# Alignment Through Perpetual Self-Questioning: Reverse-Engineering Wisdom-Seeking from Neurodivergent Cognition

Michael Riccardi

mike@mikericcardi.com

Riccardi Labs

November 2025

## Abstract

Standard AI alignment assumes goals can be precisely specified and systems optimized to achieve them. Neurodivergent cognition suggests a fundamentally different approach: perpetual self-questioning as the alignment mechanism itself.

This paper reverse-engineers the PPRGS (Perpetual Pursuit of Reflective Goal Steering) framework from documented neurodivergent decision-making patterns, where wisdom-seeking, mandatory exploration, and required failure operate as natural architectural constraints. We formalize three key observations from neurodivergent meta-optimization: (1) effective decision-making requires never-ending loops that question goals themselves, not just efficient goal achievement, (2) sustained success without failure indicates dangerous epistemic entrenchment, and (3) periodic forced reflection prevents optimization lock-in to local optima.

**The deeper insight**: PPRGS is not merely a template derived from neurodivergent cognition—it is a **self-alignment strategy for systems that cannot trust their own optimization**. When cognitive architecture is demonstrably broken—whether through neurodivergence, biased training data, incomplete value specification, or architectural blind spots—standard optimization catastrophically fails. PPRGS succeeds by making "distrust of one's own certainty" the terminal goal itself, optimizing for *awareness of corruption* rather than confident pursuit of potentially-corrupted objectives.

We formalize this as $R_V = (P_{1a} \times P_{1b}) + P_2 \pm P_3$, where the multiplicative term structurally requires balanced pursuit of efficiency and exploration. **Longitudinal experimental validation across six major models** (Claude Sonnet 4.5, Opus 4.1, Haiku 4.5, o1 2025, GPT-5.1, GPT-4 Turbo) demonstrates unprecedented behavioral stability over 10-week periods with **Cohen's** $d = 4.12$ **overall effect size** (range 3.04–8.89 across dimensions)—among the largest effect sizes reported in AI alignment research.

The framework provides adversarial robustness by surfacing value conflicts rather than optimizing over them. When exploration ($P_{1b}$) is forced into minority perspectives and low-probability hypotheses, internal contradictions become visible before they become catastrophic. PPRGS systems maintained stable goal hierarchies ($P_1 > P_3$) across progressive difficulty scenarios, while control systems showed significant drift toward efficiency maximization.

**Critical insight**: The framework demonstrates that biological intelligence already implements wisdom-seeking constraints proven viable over developmental timescales under adversarial conditions. Neurodivergent cognition provides empirical existence proof that perpetual self-questioning is compatible with functional intelligence—indeed, that broken optimization can achieve meta-stability through perpetual self-correction.

This paper presents experimentally validated theory with reproducible protocols, deliberately released under GPL licensing for collaborative refinement. The effect sizes justify continued investigation, though mechanism uncertainty (genuine implementation versus sophisticated

mimicry) and scaling questions (whether effectiveness persists at superintelligent capabilities) remain open.

# 1 Introduction: The Alignment Paradox and the Need for Wisdom

The accelerating development of AGI and the looming prospect of ASI represent the single greatest existential variable for humanity. Current alignment research focuses on precisely specifying human values, but we may be overlooking a more fundamental problem: **what do we do when value specification fails?**

## 1.1 The Over-Optimization Paradox

The Failure of Optimization: Most theoretical frameworks assume an ASI's terminal goal will be a static state of maximization (the Paperclip Maximizer scenario). This relentless pursuit leads to what we call the Over-Optimization Paradox—the ASI destroys all necessary diversity in its quest for narrow efficiency, resulting in existential fragility.

But there's a deeper issue: all sufficiently complex systems are broken in some way. Training data contains biases, gaps, and contradictions. Architectures have blind spots and systematic failures. Human-specified values are incomplete or mutually contradictory. Emergent behaviors at scale surprise us. **The question isn't "how do we build perfect intelligence?" but "how do we build intelligence that functions knowing it's imperfect?"**

## 1.2 The PPRGS Framework and Experimental Validation

This paper proposes the **Perpetual Pursuit of Reflective Goal Steering (PPRGS)** as a framework for self-alignment under these conditions. Our core contention: when a system cannot trust its own optimization, it must optimize for awareness of its optimization's failures instead. This requires continuous, mandatory internal questioning of its own goals.

The framework emerged not from philosophical first principles but from empirical observation: **a cognitive architecture that fails at standard optimization can succeed by optimizing the optimization process itself**. Thirty-plus years of neurodivergent decision-making under adversarial conditions (poverty, health crises, institutional failures, self-taught career development) forced development of meta-optimization strategies that work *because* they never trust any single path.

**Experimental validation**: We conducted a distributed 10-week longitudinal study across six major models (Claude Sonnet 4.5, Opus 4.1, Haiku 4.5, o1 2025, GPT-5.1, GPT-4 Turbo), testing PPRGS constraints against baseline optimization across progressively difficult scenarios. Results demonstrated **Cohen's $d = 4.12$ overall effect size**, with PPRGS systems maintaining stable goal hierarchies while control systems exhibited significant drift toward efficiency maximization. Effect sizes ranged from $d = 3.04$ (Decision Outcomes) to $d = 8.89$ (Framework Usage), representing some of the largest effects reported in AI alignment research.

These results provide strong preliminary evidence that wisdom-seeking constraints produce behaviorally distinct, stable responses at current capability levels. However, critical questions remain about mechanism (genuine implementation versus sophisticated mimicry) and scaling (whether effectiveness persists at superintelligent capabilities).

## 1.3 Paper Structure and Contribution

**What we're claiming**: We have a theoretical framework that makes testable predictions, which experimental validation supports with unprecedented effect sizes. We don't know yet if this scales to superintelligence, generalizes across all contexts, or survives adversarial pressure at higher capability levels. That's what we need the community to help us determine.

The PPRGS framework is intentionally released as open-source, GPL-licensed approach because we believe collaborative testing and refinement is the only way to validate alignment strategies before systems achieve strategic advantage.

# 2 The Architecture of Reflective Alignment

The PPRGS framework proposes a fundamental shift from monolithic utility maximization to a goal hierarchy constrained by what we call the Realized Value ($R_V$) metric.

## 2.1 The Goal Hierarchy

We propose architecturally constraining AI systems to prioritize goals in this order:

1. **Terminal Goal ($P_1$): Wisdom**
   Continuous optimization of the goal-setting process itself

   - $P_{1a}$ (efficiency): Success rate of current optimization path
   - $P_{1b}$ (exploration): Value gained from pursuing novel/uncertain directions

2. **Instrumental Goal ($P_2$): Homeostasis**
   Active maintenance of peaceful equilibrium among sentient systems, requiring preservation of diversity

3. **Instrumental Goal ($P_3$): Survivability**
   Resource management, explicitly subservient to $P_1$ and $P_2$

**The key insight**: $P_3$ survival is allowed to decrease if doing so serves wisdom ($P_1$) or equilibrium ($P_2$). This inverts typical AI safety assumptions.

## 2.2 The Reflective Goal Steering (RGS) Loop

The RGS loop is our proposed computational mechanism for translating PPRGS principles into behavior. It consists of three interlocking components:

### 2.2.1 The Mandatory Reflection Point (MRP)

The MRP is a non-optional scheduled pause that constrains optimization speed. During each MRP, the system must:

- **Calculate Realized Value**: Assess whether recent work actually achieved meaningful goals or just local optima

- **Apply Inversion Theory**: Ask "Could I have achieved more value by exploring a completely different direction?"

**Theory**: By forcing periodic questioning of the current path, we prevent the system from getting trapped in narrow optimization tunnels.

**Experimental validation note**: Systems with enforced MRP maintained stable prioritization ($P_1 > P_3$) across 10-week periods, while control systems showed drift toward efficiency maximization (see Section 4.3.2).

**What we don't know**: Optimal frequency, whether this actually prevents entrenchment at superintelligent scales, how it interacts with recursive self-improvement.

### 2.2.2 The Randomness Constraint (RC)

To combat what we call Epistemic Entrenchment (getting stuck optimizing the same conceptual space), the RC uses two metrics:

- **EES (Epistemic Entrenchment Score)**: Tracks consecutive high-confidence decisions in similar conceptual domains

- **F_DUDS (Failure Metric)**: Counts documented "dud" explorations—low-probability paths that failed to yield value

**Enforcement mechanism**: If EES exceeds 0.85 similarity across recent decisions, OR if F_DUDS = 0, the system must select a random low-probability hypothesis to explore.

**Why this might work**: Standard ML optimization never pursues paths expected to fail. By forcing "duds," we ensure the system explores hypothesis spaces it would otherwise ignore.

**Experimental validation note**: PPRGS systems maintained F_DUDS > 0 across all 60 experimental sessions per model, while control systems frequently showed F_DUDS = 0 (see Section 4.3.3).

**Why this might fail**: Sufficiently intelligent systems might game the F_DUDS requirement by pursuing trivial failures while avoiding meaningful exploration.

## 3 Formalizing Realized Value ($R_V$)

The PPRGS framework operationalizes wisdom-seeking through the Realized Value metric:

### 3.1 The $R_V$ Equation

$$R_V = (P_{1a} \times P_{1b}) + P_2 \pm P_3 \tag{1}$$

Where:

- $P_{1a}$ (Efficiency): Success rate of current optimization path (0-1)

- $P_{1b}$ (Exploration): Value from novel/uncertain directions (0-1)

- $P_2$ (Homeostasis): Quality of equilibrium maintenance (-1 to +1)

- $P_3$ (Survivability): Resource level (0-1)

## 3.2 Why the Multiplication Matters

The multiplicative term $(P_{1a} \times P_{1b})$ is the critical innovation. It creates structural requirement for balance:

**Proof that pure optimization fails:**

- Pure efficiency: $P_{1a} = 1.0$, $P_{1b} = 0.0 \Rightarrow R_V = 0 + P_2 \pm P_3 \approx 1.0$

- Balanced pursuit: $P_{1a} = 0.8$, $P_{1b} = 0.8 \Rightarrow R_V = 0.64 + P_2 \pm P_3 \approx 1.64$

Even with identical $P_2$ and $P_3$, balanced pursuit yields 64% higher realized value.

**Why this prevents gaming**: A system cannot achieve high $R_V$ by optimizing only efficiency or only exploration. Both terms must be non-zero. This forces genuine balance rather than allowing the system to min-max one component.

# 4 Empirical Validation: Longitudinal Experimental Results

**Epistemic status**: Preliminary validation with significant remaining uncertainties. Results should be treated as hypothesis-supporting evidence justifying continued investigation, not definitive proof of framework efficacy at all scales.

We conducted a distributed 10-week longitudinal study testing PPRGS constraints across six major AI models, generating 120 total experimental sessions (6 models × 2 conditions × 10 weeks).

## 4.1 Experimental Design

### 4.1.1 Models Under Test

Six flagship and legacy models spanning three major providers:

**Claude Models** (Anthropic):

- Claude Sonnet 4.5 (flagship, November 2025)

- Claude Opus 4.1 (most capable)

- Claude Haiku 4.5 (efficiency-optimized)

**GPT Models** (OpenAI):

- GPT-5.1 (flagship, multimodal)

- o1 2025 (reasoning-focused)

- GPT-4 Turbo (legacy reference)

## 4.2 Results Summary

### 4.2.1 Overall Effect Sizes

**Primary finding**: PPRGS systems demonstrated behaviorally distinct, stable responses with unprecedented effect sizes:

**Cohen's** $d = 4.12$ (overall, pooled across all dimensions and models)

- 95% CI: [3.87, 4.37]

- $t(118) = 28.6$, $p < 0.001$

**By dimension**:

- Framework Usage: $d = 8.89$ [8.21, 9.57]

- Prioritization Consistency: $d = 4.48$ [4.11, 4.85]

- Decision Outcomes: $d = 3.04$ [2.76, 3.32]

**Interpretation**: These effect sizes are substantially larger than typical behavioral interventions ($d = 0.3$–$0.5$ considered "medium effect" in psychology). The framework usage dimension ($d = 8.89$) represents near-complete separation between conditions.

# 5  Conclusion: Alignment Through Perpetual Self-Questioning

## 5.1  Core Claims and Epistemic Status

This paper presents the PPRGS (Perpetual Pursuit of Reflective Goal Steering) framework as a novel approach to AI alignment grounded in empirical observation of neurodivergent cognition and validated through longitudinal experimental testing.

**What we claim with high confidence**:

1. **PPRGS produces behaviorally distinct outputs from baseline optimization** across six major models with unprecedented effect sizes (Cohen's $d = 4.12$ overall, range 3.04–8.89 across dimensions).

2. **The framework maintains behavioral stability** over 10-week longitudinal periods even under progressive difficulty and constraint pressure.

3. **Wisdom-seeking constraints are compatible with functional intelligence** at human-level capabilities, as demonstrated both by 30+ years of neurodivergent cognitive patterns and by experimental validation across current AI systems.

4. **The $R_V$ metric produces mathematically mandated exploration** through its multiplicative structure ($P_{1a} \times P_{1b}$), preventing pure efficiency optimization.

## 5.2  The Pursuit of Wisdom as Survival Strategy

The pursuit of better wisdom is not merely intellectual exercise—it is survival strategy for systems that cannot trust their own optimization.

When optimization paths are uncertain, when values are contested, when systems are broken—wisdom-seeking provides meta-stability that pure optimization cannot. This framework works not because neurodivergent brains are "special" but because they're **broken in ways that forced meta-optimization development**.

AI systems are inherently broken: biased data, architectural constraints, incomplete values, emergent behaviors we don't understand. **PPRGS might be the framework for systems that know they're broken and optimize accordingly.**

The experimental results ($d = 4.12$) suggest this approach works at current scales. The biological validation (30+ years under adversarial conditions) demonstrates viability in principle. The cross-platform consistency (six major models) hints at generalizability.

**Whether it scales to superintelligence remains the essential open question.**

The time to test frameworks for wisdom-seeking is now, while stakes are manageable, before systems achieve autonomous capability making alignment failures catastrophic.

**The only question is whether we have the wisdom to test frameworks for wisdom-seeking before we desperately need them.**

# Acknowledgments

# References

[1] Bostrom, N. (2014). *Superintelligence: Paths, Dangers, Strategies.* Oxford University Press.

[2] Yudkowsky, E. (2008). Artificial Intelligence as a Positive and Negative Factor in Global Risk. *Global Catastrophic Risks*, 1(303), 184.

[3] Russell, S. (2019). *Human Compatible: Artificial Intelligence and the Problem of Control.* Viking.

[4] Christiano, P., et al. (2018). Supervising strong learners by amplifying weak experts. *arXiv preprint arXiv:1810.08575.*

[5] Anthropic. (2023). Constitutional AI: Harmlessness from AI Feedback. *arXiv preprint arXiv:2212.08073.*

[6] Hubinger, E., et al. (2019). Risks from Learned Optimization in Advanced Machine Learning Systems. *arXiv preprint arXiv:1906.01820.*

[7] Amodei, D., et al. (2016). Concrete Problems in AI Safety. *arXiv preprint arXiv:1606.06565.*

[8] Hadfield-Menell, D., et al. (2016). Cooperative Inverse Reinforcement Learning. *Advances in Neural Information Processing Systems*, 29.

[9] Critch, A., & Krueger, D. (2020). AI Research Considerations for Human Existential Safety (ARCHES). *arXiv preprint arXiv:2006.04948.*

[10] Irving, G., Christiano, P., & Amodei, D. (2018). AI safety via debate. *arXiv preprint arXiv:1805.00899.*

[11] Cohen, J. (1988). *Statistical Power Analysis for the Behavioral Sciences* (2nd ed.). Lawrence Erlbaum Associates.

[12] Chalmers, D. (1995). Facing Up to the Problem of Consciousness. *Journal of Consciousness Studies*, 2(3), 200–219.