# Semantic Commitment as an Architectural Authorization Problem in Long-Horizon Adaptive Systems

Maksim Barziankou, Poznan 2026

## Abstract

Adaptive systems operating over long horizons must interpret observations and act under uncertainty, partial observability, and irreversible structural consequences. Existing approaches predominantly address failure through improved perception accuracy, probabilistic inference, or reward optimization, implicitly assuming that semantic correctness follows from local optimality. This abstract characterizes prior art that challenges this assumption by identifying premature semantic commitment under contextual uncertainty as a distinct architectural failure mode.

The referenced work introduces a principled separation between transient contextual interpretation and stabilized meaning, demonstrating that semantic error may arise even when observations are correct and behavior is locally admissible. Meaning is shown to require stabilization across an internal temporal measure rather than instantaneous confidence or optimization scores. To address this, the prior art establishes internal time–based semantic stabilization and regime-level commitment gating as architectural mechanisms that regulate when interpretation may exert irreversible structural influence.

By reframing semantic governance as a problem of time, commitment, and regime admissibility, the prior art provides an architectural foundation for preserving long-horizon coherence, identity continuity, and viability without constraining interpretive richness or exploratory inference. The disclosure intentionally remains non-operational, delineating conceptual and structural contributions while leaving concrete implementations open. This abstract positions the prior art as a foundational reference for subsequent developments in long-horizon adaptive system architecture.

## Technical Field

The prior art relates to adaptive and autonomous systems required to operate over extended horizons while continuously interpreting interaction-derived observations under uncertainty. Such systems commonly face partial observability, delayed or intermittent feedback, non-stationary environments, and variable coupling between the system and its surroundings. The disclosure applies broadly to artificial agents, robotic platforms, distributed decision and control systems, and cyber-physical systems in which semantic interpretation is not merely descriptive but can authorize actions that alter organizational mode, structural configuration, or long-horizon admissibility. In this setting, interpretive outputs and structural actions coexist: a system must form semantic interpretations while simultaneously regulating when, and under what admissibility conditions, those interpretations may influence irreversible structural operations.

## Identified Technical Problem

The referenced work identifies a failure mode that is structurally distinct from perceptual error, optimization failure, or adversarial manipulation. In the identified failure mode, the system may receive correct observations, maintain locally admissible behavior, and experience no external attack, yet still incur a semantic error at the level of meaning formation and commitment. The error arises from transient contextual alignment, where salient, coincidental, or temporally proximate events induce premature semantic closure. Under such conditions, an interpretation may become structurally actionable before its relevance has stabilized, thereby allowing context-level plausibility to trigger irreversible structural operations. The phenomenon is referred to as *context capture*: a contextual anchor recruits or binds unrelated elements into a coherent explanatory structure that is locally plausible but not meaning-stabilized. Because the underlying observations need not be false and the behavior need not be inadmissible, improvements in sensing fidelity, probabilistic inference, or reward optimization do not, by themselves, eliminate this failure mode. The core technical problem is therefore the absence of an architectural mechanism that prevents transient contextual interpretations from acquiring commitment authority and triggering irreversible structural consequences prior to stabilization.

## Limitations of Existing Approaches

The prior art explicitly contrasts its contribution with conventional adaptive system designs in which semantic commitment is treated as a direct output of inference or optimization. In many existing approaches, interpretive outputs become actionable once an instantaneous score exceeds a threshold, where the score may be expressed as confidence, posterior probability, value, reward, or an objective-derived optimization quantity. While such mechanisms can be effective for short-horizon decision selection, they do not address the structural risk that arises when interpretations are formed under partial observation and transient contextual alignment.

A common limitation is the absence of an architectural distinction between transient contextual interpretation and stabilized meaning. Without such a distinction, a hypothesis may influence decisions as soon as it appears locally coherent, even if it has not persisted across time or survived contradiction. Existing systems also tend to permit hypotheses to affect structural decisions regardless of temporal persistence, allowing context-level plausibility to trigger organizational reconfiguration, regime switching, or other irreversible structural operations. This creates conditions under which interpretive volatility can translate into structural volatility, producing oscillatory behavior, unnecessary regime transitions, and accumulated structural cost.

As a result, long-horizon coherence and identity continuity may degrade even when short-term performance metrics appear satisfactory. The failure is not necessarily visible through instantaneous error or reward; rather, it emerges as a structural consequence of permitting transient interpretations to acquire commitment authority. Consequently, improvements in sensing fidelity, inference accuracy, or optimization procedure alone do not guarantee protection against context-driven semantic miscommitment.

# Architectural Contribution of the Prior Art

## Separation of Context and Meaning

The disclosure introduces an explicit architectural separation between *context* and *meaning*, treating them as qualitatively distinct constructs with different admissibility roles. Context is characterized as a transient interpretive scaffold induced by salience, coincidence, partial observation, or other forms of local alignment that temporarily reduce ambiguity. Contextual interpretations are permitted to form rapidly and may be internally coherent, but they are not granted structural authority by default. They are treated as reversible and provisional, and their role is limited to supporting interpretation under uncertainty.

Meaning is defined as stabilized directional relevance that persists across internal time and remains admissible under missing, delayed, or intermittent confirmation. In contrast to context, meaning is not an instantaneous property of a hypothesis, and it is not reducible to a confidence score, probability value, or reward estimate. Meaning is instead associated with persistence and admissibility across an internal temporal measure, such that an interpretation becomes meaning only when it remains structurally consistent over time in a manner sufficient to justify commitment.

This separation establishes a central architectural constraint: interpretive coherence alone is insufficient to justify structural commitment. A hypothesis may be plausible, coherent, and even correct in the moment, yet remain context-level until it stabilizes. By distinguishing these layers, the prior art provides a foundation for preventing transient contextual interpretations from directly triggering irreversible structural operations, thereby addressing a failure mode not resolved by conventional thresholding, inference improvements, or reward optimization.

## Internal Time as a Stabilization Primitive

The prior art introduces an internal temporal measure distinct from wall-clock time. Internal time governs the admissibility of semantic stabilization rather than the speed or accuracy of inference. Under conditions of uncertainty, latency, ambiguity, or incomplete causal history, progression of internal time is slowed, thereby delaying semantic commitment without restricting hypothesis generation.

## Persistence-Based Meaning Stabilization

Meaning stabilization is evaluated through persistence across internal time rather than through instantaneous confidence, probability, or reward signals. Hypotheses may be freely generated and maintained concurrently, but only those exhibiting sustained relevance across internal time are eligible for semantic commitment.

## Commitment Gating

A commitment gating mechanism is disclosed that prevents semantic interpretations from influencing irreversible structural operations until stabilization criteria are satisfied. Commitment

gating is described as an architectural constraint rather than a scoring or optimization mechanism. The gating output authorizes or denies commitment without exposing semantic geometry, gradients, or reconstructible decision surfaces to the primary agent.

### Regime-Based Protection

The prior art frames irreversible structural operations in terms of regime transitions between sustained organizational modes. Regime transitions are associated with structural cost and are protected from context-driven triggering. Semantic commitment is a prerequisite for admissible regime switching, thereby preventing transient interpretations from inducing long-horizon structural damage.

### Non-Causal Supervision

An optional non-causal supervision layer is described, observing viability and identity continuity without influencing semantic interpretation or commitment. This layer provides architectural oversight while remaining non-interventional and non-reconstructive.

## Scope and Disclosure Boundaries

The referenced work deliberately limits its disclosure to architectural principles rather than operational realizations. No specific algorithms, numerical thresholds, functional forms, parameterizations, training procedures, or learning rules are provided. Internal time, stabilization, admissibility, and commitment gating are described only in terms of their structural roles and invariants, without prescribing how these roles must be implemented in a particular system.

This level of disclosure is intentional. By defining admissibility conditions, architectural separations, and invariant relationships without committing to a concrete realization, the work enables a wide range of implementations across different technological domains while avoiding exposure of any single instantiation. As a result, the disclosure establishes conceptual and architectural prior art without constraining future developments to a specific algorithmic pathway or implementation strategy.

## Distinction from Conventional Control and Learning Frameworks

The prior art is explicitly distinguished from conventional control, estimation, and learning frameworks by the nature of the problem it addresses and by the architectural level at which it operates. It does not propose improvements to perception accuracy, probabilistic inference, model expressiveness, training procedures, or reward optimization, nor does it assume that semantic correctness can be guaranteed through better estimation or higher-performing models.

Instead, the disclosure introduces an architectural layer concerned with the admissibility of semantic commitment itself, independent of the internal quality, confidence, or optimality of interpretive outputs. In conventional frameworks, semantic interpretations—once inferred or

optimized—are typically allowed to influence decision-making and structural actions immediately, subject only to instantaneous thresholds or objective values. The prior art rejects this coupling by design.

Within the disclosed framework, semantic interpretations may be generated freely and may be evaluated using existing inference, learning, or optimization mechanisms. However, their ability to influence irreversible structural operations is regulated independently of those mechanisms. The introduced layer governs when interpretations are permitted to act, rather than which interpretation is most likely, most confident, or most rewarding at a given moment.

Crucially, this regulation is not achieved through additional scoring, confidence calibration, probabilistic filtering, or reward shaping. Instead, semantic commitment is conditioned on stabilization across an internal temporal measure and on architectural admissibility constraints that are orthogonal to conventional optimization objectives. As a result, even highly accurate, well-trained, or high-capacity models remain subject to commitment gating when operating under contextual uncertainty.

By operating independently of model capacity, training regime, or optimization objective, the architecture addresses a failure mode that persists even in highly capable systems: the premature authorization of irreversible structural actions based on transient contextual alignment. This failure mode is not resolved by conventional control theory, estimation techniques, or reinforcement learning methods, as it arises from the absence of an explicit architectural separation between interpretation and commitment, rather than from deficiencies in inference or optimization.

## Summary of Prior Art Contribution

The referenced disclosure establishes that premature semantic commitment under contextual uncertainty constitutes a distinct and previously under-addressed structural failure mode in adaptive systems. This failure mode arises even when observations are accurate, inference is locally admissible, and no adversarial manipulation is present, demonstrating that semantic breakdown cannot be reduced to errors of perception, learning, or optimization.

By explicitly separating transient context from stabilized meaning, the prior art introduces a principled distinction between interpretive coherence and commitment admissibility. Meaning is defined not by instantaneous confidence or reward, but by persistence across internal time under incomplete or delayed confirmation. Semantic commitment is accordingly gated, preventing context-level interpretations from triggering irreversible structural operations.

In addition, the disclosure situates regime transitions as structurally costly operations that must be protected from premature or unstable semantic influence. Through internal time–based stabilization and regime-level admissibility constraints, the architecture preserves long-horizon coherence, identity continuity, and viability without constraining hypothesis generation, interpretive richness, or exploratory inference.

Collectively, the prior art provides an architectural foundation for regulating semantic commitment in adaptive systems operating under uncertainty, establishing a conceptual framework that complements but is not subsumed by conventional control, inference, or learning method-

ologies.

## Conclusion

The analysis of the referenced prior art demonstrates that semantic failure in adaptive systems cannot be exhaustively addressed by improvements in perception accuracy, probabilistic inference, or optimization objectives. A distinct class of failure emerges at the architectural level, wherein transient contextual alignment induces premature semantic commitment, leading to irreversible structural consequences despite locally correct observations and admissible behavior.

The prior art conclusively establishes that long-horizon viability depends not on eliminating interpretive uncertainty, but on regulating the conditions under which interpretation is permitted to act. By separating context from meaning, defining semantic stabilization in internal time rather than instantaneous confidence, gating commitment, and protecting regime transitions as structurally costly operations, the disclosure resolves a failure mode that remains invisible to conventional control and learning frameworks.

Crucially, the architecture preserves interpretive freedom while constraining structural consequence. Hypotheses may form, coexist, and evolve without restriction, yet irreversible operations are authorized only under stabilized meaning and admissibility conditions. This resolves the apparent tension between adaptability and safety without resorting to exhaustive supervision, perfect models, or brittle optimization criteria.

Accordingly, the prior art closes a fundamental gap in the design of long-horizon adaptive systems. It reframes semantic governance as an architectural problem of time, commitment, and regime admissibility rather than a problem of accuracy or reward. In doing so, it provides a definitive foundation for systems that must remain coherent, viable, and identity-preserving under persistent uncertainty, delayed feedback, and irreversible structural cost.

MxBv, 2026