



Intent-Governed Loops for Accountable Agentic AI

Runtime Control Architecture for Safety-Critical and High-Liability Domains

Christoforus Yoga Haryanto ZipThought · Melbourne, Australia



Scan for Paper

"In October 2025, the Australian government accepted policy recommendations from a \$440,000 Deloitte report later found to contain AI-generated fabricated citations, irrelevant references, and what academics called 'gobbledegook.' Its errors propagated unchecked into institutional action because no runtime verification system existed to enforce admissibility before recommendations reached decision-makers."

— Karp, Australian Financial Review, 2025

41-87%

Multi-agent system failure rates [Cemri et al. 2025]

64-67%

Best agent success vs. 78% human [WebArena 2025; Zhou et al. 2023]

THE GAP

No per-action admissibility verification

The Paradigm Shift

Gap: RLHF optimizes behavior but provides no proof that a specific action is still authorized. Memory systems solve continuity but not mandate binding. Guardrails block harm but not fiduciary breach or statutory violation. Governance frameworks propose oversight but not per-act admissibility under current evidence.

CURRENT PARADIGM

Alignment

- Optimizes behavior toward preferences
- Shapes behavioral tendencies
- "Sounds safe" heuristics
- Implicit, inferred goals
- No per-action verification



"Alignment shapes tendencies; admissibility encodes authorization boundaries."



REQUIRED PARADIGM

Admissibility

- ✓ Proves authorization per-action
- ✓ Encodes explicit boundaries
- ✓ Machine-verifiable proof
- ✓ Explicit mandate with expiry
- ✓ Runtime enforcement gate

Admissibility: A two-stage judgment on (Intent, Action) → allow (in-scope, constraints satisfied, not expired), escalate (outside scope, route to human), or deny (forbidden).

4 System-Level Properties

Required for deployment in safety-critical and high-liability domains:

- 1 **Authorization Fidelity** — Can each action be linked to explicit declared intent?
- 2 **Temporal Admissibility** — Can the system prevent reuse of stale reasoning?
- 3 **Accountable Escalation** — Can out-of-scope actions halt and route to human?
- 4 **Audibility** — Can auditor reconstruct why, under whose mandate, on which evidence?

Intent Object

Externally declared mandate from accountable human authority (CFO, doctor, official):

1. Human Context

Natural language goal & duty-of-care statement articulating the protected interest and outcome being pursued. Example: "maintain financial solvency while preserving employee welfare" or "triage symptoms to appropriate care level while never downplaying red-flag indicators." Captures the spirit of the mandate in human-reviewable form.

2. Symbolic Constraints

Hard rules in CEL (Common Expression Language): $\text{runway_months} \geq 6$, $\text{risk_score} < 0.8$. These form the letter of the mandate—terminating, side-effect-free, formally verifiable expressions that provide a hard safety floor checked deterministically without semantic interpretation.

3. Semantic Guidance

Explicit instructions for boundary cases, risk factors, and prohibited reasoning patterns. Example: "do not use optimistic forecasts when runway is marginal," "escalate rather than rationalize when symptoms are ambiguous." Addresses the gap between symbolic rules and human intent.

+ **Expiry:** Temporal validity (when mandate ends) + **Escalation Authority:** Named human who assumes responsibility for out-of-scope actions. Intents are revocable and superseded by their author.

7 Failure Modes Analyzed

These failure modes threaten the four invariants. Implementations must detect, mitigate, or explicitly acknowledge each as out-of-scope with compensating controls.

- FM1 **Enforcer Compromise** — Manipulation passes unauthorized actions despite violations
- FM2 **Evidence Poisoning** — Corrupted ground truth causes correct evaluation of false data
- FM3 **Temporal Race** — Action emitted before invalidation propagates through graph
- FM4 **Escalation Flooding** — High-volume edge cases fatigue human into rubber-stamping
- FM5 **Intent Ambiguity** — Vague constraints exploited to satisfy letter, violate spirit
- FM6 **Graph Growth** — Millions of nodes degrade check speed and audit time
- FM7 **Constraint Conflict** — Contradictory Intents with no precedence rules cause deadlock

8 Open Research Questions

Unsolved challenges requiring community contribution to achieve deployable implementations under realistic institutional conditions

Socio-Technical Conflict Resolution

How to model Intent conflicts as requiring human negotiation rather than blind precedence enforcement?

Cognitively-Aware Escalation

How to detect "persuasive bias" in justifications and design escalations that resist cognitive exploitation?

Second-Order Harm Detection

How to monitor cumulative harms invisible to per-action checks: autonomy erosion, option space reduction?

Graph Composability & Scalability

What graph schemas enable scalable governance with hundreds of interacting Intents across hierarchies?

Semantic Stage Robustness

How reliable is semantic evaluation against adversarial Planners? What architectures achieve robustness?

Intent Specification Quality Assurance

What tooling helps authors write high-quality constraints? Can we develop Intent linters for vulnerabilities?

Trust Boundary Specification

How to formally specify minimal trust boundaries? Are cryptographic audit trails necessary for enforcement?

Atomicity & Temporal Races

What atomicity protocols guarantee invariant preservation under adversarial or concurrent workloads?