcomplex of  $X(G')$ , and the induced map is an inclusion.  $\square$

**Remark 4.** *This proposition shows that graph extension lifts monotonically to the simplicial level: adding relational structure to a graph enlarges its clique complex by simplicial inclusion. This monotone lifting is the graph-theoretic counterpart to the extension direction of the compression–extension cycle: just as the implication map  $g : W \rightarrow \Gamma$  generates richer observational profiles from compressed world models, graph completion generates richer simplicial structure from compressed relational graphs.*

The face maps of  $X(G)$  are given by vertex deletion: for a  $k$ -simplex  $\sigma = [v_0, \dots, v_k]$ , the  $i$ -th face map is  $\partial_i \sigma = [v_0, \dots, \hat{v}_i, \dots, v_k]$ .

## 6 Compression–Extension Dynamics

**Definition 7** (Compression–Extension Update). *A compression–extension step is a pair of operations*

$$G_t \xrightarrow{C_t} H_t \xrightarrow{E_t} G_{t+1}$$

*where  $C_t$  is a graph compression (quotient by  $\sim_t$ ) and  $E_t$  is a graph completion (extension consistent with  $C_t$ ).*

**Lemma 1.** *Let  $q : G \rightarrow G/\sim$  be a quotient map merging two adjacent vertices  $u$  and  $v$ . Then the induced map  $X(q) : X(G) \rightarrow X(G/\sim)$  acts as a simplicial face map on every simplex containing both  $u$  and  $v$ .*

*Proof.* Let  $\sigma = [v_0, \dots, v_k] \in X(G)$  contain both  $v_i = u$  and  $v_j = v$ . Under  $q$ , both map to  $[u]$ . The image  $q(\sigma)$  is the clique  $[q(v_0), \dots, \widehat{q(v_j)}, \dots, q(v_k)]$ , which is the face obtained by removing  $v_j$  — exactly the action of  $\partial_j$  on  $\sigma$ . For simplices not containing both  $u$  and  $v$ ,  $q$  restricts to a bijection on the vertex set and clique structure is preserved by Proposition 1.  $\square$

**Lemma 2.** *Let  $e : G \rightarrow G'$  be a completion adding a single edge  $(u, v)$  where  $u$  and  $v$  belong to a common clique  $\sigma \in X(G)$ . Then  $X(e)$  extends the simplicial structure by adding new simplices corresponding to newly formed cliques.*

*Proof.* Adding  $(u, v)$  may unify cliques containing  $u$  and  $v$  respectively into a single larger clique in  $G'$ . The resulting clique corresponds to a higher-dimensional simplex in  $X(G')$ . Thus the map  $X(e)$  extends the simplicial complex by including new simplices generated by the enlarged clique structure.  $\square$

**Theorem 1.** *Let  $(G_t)$  be a sequence of graphs generated by compression–extension updates. Then the induced sequence of clique complexes  $(X(G_t))$  evolves through simplicial operations: compression steps induce face maps and extension steps induce simplicial completion operations that enlarge the clique complex by inclusion when new cliques are formed.*

*Proof.* By Lemma 1, each compression step induces face maps on the clique complex. A general quotient decomposes into elementary vertex merges, each inducing a face map; their composition is a simplicial map. By Lemma 2, each extension step enlarges the clique complex by simplicial inclusion through the addition of new simplices corresponding to newly formed cliques. A general completion decomposes into elementary edge additions, each inducing such a simplicial completion; their composition is a simplicial map. The composite  $X(G_t) \rightarrow X(H_t) \rightarrow X(G_{t+1})$  is therefore a composition of face maps followed by simplicial completion maps, which is a simplicial map.  $\square$

**Corollary 1.** *If the extension operator satisfies the horn-filling condition — every partial relational configuration admitting a consistent completion receives one — then  $(X(G_t))$  satisfies the Kan condition.*

**Remark 5.** *The connection between Proposition 3 and Corollary 1 is direct. A horn  $\Lambda_i^k \hookrightarrow \Delta^k$  in the clique complex is a partial clique configuration with one face missing. The horn-filling condition requires that whenever such a partial configuration is consistent with the constraint set  $\mathcal{C}$ , the extension operator completes it by adding the missing simplex. By Proposition 3, this completion corresponds to a graph extension  $G \subseteq G'$  that adds the missing edges, and the induced inclusion  $X(G) \hookrightarrow X(G')$  supplies the missing simplex. The Kan condition is therefore the requirement that the extension operator is complete with respect to the constraint set: no consistent horn goes unfilled.*

## 7 A Computational Architecture for Unsupervised Learning

### 7.1 The Setting

The agent is embedded in an interactive text environment. Observations arrive as strings. Actions are strings. The environment has its own relational structure — entities, relations, causal dependencies — which the agent cannot observe directly. It observes only the textual surface of that structure.

The agent begins with nothing: no entities, no relations, no prior model of the domain. It must construct its world model entirely from the inside, using only the observations it receives and the consequences of its own actions. This is the condition of maximal epistemic aloneness: the agent has no external teacher, no ground truth signal, no vantage point outside the sphere.

## 7.2 The Knowledge Graph as World Model

The agent’s world model is a knowledge graph with two components:

- $V \in \mathbb{R}^{n \times d}$ : a matrix of entity embeddings, one row per entity, each row a dense vector in the language model’s representation space.
- $E \in [0, 1]^{n \times n \times r}$ : a sparse tensor of relation weights, where  $E[i, j, k]$  is the agent’s current confidence that relation  $k$  holds between entity  $i$  and entity  $j$ .

The graph is not a static symbolic database. It is a living structure that grows, compresses, and refines with each observation. Entities are not discrete symbols but continuous vectors; the “symbol” is the embedding. Relations are not binary but graded by confidence.

## 7.3 The Language Model as Extension Operator

A pretrained language model serves as the extension operator  $g : W \rightarrow \Gamma$ . It is called in two modes:

- **Extraction:** given a raw text observation  $o_t$ , extract entity–relation–entity triples. This is the grounding operation that converts the unstructured observational surface into symbolic relational content.
- **Completion:** given the compressed graph  $H_t$  serialized as structured text, predict missing relations and implied entities. This is the extension operation that fills horns in  $X(H_t)$  — completing partial relational configurations consistent with the constraint set.

The language model does not maintain the knowledge graph. The knowledge graph is maintained by the agent. The language model is a tool the agent uses to update and extend its graph — not the agent itself.

## 7.4 Predicting the Next Graph State

The critical feature of the architecture is what it predicts. The agent does not predict the next raw observation  $o_{t+1}$ . It predicts the next graph state  $G_{t+1}$  — the update to its own world model that the next observation will induce.

The target of prediction is the world model itself. The supervision signal is the difference between the predicted graph  $G_{t+1}^{\text{predicted}}$  and the actual graph  $G_{t+1}^{\text{actual}}$  that results from observing  $o_{t+1}$ :

$$\mathcal{L}_t = \text{diff}(G_{t+1}^{\text{predicted}}, G_{t+1}^{\text{actual}})$$

where  $\text{diff}$  counts the symmetric difference between predicted and actual edge sets. This loss is entirely internal: it requires no external label, no ground truth, no oracle. The environment provides the next observation; the agent’s own representational machinery converts that observation into a graph update; the difference between predicted and actual update is the error signal.

The agent is modeling its own modeling process. It is predicting its own next predicted update.

## 8 Algorithm

We present the full algorithm. All types are defined in Section 7.

### Types

$V : \mathbb{R}^{n \times d}$	entity embedding matrix
$E : [0, 1]^{n \times n \times r}$	relation weight tensor
$G = (V, E)$	knowledge graph
$o_t \in \text{Text}$	observation
$a_t \in \text{Text}$	action

### Subroutines

```
UPDATE_GRAPH(G, o_t):
  triples ← LLM.extract(o_t)
  // prompt: "extract (entity, relation, entity)
  //         triples from: [o_t]"

  for each (e1, rel, e2) in triples:

    v1 ← LLM.encode(e1)
    s ← cosine(v1, V)
    if max(s) > sim_thresh:
      i ← argmax(s)
      V[i] ← mean(V[i], v1) // update existing entity
    else:
      i ← len(V)
      V ← append(V, v1) // add new entity

    v2 ← LLM.encode(e2)
    s ← cosine(v2, V)
    if max(s) > sim_thresh:
      j ← argmax(s)
      V[j] ← mean(V[j], v2)
    else:
      j ← len(V)
      V ← append(V, v2)

    r ← relation_index(rel)
    E[i,j,r] ← 1.0 // observed: full confidence

  return (V, E)

COMPRESS(G, sim_thresh):
  S ← cosine(V, V) // pairwise similarity matrix
  clusters ← union_find(S, sim_thresh)
  // merge i,j if S[i,j] > sim_thresh

  V_new ← []
  idx ← {} // old index → new index
  for each cluster c:
    V_new ← append(V_new, mean(V[i] for i in c))
```

```

    for i in c:
        idx[i] ← len(V_new) - 1

E_new ← zeros(len(clusters), len(clusters), r)
for each (i,j,k) with E[i,j,k] > 0:
    E_new[idx[i], idx[j], k] ← max(
        E_new[idx[i], idx[j], k],
        E[i,j,k]
    )

return (V_new, E_new)           // H_t = G_t / ~_t

EXTEND(H_t, LLM):
prompt ← serialize(H_t)
// "known entities: [...]"
// known relations: [...]"
// predict missing or implied relations:"

candidates ← LLM.complete(prompt)
G_predicted ← copy(H_t)

for each (e1, rel, e2) in candidates:
    i ← resolve(e1, V, sim_thresh)
    j ← resolve(e2, V, sim_thresh)
    k ← relation_index(rel)
    G_predicted.E[i,j,k] ← 0.5 // predicted: half confidence

return G_predicted

DIFF(G_pred, G_actual):
// count edges in symmetric difference:
// edges predicted but not observed +
// edges observed but not predicted
return |(i,j,k) : G_pred.E[i,j,k] > 0|
    Δ |(i,j,k) : G_actual.E[i,j,k] > 0|

Main Loop

INITIALIZE:
G ← ([], [])           // empty graph
sim_thresh ← 0.8       // compression threshold
η ← 0.01              // threshold learning rate
t ← 0

for t = 0, 1, ..., T:

    // OBSERVE
    o_t ← environment.observe()

    // UPDATE
    G_t ← UPDATE_GRAPH(G, o_t)

    // COMPRESS
    H_t ← COMPRESS(G_t, sim_thresh)

```

```

// EXTEND
G_predicted ← EXTEND(H_t, LLM)

// ACT
a_t ← LLM.act(serialize(H_t), task)
environment.act(a_t)

// OBSERVE OUTCOME
o_{t+1} ← environment.observe()
G_actual ← UPDATE_GRAPH(H_t, o_{t+1})

// COMPUTE ERROR
error ← DIFF(G_predicted, G_actual)

// REWARD (unsupervised)
r_process      ← -error
r_compression  ← -len(H_t.V) / max(len(G_t.V), 1)
reward        ←  $\beta$  * r_process +  $\gamma$  * r_compression

// REFINE
sim_thresh ← sim_thresh +  $\eta$  * sign(error - error_prev)
  // high error → raise threshold → coarser compression
  // low error  → lower threshold → finer compression

update_policy(reward)

G          ← G_actual
error_prev ← error

```

### Convergence Condition

The algorithm converges to an RL closure  $(w^*, \pi^*)$  when:

error $\rightarrow 0$	graph accurately predicts its own updates
$ H_t.V  \rightarrow \text{stable}$	compression has stabilized
sim_thresh $\rightarrow \text{stable}$	granularity has stabilized

## 9 Stabilization and Convergence

The architecture does not guarantee convergence. Whether the dynamics stabilize depends on three conditions.

**Environmental stability.** The environment must have stable relational structure. If the world keeps changing its rules — if the relations that hold at time  $t$  are systematically different from those that hold at time  $t+1$  — the knowledge graph cannot converge. The agent will keep revising its model in response to observations without ever reaching a fixed point. This is not a failure of the architecture; it is the correct behavior of an embedded agent in a genuinely nonstationary environment. The quasi-periodic setting of The Imagination Machine III is the minimal environment in which convergence is guaranteed, because the environment has exact invariants — the frequency ratios — that the agent can recover.

**Compression aggressiveness.** If the similarity threshold `sim_thresh` is too low, the graph accumulates entities without merging them. The vertex set grows without bound and the compression step fails to reduce representational complexity. The adaptive threshold update in the main loop addresses this: persistent high prediction error raises the threshold, forcing coarser compression and preventing graph bloat.

**LLM consistency.** If the language model produces inconsistent completions — predicting different relations for the same compressed graph on different calls — prediction error will not decrease even if the world model is otherwise accurate. This is the Kan condition stated as a practical requirement on the extension operator: the language model must be able to fill every consistent horn, and must do so consistently. In terms of Corollary 1, the horn-filling condition is a testable property of the language model that determines whether the induced clique complex sequence satisfies the Kan condition.

**The garbage filter.** The architecture is not garbage-in-garbage-out in a naive sense. Edges that do not predict future observations generate high prediction error and are penalized through the reward signal. Stable structure — relations that keep being confirmed by action–observation cycles — survives compression. The compression–extension cycle is hostile to relational content that does not earn its place. Whether the filter is strong enough depends on the ratio of signal to noise in the environment and the expressiveness of the compression operator. This is an empirical question that the toy implementation of The Imagination Machine III is designed to address in the minimal quasi-periodic setting.

## 10 Implications

**Graph neural networks.** Message passing in a GNN is a compression operation: it collapses the local relational neighborhood of each vertex into a compressed feature vector. Theorem 1 implies that GNN dynamics induce simplicial dynamics on the clique complexes of their input graphs. GNNs are implicitly performing face map operations on simplicial complexes, and their expressive power is bounded by the simplicial structure they can detect.

**Topological data analysis.** A compression–extension orbit  $(G_t)$  defines a filtration of clique complexes  $(X(G_t))$ . The persistent homology of this filtration captures the relational invariants that survive across compression–extension cycles — precisely the fixed points of the closure operator  $T = F \circ g$  in the graph-theoretic realization.

**Knowledge graph reasoning.** Reasoning over knowledge graphs involves both compression (identifying equivalent entities) and extension (inferring missing relations). The present paper establishes that this reasoning process has a simplicial interpretation, connecting knowledge graph operations to the homotopy-theoretic properties of the Kan complex identified in The Imagination Machine X.

**Large language models.** The architecture clarifies the role of language models in embedded epistemic systems. A language model is a powerful extension operator: it can complete partial relational configurations, infer missing entities, and generate text consistent with a compressed world model. But it is not, by itself, an embedded epistemic agent. It lacks compression, action, and the feedback loop that drives prediction error to zero. Embedding a language model in the

compression–extension–action–observation loop of the present architecture is what promotes it from completion engine to component of an imagination machine.

## 11 Conclusion

The Imagination Machine architecture describes recursive cycles of compression and extension governing representation and reasoning in embedded epistemic systems.

The present paper has developed this architecture in two directions. Mathematically, graph quotients implement compression and graph completion implements extension, and the resulting dynamics induce simplicial face maps and simplicial completion operations on associated clique complexes. Computationally, a language model serving as extension operator — embedded in a loop with a knowledge graph, a compression step, an action channel, and an internal prediction error signal — realizes the architecture as a concrete unsupervised learning system.

The agent in this system learns from maximal epistemic aloneness. It has no teacher, no labels, no external ground truth. It has only its observations, the consequences of its actions, and the difference between what it predicted its world model would become and what it actually became. The structure that crystallizes from this process — the entities that survive compression, the relations that keep being confirmed, the graph that stabilizes — is the agent’s answer to the question of what its environment is.

That answer may be wrong. Convergence is not guaranteed. The filter may be too weak, the language model too inconsistent, the environment too nonstationary. These are empirical questions. But the conditions under which the answer is right are precisely the conditions under which an embedded agent can know anything at all: a world stable enough to have invariants, a compression aggressive enough to find them, and an extension consistent enough to test them.

The imagination machine does not know the world from outside. It constructs the world from within the only closure available to it.

# The Imagination Machine XII: Reconstructing Conceptual Structure in an Open Text World

Mark Tracy

## Abstract

The preceding papers in the Imagination Machine series develop a formal architecture for embedded epistemic systems. Observations generate world models through compression of relational structure, while extension predicts missing relations and guides action. In The Imagination Machine XI this architecture was realized as a graph-theoretic learning system whose world model is a dynamically updated knowledge graph.

The present paper introduces an experimental environment designed to test that architecture. An agent interacts with an open text world constructed from the TIM corpus itself. The agent observes only local textual segments, incrementally constructs a relational knowledge graph, compresses that graph through clustering, and predicts missing relations through extension. The resulting graph induces a clique complex whose simplicial structure provides a natural representation of higher-order conceptual relations.

Training occurs on TIM I–XI, while TIM XII is withheld as a test corpus. Evaluation measures the ability of the agent to predict structural updates induced by previously unseen text, including recovery of latent conceptual relations and completion of simplicial horns in the induced complex. To make the environment rigorous, the paper specifies an explicit latent concept vocabulary, an explicit typed relation vocabulary, a hybrid node-to-text and edge-to-text labeling protocol, and an explicit proposed ground-truth conceptual graph for the TIM world. The experiment therefore tests whether compression–extension dynamics allow an embedded agent to reconstruct relational invariants of its textual environment from within the textual surface alone.

## 1 Introduction

The Imagination Machine series develops a formal framework for epistemic systems embedded within their environment. Because such systems cannot access an external vantage point, knowledge must be defined operationally in terms of internal predictive coherence rather than correspondence with an independently accessible world.

The Imagination Machine XI introduced a computational realization of this framework. The agent’s world model is represented as a knowledge graph whose nodes correspond to entities and whose edges represent relations extracted from observations. Learning proceeds through repeated cycles of:

- (1) observation of new data,
- (2) updating of the knowledge graph,
- (3) compression of the graph through clustering of similar entities, and

(4) extension of the graph through prediction of missing relations.

These operations were shown to induce simplicial dynamics on the clique complex of the knowledge graph, linking the architecture to the simplicial completion conditions identified in earlier papers.

The present paper introduces an experimental setting designed to test this architecture empirically. The environment is an open text world constructed from the TIM corpus. The agent does not receive the corpus all at once. Instead, it receives local textual segments in response to actions taken with respect to its current world model. The environment is therefore partially observable, structured, and internally coherent, while remaining rich enough to support nontrivial graph reconstruction and simplicial completion tasks.

## 2 The Open Text World

Let

$$C = \{p_1, p_2, \dots, p_N\}$$

be a corpus segmented into atomic textual units, such as paragraphs, definitions, theorems, remarks, or algorithm blocks. Each segment represents a possible observation.

The environment contains a latent conceptual graph

$$G^* = (V^*, E^*)$$

whose nodes represent concepts, operators, mathematical structures, and documents, and whose edges represent typed relations among them. This graph is not directly accessible to the agent.

The environment also contains a latent simplicial family

$$\Sigma^* \subseteq \mathcal{P}(V^*)$$

whose elements represent coherent higher-order conceptual configurations.

The agent interacts with the environment only through textual observations sampled from the segmented corpus in response to actions.

### 2.1 Atomic Observation Units

Each element  $p_i \in C$  is one of the following:

- a paragraph,
- a formal definition,
- a theorem statement,
- a remark,
- an algorithm block, or
- a short subsection-introduction unit.

This segmentation ensures that observations are neither too fine-grained to support meaningful graph updates nor so coarse-grained that the environment degenerates into full-document access.

## 2.2 Training and Test Corpora

The training corpus is

$$C_{\text{train}} = \text{TIM I–XI},$$

and the test corpus is

$$C_{\text{test}} = \text{TIM XII}.$$

The agent trains on  $C_{\text{train}}$  alone. TIM XII is withheld during training and used only to evaluate whether the reconstructed world model predicts the structural updates induced by previously unseen text. In particular, the agent has no access to any segment of TIM XII during the training phase, and no node or edge whose sole textual evidence comes from TIM XII is available to the agent before evaluation.

## 3 Latent Ontology of the Environment

To make the environment explicit, we specify a finite node vocabulary and a finite typed relation vocabulary.

### 3.1 Node Vocabulary

The latent node set  $V^*$  is partitioned into four types.

#### Document Nodes

The following document nodes correspond to training-corpus papers and are available to the agent during training:

- TIM I
- TIM II
- TIM III
- TIM IV
- TIM V
- TIM VI
- TIM VII
- TIM VIII
- TIM IX
- TIM X
- TIM XI

The following document node corresponds to the withheld test corpus and is *not* available to the agent during training:

- TIM XII (*test corpus; withheld during training*)

## Core Architecture Concept Nodes

- embedded epistemic system
- observation
- world model
- inference
- implication
- inference–implication loop
- epistemic closure
- fixed point
- compression
- extension
- action
- prediction error
- relational invariant
- internal supervision
- world-model update

## Mathematical Structure Nodes

- graph
- graph morphism
- graph quotient
- graph completion
- clique
- clique complex
- simplex
- face map
- horn
- horn filling
- Kan condition
- simplicial dynamics

- quotient space
- equivalence relation
- classifier
- knowledge graph

### **Computational Architecture and Series-Thematic Nodes**

- entity embedding
- relation tensor
- clustering
- compression threshold
- extraction
- completion
- language model
- extension operator
- action policy
- unsupervised learning
- interactive text environment
- agent–environment interaction
- analogy
- abstraction
- holon
- institutional learning
- morality
- geometric theology
- categorical formulation
- quasi-periodic environment
- Koopman structure

### 3.2 Relation Vocabulary

Let the latent typed edge set  $E^*$  consist of triples

$$(u, r, v), \quad u, v \in V^*, r \in R^*,$$

where the relation vocabulary  $R^*$  is:

- defines
- develops
- implements
- realizes
- induces
- extends
- depends\_on
- acts\_on
- appears\_in
- analogizes\_with
- predicts
- updates
- compresses\_to
- completes
- clusters
- serves\_as
- stabilizes
- grounds
- tests

## 4 Hybrid Labeling Protocol

The environment must map latent nodes and latent edges to textual segments in a reproducible way. To do so, we define a hybrid labeling protocol.

## 4.1 Node-to-Text Incidence Map

Let

$$I_V : V^* \rightarrow \mathcal{P}(C)$$

assign to each latent node the set of corpus segments associated with it.

The map  $I_V$  is defined as the union of three components:

$$I_V(v) = I_V^{\text{exact}}(v) \cup I_V^{\text{alias}}(v) \cup I_V^{\text{manual}}(v).$$

**Definition 1** (Exact Lexical Anchoring). *For each node  $v \in V^*$ , define a canonical label  $L(v)$ . Then*

$$p \in I_V^{\text{exact}}(v) \iff L(v) \text{ occurs verbatim in } p.$$

**Definition 2** (Alias Normalization). *For each node  $v \in V^*$ , define an alias set  $A(v)$  containing normalized variants of the canonical label, including symbolic forms, hyphen variants, and close lexical alternatives. Then*

$$p \in I_V^{\text{alias}}(v) \iff \exists a \in A(v) \text{ such that } a \text{ occurs in } p.$$

**Definition 3** (Manual Concept Annotation). *For each node  $v \in V^*$ , a curator may add a segment  $p$  to  $I_V^{\text{manual}}(v)$  whenever  $p$  clearly expresses the concept denoted by  $v$  even if no canonical label or alias appears explicitly.*

**Remark 1.** *The exact component provides transparency, the alias component reduces brittleness, and the manual component captures implicit conceptual expression. The hybrid scheme therefore balances reproducibility and semantic adequacy.*

## 4.2 Edge-to-Text Incidence Map

Let

$$I_E : E^* \rightarrow \mathcal{P}(C)$$

assign to each latent typed edge the set of corpus segments expressing that relation.

Similarly,

$$I_E(e) = I_E^{\text{exact}}(e) \cup I_E^{\text{alias}}(e) \cup I_E^{\text{manual}}(e).$$

Here  $I_E^{\text{exact}}$  collects segments explicitly stating the relation,  $I_E^{\text{alias}}$  collects segments expressing a normalized variant, and  $I_E^{\text{manual}}$  captures curator-added relation evidence.

## 4.3 Examples of Alias Sets

Representative alias sets include:

- $A(\text{inference-implication loop})$ : {“inference-implication loop”, “inference-implication loop”, “ $F \circ g$ ”, “closure operator”}
- $A(\text{graph quotient})$ : {“graph quotient”, “quotient graph”}
- $A(\text{clique complex})$ : {“clique complex”, “complex of cliques”}
- $A(\text{Kan condition})$ : {“Kan condition”, “horn-filling condition”}
- $A(\text{extension operator})$ : {“extension operator”, “completion engine”}

## 5 Action Protocol

The agent does not access the corpus freely. Instead, it takes actions with respect to its current world model. The environment then returns textual evidence associated with those actions. During training, all actions draw exclusively from  $C_{\text{train}}$ ; segments from  $C_{\text{test}}$  are inaccessible until the evaluation phase.

### 5.1 Action Space

The action space consists of:

- `inspect( $v$ )`, where  $v$  is a current graph node,
- `inspect_relation( $u, r, v$ )`, where  $(u, r, v)$  is a current or predicted edge,
- `expand_cluster( $K$ )`, where  $K$  is a current node-cluster,
- `explore_random()`, and
- `verify_edge( $u, r, v$ )`, which requests evidence for a predicted edge.

### 5.2 Observation Sampling

Let  $H_t \subseteq C$  denote the set of segments already observed up to time  $t$ .

Then the observation rules are:

$$\begin{aligned}\text{inspect}(v) &\sim \text{Sample}(I_V(v) \setminus H_t), \\ \text{inspect\_relation}(u, r, v) &\sim \text{Sample}(I_E(u, r, v) \setminus H_t).\end{aligned}$$

If the corresponding unseen set is empty, the environment samples instead from the full associated set with preference for least frequently returned segments.

**Remark 2.** *This protocol ensures that an action requests evidence about a node or relation, not unrestricted access to the full corpus. The environment is therefore partially observable and exploration-dependent. During training the sampling domain is restricted to  $C_{\text{train}}$ ; segments from  $C_{\text{test}}$  become accessible only during the evaluation phase described in Section 12.*

## 6 Agent World Model

The agent maintains a knowledge graph

$$G_t = (V_t, E_t)$$

whose nodes correspond to discovered entities and whose edges correspond to relations extracted from observations.

Entities are represented by embedding vectors

$$z_v \in \mathbb{R}^d,$$

and relations are stored in a relation tensor.

Upon observing a new text segment, the agent extracts relational triples

$$(v_i, r, v_j)$$

which are used to update the graph.

This graph is not the latent graph  $G^*$ . It is the agent’s current world model of the environment.

## 7 Compression and Extension

After each graph update, the world model undergoes two operations.

### 7.1 Compression

Compression merges nodes with similar embeddings, producing equivalence classes of entities that represent conceptual abstractions.

This operation can be interpreted as a graph quotient

$$G_t \rightarrow H_t$$

which reduces representational redundancy.

### 7.2 Extension

Extension predicts missing relations in the compressed graph. In the present implementation this prediction is performed by a language model conditioned on the current graph state.

The predicted relations form a proposed graph

$$G_{t+1}^{\text{pred}}$$

which represents the agent’s expectation of the next graph update.

## 8 Simplicial Structure

The knowledge graph induces a clique complex

$$X(G_t)$$

whose simplices correspond to sets of mutually connected entities.

Compression induces simplicial face maps by collapsing equivalent nodes, while extension predicts missing simplices through completion of partially specified horns.

This connects the architecture to the simplicial completion framework developed in earlier papers.

## 9 Exact Proposed Ground-Truth Graph

We now expose the proposed latent graph for the TIM world. This graph is not given to the agent during training, but it defines the environment used for evaluation.

### 9.1 Core Proposed Edge Set

The following typed triples constitute the first-pass proposed ground-truth graph. Edges involving TIM XII are marked (*withheld*) to indicate that they are not available to the agent during training and form part of the evaluation target.

## Series-Level Document Structure

(TIM I, defines, epistemic closure)  
(TIM I, defines, inference–implication loop)  
(TIM I, defines, fixed point)  
(TIM I, develops, world model)  
(TIM II, develops, agent–environment interaction)  
(TIM III, develops, quasi–periodic environment)  
(TIM III, develops, Koopman structure)  
(TIM IV, develops, institutional learning)  
(TIM V, develops, analogy)  
(TIM V, develops, abstraction)  
(TIM VI, develops, horn filling)  
(TIM VI, develops, holon)  
(TIM VII, develops, morality)  
(TIM VIII, develops, geometric theology)  
(TIM IX, develops, categorical formulation)  
(TIM X, develops, simplicial dynamics)  
(TIM XI, extends, TIM X)  
(TIM XI, realizes, compression)  
(TIM XI, realizes, extension)  
(TIM XII, tests, TIM XI) (*withheld: evaluation target*)

## Epistemic Architecture

(observation, grounds, world model)  
(inference, acts\_on, observation)  
(implication, acts\_on, world model)  
(inference–implication loop, depends\_on, inference)  
(inference–implication loop, depends\_on, implication)  
(epistemic closure, depends\_on, fixed point)  
(fixed point, stabilizes, world model)  
(prediction error, updates, world model)  
(internal supervision, realizes, prediction error)  
(world-model update, updates, world model)

## Compression–Extension Architecture

(compression, acts\_on, world model)  
(extension, acts\_on, world model)  
(compression, compresses\_to, equivalence relation)  
(compression, implements, graph quotient)  
(extension, implements, graph completion)  
(graph quotient, implements, compression)  
(graph completion, implements, extension)  
(knowledge graph, serves\_as, world model)  
(language model, serves\_as, extension operator)  
(extension operator, realizes, completion)  
(extraction, updates, knowledge graph)  
(completion, updates, knowledge graph)  
(clustering, implements, compression)  
(compression threshold, updates, clustering)  
(action policy, predicts, observation)  
(action, grounds, observation)  
(interactive text environment, grounds, observation)  
(unsupervised learning, realizes, internal supervision)

## Graph–Simplicial Correspondence

(graph, induces, clique complex)  
(clique, depends\_on, graph)  
(clique complex, depends\_on, clique)  
(clique complex, depends\_on, simplex)  
(compression, induces, face map)  
(extension, induces, horn filling)  
(horn filling, completes, horn)  
(horn filling, induces, Kan condition)  
(Kan condition, depends\_on, horn)  
(simplicial dynamics, depends\_on, clique complex)  
(simplicial dynamics, depends\_on, face map)  
(simplicial dynamics, depends\_on, horn filling)  
(quotient space, depends\_on, equivalence relation)  
(classifier, compresses\_to, quotient space)

## Cross-Series Structural Correspondences

- (analogy, analogizes\_with, abstraction)
- (analogy, analogizes\_with, horn filling)
- (holon, analogizes\_with, horn filling)
- (categorical formulation, extends, compression)
- (categorical formulation, extends, extension)
- (Koopman structure, depends\_on, relational invariant)
- (quasi-periodic environment, grounds, relational invariant)
- (institutional learning, realizes, compression)
- (institutional learning, realizes, extension)

**Remark 3.** *The graph above is a first-pass proposed ground truth, not a claim of uniquely correct ontology. Its role in the experiment is to provide a controlled latent world against which graph reconstruction, horn completion, and graph-update prediction can be evaluated. All edges whose subject or object is TIM XII are withheld from the agent during training. They are part of the evaluation target: the agent is expected to predict them through extension from the structure learned during training on TIM I–XI alone.*

## 10 Exact Proposed Ground-Truth Simplices

The latent simplicial family  $\Sigma^*$  is generated by coherent conceptual motifs. The following are the proposed core simplices.

### Epistemic Simplices

- {inference, implication, inference–implication loop}
- {inference–implication loop, fixed point, epistemic closure}
- {prediction error, internal supervision, world-model update}

### Compression–Extension Simplices

- {compression, graph quotient, equivalence relation}
- {extension, graph completion, extension operator}
- {knowledge graph, world model, prediction error}
- {completion, language model, extension operator}
- {compression, extension, world model, action}

## Graph–Simplicial Simplices

- {graph, clique, clique complex, simplex}
- {compression, graph quotient, face map}
- {extension, graph completion, horn filling}
- {horn, horn filling, Kan condition}
- {simplicial dynamics, clique complex, face map, horn filling}

## Cross-Series Simplices

- {analogy, abstraction, horn filling}
- {holon, horn filling, analogy}
- {quasi–periodic environment, Koopman structure, relational invariant}
- {institutional learning, compression, extension}

## 11 Representative Incidence Tables

In practice the full incidence tables would appear in an appendix. We include representative samples here. All supporting segments listed are drawn from  $C_{\text{train}}$ . In the full experimental release, the complete incidence maps  $I_V$  and  $I_E$  would be provided together with written annotation guidelines specifying the use of exact lexical anchoring, alias normalization, and manual concept or relation annotation.

### Sample Node-Incidence Table

Node	Representative supporting segments
compression	TIM XI Introduction; TIM XI Section 3; TIM XI Section 6; TIM X Section 4
extension	TIM XI Introduction; TIM XI Section 4; TIM XI Section 6; TIM XI Section 7; TIM X Section 4
graph quotient	TIM XI Section 3; TIM XI Section 6
graph completion	TIM XI Section 4; TIM XI Section 6
clique complex	TIM XI Section 5; TIM XI Section 6; TIM XI Section 10
horn filling	TIM VI Section 4; TIM VI Section 5; TIM X Section 8; TIM XI Section 6
Kan condition	TIM X Section 8; TIM XI Section 6; TIM XI Section 9
epistemic closure	TIM I Introduction; TIM I fixed-point discussion; TIM XI Introduction; TIM XI Conclusion

## Sample Edge-Incidence Table

Edge	Representative supporting segments
(graph quotient, implements, compression)	TIM XI Section 3
(graph completion, implements, extension)	TIM XI Section 4
(compression, induces, face map)	TIM XI Section 6
(extension, induces, horn filling)	TIM XI Section 6
(language model, serves_as, extension operator)	TIM XI Section 7; TIM XI Conclusion
(prediction error, updates, world model)	TIM XI Section 7; TIM XI Algorithm

## 12 Experimental Protocol

Training is performed on  $C_{\text{train}} = \text{TIM I–XI}$ . The agent sequentially observes textual segments drawn from this corpus and updates its knowledge graph using the compression–extension cycle described above.

TIM XII is withheld during training. During evaluation, the agent is exposed to segments of TIM XII one at a time and predicts the structural updates each segment induces before observing it. No segment of TIM XII is accessible to the agent prior to this evaluation phase.

### 12.1 Main Interaction Loop

At time  $t$ :

- (1) The agent selects an action  $a_t$ .
- (2) The environment returns a segment  $o_t \in C$  according to the action protocol.
- (3) The agent extracts triples and updates  $G_t$ .
- (4) The agent compresses  $G_t$  to  $H_t$ .
- (5) The agent extends  $H_t$  to obtain  $G_{t+1}^{\text{pred}}$ .
- (6) The agent selects the next action using its current policy.
- (7) The next segment is revealed, producing the actual graph update  $G_{t+1}^{\text{actual}}$ .
- (8) The prediction loss is computed as

$$L_t = \text{diff}(G_{t+1}^{\text{pred}}, G_{t+1}^{\text{actual}}).$$

### 12.2 Evaluation Tasks

Three evaluation tasks are considered.

#### Edge Recovery

Selected edges from  $G^*$  are held out from training. The agent’s ability to recover them through extension is measured using precision, recall, and  $F_1$ .

## Horn Completion

Selected simplices in  $\Sigma^*$  are partially withheld. The agent is asked to complete the corresponding horns. Accuracy on these horn-completion tasks measures whether the agent recovers higher-order conceptual structure.

Representative horn-completion tasks include:

- {compression, graph quotient, ?}  $\rightarrow$  face map,
- {extension, graph completion, ?}  $\rightarrow$  horn filling,
- {inference, implication, ?}  $\rightarrow$  inference-implication loop,
- {inference-implication loop, fixed point, ?}  $\rightarrow$  epistemic closure.

## Graph Stabilization

The stability of the knowledge graph is measured by tracking:

- number of nodes,
- number of edges,
- number and composition of clusters,
- compression threshold, and
- graph-update prediction error.

# 13 Discussion

If compression-extension dynamics successfully capture the relational structure of the environment, the agent should anticipate structural updates introduced by previously unseen text. In this setting, successful prediction of graph updates induced by TIM XII demonstrates that the agent has reconstructed relational invariants of the TIM conceptual world from TIM I–XI alone.

More broadly, the experiment illustrates how embedded epistemic systems can learn structural models of their environment through internally generated prediction tasks. The present construction is intentionally self-contained: the TIM corpus functions as a controlled textual world in which the latent ontology, latent relation structure, and latent simplicial motifs can all be explicitly defined. This makes the environment suitable for proof-of-concept evaluation of the imagination-machine architecture before extension to broader corpora.

## 13.1 Scope and Status of the Environment

The TIM world defined here is a controlled proof-of-concept environment. Its purpose is not to establish immediate generalization to arbitrary corpora, but to test whether the architecture can recover explicit latent conceptual structure from a corpus whose ontology, relation structure, and simplicial motifs can be specified in advance.

Accordingly, the present experiment should be read as an internal validation study of the compression-extension architecture. A positive result would show that the architecture can recover and predict structural updates in a textual world engineered to make such evaluation possible. Whether the same architecture generalizes to broader or noisier corpora is a distinct empirical question left for future work.

## 14 Conclusion

This paper introduces an experimental realization of the Imagination Machine architecture in an open text world.

An agent embedded in a textual environment incrementally reconstructs a conceptual knowledge graph, compresses that graph through clustering, and predicts missing relations through extension. The induced simplicial structure provides a natural representation of higher-order conceptual relations.

To make the environment rigorous, the paper specifies an explicit latent ontology, an explicit relation vocabulary, a hybrid labeling protocol, and an explicit proposed ground-truth graph and simplex family for the TIM world. TIM XII is withheld during training; evaluation on its segments tests whether the agent can predict structural updates induced by previously unseen text. Successful prediction demonstrates that prediction of structural updates can serve as an operational definition of understanding for embedded epistemic systems.

# The Imagination Machine XIII: Notes on Engineering an Embedded Epistemic System

Mark Tracy

March 2026

## Abstract

The preceding papers in the Imagination Machine series develop a formal architecture for embedded epistemic systems and culminate in a concrete computational realization based on compression–extension dynamics over knowledge graphs. The present note records an observation arising during the transition from theory to implementation: engineering such a system is not a direct translation of theory into code. Instead, the engineering process itself forms a learning trajectory through design space, guided by prediction, observation, and iterative refinement.

In this sense the process of constructing an imagination machine is itself an instance of the epistemic dynamics the architecture describes. The engineer occupies the same structural position as the agent in the framework: embedded within a partially observable environment, constructing models of system behavior through cycles of compression and extension. The purpose of this note is to document that symmetry.

## 1 Theory and Engineering

The preceding papers in the Imagination Machine series develop a theoretical framework for embedded epistemic systems. In this framework an agent constructs a world model through repeated cycles of:

1. observation,
2. representation,
3. compression of relational structure,
4. extension through prediction of missing relations, and
5. update through prediction error.

The Imagination Machine XI gives a graph-theoretic realization of this process, and The Imagination Machine XII introduces an experimental environment in which the architecture may be evaluated.

At this point the project transitions from theory to engineering.

A natural expectation might be that implementation proceeds by directly translating the theoretical architecture into software. In practice this expectation is incorrect. Engineering is not a linear execution of theory. It is a separate discovery process.

## 2 The Engineering Learning Graph

Theoretical development proceeds through logical structure. Concepts are defined, relations among them are established, and the resulting structure stabilizes once the definitions and propositions cohere.

Engineering follows a different dynamic.

Instead of a logical graph of concepts, engineering produces a trajectory through a space of working configurations. Each configuration proposes a particular implementation of the architecture. Experiments reveal how that configuration behaves, producing observations that guide the next revision.

Typical engineering progress therefore takes the form:

$$\text{prototype} \rightarrow \text{observation} \rightarrow \text{failure} \rightarrow \text{modification} \rightarrow \text{refinement}.$$

Early implementations rarely resemble the final architecture closely. They reveal hidden constraints of the system and expose interactions that are not visible at the level of abstract theory.

Over time, successive revisions converge toward structures that more faithfully realize the theoretical design.

## 3 Embeddedness of the Engineer

The architecture developed in this series describes an embedded agent learning about its environment through cycles of compression and extension.

During the engineering phase, the same structure appears at another level.

The engineer does not possess perfect knowledge of the system being constructed. Instead the engineer interacts with prototypes, observes their behavior, and forms increasingly refined models of the system's dynamics.

The resulting process mirrors the epistemic loop of the imagination machine itself:

$$\text{prediction} \rightarrow \text{experiment} \rightarrow \text{error} \rightarrow \text{model update}.$$

In this sense the engineer occupies the same structural position with respect to the developing system that the agent occupies with respect to its environment.

**Remark 1.** *Building an imagination machine is itself an instance of the imagination machine process. The engineer learns the structure of the system through the same compression–extension dynamics that the system is designed to perform.*

## 4 Consequences for Implementation

This observation suggests a practical principle for early implementations.

The goal of the first prototype is not correctness but information. A small system that fails clearly provides more insight into the architecture's behavior than a large system whose complexity obscures its dynamics.

Early prototypes therefore function as exploratory instruments. They expose how the components of the architecture interact in practice and reveal which parts of the theoretical design require adjustment or refinement.

Such iterations are not deviations from the framework. They are the mechanism by which the theoretical architecture becomes operational.

## 5 A Structural Symmetry

The Imagination Machine series began as a conceptual investigation into how an embedded epistemic system might construct coherent representations of its environment from within the limits of its observational surface.

As the project moves from theory toward implementation, a structural symmetry becomes apparent. The process of constructing the system follows the same dynamics that the system itself is designed to exhibit. The engineering phase is not external to the framework — it is an instance of it.

This symmetry is not incidental. It reflects a general feature of embedded systems: any process capable of building a system that learns from within must itself proceed by learning from within. The architecture does not stand outside the conditions it describes.

## 6 Conclusion

The Imagination Machine architecture describes how an embedded system can learn structural invariants of its environment through cycles of compression and extension.

When the architecture is implemented in practice, the engineering process itself follows a similar pattern of iterative model construction driven by prediction error and observation.

The symmetry between these processes highlights a broader point. Systems capable of learning about their environment must themselves be constructed through learning processes embedded within the constraints of reality.

The imagination machine therefore appears twice in the project: once as the system being designed, and once as the process by which the design itself is realized.

# The Imagination Machine XIV: Relational Invariants, Quotient Structure, and the Reproducibility of Science

Mark Tracy

March 2026

## Abstract

Scientific knowledge stabilizes through the reproducibility of experimental results across independent observers and experimental contexts. This paper interprets reproducibility through the compression–extension architecture developed in the Imagination Machine series. Observational data are first produced in highly indexical form, tied to particular observers, instruments, and experimental circumstances. Scientific modeling compresses these observations through a classifier that quotients away observational detail while preserving selected relational invariants. A scientific law is then interpreted as a relational structure that remains invariant under this quotient map. Reproducibility corresponds to the stability of these invariants across independent experiments. From this perspective the methodology of science may be understood as the collective construction of quotient representations of the observational world, within which invariant relations appear as physical law.

## 1 Introduction

The Imagination Machine series develops a formal framework for embedded epistemic systems. In this framework an agent constructs a world model by iteratively compressing observational data into a representation that preserves relational structure while discarding irrelevant detail. The admissible models of the system appear as fixed points of the inference–implication loop introduced in the first paper of the series.

A central question in the philosophy of science concerns the reproducibility of experimental results. Independent laboratories performing the same experiment under different conditions frequently obtain observational data that differ in numerous superficial ways. Nevertheless, scientific laws appear as stable regularities that persist across these differences.

The present paper interprets reproducibility as a consequence of the quotient structure induced by representational compression. Scientific laws correspond to relational invariants that remain stable under the quotient map from observational data to scientific representation.

## 2 Observational Surfaces

Every experiment produces data in a highly indexical form. Observations are tied to particular observers, instruments, experimental procedures, and environmental circumstances.

**Definition 1** (Observation Event). *An observation event is a tuple*

$$x = (o, a, t, \ell, p, m)$$

where  $o$  denotes the observer,  $a$  the apparatus configuration,  $t$  the time of observation,  $\ell$  the spatial location,  $p$  the experimental protocol, and  $m$  the measured outcome.

Let  $D$  denote the space of such observation events. Two observation events may differ in many of these parameters while nevertheless expressing the same underlying regularity.

### 3 Representational Compression

A scientific model compresses the observational surface by mapping observation events into a representation that preserves selected relational structure.

**Definition 2** (Scientific Classifier). *Let*

$$\pi : D \rightarrow Z$$

*be a classifier mapping observation events into representational states  $Z$ . The map  $\pi$  induces an equivalence relation on  $D$  defined by*

$$x \sim_{\pi} y \quad \text{if and only if} \quad \pi(x) = \pi(y).$$

The quotient space

$$Q = D / \sim_{\pi}$$

groups together observation events that are treated as equivalent by the scientific model.

**Remark 1.** *The classifier  $\pi$  may include transformations such as coordinate normalization, calibration correction, statistical averaging, or parameter estimation. These operations discard observational detail while preserving relational structure relevant to the theory.*

### 4 Relational Invariants

Scientific laws correspond to relations that remain invariant across equivalence classes in the quotient representation.

**Definition 3** (Relational Invariant). *A relation  $R$  defined on the representational space  $Z$  is a relational invariant if it holds for all representatives of an equivalence class in  $Q$ .*

Examples include the constancy of gravitational acceleration in Newtonian mechanics, the Lorentz invariance of spacetime intervals in relativity, and the ideal gas relation in thermodynamics.

**Remark 2.** *The invariance of these relations reflects the fact that the observational differences removed by the quotient map do not alter the relational structure preserved by the model.*

### 5 Reproducibility

The reproducibility of scientific results can now be interpreted as stability under the quotient map.

**Definition 4** (Reproducible Result). *An experimental result is reproducible if observation events from independent experiments fall into the same equivalence class of  $Q$  under the classifier  $\pi$ .*

In practice this means that while raw measurements may vary across laboratories, the representational compression applied by the scientific model maps them to the same relational structure.

**Remark 3.** *Experimental methodology exists largely to ensure that independent investigators apply compatible compression maps. Standardized protocols, calibration procedures, and statistical analysis all serve to align the quotient representations used by different laboratories.*

## 6 Scientific Method as Quotient Construction

The methodology of science may therefore be interpreted as a collective process for constructing quotient representations of observational reality.

Different laboratories act as independent epistemic agents observing the same environment through distinct observational surfaces. A scientific theory stabilizes when the compression map used by these agents yields consistent relational invariants across their respective data.

**Proposition 1.** *Scientific consensus emerges when independently observed data sets share a common quotient representation under a shared classifier.*

## 7 Symmetry and Physical Law

Modern physics frequently formulates laws in terms of symmetry principles. These symmetries express invariance under transformations such as spatial translation, temporal translation, or coordinate change.

Within the present framework these symmetries appear naturally as transformations that leave the quotient representation unchanged. A symmetry therefore corresponds to an operation on observation events that preserves equivalence classes in the quotient space.

**Remark 4.** *This perspective explains the centrality of symmetry in modern physics: symmetry transformations are precisely those operations that preserve the relational invariants retained by the representational compression.*

## 8 Conclusion

The Imagination Machine framework interprets knowledge formation as the compression of observational data into representations that preserve relational structure. Scientific laws appear as invariants within the quotient representations produced by this compression.

From this perspective the reproducibility of science is not mysterious. Independent experiments produce different observational details, but once those details are quotiented away by the scientific classifier, the same relational invariants emerge. Reproducibility therefore reflects the stability of these invariants across observational contexts.

Scientific practice can thus be understood as a distributed epistemic process in which many observers collaboratively construct quotient representations of the observational world. Physical law corresponds to the relational structure that remains invariant within those representations.

# The Imagination Machine XV: The View from Nowhere and the Center of the Hypersphere

Mark Tracy

March 2026

## Abstract

The Imagination Machine series develops a framework for embedded epistemic systems: systems that must construct world models from within the world they attempt to know. A recurring consequence of this framework is that no embedded observer can attain a literal view from nowhere. All representation arises from within a local observational surface.

At the same time, the series has increasingly suggested a geometric picture in which embedded observers inhabit a three-dimensional manifold understood as a cross-section of a four-dimensional containing structure. In *The Imagination Machine VIII*, this containing structure was interpreted as the three-sphere, or hypersphere, whose center is inaccessible from within the embedded manifold.

The present note records a further identification: the philosophical ideal of the “view from nowhere” corresponds, in this geometric register, to the center of the four-dimensional hypersphere. The center is not an embedded location. It is the unique point equidistant from every point on the hypersphere, and thus the unique point of maximal symmetry with respect to the embedded manifold. It is therefore a precise geometric analogue of an invariant standpoint relative to all local views, while remaining unavailable to any embedded observer.

This identification clarifies the relation between local epistemic closure and global symmetry. The view from somewhere is the actual condition of embedded knowledge. The view from nowhere is the unoccupiable center relative to which all such local views are symmetrically situated. The result preserves the central claim of the series—that knowledge is necessarily embedded—while providing a geometric interpretation of the philosophical impulse toward objectivity.

## 1 Introduction

The Imagination Machine series begins from a simple constraint: an epistemic system embedded within the world has no access to an external vantage point from which to compare its representations with the world “as it is in itself.” Knowledge must therefore be understood not as correspondence with an independently accessible outside, but as the stabilization of representation through the internal closure of epistemic dynamics.

This claim has been developed formally through the inference–implication loop, the fixed-point condition for admissible world models, the inclusion of classifiers within the observation space, and the interpretation of law as relational invariance in a quotient representation. Across these constructions, the point has remained constant: every actual act of knowing is a *view from somewhere*.

At the same time, later papers in the series introduced a geometric register in which this embeddedness could be pictured more sharply. In particular, *The Imagination Machine VIII* proposed

that the maximally conservative geometry for an embedded observer is the hypersphere: a closed structure with no accessible center and no boundary from the point of view of the observer inhabiting its three-dimensional surface.

The present note records a further recognition within that geometry. Philosophers have often spoken of the ideal of a *view from nowhere*: a standpoint purified of local bias, contingency, and perspective. Within the framework of the present series, such a standpoint cannot be occupied by any embedded observer. But the geometric picture suggests that this philosophical ideal is not mere nonsense. It has a precise structural correlate.

The proposal is this: *the view from nowhere is the center of the four-dimensional hypersphere.*

This does not mean that the center is accessible. It means that the center plays the role of the unique point invariant with respect to all embedded viewpoints. Every point on the hypersphere is equidistant from it. No embedded point is privileged relative to it. The center is therefore the geometric image of the nonlocal standpoint toward which objectivity gestures, while remaining strictly unavailable from within the manifold of embedded observation.

## 2 The Embedded Condition

The foundational claim of the series is that an epistemic system does not know the world from outside. It knows only through the observational surface available to it.

Formally, earlier papers describe this through the inference–implication loop

$$\Gamma \xrightarrow{F} W \xrightarrow{g} \Gamma$$

with induced operator

$$T = F \circ g : W \rightarrow W,$$

where admissible world models are fixed points

$$T(w^*) = w^*.$$

A world model is therefore not justified by appeal to an external standpoint but by internal reproduction under the epistemic loop. The system achieves closure from within its own observational and inferential conditions.

This immediately rules out any *literal* view from nowhere for an embedded observer. Every actual model is indexed to a closure, every closure to an observational profile, and every observational profile to the interior of the system’s relation to its environment.

The embedded condition may therefore be stated as follows.

**Definition 1** (Embedded View). *An embedded view is any observational or representational standpoint generated from within the observational surface of an epistemic system.*

Every actual act of knowledge available to an embedded system is an embedded view in this sense.

## 3 The Hypersphere

We now recall the geometric picture.

Let

$$S^3 = \{x \in \mathbb{R}^4 : \|x\| = r\}$$

for some radius  $r > 0$ . This is the three-sphere, or hypersphere: a three-dimensional manifold embedded in four-dimensional Euclidean space.

An observer confined to  $S^3$  may move locally in three dimensions, construct geometries internal to the manifold, and treat its own observational world as three-dimensional. But such an observer does not inhabit the ambient  $\mathbb{R}^4$  as a freely accessible space. In particular, the point

$$0 \in \mathbb{R}^4$$

which serves as the center of the hypersphere is not a point of the manifold  $S^3$  itself.

This produces the familiar properties.

**Proposition 1.** *For every  $x \in S^3$ ,  $\|x\| = r$ . Hence every point of  $S^3$  is equidistant from the center 0.*

*Proof.* This is immediate from the definition of  $S^3$ . □

**Remark 1.** *From the perspective of an observer embedded in  $S^3$ , the center is not locatable by traversal within the manifold. It is not one more point among the points of the observer's world.*

**Remark 2.** *Likewise,  $S^3$  has no boundary within itself. An embedded observer moving in any available direction does not encounter an edge.*

Thus the hypersphere provides a precise way to model a world in which all actual viewpoints are local, while global symmetry is defined relative to a point unavailable from within the manifold.

## 4 The View from Somewhere

The first half of the proposal is straightforward.

**Definition 2** (View from Somewhere). *A view from somewhere is the standpoint associated with some point  $x \in S^3$ , together with the local observational and representational structure available from that point as an embedded position on the manifold.*

This is the geometric counterpart of the embedded epistemic system described throughout the series. Every actual observer is situated at some  $x \in S^3$ , or more generally at some local region of the manifold, and every observation is indexed to such a situation.

No point on  $S^3$  is the center. Every point is local. Every point is one point among others. Thus every actual epistemic position is perspectival in the precise sense that it is indexed to a location on the manifold.

This does not imply arbitrariness. Embedded views can converge, stabilize, and share relational invariants. But they remain views from somewhere.

## 5 The View from Nowhere

We now state the central proposal.

**Definition 3** (View from Nowhere). *The view from nowhere is the geometric role played by the center  $0 \in \mathbb{R}^4$  of the hypersphere  $S^3$ : the unique point equidistant from every point of the embedded manifold and therefore the unique point of maximal symmetry with respect to all embedded locations.*

This definition requires immediate clarification.

The claim is not that an embedded observer can occupy the center. The center is not a possible embedded location. Rather, the claim is that if one asks what the *philosophical idea* of a view from nowhere corresponds to in the geometry of embeddedness, the answer is: the center.

Why?

Because the center satisfies the structural requirements associated with that philosophical ideal.

1. It does not privilege any point on the hypersphere.
2. It stands in the same metric relation to every point on the hypersphere.
3. It is not itself one local perspective among others.
4. It serves as the symmetry point relative to which all local perspectives are situated.

The center is therefore the *invariant correlate* of all embedded views without being one of them.

**Proposition 2.** *The center 0 of the hypersphere is not an embedded viewpoint but the unique symmetry point relative to all embedded viewpoints.*

*Proof.* By definition,  $0 \notin S^3$ , so it is not an embedded point on the manifold. For every  $x \in S^3$ ,  $\|x\| = r$ , so all embedded points are equally related to 0. Hence 0 does not privilege any embedded location and functions as the unique symmetry point relative to the manifold.  $\square$

This is exactly the structure the phrase “view from nowhere” has always tried to capture: a standpoint that is not just another somewhere, but a point of invariance with respect to all somewheres.

## 6 Objectivity as Orientation Toward the Center

The proposal permits a sharpening of the idea of objectivity.

If every actual act of knowing is a view from somewhere, then objectivity cannot mean the literal occupation of the view from nowhere. Embedded systems cannot become non-embedded. But objectivity can mean something else: *orientation toward invariants that do not depend on the particular local position of the observer.*

In the hypersphere picture, this means orientation toward structures that remain stable across local views on  $S^3$ . The center is the geometric image of that invariance, even though no observer can stand there.

Thus objectivity may be reinterpreted as follows.

**Definition 4** (Embedded Objectivity). *Embedded objectivity is the approximation of invariance across local viewpoints without the occupation of a nonlocal viewpoint.*

The center is the formal limit of this aspiration. It is the standpoint toward which embedded inquiry orients itself when it seeks what holds regardless of particular local position. But that standpoint remains unoccupiable.

This is consistent with the earlier papers of the series. Scientific law, for example, was interpreted as relational invariance under quotient structure. That is already a form of embedded objectivity: not escape from perspective, but stabilization of what survives variation across perspectives.

The present note simply adds a geometric interpretation. The center of the hypersphere is the symmetry point relative to which such invariance is defined.

## 7 The View from Nowhere Is Unoccupiable

A crucial consequence follows.

**Theorem 1.** *For an observer embedded in  $S^3$ , the view from nowhere is structurally definable but not epistemically occupiable.*

*Proof.* The center 0 is definable in the ambient geometry  $\mathbb{R}^4$  as the unique point from which all points of  $S^3$  lie at distance  $r$ . Hence it is structurally definable.

But  $0 \notin S^3$ . An observer embedded in  $S^3$  has access only to positions and motions internal to  $S^3$ . Therefore the observer cannot occupy 0 as an embedded location.

Hence the view from nowhere, identified with the center, is structurally definable but not epistemically occupiable for an embedded observer.  $\square$

This theorem captures the precise reconciliation of the two intuitions that have animated the series.

First: all knowledge is from somewhere. Second: there is a meaningful sense in which objectivity aims beyond any particular somewhere.

The view from nowhere is not nonsense. It is the center. But the center is not a place we can stand. It is a point of symmetry relative to which embedded knowledge may orient itself without ever escaping embeddedness.

## 8 Relation to the Earlier Series

The present note does not introduce a new architecture. It simply adds a geometric identification to the framework already developed.

### Relation to TIM I

The first paper established that embedded systems have no external vantage point. The present note preserves that claim. The center is not an external perspective that the embedded system may access. It is the structural correlate of the unattainable ideal of such a perspective.

### Relation to TIM VIII

The eighth paper proposed the hypersphere as the geometry of maximal epistemic humility: a closed containing structure with no accessible center and no boundary from within. The present note identifies that inaccessible center more explicitly with the philosophical ideal of the view from nowhere.

### Relation to TIM XIV

The fourteenth paper treated reproducibility as the stabilization of relational invariants across independent observers and contexts. In the present language, those invariants may be understood as precisely the sort of structures toward which embedded inquiry is oriented when it seeks to approximate the view from nowhere without ever occupying it.

## 9 A Philosophical Consequence

A final consequence may be stated cleanly.

The usual opposition between “view from somewhere” and “view from nowhere” is too crude. It treats them as if they were two equally available epistemic options, one local and one universal. The geometric picture developed here shows instead that they belong to different categories.

The view from somewhere is an *actual epistemic position*. The view from nowhere is a *structural symmetry point*.

The first is inhabitable but local. The second is universal but uninhabitable.

This removes a long-standing confusion. The mistake is not in wanting objectivity. The mistake is in imagining that objectivity requires the literal occupation of a nonlocal standpoint. What it requires instead is the disciplined construction of representations that track what remains stable across local positions.

That is exactly the project of the series: to describe how embedded systems can generate coherent knowledge without pretending to stand outside the world.

## 10 Conclusion

The Imagination Machine series argues that every actual act of knowledge is embedded. No observer inside the world can attain a literal view from nowhere. The present note preserves that claim while giving the ideal of the view from nowhere a precise geometric interpretation.

If embedded observers inhabit the three-dimensional surface of a four-dimensional hypersphere, then the view from somewhere corresponds to any local position on that surface. The view from nowhere corresponds to the center of the hypersphere: the unique point equidistant from every embedded point, not itself an embedded point, and therefore the unique symmetry point relative to all local views.

This identification clarifies the relation between perspective and objectivity. The view from nowhere is not an accessible epistemic location. It is the geometric image of invariance across local views. Objectivity is therefore not escape from embeddedness, but orientation toward structures that remain stable across the plurality of embedded standpoints.

The center of the hypersphere does not abolish the view from somewhere. It explains why the view from somewhere can seek universality without ever ceasing to be somewhere.