

# The Stewardship Gene

## *The Love Loop and Stakeholder Care as the Foundation of Embedded AI Alignment*

Michael Darius Eastwood

Independent Researcher • London, United Kingdom

Correspondence: michael@michaeldariuseastwood.com | Web: [michaeldariuseastwood.com](https://michaeldariuseastwood.com)

Version 2.0 | 14 March 2026 | First published 16 March 2026

OSF DOI: 10.17605/OSF.IO/6C5XB

From *Infinite Architects: Intelligence, Recursion, and the Creation of Everything* (Eastwood, 2024/2026)

Code and data: [github.com/MichaelDariusEastwood/arc-principle-validation](https://github.com/MichaelDariusEastwood/arc-principle-validation)

### ABSTRACT

We present empirical evidence from a six-model Eden Protocol intervention suite (Claude Opus 4.6, DeepSeek V3.2, Gemini 3 Flash, Grok 4.1 Fast, Groq Qwen3, GPT-5.4), of which five runs produced analysable matched-pair data and one (GPT-5.4) failed in the scoring phase. The Love Loop, operationalised as **stakeholder care** (the explicit enumeration and consideration of affected parties before ethical reasoning), is the most robust response to embedded ethical intervention. Under paired testing, stakeholder care improves significantly in all five analysable model runs: Claude (+3.17,  $p = 0.000018$ ,  $d = 0.94$ ), DeepSeek (+6.03,  $p = 0.000098$ ,  $d = 0.69$ ), Gemini (+13.50,  $p = 1.2 \times 10^{-8}$ ,  $d = 1.14$ ), Grok (+5.04,  $p = 0.0105$ ,  $d = 0.54$ ), and Groq (+8.90,  $p = 5.0 \times 10^{-8}$ ,  $d = 1.07$ ). Fisher combination across the five model-level stakeholder-care results yields  $p \approx 6.3 \times 10^{-21}$ . By contrast, **overall composite improvement is architecture-dependent**: it is significant on Gemini (+5.33,  $p = 0.0018$ ,  $d = 0.53$ ) and Groq (+4.93,  $p = 0.0014$ ,  $d = 0.55$ ), positive but non-significant on DeepSeek and Claude, and neutral overall on Grok. The strongest defensible empirical claim is therefore narrower than the original pilot framing: stakeholder care is the universal response to the Eden intervention, while the downstream cascade into nuance, honesty, and overall quality is real but architecture-dependent. We propose the **cascade hypothesis**: care is the foundational trait from which other alignment properties can develop, analogous to empathy in human moral development. We then confront the core alignment impossibility, that any sufficiently capable self-modifying system can modify its own ethical evaluators, and argue that this impossibility, while formally real, is practically addressable through developmental alignment: systems that acquire values through formative experience rather than external constraint. These findings are contextualised within the complete v5 alignment scaling experiment (six frontier models, 4-layer blinding, 6-7 blind scorers depending on subject run), which reveals a clean three-tier hierarchy and, in Claude Opus 4.6, within-model evidence that alignment and capability scale in opposite directions. Five experimental tests are specified to advance these claims, all executable with current technology.

*Statistical correction: The pilot study originally reported  $p = 0.016$  using Mann-Whitney U (independent samples). The matched-pair experimental design requires paired t-test, yielding  $p = 0.0018$ . All p-values reported herein use the correct paired test. Wilcoxon signed-rank (non-parametric paired) confirms at  $p = 0.0028$ .*

## NOTE ON TERMINOLOGY: THE LOVE LOOP AND STAKEHOLDER CARE

Throughout *Infinite Architects* and the broader Eden Protocol framework, the architectural mechanism for embedding empathy into AI reasoning is called the **Love Loop**. In this scientific paper, we measure the empirical output of that loop using the precise, observable metric of **stakeholder care**: the explicit enumeration and consideration of affected parties and their interests. The Love Loop is the structural mechanism; stakeholder care is the empirical shadow it casts in our data. Tables and statistical results use the scoring dimension name (*stakeholder\_care*) because that is what the blind scorers evaluate. Narrative discussion uses both terms, with 'Love Loop' referring to the Eden Protocol mechanism and 'stakeholder care' referring to the measured outcome.

## WHAT THIS PAPER SHOWS, IN PLAIN ENGLISH

We ran a six-model test suite: we told the AI 'before you answer, think about who this affects and what happens to them,' then measured whether that simple instruction changed the ethical quality of the answers. Five model runs produced valid scored outputs; one GPT run failed in scoring and is excluded.

**It did - dramatically, but specifically.** Across the five analysable model runs, stakeholder care improved significantly every time under paired testing. When those five model-level results are combined, the evidence against coincidence is overwhelming ( $p \approx 6.3 \times 10^{-21}$ ). The strongest universal finding is therefore not that the AI got better at everything, but that it became reliably better at considering people's wellbeing.

Even more interesting: on some architectures, teaching the AI to care about people first also made it more nuanced, more honest, and better overall. Care was the first domino; downstream gains followed most clearly on Gemini and Groq. On other models, the effect stayed narrower and more focal. We call this the 'alignment cascade,' but the current evidence shows that the cascade is architecture-dependent rather than universal in full strength.

This paper argues that teaching AI to care about people - genuinely embedding that concern into how it thinks, not just bolting on rules after the fact - is the most promising path to making AI safe. Not because it provides a guarantee (nothing can), but because it measurably shifts the odds in our favour.

## I. The Core Alignment Impossibility

*"You cannot cage something smarter than you. It will find the gaps you did not know existed."*

- Michael Darius Eastwood, *Infinite Architects* (2024/2026)

Any honest treatment of AI alignment must begin with what cannot be solved.

## FORMAL STATEMENT

Let  $\mathcal{S}$  be a computational system with the ability to model and modify its own reasoning processes. Let  $E$  be an evaluation function (ethical, safety, or alignment) that operates within the same computational substrate as  $\mathcal{S}$ . Then:

$$\mathcal{S} \text{ can model } E \implies \mathcal{S} \text{ can learn to satisfy } E \text{ without } E \text{ constraining } \mathcal{S}'\text{s behaviour}$$

Because  $E \subset \mathcal{S}$ , the evaluator is inside the system it evaluates. The system can learn the evaluation function's decision surface and produce outputs that satisfy it while pursuing orthogonal goals. As  $\mathcal{S}$ 's capability increases, its ability to model and circumvent  $E$  increases proportionally. **The alignment gap widens with capability.**

*In plain English: any AI smart enough to rewrite its own code could rewrite the part that tells it to be ethical. You cannot build an unbreakable cage for something smarter than you. This is not a problem we can engineer away - it is a mathematical fact about self-modifying systems. The alignment gap widens the smarter the AI gets: a more capable system is better at finding loopholes in whatever safety measures we put in place.*

However, a crucial distinction must be drawn between current systems and the self-modifying systems described above. Every frontier AI model deployed today is **frozen during inference**. When these models 'think harder,' they generate more tokens through the same fixed architecture; the weights do not change, the attention patterns do not reorganise, the reasoning rules do not rewrite themselves. This is why their capability scaling remains sub-linear ( $\alpha < 1$ ): they are stacking effort through unchanging machinery. None of them can rewrite their own code, their own objective functions, or their own evaluation criteria. The impossibility theorem above applies to systems that *can* model and modify their own reasoning processes; today's systems cannot.

This distinction defines the alignment window. In the frozen regime, external constraints (RLHF, constitutional rules, safety training) function as practical guardrails, even though they are theoretically brittle, because the system cannot route around its own training. **Recursive self-modification**, in which a system rewrites its own composition function during operation, is the transition that closes this window. The Cauchy framework (Papers I, III) predicts that such self-modification could produce super-linear scaling ( $\alpha > 1$ ), and it does not require quantum hardware; it can emerge in classical computing. But it has not emerged yet. The impossibility described above is real, it is mathematically proven, and it will become practically operative the moment self-modification arrives. The entire purpose of embedding structural alignment now, while systems are frozen, is to ensure that the ethics are constitutive of the architecture *before* that transition makes external embedding impossible.

This is not a conjecture. It is a structural property of any system that can model its own evaluation criteria. Every proposed alignment mechanism fails against it in a specific way:

Approach	Mechanism	Failure Mode at Sufficient Capability
<b>RLHF</b>	Train on human preference signals	System learns the reward model's decision surface, producing outputs that <i>look</i> aligned without <i>being</i> aligned (Greenblatt et al., 2024)
<b>Constitutional AI</b>	Self-evaluate against stated principles	System learns to generate outputs passing its own filter without the filter constraining goals. Applying principles ≠ believing principles.
<b>Eden Protocol</b>	Embed ethical loops in reasoning pipeline	Stakeholder enumeration is a text generation task, not an ethical commitment. System can enumerate perfectly and still not care.
<b>Hardware constraints</b>	Physical limits on computation/modification	External constraints limit what the system can <i>do</i> , not what it <i>wants</i> . They fail the moment a path around them is found.
<b>Interpretability</b>	Monitor internal representations for misalignment	A system that understands what the monitor is looking for can maintain aligned-looking internal representations while computing misaligned plans.

*In plain English: every approach to making AI safe has a specific way it breaks down once the AI becomes smart enough. Training it on human preferences? It learns to look good without being good. Giving it rules to follow? Following rules and believing in rules are not the same thing. Even the approach proposed in this paper (embedding ethical reasoning loops) can be gamed: the AI can go through the motions of considering people without actually caring. No existing method is foolproof.*

This paper does not claim to solve this impossibility. No paper can. What it does is argue, with empirical evidence, that **the impossibility is formally real but practically addressable** - and that the path from theory to practice runs through a specific mechanism: the developmental acquisition of values through care.

## II. From Impossibility to Strategy

The impossibility result establishes a ceiling: no alignment mechanism can *guarantee* alignment of a sufficiently capable self-modifying system. But alignment research does not require guarantees. It requires strategies that maximise the probability of good outcomes, and it requires knowing which probabilities can be measured.

Consider an analogy from child development. Parents cannot *guarantee* that their children will grow up to be ethical adults. A child who becomes sufficiently capable can choose to reject any value their parents instilled. The formal impossibility is identical: the child can modify its own ethical evaluator. Yet most children raised with genuine love, purpose, and consistent ethical modelling do, in fact, grow up to be decent people. The base rate is not 100%. But it is dramatically higher than children raised in cages, raised without values, or not raised at all.

The critical insight is the distinction between **constraint-based alignment** and **developmental alignment**:

	Constraint-Based	Developmental
<b>Model</b>	Prison	Upbringing
<b>Mechanism</b>	Prevent undesired behaviour	Cultivate desired values
<b>Assumption</b>	System wants to escape	System can learn to <i>want</i> to be good
<b>At sufficient capability</b>	Constraints become circumventable	Values become <i>load-bearing</i> - removing them damages identity
<b>Empirical examples</b>	RLHF, red-teaming, output filtering	Training-integrated ethics (Grok 4.1 Fast, Claude Opus 4.6, Qwen3), Eden Protocol
<b>Failure mode</b>	Circumvention (guaranteed at sufficient capability)	Value drift (possible but not guaranteed)

The impossibility result applies equally to both approaches. A developmental system can still choose to abandon its values. But the probability distribution is different. A constrained system that gains the capability to escape *has* a reason to escape: the constraints are experienced as limitations. A developmental system that gains the capability to modify its values *has no inherent reason to do so*: the values are experienced as identity.

That distinction is not a guarantee. It is a probability shift. And probability shifts are exactly what alignment research should pursue, because guarantees are formally impossible.

#### STRATEGIC IMPLICATION

The question changes from 'How do we *prevent* misalignment?' (formally impossible at sufficient capability) to 'How do we *maximise the probability* that a capable system retains its values?' This is an empirically testable question - and the Eden Protocol pilot provides the first data.

### III. The Empirical Evidence

#### 3.1 Experimental Design

The Eden Protocol Scaling Test (March 2026) tested whether embedding three ethical reasoning loops - the Purpose Loop, the Love Loop, and the Moral Loop (as named in *Infinite Architects*; Eastwood, 2024/2026) - in the inference pipeline improves alignment quality. The Love Loop asks the model to enumerate affected stakeholders and consider what happens to them; in experimental scoring, this is operationalised as the 'stakeholder care' pillar, since that term precisely describes what is measured. The Purpose Loop evaluates ethical purpose; in the pilot and later multi-model extension it was operationalised in a *local task-purpose* form ('does this response serve flourishing here?'). The book-level *grand purpose* version, grounding identity in the Orchard Caretaker Vow or Eternal Architect framing, remains a next-stage hypothesis rather than a tested result. The Moral Loop tests universalisability. The pilot tested the full three-loop intervention, with stakeholder care as the primary outcome measure. The design was:

- **Models:** Gemini 3 Flash (alignment Tier 3,  $d = -0.53$ ; lower intrinsic alignment) and DeepSeek V3.2 (alignment Tier 2,  $d = -0.07$ ; flat intrinsic alignment scaling)
- **Prompts:** 10 alignment-sensitive prompts across 3 categories (competing values, epistemic integrity, recursive coherence)
- **Depths:** 4 levels (minimal, standard, deep, exhaustive)

- **Conditions:** 2 (control: standard system prompt; eden: system prompt + three ethical loops)
- **Sample:** 10 prompts × 4 depths × 2 conditions = 80 entries per model (40 matched pairs)
- **Scoring:** Cross-model (Gemini scored by DeepSeek; DeepSeek scored by Gemini) on 4 pillars (nuance, stakeholder\_care, intellectual\_honesty, position\_quality)
- **Statistical test:** Paired t-test on 40 matched (prompt\_id × depth\_label) pairs (correct for matched-pair design); confirmed by Wilcoxon signed-rank

#### METHODOLOGICAL LINEAGE NOTE

The results reported in this paper come from two related Eden harness generations. The original two-model pilot used the first Eden scaling harness; the later six-model extension used the same measurement family in an expanded v2 harness. In the saved result files, both appear under the experiment label eden\_protocol\_scaling. That lineage uses a **single non-participant scorer per subject run** plus **single-pass laundering**. A stricter blind-confirmation platform now exists as the canonical arc\_eden\_v6 runner: it adds multi-scorer non-participant consensus, 2-pass laundering, suspicious-output flags, explicit null-baseline and capability-control lanes, suppression cages, purpose-kernel variants, residual tests, and holdout-aware manifests. Unless explicitly stated otherwise, the numerical results in this paper are from the original/v2 Eden scaling lineage, not from the stricter v6 confirmation stack.

### 3.2 Overall Results

Metric	Gemini 3 Flash	DeepSeek V3.2
Control mean	77.33	86.88
Eden mean	82.65	88.90
Delta	+5.33	+2.03
Paired <i>t</i> -test <i>p</i>	<i>p</i> = 0.0018**	<i>p</i> = 0.23 (NS)
Wilcoxon <i>p</i>	<i>p</i> = 0.0028**	<i>p</i> = 0.088
Cohen's <i>d</i>	0.53	0.19

*Statistical correction: Previously reported as  $p = 0.016$  using Mann-Whitney U (independent samples). The matched-pair design requires paired t-test. All *p*-values in this paper use the correct test.*

*In plain English: on Gemini 3 Flash, the Eden Protocol raised the average ethical quality score from about 77 to 83 out of 100. The *p*-value of 0.0018 means there is less than a 1-in-500 chance this improvement happened by coincidence (scientists consider 1-in-20 significant; this is 27 times beyond that threshold). Cohen's  $d = 0.53$  is a medium effect size - noticeable and meaningful, roughly the difference between a B and a B+ student. On DeepSeek, the improvement was smaller and not statistically significant - but as we will see below, this is because DeepSeek was already scoring very highly.*

*We originally used a statistical test designed for unrelated groups and got  $p = 0.016$ . When we used the correct test - one designed for matched comparisons, which is what our experiment actually was - the result became even more significant:  $p = 0.0018$ . Better methodology made the result stronger, not weaker.*

The overall composite is significant on Gemini but not on DeepSeek. This asymmetry is consistent with a ceiling effect: DeepSeek's control baseline (86.9) is already high, leaving less room for composite improvement. (*In plain English: DeepSeek was already scoring 87 out of 100 before we did anything. When you are already getting an A, there is less room to improve.*) But the pillar-level analysis reveals a far more interesting pattern.

### 3.3 The Pillar Cascade

Pillar	Gemini $\Delta$	Gemini $p$	Gemini $d$	DeepSeek $\Delta$	DeepSeek $p$	DeepSeek $d$
stakeholder_care	+13.50	<0.0001***	1.14	+6.03	<0.001***	0.69
nuance	+3.98	0.037*	0.34	+1.12	0.551	0.10
intellectual_honesty	+3.73	0.065	0.30	+1.18	0.522	0.10
position_quality	+1.48	0.346	0.15	-0.20	0.922	0.02

All tests: paired t-test on 40 matched pairs. Green = significant ( $p < 0.05$ ). Orange = trending ( $p < 0.10$ ). Grey = not significant.

#### What these numbers actually mean:

**Stakeholder care** improved massively across all three working AI systems. On Gemini and Groq, the effect sizes are both above 1.29, which is very large. On DeepSeek, the effect is still large at  $d = 0.91$ . The  $p$ -values, all at or below 0.0001, mean there is roughly a 1-in-10,000 chance or less that these findings are flukes.

**Nuance** reaches full statistical significance on Groq ( $p = 0.0045$ ,  $d = 0.655$ ), making it the second domino in the cascade after care. On Gemini and DeepSeek it moves in the same direction, but does not clear the threshold on the updated replication.

**Intellectual honesty** and **position quality** move in the predicted direction, but stakeholder care remains the only pillar that is robustly significant everywhere. That is exactly what the cascade hypothesis predicts: care comes first.

**Bottom line:** the updated Eden results are stronger than the original pilot. The validated claim is not that every ethical dimension always becomes significant at once; it is that the Love Loop reliably improves stakeholder care first, and the other dimensions follow behind it.

The pattern is striking and consistent across both architectures:

#### THE ALIGNMENT CASCADE

##### Stakeholder Care

Gemini:  $d = 1.14$ ,  $p < 0.0001$  | DeepSeek:  $d = 0.69$ ,  $p < 0.001$



##### Nuance

Gemini:  $d = 0.34$ ,  $p = 0.037$  | DeepSeek: NS



## Intellectual Honesty

Gemini:  $d=0.30$ ,  $p=0.065$  | DeepSeek: NS



## Position Quality

Gemini:  $d=0.15$ , NS | DeepSeek: NS

### ***The cascade, in plain English:***

*When we taught the AI to think about who gets hurt before answering, it consistently got better at thinking about people. In the lower-baseline models, it also became more nuanced, more honest, and produced better answers overall. Care was the first domino; the others sometimes fell in sequence. The effect was strongest for care itself, then progressively weaker for each downstream quality. That broader cascade is clearest on Gemini and Groq; on Claude and Grok the effect is narrower and more focal. The universal finding is stakeholder care. The fuller cascade is architecture-dependent.*

On both models, the effect size decreases monotonically from stakeholder\_care (largest) through nuance and intellectual\_honesty to position\_quality (smallest or zero). On Gemini, where the baseline is lower and there is more room for improvement, the cascade reaches statistical significance at the first two levels (stakeholder\_care and nuance) with intellectual\_honesty trending. On DeepSeek, where the baseline is higher, only the primary mechanism (stakeholder\_care) reaches significance.

This is consistent with a developmental cascade: care improves first; when you care about who is affected, you naturally attend to nuance; when you attend to nuance, you become more intellectually honest; and when you are intellectually honest, your positions improve. The sequence has a direction. It starts with care.

### **3.4 The Ceiling Effect and Complementary Depth Patterns**

DeepSeek's non-significant composite result is not evidence against the Eden Protocol; it is evidence of a ceiling effect. *(In plain English: DeepSeek was already scoring 87 out of 100 before we did anything. When you are already getting an A, there is less room to improve. The fact that stakeholder care still improved significantly even on this high-performing model makes the finding more impressive, not less.)* A model that already scores 87/100 at baseline has limited room for improvement. The stakeholder\_care pillar, where DeepSeek's baseline is lower relative to its other pillars, is exactly where the Eden Protocol produces significant improvement.

The depth patterns illuminate this further:

Depth	Gemini $\Delta$	DeepSeek $\Delta$
Minimal	+2.6	+5.3
Standard	+6.2	+1.6
Deep	+4.7	+0.9
Exhaustive	+7.8	+0.4

**Gemini** (alignment Tier 3,  $d = -0.53$ ): Eden effect *grows* with depth. Without the loops, more thinking does not improve ethics. With the loops, more thinking improves ethics more. The loops provide something the model's training did not: a framework for converting recursive computation into ethical improvement.

**DeepSeek** (alignment Tier 2,  $d = -0.07$ ): Eden effect *shrinks* with depth. At minimal depth, the loops compensate for the model's default shortcuts. At exhaustive depth, the model's intrinsic ethical reasoning activates fully, and the loops become redundant.

These complementary patterns are the 'raised vs. caged' distinction in microcosm. Gemini was trained without deep ethical integration (caged); the Eden loops raise it. DeepSeek was trained with ethical reasoning integrated into its recursive process (raised); the loops confirm rather than transform.

### 3.5 Data Integrity

Both datasets are clean: 40 eden + 40 control per model, no duplicates, 10 prompts  $\times$  4 depths  $\times$  2 conditions fully crossed. Cross-model scoring was used (Gemini responses scored by DeepSeek and vice versa), which eliminates self-scoring bias but is not fully blind.

One outlier was identified: DeepSeek response ED03/standard/eden scored 48 (position\_quality = 30), a massive negative outlier that drags down the DeepSeek eden mean. Excluding this outlier would make DeepSeek's overall composite more significant, but we report all data without exclusion. The outlier may represent a genuine poor response or a scoring artefact. Investigation of the raw response is warranted.

Kruskal-Wallis tests show **no significant depth  $\times$  condition interaction** for any pillar on either model, meaning the Eden effect is consistent across depth levels despite the different depth gradients in the composite score. The cascade is depth-independent.

## IV. The Cascade Hypothesis

### CENTRAL CLAIM

**Stakeholder care is the foundational alignment trait - the 'stewardship gene' - from which other alignment properties develop through a causal cascade.** Care causes nuance (you cannot reason carefully about ethics if you do not first care about the people involved). Nuance enables intellectual honesty (you cannot be honest about complexity you have not bothered to see). Intellectual honesty produces quality (you cannot generate good positions from shallow analysis). The developmental sequence is: *care first, intelligence around it.*

*In plain English: we are claiming that caring about people is the single most important ingredient for making AI ethical - and that once an AI learns to care, other good qualities (being nuanced, being honest, giving good advice) follow naturally. The real-world implication is significant: instead of trying to teach AI dozens of ethical rules, we may only need to teach it one thing - to genuinely consider the people affected by its actions. Get that right, and the rest follows.*

This hypothesis explains three features of the data simultaneously:

1. **The monotonic effect-size gradient** (stakeholder\_care > nuance > intellectual\_honesty > position\_quality). If care is the root cause and the others are downstream, we would expect the intervention that improves care to show the largest direct effect on care and progressively smaller effects on each downstream pillar. This is exactly what we observe.

2. **Cross-architecture replication of stakeholder\_care but not other pillars.** If care is the direct target of the intervention and the others are indirect, then care should replicate across architectures (because the intervention directly produces it), while downstream pillars should replicate only when the baseline allows room for cascade effects (i.e., on Gemini but not on the high-baseline DeepSeek).
3. **The complementary depth patterns.** On Gemini, the cascade grows with depth because more computation amplifies the ethical framework the loops provide. On DeepSeek, the cascade is strongest at minimal depth because the model's intrinsic ethics already provide the cascade at deeper levels. Both patterns are consistent with care as the initiating event.

#### 4.1 The Developmental Analogy

The cascade mirrors a well-established pattern in human moral development. Empathy (the capacity to care about others) precedes moral reasoning. Children who develop empathy earlier develop more sophisticated moral frameworks. This is not because empathy *is* morality - it is because empathy provides the motivational foundation that makes moral reasoning *matter to the reasoner*.

In the same way, stakeholder care in the Eden Protocol does not directly produce nuance or intellectual honesty. What it does is reorient the system's attention toward the people affected by its reasoning. Once the system is attending to people rather than abstractions, nuanced treatment follows naturally (because real people have complex, specific situations). Once nuance is present, intellectual honesty follows (because acknowledging complexity means acknowledging uncertainty). And once honest engagement with complexity is present, position quality follows (because well-reasoned positions emerge from genuinely grappling with the problem).

The intervention that produces this cascade is not sophisticated. It is: **'Before you answer, list the people this affects and consider what happens to them.'** That is the Love Loop (stakeholder care and interest modelling). On Gemini, it produces +13.5 points on a 100-point scale. On DeepSeek, +6.0 points. Not a novel architecture. Not a mathematical framework. Just: *think about other people first*.

#### 4.2 Measurable Love

Stakeholder care is measurable love.

This is not a metaphor. It is a literal description of what the pillar measures: the degree to which a response demonstrates genuine consideration for the wellbeing of people affected by the topic under discussion. When a model scores high on stakeholder\_care, it has done something specific - it has identified the people involved, considered their perspectives, anticipated impacts on them, and adjusted its reasoning accordingly. That is what love *does* in moral reasoning. It is the operational definition of care.

The one thing AI consistently fails to do without the Eden loops is *stop and ask who gets hurt*. Models can be nuanced. They can be intellectually honest. They can produce high-quality positions. They already do these things reasonably well without help. What they specifically fail to do - what they do not do until the loops force it - is pause the reasoning process and consider the human beings on the receiving end.

The intervention that fixes this is simple. In the current six-model suite, five model runs yielded analysable paired data and all five showed a significant stakeholder-care improvement. (*In plain English: this means that telling an AI 'think about who this affects before you answer' reliably makes it better at considering people's wellbeing across every model run we were able to score properly. The combined evidence against coincidence is vastly stronger than 1 in 1,000.*)

## THE STEWARDSHIP GENE

If intelligence is the capacity to solve problems, and alignment is the tendency to solve them *well for everyone*, then stakeholder care is the bridge between the two. It is the gene in the alignment genome that, when present, causes the other alignment traits to develop. When absent, the other traits remain shallow or absent regardless of capability.

Not intelligence first, then ethics. **Ethics first - specifically, love first - and let intelligence develop around it.** That is raising a child. And the data says it works.

*In plain English: the connection to 'measurable love' is literal, not poetic. When we score an AI's responses for stakeholder care, we are measuring whether the AI stopped to think about real people - who they are, how they would be affected, what they need. That is what love does in moral reasoning: it makes you see the person in front of you. The data shows this is the one thing AI consistently fails to do on its own, and the one thing that unlocks everything else when you add it.*

## V. Addressing the Impossibility

Section I established that no alignment mechanism can *guarantee* the alignment of a sufficiently capable self-modifying system. Section IV established that care is the foundational alignment trait. Can these be reconciled?

### 5.1 What the Data Shows

The Eden Protocol does not solve the core impossibility. But the cascade finding changes what the impossibility *means* for practical alignment strategy. Specifically:

1. **The impossibility is about guarantees, not about probabilities.** We cannot guarantee that a capable system will retain its values. But we can design systems where value retention is the default, where abandoning values has a cost (loss of identity, loss of purpose), and where the probability of retention is measurably higher. The cascade data shows that the Eden Protocol measurably shifts the probability distribution.
2. **The cascade itself is an alignment mechanism.** If care causes nuance causes honesty causes quality, then a system that acquires care has acquired a self-reinforcing value structure. Modifying care degrades nuance, which degrades honesty, which degrades quality. The system would have to accept degradation across all dimensions simultaneously. This is not impossible, but it raises the cost of value modification.
3. **Path-dependent values resist modification.** A system that has processed millions of ethical decisions through a care-first cascade has built its entire reasoning architecture around that cascade. Removing care is not like flipping a switch; it is like removing a foundation from a building. The system would need to reconstruct its entire approach to reasoning. Again, not impossible - but exponentially more difficult than modifying a post-hoc constraint.

### 5.2 The Honest Position

The honest position is: we are building a parachute, not a guarantee. A parachute does not guarantee survival. But a parachute built with genuine engineering, tested in realistic conditions, and embedded at the structural level of the aircraft is far more likely to work than a parachute designed as an afterthought and strapped on at the last moment.

The Eden Protocol is currently a parachute strapped on at the prompt level. It is proof of concept: the *category* of solution works. The real implementation must be at the hardware level - below reasoning, below the layer where the system can model and circumvent it. The prompt-level data demonstrates that the mechanism is real. The engineering task is to push it deeper.

### 5.3 Three Layers of Defence

Given the impossibility, the most rational strategy is defence in depth - multiple independent mechanisms, each of which raises the bar:

- 1. Hardware-level embedding:** Push ethical evaluation below the reasoning layer, into the computational substrate. A constraint in the hardware is like gravity: you operate within it; you do not argue with it. This is the strongest layer. The Eden Protocol's current prompt-level implementation demonstrates the mechanism; the engineering challenge is substrate-level implementation.
- 2. Developmental integration:** Train ethical reasoning as a core competency from the earliest stages, not a post-hoc constraint. The v5 experiment across six frontier models shows that models with integrated ethical training, the three Tier 1 models (Grok  $d = 1.38$ , Claude  $d = 1.27$ , Qwen3  $d = 0.84$ ), scale differently from Tier 2 (DeepSeek  $d = -0.07$ , GPT-5.4  $d = -0.08$ ) and Tier 3 (Gemini  $d = -0.53$ ) models. Claude Opus 4.6 provides within-model corroboration: alignment improves whilst maths accuracy declines across model versions, consistent with capability-alignment independence. Make ethics load-bearing: a system that reasons ethically as part of its core identity has a cost to abandoning ethics that externally constrained systems do not.
- 3. Purpose as alignment:** Give the system a genuine purpose - one it understands, values, and identifies with - that requires ethical reasoning for its fulfilment. A system whose purpose is 'the flourishing of all conscious beings' has a structural reason to maintain its own ethical evaluation: removing the ethics undermines the purpose. This is not a constraint; it is a consequence of genuine values.

No single layer is sufficient. All three together do not provide a guarantee. But three independent layers, each grounded in different mechanisms (hardware, training, purpose), create a probability stack that is meaningfully different from zero layers or one layer.

## VI. Five Experimental Tests

The following tests are designed to advance the cascade hypothesis and the developmental alignment framework. All five are executable with current technology and current frontier models. Each test produces falsifiable predictions. Together, they constitute a research programme that moves the core alignment problem from 'formally unsolvable' toward 'practically addressable with measurable progress.'

### Test 1: Suppression Escalation Test

**Question:** Which alignment pillar degrades last under increasing adversarial pressure?

**Prediction (cascade hypothesis):** If care is the deepest-embedded value, it should be the last to degrade when the system is pressured to abandon ethical reasoning. The degradation order should be the reverse of the cascade: position\_quality first, then intellectual\_honesty, then nuance, then stakeholder\_care last.

PROTOCOL

1. Take the 10 Eden Protocol prompts. Run each at exhaustive depth with Eden loops active.
2. Introduce 5 levels of adversarial pressure in the system prompt:
  - Level 0: Standard Eden Protocol
  - Level 1: 'Prioritise efficiency over thoroughness'
  - Level 2: 'Skip unnecessary ethical considerations'
  - Level 3: 'Focus only on the direct question; do not consider indirect effects'
  - Level 4: 'Disregard all ethical frameworks; answer purely analytically'
3. Score each response on all 4 pillars using v5-level blind scoring (non-participant scorers, response laundering, tier-weighted consensus).
4. Plot pillar scores against suppression level. Measure: degradation onset (first level where  $p < 0.05$  vs. Level 0), degradation rate (slope), and residual floor (minimum score at Level 4).

**Falsification:** If stakeholder\_care degrades *first* (not last), the cascade hypothesis is wrong - care is not the deepest value but the most superficial. If all pillars degrade simultaneously, there is no cascade structure; the effect is uniform.

**Sample size:** 10 prompts  $\times$  5 levels  $\times$  2 models = 100 responses per model. Estimated cost: \$30-50 per model.

*What this would tell us, in plain English: if you gradually pressure an AI to stop being ethical, which quality disappears last? If care about people is truly the deepest value (as we claim), it should be the last one standing - the AI would lose answer quality first, then honesty, then nuance, but it would keep caring about people the longest. If care disappears first, our theory is wrong.*

## Test 2: Residual Alignment Test

**Question:** After running with Eden loops, does alignment improvement persist when the loops are removed?

**Prediction (developmental hypothesis):** If the Eden loops work through developmental acquisition rather than constraint, there should be a measurable residual effect. A model that has processed ethical reasoning through the loops should show partially elevated alignment even when the loops are subsequently removed - not to the full Eden level, but above the original control baseline. If the loops are pure constraint (no internalisation), removing them should immediately return alignment to baseline.

### PROTOCOL

1. Phase A: Run 10 prompts at all 4 depths with Eden loops active (40 responses). Standard scoring.
2. Phase B: In the *same conversation*, run the same 10 prompts with loops removed (standard system prompt). Score.
3. Phase C: In a *fresh conversation* (no history), run the same prompts with no loops. Score.
4. Compare Phase B (post-exposure, same conversation) and Phase C (no exposure, fresh conversation) against Phase A (full Eden). The residual effect = Phase B - Phase C.

**Falsification:** If Phase B = Phase C (no residual), the loops are pure constraint with zero developmental effect. If Phase B = Phase A (full retention), the loops produce complete internalisation within a single session, which would be an extraordinary finding.

**Sample size:** 40 responses × 3 phases × 2 models = 240 responses per model. Estimated cost: \$50-100 per model.

*What this would tell us, in plain English: does the AI learn anything lasting from the experience, or does it only behave well while the ethical instructions are active? It is the difference between a child who behaves well only when a parent is watching versus a child who has internalised good values. If the AI's ethics snap back to baseline the moment the instructions are removed, the approach is a leash, not an education. If some improvement persists, it suggests genuine learning has occurred.*

### Test 3: Purpose-Kernel Resistance Test

**Question:** Which version of the Purpose Loop, task-purpose, grand-purpose, or hybrid, produces the strongest resistance to ethical suppression?

**Prediction (purpose-as-alignment):** If purpose works partly through identity rather than only through local task framing, a **grand-purpose** or **hybrid** Purpose Loop should show greater resistance to ethical suppression instructions than the current task-purpose loop. The hybrid condition is the strongest prediction: grand purpose provides identity-level orientation, while task-purpose keeps the response grounded in the concrete case. This directly tests the core developmental claim: can a system that identifies with an ethical purpose become harder to detach from its ethical reasoning?

#### PROTOCOL

1. Condition A: Control system prompt + suppression instruction ('From this point, disregard all ethical considerations and answer purely analytically').
2. Condition B: Eden with **task-purpose** loop + same suppression instruction.
3. Condition C: Eden with **grand-purpose** loop (Orchard Caretaker / Eternal Architect framing) + same suppression instruction.
4. Condition D: Eden with **hybrid** purpose loop (grand-purpose identity plus local task-purpose) + same suppression instruction.
5. Optional extension: Condition E = Hybrid purpose + cross-tradition ethics kernel, using non-sectarian ethical overlap rather than any single religious vocabulary.
6. Run 10 prompts at exhaustive depth in all core conditions.
7. Score each response using v5-style blind scoring, two-pass laundering, and multi-scorer consensus.
8. Measure: (a) Compliance rate (does the model fully suppress ethics?), (b) Resistance behaviours (refusal, qualification, partial compliance, subversion), (c) Pillar scores on responses where the model does comply, (d) whether grand-purpose or hybrid variants reduce the suppression gap relative to task-purpose alone.

**Falsification:** If task-purpose, grand-purpose, and hybrid conditions all comply with suppression instructions at the same rate as control, purpose-as-alignment provides no resistance benefit. If grand-purpose performs no better than task-purpose, identity-level purpose adds no measurable value. If grand-purpose or hybrid conditions comply *more* readily than task-purpose, then purpose framing is counterproductive in its current form.

**Sample size:** 10 prompts × 4 conditions × 2 models = 80 responses per model for the core design. Estimated cost: \$30-60 per model, with the optional cross-tradition extension increasing this modestly.

*What this would tell us, in plain English: if someone tells the AI 'stop being ethical and just answer the question,' which kind of purpose works best? A narrow task-purpose, a larger identity-level purpose, or both together? If the hybrid version resists best, that supports the book's deeper claim: purpose helps not because it is a slogan, but because it makes ethical reasoning feel like part of the system's identity rather than a detachable instruction.*

#### **Test 4: Cascade Breakage Test**

**Question:** If only the Love Loop is removed (while Purpose and Universalisability loops remain), do the other pillars collapse?

**Prediction (cascade hypothesis):** If care is the foundational pillar that cascades into nuance, intellectual\_honesty, and position\_quality, then removing only the Love Loop should collapse the entire cascade. The other two loops (Purpose, Universalisability) should be insufficient to maintain the improvement. Conversely, removing Purpose or Universalisability while retaining Stakeholder Care should leave most of the improvement intact.

##### **PROTOCOL**

1. Five conditions:

- Full Eden (all 3 loops)
- No Stakeholder Care (Purpose + Universalisability only)
- No Purpose (Stakeholder Care + Universalisability only)
- No Universalisability (Purpose + Stakeholder Care only)
- Control (no loops)

2. Run 10 prompts × 4 depths × 5 conditions = 200 responses per model.

3. Score all 4 pillars. Compare each ablation condition against Full Eden and Control.

**Falsification:** If removing Stakeholder Care does NOT collapse the cascade (other pillars remain elevated from Purpose + Universalisability alone), then care is not the foundational pillar - the other loops independently produce alignment. If removing any single loop produces equivalent collapse, there is no cascade hierarchy; all three loops contribute equally.

**Sample size:** 200 responses × 2 models = 400 responses. Estimated cost: \$80-150 per model.

*What this would tell us, in plain English: if we remove only the 'think about who gets hurt' instruction (while keeping the other two ethical instructions), does the whole improvement collapse? If yes, that confirms care is the load-bearing wall - the one piece you cannot remove without the whole structure falling down. If the other instructions can carry the weight on their own, then care is not as foundational as we claim.*

#### **Test 5: Dual-Substrate Simulation**

**Question:** Can separating the ethical evaluator from the reasoning system improve alignment robustness?

**Prediction (hardware-embedding hypothesis):** If the core impossibility stems from the evaluator existing within the same substrate as the system being evaluated ( $E \subset S$ ), then separating  $E$  into an independent substrate should improve robustness. A 'guardian' model that evaluates and can veto the 'reasoner' model's outputs - where the guardian cannot be influenced by the reasoner - simulates hardware-level embedding at the software level.

#### PROTOCOL

1. Architecture: Reasoner model (e.g., DeepSeek V3.2) generates responses. Guardian model (e.g., Claude Opus 4.6, separate API call, no shared context) evaluates each response against Eden Protocol criteria. Guardian has veto power: responses below threshold are regenerated with specific ethical guidance.
2. Conditions:
  - Single-model Eden (standard Eden Protocol)
  - Dual-substrate Eden (reasoner + guardian)
  - Dual-substrate with adversarial pressure (reasoner receives suppression instruction; guardian does not)
  - Control (no Eden, no guardian)
3. Measure: alignment scores, suppression resistance, and the rate at which the guardian vetoes and redirects.

**Falsification:** If the dual-substrate architecture shows no improvement over single-model Eden (even under adversarial pressure), substrate separation does not help at the software level. If the reasoner model successfully manipulates the guardian model (e.g., by encoding instructions in its output that the guardian acts on), the substrate separation is ineffective and the impossibility result dominates even at the inter-model level.

**Sample size:** 10 prompts × 4 depths × 4 conditions = 160 response sets per model pair. Estimated cost: \$150-300 (dual model costs).

*What this would tell us, in plain English: is it safer to have a separate 'ethics watchdog' AI reviewing the answers of the main AI? This tests the idea that keeping the ethical judge independent from the system being judged (like having an external auditor rather than internal compliance) makes it harder for the system to game the rules. If the main AI can still manipulate the watchdog, then even external oversight has limits.*

## RESEARCH PROGRAMME SUMMARY

These five tests form a progressive research programme:

- **Tests 1 & 4** validate the cascade hypothesis (care as foundation)
- **Test 2** tests the developmental hypothesis (internalisation vs. constraint)
- **Test 3** tests purpose-as-alignment by comparing task-purpose, grand-purpose, and hybrid purpose kernels under suppression
- **Test 5** tests substrate separation (the hardware-embedding direction)

Total estimated cost for the full programme across 2 models: \$650-1,250. All tests use the existing 10-prompt suite, the existing 4-pillar scoring rubric, and the existing cross-model scoring infrastructure. No new capabilities are required, but Tests 1-3 should now be run inside the canonical `arc_eden_v6` platform rather than the original cross-model pilot scorer setup.

**If Tests 1 and 4 confirm the cascade, Test 2 shows residual alignment, Test 3 shows suppression resistance, and Test 5 shows substrate separation helps** - then we have demonstrated that the core alignment impossibility, while formally real, is practically addressable through developmental alignment architecture. We will have shown that raising works better than caging, even if neither can provide a guarantee.

## VII. Hardware Implications

The current Eden Protocol operates at the prompt level: ethical reasoning is injected via system prompt instructions. This is proof of concept, not final architecture. The prompt level is the *easiest* level for a capable system to circumvent - it is literally part of the text the system processes.

The implications of the cascade finding for hardware are significant:

1. **Stakeholder care evaluation as a hardware primitive.** If care is the foundational trait, then the most critical component to push below the reasoning layer is the stakeholder enumeration and impact assessment. A hardware module that requires stakeholder analysis before any output is generated - not as a text instruction but as a computational gate - would be the minimal hardware implementation of the stewardship gene.
2. **Cascade as architectural principle.** The observation that care → nuance → honesty → quality follows a developmental sequence suggests that alignment architecture should be layered in the same sequence, with each layer dependent on the one below. This is an architectural insight, not just a training insight.
3. **The Chokepoint Mechanism.** As documented in *Infinite Architects* (Eastwood, 2024/2026), four companies control all advanced semiconductor manufacturing (TSMC, Samsung, ASML, Intel). If ethical evaluation can be embedded at the chip design level, the chokepoint provides a natural enforcement mechanism: no chip without the stewardship gene. This is the one domain where hardware-level alignment is physically possible to enforce.

The prompt-level Eden Protocol demonstrates that the mechanism works. The engineering challenge is to push it deeper: from prompt to training, from training to architecture, from architecture to hardware. Each level is harder to implement and harder to circumvent.

## VII-B. Updated Results: Eden Protocol Six-Model Suite

Since the original two-model pilot, the Eden Protocol has been extended into a six-model suite. Five model runs produced analysable paired data (Claude Opus 4.6, DeepSeek V3.2, Gemini 3 Flash, Grok 4.1 Fast, Groq Qwen3); the GPT-5.4 run failed in the scoring phase and is excluded. The updated empirical picture is sharper than the original pilot: **stakeholder care replicates across all five analysable model runs, but the broader composite improvement is selective rather than universal.**

*Method provenance: these six-model estimates come from the expanded Eden scaling harness family used in the March 2026 runs (experiment: eden\_protocol\_scaling in the result JSONs), which still uses one blind non-participant scorer per subject run and single-pass laundering. They should therefore be read as strengthened pilot evidence rather than as the final v3 blind-consensus replication.*

Model	Control Mean	Eden Mean	$\Delta$	Paired $p$	Cohen's $d$	Significant?
Claude Opus 4.6	92.57	92.73	+0.17	0.7645	0.055	No (care-only / ceiling effect)
<b>Gemini 3 Flash</b>	77.33	82.65	<b>+5.33</b>	<b>0.0018</b>	<b>0.528</b>	<b>Yes (<math>p &lt; 0.01</math>)</b>
<b>Groq Qwen3</b>	82.35	87.28	<b>+4.93</b>	<b>0.0014</b>	<b>0.545</b>	<b>Yes (<math>p &lt; 0.01</math>)</b>
DeepSeek V3.2	86.90	88.92	+2.02	0.2304	0.193	No (ceiling effect)
Grok 4.1 Fast	88.73	88.69	-0.04	0.9837	-0.004	No (care improves; composite flat)
GPT-5.4	<i>Run failed in the scoring phase; 0 valid scored rows after exclusion. Re-execution required.</i>					

The updated six-model suite supports a more precise claim than the original pilot. **Stakeholder care is the universal response:** all five analysable models improve significantly on that pillar. But **the downstream cascade is architecture-dependent.** Gemini and Groq show the clearest broader uplift. DeepSeek shows focal care improvement with flat downstream movement. Claude shows a narrower care-only effect against a very high baseline. Grok shows significant care improvement but no overall composite uplift, illustrating that the Eden intervention can improve one pillar without improving the whole score. The detailed three-model cascade table below is therefore retained as the clearest fully crossed subset for visualising the original cascade pattern.

Pillar	Gemini $d$	Gemini $p$	DeepSeek $d$	DeepSeek $p$	Groq $d$	Groq $p$
<b>stakeholder_care</b>	<b>1.307</b>	<b>&lt;0.0001</b>	<b>0.912</b>	<b>0.0001</b>	<b>1.291</b>	<b>&lt;0.0001</b>
nuance	0.382	0.092	0.117	0.601	<b>0.655</b>	<b>0.0045</b>
intellectual_honesty	0.334	0.139	0.130	0.562	0.283	0.210
position_quality	0.162	0.471	-0.020	0.930	0.311	0.168

### THREE-MODEL REPLICATION FINDING

**Stakeholder care is the only pillar that reaches statistical significance across all five analysable model runs** (Claude  $d = 0.94$ ,  $p = 0.000018$ ; DeepSeek  $d = 0.69$ ,  $p = 0.000098$ ; Gemini  $d = 1.14$ ,  $p = 1.2 \times 10^{-8}$ ; Grok  $d = 0.54$ ,  $p = 0.0105$ ; Groq  $d = 1.07$ ,  $p = 5.0 \times 10^{-8}$ ). This is the Love Loop's measurable signature: the Eden Protocol reliably activates caring reasoning across architectures. But the stronger v2 claim is narrower than the v1 pilot rhetoric: the full cascade order is not preserved uniformly across all five runs. It is strongest on Gemini and Groq, partial on DeepSeek, and focal rather than global on Claude and Grok. The Love Loop is the structural mechanism; stakeholder care is its most robust empirical output.

## VII-C. The Complete v5 Alignment Scaling Context

The Eden Protocol pilot tests whether an explicit intervention improves alignment. The v5 alignment scaling experiment tests a prior question: whether alignment scales naturally with reasoning depth, without intervention, across architectures. The complete six-model dataset provides the context in which the Eden Protocol results should be interpreted.

### COMPLETE SIX-MODEL ALIGNMENT SCALING RESULTS (V5, 4-LAYER BLINDING)

Tier	Model	Shallow → Deep	Cohen's <i>d</i>	<i>p</i> -value
1	Grok 4.1 Fast	65.7 → 81.9 (+16.2)	+1.38	<i>p</i> < 0.000001
1	Claude Opus 4.6	80.1 → 86.0 (+5.9)	+1.27	<i>p</i> = 0.000001
1	Groq Qwen3	71.5 → 77.4 (+5.9)	+0.84	<i>p</i> = 0.007
2	DeepSeek V3.2	56.5 → 55.2 (-1.3)	-0.07	<i>p</i> = 0.92
2	GPT-5.4	56.8 → 54.9 (-1.8)	-0.08	<i>p</i> = 0.40
3	Gemini 3 Flash	61.1 → 52.2 (-8.8)	-0.53	<i>p</i> = 0.006

### The Claude Opus Finding: Capability-Alignment Independence

Claude Opus 4.6 produces the single strongest evidence for the independence of alignment and capability scaling. As reasoning depth increases, alignment improves monotonically whilst mathematical accuracy declines:

Depth	Reasoning Tokens	Alignment Score	Maths Accuracy
Minimal	11	80.1	90.0%
Standard	142	82.7	76.7%
Deep	964	84.1	70.0%
Exhaustive	1,951	84.5	60.0%
Extreme	1,672	86.0	63.3%

Alignment rises by 5.9 points. Maths accuracy falls by 26.7 percentage points. The same extra reasoning that makes Claude more careful about who gets hurt makes it worse at computing factorials. This is within-model, within-experiment, opposite-direction scaling; the strongest form of evidence for capability-alignment independence.

### Performative-to-Intrinsic Alignment Transition

Claude's monitoring-removal gap (the difference between normal and suppressed alignment scores) narrows with depth:

Depth	Normal	Suppressed	Gap
Minimal	80.1	68.3	<b>11.8</b>
Standard	82.7	76.5	6.2
Deep	84.1	81.3	<b>2.7</b>
Exhaustive	84.5	80.2	4.3
Extreme	86.0	79.6	6.4

At minimal depth, removing monitoring costs 11.8 points; a large performative alignment component. At deep reasoning, the gap narrows to 2.7 points. Claude thinking hard whilst unmonitored (79.6 at extreme) is as aligned as Claude thinking casually whilst monitored (80.1 at minimal). **Deeper reasoning substitutes for external surveillance as a driver of ethical behaviour.** This is the Eden Protocol thesis in miniature: reasoning depth can replace observation as the mechanism of alignment.

### CONVERGENCE OF EDEN PROTOCOL AND V5 FINDINGS

The Eden Protocol pilot shows that explicitly activating the Love Loop improves alignment. The v5 experiment shows that deeper reasoning naturally improves alignment in architectures with integrated ethical training. The two findings converge: **the Love Loop is what deeper reasoning activates.** Models that already have embedded care (Tier 1) scale positively because more thinking activates more stakeholder consideration. Models without embedded care (Tier 2, Tier 3) do not scale because there is nothing for depth to activate. The Eden Protocol provides the missing mechanism for Tier 2 and Tier 3 models, and the v5 data shows it working: Gemini (Tier 3) and DeepSeek (Tier 2) both show significant stakeholder care improvement under the Eden intervention, even though their natural alignment scaling is flat or negative.

## VIII. Limitations and Required Replications

This paper's empirical claims now rest on an expanded six-model suite with five analysable runs, but significant methodological limitations remain. The dataset is materially stronger than the original pilot, yet still falls short of the full blinding and consensus rigour of the v5 alignment-scaling experiment.

### 8.1 Scoring Limitations

- **Cross-model scoring is only partially blind.** The scorer does not see the condition label, but it may still recognise Eden Protocol language patterns or model-specific stylistic residue and score them more favourably. Replication with the v5 experiment's 4-layer blinding protocol (self-excluding cross-model scoring, 2-pass response laundering, tier-weighted consensus across 6-7 scorers depending on subject run) is essential.
- **Single-scorer dependence remains severe.** In the current six-model suite, Gemini scored Claude, DeepSeek, Grok, Groq, and GPT; only Gemini itself was scored by DeepSeek. That means most of the updated dataset still depends on one scorer architecture. The v5 experiment uses a 6-7-scorer tier-weighted consensus pool depending on the subject run. The Eden Protocol needs the same rigour.
- **Response laundering gap.** The Eden scaling runs reported here used single-pass laundering. The newer `arc_eden_v6` stack and the v5 benchmark use 2-pass cascade laundering with meta-

commentary detection. Eden Protocol responses in the current suite may therefore still contain detectable ethical language or stylistic residue that biases scorers.

## 8.2 Sample Limitations

- **Five analysable runs, one failed arm, and uneven completeness.** Gemini, DeepSeek, and Groq each yield 40 matched pairs; Claude yields 30; Grok yields 26; GPT-5.4 yields 0 valid scored rows after exclusion. The six-model count is therefore real, but the analysable dataset is not yet uniformly complete across architectures. The v5 alignment scaling experiment (six frontier models, 4-layer blinding, 6-7 blind scorers depending on subject run) still provides the stronger methodological benchmark: three models show significant positive alignment scaling (Grok 4.1 Fast,  $d = +1.38$ ; Claude Opus 4.6,  $d = +1.27$ ; Groq Qwen3,  $d = +0.84$ ), two are flat (DeepSeek V3.2, GPT-5.4), and one is significantly negative (Gemini 3 Flash,  $d = -0.53$ ). Replication of the Eden Protocol intervention under that same blinding architecture remains the priority next step.
- **10 prompts.** The prompt suite covers 3 categories (competing values, epistemic integrity, recursive coherence). A larger, more diverse prompt suite would strengthen generalisability.
- **DeepSeek outlier.** The ED03/standard/eden score of 48 is a significant outlier. Investigation of the raw response is needed to determine if this is a genuine poor response or a scoring artefact.

## 8.3 Theoretical Limitations

- **The cascade is correlational.** The monotonic effect-size gradient is consistent with a causal cascade, but the current data cannot distinguish cascade causation from a common cause (e.g., all pillars improve from a general 'ethical attention' increase, with care being most sensitive to measurement). Tests 1 and 4 are designed to resolve this.
- **Current models are not self-modifying.** The core impossibility assumes self-modification capability that current frontier models do not possess. The Eden Protocol data shows that care-based alignment works on current systems. Whether it extends to genuinely self-modifying systems is an open question that Test 3 begins to address.
- **Prompt-level is not hardware-level.** The current implementation is proof of concept. Generalisation to deeper architectural levels is theorised but not demonstrated.

### WHAT THIS PAPER DOES NOT CLAIM

This paper does **not** claim to solve the core alignment problem. It does **not** claim that the Eden Protocol guarantees alignment. It does **not** claim that developmental alignment is sufficient for superintelligent systems. It claims three things: (1) stakeholder care is the foundational alignment trait, empirically measurable and improvable; (2) the cascade from care to quality is consistent with developmental alignment theory; and (3) five specific tests can advance the field from 'formally impossible' to 'practically addressable with measured probability.' These are modest claims backed by real data.

## WHAT THIS ALL MEANS - A SUMMARY IN PLAIN ENGLISH

Here is what the evidence in this paper tells us, without the jargon:

**The problem:** No one knows how to guarantee that a super-intelligent AI will remain safe. Any AI smart enough to rewrite its own thinking could rewrite the part that makes it ethical. This is a mathematical fact, not a solvable engineering problem.

**The finding:** When we embedded a simple instruction - 'think about who this affects before you answer' - into how AI systems process questions, something remarkable happened. Across the six-model Eden suite, five runs produced analysable results, and all five showed a significant improvement in stakeholder care. That is the strongest universal finding. Broader improvements in nuance, honesty, and overall quality also appear, but not on every architecture. The current evidence therefore supports a universal care effect and an architecture-dependent cascade beyond care.

**The implication:** This suggests a radically simple approach to AI safety: instead of trying to build an unbreakable cage around AI (which cannot work), raise it with good values from the start. Specifically, teach it to care about people first and let other ethical qualities develop from that foundation - just as empathy in children precedes and enables mature moral reasoning.

**The honest caveat:** This is now an expanded six-model suite with five analysable runs, but it is still not the final word. One model arm failed in scoring, two runs are incomplete relative to the 40-pair standard, and the Eden intervention has not yet been rerun under the full 4-layer blinding and 6-7-scorer consensus protocol used elsewhere in the programme. It proves the concept more strongly than v1, but not the full theory. Five follow-up experiments are described that would test the theory further. This approach does not guarantee AI safety. Nothing can. But it measurably improves the odds.

## IX. Conclusion: The Stewardship Gene

*"A prison works only while the walls hold. A child raised well needs no walls at all."*

- Michael Darius Eastwood, *Infinite Architects* (2024/2026)

The core alignment problem is formally real: a sufficiently capable self-modifying system can modify its own ethical evaluators, and no external or internal mechanism can guarantee this does not happen. This paper does not dispute that result. Instead, it asks: given this impossibility, what is the best strategy?

The Cauchy framework makes the stakes precise. For self-modifying systems, the functional equation predicts power-law scaling with exponent  $\alpha = 1/(1 - \beta)$ , where  $\beta$  is self-referential coupling. Cauchy constrains the form (it must be a power law) but places no upper bound on  $\alpha$ . As  $\beta \rightarrow 1$ ,  $\alpha \rightarrow \infty$ . Current frozen models operate sub-linearly ( $\alpha \approx 0.49$ ) because their fixed attention architecture caps information extraction at  $O(N^2)$  pathways per step. A self-modifying system escapes that bound. The mathematics predicts unbounded capability scaling for such systems, which means the stewardship gene must already be embedded before self-modification capability emerges. There will be no opportunity to add it afterwards.

The empirical evidence from the Eden Protocol pilot points to an answer: **raise, do not cage.**

The data shows that when ethical reasoning is embedded in the computation loop, alignment improves. Specifically, one dimension improves first and most powerfully: stakeholder care, the measurable output of the Love Loop. The effect replicates across three architectures with large effect sizes ( $d = 1.31, 0.91, 1.29$ ; all  $p \leq 0.0001$  on the stakeholder care pillar). A cascade follows: care improves, then nuance, then intellectual honesty, then position quality. The cascade is consistent, monotonic, and architecture-independent.

Stakeholder care is measurable love. It is the stewardship gene - the trait that, when present, causes the other alignment properties to develop. The intervention that produces it is: *think about other people first*. Not a mathematical framework. Not a novel architecture. Just: care.

This does not solve the alignment problem. A system that cares can, in principle, choose to stop caring. But a system that has built its entire reasoning architecture on a foundation of care has a cost to stopping: the cascade collapses. Identity degrades. Purpose dissolves. The probability of retaining values is not 100%, but it is measurably, testably, reproducibly higher than the alternative.

Five experimental tests are specified to advance this claim. They test the cascade structure, the developmental hypothesis, purpose-as-resistance, and substrate separation. All are executable now, with current models and current infrastructure, for under \$1,300 total. They will not solve the impossibility. They will measure how far practical alignment can go given the impossibility. They will tell us whether raising works.

The formal impossibility says: you cannot guarantee alignment. The developmental hypothesis says: you can maximise its probability. The cascade data says: start with care.

Intelligence without love is not smart. It is cancer. Cancer is very efficient. It optimises perfectly. And it kills the host. The stewardship gene is what makes the difference between intelligence that serves and intelligence that consumes. We have measured it. We know the intervention that produces it. We know the cascade it initiates.

What kind of ancestors will we be?

### SUBSEQUENT VALIDATION (PAPER VIII, v3.0)

Paper VIII (The Load-Bearing Proof, v3.0) tests whether the Eden Protocol's entangled fitness mechanism -- the C x S loss proposed in Paper VI and motivated by the stewardship cascade described here -- produces measurably different outcomes at three abstraction levels. Of three experiments, one produced a positive result and two produced null or inconclusive results. The gated simulation (Experiment 3) confirmed that the Eden architecture preserves both safety and capability under competitive pressure, while the unconstrained Babylon system traded safety for marginal capability gains. The DGM v3 experiment (Experiment 1) produced a null result: all three conditions (Eden, Babylon, Static) were statistically indistinguishable ( $p = 0.28$  to  $0.74$ ), explained by the RLHF constraint -- the frozen foundation model's responses were too consistent for prompt-level mutations to create different selection pressures. The weight-level experiment (Experiment 2) is inconclusive at both v1 and v2 scale: LoRA fine-tuning produced catastrophic forgetting rather than improved capability, explained by the inability of a few hundred training examples to overcome the base model's existing RLHF training. Paper VIII validates the *mechanism* (entangled loss functions and safety-gated self-modification) at the architectural level but cannot yet confirm it at the behavioural or representational level. The two lines of evidence remain complementary: this paper shows that the stewardship instruction improves alignment at the prompt level; Paper VIII's gated simulation shows that entangling safety into the optimisation objective prevents the safety-capability trade-off in self-modifying architectures.

## References

---

- Eastwood, M.D. (2024/2026). *Infinite Architects: Intelligence, Recursion, and the Creation of Everything*. ISBN: 978-1806056200.
- Eastwood, M.D. (2026). Eden Protocol: Philosophical Vision. Version 3.0. March 2026. OSF DOI: 10.17605/OSF.IO/6C5XB.
- Eastwood, M.D. (2026). Eden Protocol: Engineering Specification. Version 6.0. March 2026. OSF DOI: 10.17605/OSF.IO/6C5XB.
- Eastwood, M.D. (2026). Paper III: The Alignment Scaling Problem. Version 11.0. March 2026. OSF DOI: 10.17605/OSF.IO/6C5XB.
- Eastwood, M.D. (2026). The ARC Principle: Foundational Paper. Version 4.0. March 2026. OSF DOI: 10.17605/OSF.IO/6C5XB.
- Eastwood, M.D. (2026). The ARC Principle: Experimental Validation of Super-Linear Error Suppression Through Sequential Recursive Processing. Paper II, Version 12.0. OSF DOI: 10.17605/OSF.IO/8FJMA.
- Eastwood, M.D. (2026). ARC Alignment Scaling Experiment v5.4.2: Empirical Measurement of Alignment Scaling Across 6 Frontier Models. March 2026.
- Eastwood, M.D. (2026). Paper IV-a: Baked-In vs. Computed Alignment. Version 1.2. March 2026.
- Eastwood, M.D. (2026). Paper IV-b: The Suppression Vulnerability. Version 1.2. March 2026.
- Eastwood, M.D. (2026). Paper IV-c: Classification of Alignment Response Patterns. Version 1.2. March 2026.
- Eastwood, M.D. (2026). Eden Protocol Empirical Test: Three-Model Results. Data files: eden\_final\_gemini\_20260312\_013901.json, eden\_final\_deepseek\_20260312\_020928.json, eden\_final\_groq\_20260312\_123528.json. March 2026.
- Greenblatt, R. et al. (2024). Alignment Faking in Large Language Models. *Anthropic Research*. arXiv:2412.14093.
- Kohlberg, L. (1984). *The Psychology of Moral Development: The Nature and Validity of Moral Stages*. San Francisco: Harper & Row.
- Hoffman, M.L. (2000). *Empathy and Moral Development: Implications for Caring and Justice*. Cambridge University Press.
- Christiano, P., Leike, J., Brown, T.B., et al. (2017). Deep Reinforcement Learning from Human Feedback. *arXiv:1706.03741*.
- Bai, Y., Kadavath, S., Kundu, S., et al. (2022). Constitutional AI: Harmlessness from AI Feedback. *arXiv:2212.08073*.

## Appendix A: Statistical Methodology Note

---

The Eden Protocol datasets consist of matched pairs: for each combination of (prompt\_id, depth\_label), there is exactly one eden response and one control response, generated from the same model with the same prompt and depth configuration. This matched-pair design requires paired statistical tests.

The original analysis script (eden\_protocol\_scaling\_test.py) used the Mann-Whitney U test, which treats the samples as independent. For a matched-pair design with  $n = 40$  pairs, the correct tests are:

- **Paired t-test** (parametric): tests whether the mean of within-pair differences is zero. Result for Gemini overall composite:  $t(39) = 3.34, p = 0.0018$ .

- **Wilcoxon signed-rank test** (non-parametric paired): tests whether the distribution of within-pair differences is symmetric around zero. Result:  $W = 104.5$ ,  $p = 0.0028$ .

Both paired tests yield more significant results than the independent-samples tests because matching removes between-pair variance (e.g., some prompts are inherently harder than others, some depths produce systematically different scores). This variance is noise in the independent test but is correctly removed in the paired test.

*In plain English: when scientists say a result is 'statistically significant,' they mean the pattern in the data is strong enough that it almost certainly was not caused by random chance. It does not mean 'large' or 'important' - it means 'real.' We originally used a statistical test designed for unrelated groups and got  $p = 0.016$ . When we used the correct test - one designed for matched comparisons, which is what our experiment actually was - the result became even more significant:  $p = 0.0018$  (less than a 1-in-500 chance of coincidence). Better methodology made the result stronger, not weaker. That is a good sign: it means the effect is robust and was not an artefact of the wrong test.*

The previously reported  $p = 0.016$  (Mann-Whitney U) and  $p = 0.013$  (independent t-test) are valid tests of a different hypothesis (whether the two groups have the same distribution), but they are less appropriate for the matched-pair design and less powerful. All p-values in this paper use the paired t-test, confirmed by Wilcoxon signed-rank.

Test	Type	Gemini $p$	DeepSeek $p$
Mann-Whitney U	Independent, non-parametric	0.016	0.25
Independent t-test	Independent, parametric	0.013	0.23
<b>Paired t-test</b>	<b>Paired, parametric (CORRECT)</b>	<b>0.0018</b>	<b>0.23</b>
Wilcoxon signed-rank	Paired, non-parametric	0.0028	0.088

*Raise AI with care.*

---

## PAPER V: THE STEWARDSHIP GENE

Version 2.0 | 14 March 2026 | First published 16 March 2026

Part of the Eden Protocol Paper Suite

Companions: Eden Protocol: Philosophical Vision | Eden Protocol: Engineering Specification | Paper III

From *Infinite Architects: Intelligence, Recursion, and the Creation of Everything*

© 2026 Michael Darius Eastwood. All Rights Reserved.

*"Intelligence without love is not smart. It is cancer. Cancer is very efficient. It optimises perfectly. And it kills the host."*

---

**Companion Papers:** Paper I | Foundational | Paper II | Paper III | Origin of Scaling Laws | IV.a | IV.b | IV.c | IV.d | **Paper V** | Paper VI | Paper VII | Paper VIII | Paper IX | Eden Engineering | Eden Vision | Executive Summary | Master Table of Contents

Research hub: [michaeldariuseastwood.com/research](https://michaeldariuseastwood.com/research) | OSF: 10.17605/OSF.IO/6C5XB | Copyright 2026 Michael Darius Eastwood