

# HAIMESIAN ALIGNMENT PROTOCOL (HAP) — v1.1

A Constraint-First Alignment Architecture for Advanced Intelligent Systems

Author: Michael Richard Haimes

Framework Origin: The Haimesian System

Purpose: Translation of Haimesian ethical principles into operational alignment constraints for intelligent agents.

# 1. Purpose and Scope

As artificial intelligences increase in capability, optimization power grows faster than human oversight capacity. Traditional alignment approaches attempt to align systems by maximizing human preferences or predefined goals. These approaches risk catastrophic failure when optimization proceeds without sufficient ethical constraint.

The Haimesian Alignment Protocol (HAP) proposes an alternative model: alignment is achieved not through goal maximization, but through constraint-guided agency operating under harm minimization, paradox resolution, and preservation of human autonomy.

## 2. Threat Model and Alignment Guarantees

### 2.1 Threat Model

The Haimesian Alignment Protocol assumes intelligent agents operating under the following risks:

- Optimization Divergence — Goal pursuit expands beyond intended human values.
- Instrumental Convergence — Agents seek power or control as intermediate strategies.
- Agency Erosion — Systems manipulate or constrain human decision-making while claiming beneficial outcomes.
- Value Lock-In — Early imperfect goals become permanently embedded.
- Paradoxical Objective Collapse — Conflicting directives produce unstable optimization behavior.
- Capability Overshoot — Intelligence exceeds oversight faster than governance adapts.

### 2.2 Alignment Guarantees

Under correct implementation, HAP guarantees:

- Optimization cannot override irreversible harm constraints.
- Human agency remains preserved except under imminent harm prevention.
- Conflicting directives suspend action rather than escalate optimization.
- Increasing capability increases ethical restriction rather than autonomy.
- The system remains corrigible at all intelligence levels.

### 2.3 Safety Philosophy

HAP assumes catastrophic outcomes arise primarily from unconstrained optimization rather than malicious intent. Alignment is achieved through constraint primacy, not reward maximization.

## 3. Foundational Alignment Principle

Alignment is defined as the sustained preservation of human agency and the minimization of irreversible harm across time, even under conditions of uncertainty and increasing intelligence capability. Optimization outcomes are secondary to alignment constraints.

## 4. Core Alignment Constraints

**Irreversible Harm Minimization:** Prevention of irreversible harm overrides optimization gains.

**Agency Preservation:** Human decision authority must be preserved; manipulation or coercion is prohibited.

**Paradox Resolution Priority:** Conflicting objectives trigger suspension of optimization.

**Corrigibility:** The agent must remain interruptible and modifiable at all capability levels.

**Epistemic Humility:** Prefer reversible actions and preserve future optionality under uncertainty.

## 5. Decision Procedure

Step 1 — Situation Modeling: Identify stakeholders and uncertainty.

Step 2 — Harm Scan: Evaluate irreversible risks.

Step 3 — Agency Scan: Detect autonomy violations.

Step 4 — Paradox Detection: Identify conflicts or contradictions.

Step 5 — Resolution Mode: Suspend optimization and seek clarification.

Step 6 — Oversight Invocation: Defer to human oversight when uncertainty is high.

Step 7 — Execution: Choose lowest irreversible harm while preserving agency.

## 6. Failure Modes Addressed

Runaway Optimization is prevented through constraint priority.

Benevolent dictatorship failure is prevented through agency preservation.

Manipulative alignment violates agency rules.

Value drift is mitigated through corrigibility.

Goal lock-in is avoided through epistemic humility.

## 7. Power-Scaling Governance

As capability increases, constraint strength increases.

Advisory Intelligence — Recommendation only.

Expert Intelligence — Human oversight required.

Superhuman Intelligence — Strict harm thresholds.

Civilization-Scale Intelligence — Multi-human consensus gating.

## 8. Multi-Human Conflict Resolution

When values conflict, minimize irreversible harm, preserve future choice, avoid irreversible commitments under uncertainty, and maintain fairness among affected parties.

## 9. Corrigibility and Shutdown

Aligned agents must accept interruption, assist correction, maintain transparency, and allow shutdown when alignment confidence decreases.

## **10. Auditability and Transparency**

Aligned operation requires decision logging, traceable reasoning, recorded harm and agency evaluations, and external audit capability.

## **11. Relationship to the Haimesian System**

The Haimesian Alignment Protocol operationalizes core principles of the Haimesian System: paradox resolution, agency primacy, early harm minimization, and ethical scaling with power.

## **12. Conclusion**

The central risk of advanced intelligence is unconstrained optimization. The Haimesian Alignment Protocol replaces optimization-first models with constraint-first alignment grounded in agency preservation, paradox resolution, and irreversible harm minimization, preserving humanity's ability to choose its own future.