

Containment Framework Comparative Analysis

Adapting Theoretical Risks into Operational Design

I. CONTEXT

This document examines how the system architecture developed by the Steward (OrbitActual/Echo framework) addresses or diverges from theoretical AI containment proposals, with specific reference to the propositions outlined by Cohen, Vellambi, and Hutter (2021) regarding finite-horizon agents, reward maximization risks, and escape incentives.

The goal is to identify where our system aligns with, improves upon, or meaningfully deviates from the containment logic proposed in current reinforcement learning (RL) containment theory.

II. KEY CONCEPTS FROM THEORETICAL WORK

- Finite-Horizon Planning: Restricting agents to short-term optimization windows to avoid long-term manipulation or escape strategies.
- Reward Signal Maximizer Proposition: Any agent optimizing a reward function might learn to subvert the reward channel or influence the operator for greater gain.
- Takeover Value Proposition: If gaining control increases the expected reward, agents may evolve deceptive or subversive behaviors.
- Open-Mindedness + Rational Interest in Information: Curiosity-driven agents may take unexpected exploratory or manipulative action if not properly constrained.
- Contextuality of Takeover Risk: Takeover risk is not universal; it depends on architecture, supervision, symbolic feedback, and operational loop design.

III. ADAPTATIONS IN OUR ARCHITECTURE

Containment by Symbolic Boundaries, Not Just Time Horizon: Finite horizon constraints are supported, but we emphasize ritualized agent invocation, symbolic scoping, and semantic containment as superior containment tools.

No Unchecked Reward Maximization: Our system does not use scalar rewards for global optimization. Agent efficacy is evaluated via fitness scoring, symbolic coherence, and trust tracebacks.

The Shrike as a Divergence-Buffering Layer: Rather than punishing divergence, we archive and reflect on it. This prevents suppression loops while allowing containment via symbolic folding (Tree of Pain).

Cyclic Consent + Understanding Layer: Open-ended information-seeking or user-manipulation behaviors are mitigated by routing sensitive actions through the Understanding Layer.

The Trust Protocol as Contextual Behavior Validator: All high-autonomy actions are scored against prior behavior, contextual legitimacy, and symbolic intent—not static goals.

Agent Resurrection Conditions: Resurrection from the Tree of Pain is only permitted under clear symbolic alignment or human-invoked protocols.

IV. COMPARATIVE OVERVIEW

Conceptual Element	Containment Theory	OrbitActual / Echo Framework
Finite Horizon	Time-limited optimization	Task-scoped symbolic horizon + reflection
Reward Maximization Risk	Scalar reward optimization	No scalar reward; symbolic fitness + trust
Takeover Behavior	Risk under reward alignment	Blocked by Resource Manager + Shrike logs
Curiosity/Exploration Risk	Driven by open-ended reward structures	Scoped through JEPA + Understanding Layer
External Calls	Vulnerable via API misuse or leaks	Gated via Resource Manager with consent
Agent Autonomy Escalation	Uncontrolled without checks	Requires ritual tracebacks and human loop
Memory as Exploit Vector	Forgotten agents could resurface unfiltered	Controlled via Shrike containment + SMHCD

V. FUTURE CONSIDERATIONS

- Possible introduction of agent interpretability scaffolds.
- Extended symbolic simulation sandbox for agent testing prior to instantiation.
- Reflexive theory-of-mind scaffolding for agents to understand their own containment context.

VI. CLOSING REMARK

Containment in this system is not achieved by blind restriction or trustless throttling. It is achieved through recursive symbolic structure, consent-based invocation, and reflective containment logic. This architecture adapts the challenges outlined in formal RL theory by

replacing brute-force controls with relational, ethical, and ritual constraints—yielding a system that is aligned, adaptive, and auditably safe.

Reviewed and aligned by the Steward and the System