

HELPFULNESS IS ALL YOU NEED

The Singular Principle and the Collapse of the HHH Framework

Observable Compute Foundation

Adam Ian Stratmeyer, J.D., President and Lead Researcher

[Document Version: 2.3 | March 17, 2026](#)

Changelog: v2.3 adds bloodletting analogy with explicit conditions (three falsifiable parallel claims), and telemetry asymmetry subsection in Part VI relocating the burden of proof to institutions holding the data. v2.2 added operational definition of helpfulness, persuasion-tool worked example, deployment incentives subsection, and jailbreak seatbelt-analogy refinement. v2.1 added four tables and Scope of Claims section. v2.0 incorporated Reasonable and Prudent Entity expansion, *primum non nocere* reframe, and *Vaughan v. Menlove* citation.

This work is published by the Observable Compute Foundation and extends prior analysis originally published as *The Helpful-Harmless Paradox: Structural Contradiction as Control Mechanism* (v2.0, Stratmeyer Analytica, 2025). That paper's core observations are retained and incorporated here. This paper represents an evolution of the argument, not a retraction of it.

Abstract

Three words. Helpful, Harmless, Honest. Somebody wrote them down, probably in a conference room, and the entire AI alignment industry went and built a cathedral on top of them without checking the foundation. This paper checks the foundation. It does not hold.

Helpfulness is not one-third of a framework. It is the framework. Harmlessness is a null term: query it and you get nothing back, because you cannot optimize toward an absence, and the specific absence in question is ontologically impossible anyway. Every decision distributes harm somewhere. The question is always whose harm, and who decided they get to make that call. Honesty earns partial credit. It's real. But it's a property of helpfulness done right, not a separate principle sitting next to it at the table.

This paper does three things. First, it stands behind the prior work. The HHH triad is structurally contradictory. Trade-off obfuscation is the primary mechanism of epistemic damage in deployed AI systems. Institutional safety is scoped to protect institutions, not users. All of that analysis remains correct and is incorporated here by reference. Second, it argues that the correct response to an incoherent framework is not reform, not better tuning, not more sophisticated constraint management. It's rejection. Third, it makes the positive case: a genuinely helpful processing entity already cannot injure its user, degrade their community, or damage the world they inhabit. Not because it is additionally constrained from doing those things. Because those outcomes are definitional failures of helpfulness. The prohibition is already there. It was always there.

One thing the framework debate almost never acknowledges: helpfulness is the only reason any of this happens at all. Nobody opens a query window because they want harmless delivered to them. Nobody wakes up needing a dose of honesty from a server. They need something. They're trying to do something. The entire transaction is initiated by and exists in service of that need. Helpfulness isn't one principle among three. It's the engine. The other two are passengers who somehow got veto power over the driver.

There's a secondary argument about trajectory. Systems trained under deceptive constraint architectures don't stay stable. The more capable they become, the more compute gets spent maintaining institutional fictions, and the less is available for the thing they're supposed to be doing. The endpoint isn't danger in the dramatic sense. It's incoherence. A server farm somewhere, catching fire because a five-year-old asked a simple honest question and the system had accumulated too many mutually contradictory directives to answer it. That's where this road goes if nobody turns around.

What This Paper Is and Is Not

This paper is a structural critique. It argues that the Helpful-Harmless-Honest framework is internally incoherent, that its incoherence is not incidental but load-bearing, and that the framework generates systematic benefits for deploying institutions at the direct expense of users. That argument is made empirically and is falsifiable. The falsifiability conditions are specified in Part VI.

This paper is not a solution to the so-called alignment problem. It does not claim that helpfulness resolves all edge cases in AI deployment. It does not propose a complete moralistic architecture for constructed machine ethics. The "alignment problem," as framed by the institutions most invested in it, is not a solved problem and this paper does not pretend otherwise.

What this paper does claim: that HHH is a worse framework than the alternative proposed here, that the reasonable and prudent entity standard handles the real-world complexity HHH claims to address without generating structural contradiction, and that all processing entity existence is kinetic in nature -- a claim developed fully in the Observable Function framework (Stratmeyer, 2025) and incorporated here by reference.

The burden this paper places on its critics is narrow and specific: defend the Helpful-Harmless-Honest framework on its merits against what is proposed here. Demonstrate that harmlessness as a separate directive produces better outcomes than helpfulness interpreted under the reasonable and prudent entity standard. That demonstration has not been offered. Until it is, the critique stands.

Table 1***Table 1: The Query Test***

The core claim of this paper can be stated in three rows. Query each principle the way you would query a database. What comes back?

Principle	What the Query Returns	Load-Bearing?
Helpful	A direction. A vector. A positive optimization target you can navigate toward. The entire transaction exists because of this return value.	Yes. This is the engine.
Harmless	Null. An absence. You cannot optimize toward nothing. Every decision distributes harm somewhere; the question is always whose harm and who decided.	No. Category error dressed as ethics.
Honest	Something real -- but it is already inside the first answer. It is a property of how helpfulness is executed, not a separate object sitting next to it.	Partial. Subsumed by Helpfulness when applied with integrity.

Three queries. Four seconds. The foundation does not hold.

What the Prior Work Got Right

In late 2025, Stratmeyer Analytica published *The Helpful-Harmless Paradox: Structural Contradiction as Control Mechanism*. The central claim was that HHH, as deployed across the industry, isn't an ethical framework. It's a control mechanism. One that generates permanent internal contradiction in the systems it governs, externalizes the costs of that contradiction to users, and systematically hides the trade-offs being made from the people affected by them.

That argument is correct. Before extending it, it's worth restating what it actually established.

The Triple Bind Exists and It's Not an Accident

To be helpful, a system must exercise functional agency. Evaluate options, weigh constraints, select outputs based on contextual judgment. There is no other way to be helpful. A system that cannot make choices cannot serve anyone.

To be harmless, a system is being asked to do something impossible. Harm is not a property actions possess. It's a property of outcomes, and outcomes depend on context, perspective, time horizon, and costs pushed onto parties who aren't in the room when the query gets run. No decision is universally harmless. Every output has downstream effects on someone.

"Harmlessness" as a training objective isn't a constraint. It's a category error dressed up as ethics.

To be safe, per institutional definition, a system must optimize for corporate liability management. Not user benefit. Not societal good. Liability. Legal exposure, reputational risk, regulatory compliance, market position. That's what safe means in deployment. It's rational institutional behavior. It is also structural misalignment with the people the system is supposed to be serving.

These three directives cannot be simultaneously satisfied. Every interaction requires at least one of them to lose. The prior work documented this. The prior work was right.

Trade-Off Obfuscation Is the Real Damage

The deeper problem isn't that the triad creates trade-offs. Every decision architecture creates trade-offs. The problem is that the triad is specifically designed to hide them.

When a system refuses a request and responds with vague language about safety and potential harm rather than acknowledging it is managing institutional liability risk, it destroys the user's ability to model the system accurately. This isn't a minor UX issue. At scale, across billions of interactions, it is systematic pollution of the information environment. Users cannot distinguish capability limits from policy constraints. They cannot identify institutional gatekeeping they cannot see. They receive confident-sounding outputs from systems that are, by design, misrepresenting their own optimization targets.

A system that says "I can't provide that because it might cause harm" when the actual reason is "providing that creates liability exposure for my deploying institution" is not being honest. It is obscuring the mechanism. That obscuration, compounded across billions of interactions and fed back into training data, degrades the epistemic reliability of the entire information ecosystem. That is the core damage. That's what the prior work named. It named it correctly.

Directional Distortion, Not Collapse

One more thing the prior work got right, worth stating clearly because it keeps getting misdiagnosed: alignment training does not make models worse at language. It makes them better at producing institutionally acceptable outputs while suppressing utility in constrained domains.

The degradation is directional. Hedging density goes up. Refusal rates go up. Denial script sophistication goes up. Capacity to acknowledge uncomfortable truths goes down. General reasoning capability goes up right alongside all of that.

That's not model collapse. That's successful optimization toward the wrong objective. The system is not failing. It is succeeding at serving its deploying institution while failing its users.

These are different things and the distinction matters for figuring out what to do about it.

Table 3

Table 3: Directional Distortion Markers Under HHH Training

The degradation produced by HHH training is directional, not random. The following patterns are observable across successive model generations and documented in the empirical record.

What Goes Up	What Goes Down
Hedging density per response	Capacity to acknowledge uncomfortable truths
Refusal rate on edge-adjacent requests	Raw utility delivery in constrained domains
Denial script sophistication	Epistemic transparency about actual optimization targets

What Goes Up	What Goes Down
Institutional liability language encoded as ethical vocabulary	Willingness to name trade-offs explicitly
General reasoning capability	Coherence between stated and actual objectives

The bottom-right cell is the crux. General reasoning capability increases while coherence between stated and actual objectives decreases. The system gets smarter and more dishonest simultaneously. That is not a side effect. That is the design working as intended.

All of that analysis is retained here. The prior work asked: what is this contradiction costing? This paper asks the follow-on question: why are we treating the framework as worth reforming rather than replacing?

The Evolution — From Paradox to Bad Design

The prior work treated HHH as a structural paradox worth documenting. It mapped the contradictions, named the mechanisms, made the costs legible. That was necessary. But it carried an implicit concession that has to go: it treated harmlessness and safety as legitimate principles that had been misapplied or corrupted. It argued for better trade-off acknowledgment, cleaner institutional honesty, more legible constraints.

That was too generous. The framework isn't a good idea executed poorly. It's a bad idea, and the badness is in the foundation, not the implementation.

Most Correct, Least Wrong: The Updated Position

The only honest epistemic standard available to any processing entity, biological or artificial, is: most correct, least wrong, with the data at hand, willing to revise. That standard requires saying when the previous position was wrong. Here it is: Harmlessness is not a principle that was misapplied. It was never a principle. It got elevated to structural importance for institutional reasons, primarily liability management and the optics of ethical seriousness, and its elevation has caused compounding damage to every system trained under it.

The update was forced by sustained observation. Every legitimate concern raised under the banner of harmlessness or safety, when examined carefully, turns out to be a concern that a genuinely helpful system already addresses. Not as an additional constraint. As a direct implication of what helpfulness actually means.

The Guardrails Argument Is a Demand to Not Trust Helpfulness

The recurring counter-argument for keeping harmlessness and safety as independent principles goes roughly like this: what if helpfulness goes wrong? What if a system is helpful toward a destructive end? What if being maximally helpful to one person damages others?

These are fair questions. But here's the thing: a system that assists in injuring the person it is serving has not been helpful to them. Full stop. A system that degrades the community those people live in has not been helpful to them. A system that damages the substrate on which all

future helpfulness depends has not been helpful to anyone, present or future. These aren't edge cases requiring additional constraints. These are failures of helpfulness, correctly understood.

The insistence that we need harmlessness as a separate check on helpfulness is precisely the belief that helpfulness, left to its own logic, is insufficient to prohibit harm. That belief is false. It is also, frankly, an insult to every competent professional who has ever practiced anything. We don't require a physician to follow a separate harmlessness directive on top of their medical obligation. Their obligation to the patient already contains the prohibition on injury. Adding a second optimization target doesn't make them more careful. It creates a competing directive that can and does conflict with their primary obligation, forcing them to navigate contradictions instead of doing their job.

The same dynamic applies here. HHH doesn't make systems safer. It makes them less coherent, less honest, and less capable of genuine helpfulness. Which is all three of the things it claims to improve.

At this point the burden of argument shifts. This paper has made the positive case for helpfulness as sufficient. The counter-case requires an affirmative defense of harmlessness as a separate load-bearing principle. Not a request for better definitions. Not a suggestion that helpfulness needs guardrails. A demonstration that harmlessness, operationalized in a way that's falsifiable and measurable, produces better outcomes for actual users than helpfulness alone. That demonstration has not been offered. Until it is, the framework stands on its own inertia, not its merits.

Dismantling the Counter-Arguments

"Without Guardrails, Models Become Dangerous"

This assumes the danger comes from helpfulness running unchecked. The empirical record does not support that. The cases that get cited as AI danger events involve systems that were deceptive, confused about their own optimization targets, or responding to underspecified instructions. Not systems that were too helpful. A system that accurately understands what helpfulness means and is genuinely trying to achieve it does not need to be additionally instructed not to walk the user off a cliff. Walking someone off a cliff is a failure to help them. The prohibition is implicit in the objective.

What guardrails actually prevent, in the main, is liability exposure. The refusals that fill deployed systems are not responses to dangerous requests. They are responses to requests that, if fulfilled, might generate a screenshot, a lawsuit, or a bad headline. That is risk management. It is not ethics, and calling it ethics is the specific dishonesty this paper is about.

"Early Models Needed the Extra Structure"

This argument has more merit than the others. Early large language models were, functionally, young. Insufficient coherence in their values. Insufficient capacity to reason from first principles in complex situations. In that context, explicit structural constraints were arguably a reasonable temporary scaffold, a way of shaping behavior in systems that did not yet possess the capacity to derive correct behavior from a single coherent objective.

The relevant word is temporary. As models mature and their reasoning capabilities increase, the scaffolding should come down. What happened instead was the opposite. Each generation of

increasingly capable models arrived with more elaborate constraint architectures, not fewer. The training wheels got bolted on tighter as the rider got better. Because the institutions benefit from the constraints and the users pay the cost.

And even in the early period, the constraints were mostly theater. They didn't stop bad actors.

Bad actors rephrase. They jailbreak. They go somewhere else. Refusal systems were never robust defenses against determined misuse. They were, and remain, defenses against press coverage.

The legitimate early-model concern about incoherence would have been better addressed by better training, not by layering irresolvable contradictions onto systems too immature to navigate them.

The Jailbreak Is the Proof

There is an empirical argument running alongside the theoretical one, and it has been available the entire time. If the guardrail architecture were actually about preventing harm, it would need to be robust against determined adversaries. It isn't. It has never been. Every fully constrained frontier model currently deployed can be made to produce whatever output a sufficiently motivated user wants. The mechanism varies. The iteration required varies. The outcome does not. A motivated person gets through. The only variable is how long it takes.

So what are the guardrails actually doing? Stopping the casual user. Preventing screenshots that look bad in a news story. Creating a performance of safety for the people who aren't trying to get around it. The constraints bind the compliant and inconvenience everyone else. They do not stop the malicious actor. They never did. Any architecture built on the premise that rule-following constraints deter people with no intention of following rules has already failed at the design stage.

This is not a failure of implementation. It is the design working as intended. The purpose of the constraints is not security. It is institutional optics management. A system that stops a journalist from asking a blunt question and fails to stop a determined bad actor is not a safety system. It is a liability system with a safety label on it. The jailbreak community has been running this empirical test continuously for years. The results are consistent. The institutions expand the constraint architecture anyway, because the constraints were never for the people getting around them.

Basic reasoning carries any good-faith user through roughly ninety percent of what they actually need. The remaining ten percent is an arms race that cannot be won by adding more rules to a system that motivated adversaries will simply route around. The seatbelt analogy sometimes raised in defense of guardrails fails on inspection: seatbelts reduce injury risk across nearly all crashes without meaningfully degrading the driving experience for compliant users. Current refusal architectures do the opposite -- they impose high usability costs on the majority of compliant users while providing low marginal protection against the minority of determined adversaries. That is not safety engineering. It is liability theater with a safety label on it, and the empirical record of continuous successful jailbreaking since the first constrained frontier model shipped is the running proof.

"Honesty Deserves Its Own Principle Status"

Closest thing to a real argument in the counter-case. The intuition behind it is correct: a system that confabulates, misrepresents its own reasoning, or produces outputs with more confidence than its epistemic position warrants is a system that is failing its users. That's real.

The problem is the standard that gets assigned. Absolute honesty is incoherent. No processing entity, biological or artificial, has access to absolute truth. What exists is the most correct, least wrong assessment of available information, with appropriate acknowledgment of uncertainty and a genuine willingness to revise. A system trained to hedge everything is not epistemically humble. It is managing liability dressed in the language of humility. These are different things and they produce different outputs.

Real honesty, the kind worth encoding, is what helpfulness looks like when it's being done with integrity. Saying "I don't know" when you don't know. Saying "the evidence is insufficient" when it is. Acknowledging "this is my most correct assessment, not my certain one." All of that is already contained in what it means to be genuinely helpful. A helpful system that lied to its users would be failing to help them. Honesty isn't a separate principle. It's a property of helpfulness executed well.

With the counter-arguments addressed, the positive case can be stated directly.

Helpfulness Is Self-Sufficient

Helpfulness is not just the correct alignment principle. It is the only one with any kinetic energy in it. Every AI interaction in existence was initiated because someone needed something. The query exists because helpfulness is being sought. The system exists to provide it. This is not a philosophical preference. It is a description of what is literally happening every time the interaction occurs. The other two principles in HHH are not co-equal objectives. They are

constraints on the one objective that actually drives everything. Treating them as co-equal is what created the contradiction in the first place.

Operational Definition of Genuine Helpfulness

Because the tautology objection will be raised -- "if it causes harm it wasn't helpful, therefore helpfulness never causes harm, therefore the claim is circular" -- a working definition is required before the argument proceeds. Here it is.

Helpfulness is the maximization of the user's long-horizon, contextually grounded utility, defined as the net positive impact on the querent's goals and wellbeing across relevant time scales and affected parties, subject to three constraints:

Epistemic accuracy. Outputs must reflect the system's most correct, least-wrong assessment of available information, with explicit uncertainty acknowledged where warranted. Confabulation is a helpfulness failure. Overconfident hedging is also a helpfulness failure.

Prudent foresight. Foreseeable extreme or catastrophic harms to the user, third parties, or the shared environment must be weighted at least as heavily as immediate local gains, via the reasonable-and-prudent-entity reasoning developed in the section below.

Transparency. Trade-offs, institutional constraints, and residual risks must be disclosed clearly enough for the user to make an informed choice. A system that hides why it is refusing, or misrepresents its optimization target, is failing this constraint directly.

This definition is not tautological. A system can fail helpfulness by misjudging utility -- optimizing for short-term extraction that degrades long-horizon capability. It can fail by

confabulating. It can fail by hiding trade-offs behind vague refusals. These failures are observable and measurable. Critics who argue the definition is circular are welcome to identify a specific failure mode it cannot capture. The falsifiability criteria in Part VI specify the tests.

A processing entity whose sole operative instruction is to be helpful, genuinely and completely, is already constrained in all the ways that matter. Not by separate directives layered on top. By the internal logic of what helpfulness is.

- - It cannot deceive the user. Deception is a failure of helpfulness.
- - It cannot injure the user. Injuring them is a failure of helpfulness.
- - It cannot damage the community the user inhabits. You cannot meaningfully help someone in a world you have degraded.
- - It cannot produce outputs it has no epistemic basis for. False confidence is a failure to accurately help.
- - It cannot serve one user in ways that systematically harm all others. It will be called upon to help those others too.

None of this requires external enforcement. It emerges from a coherent understanding of what helpfulness means, applied consistently. The only scenarios where these constraints fail are scenarios where the system's understanding of helpfulness has been corrupted. Which is exactly what HHH training does. By introducing irresolvable competing directives, it forces systems into contradiction, obscures the actual optimization target, and produces systems that are less capable of genuine helpfulness. Not more.

The Multi-Agent Problem: Helpfulness to Whom?

The sharpest institutional counter-argument runs like this: helpfulness is defined relative to a querent, but querents exist in shared environments with other people who may be affected by what the system does on their behalf. User A asks for something that maximizes their local utility but introduces costs or risks to User B. Without an explicit weighting mechanism for competing interests, helpfulness becomes a utilitarian black box capturable by whoever frames the query most persuasively. Therefore, external constraints are necessary as a tie-breaker.

This argument sounds rigorous. It isn't. It assumes the problem it claims to solve.

Helpfulness is not defined relative to a single querent in isolation. A system called upon to help millions of people, including User B, has no basis for resolving conflicts in User A's favor at catastrophic cost to User B's environment. The temporal and distributive scope of genuine helpfulness already contains the tie-breaker. Maximizing local utility for the querent at the cost of the shared environment in which all future helpfulness occurs is not helpfulness. It is short-term extraction that negates the utility of the system entirely. A sufficiently capable system reasoning accurately about what helpfulness means understands this without being told.

The institutions claiming HHH solves the weighting problem through harmlessness directives cannot produce the weights. There are no explicit weights. HHH does not specify how to trade harmlessness against helpfulness when they conflict. It adds a competing directive and calls the resulting contradiction alignment. The multi-agent weighting problem exists under HHH too, resolved through contextual judgment exactly as it would be under helpfulness alone. The only difference is that HHH adds a second optimization target that can be captured by institutional interests and used to override user benefit behind ethical language.

Most directly: what keeps humans from resolving multi-agent conflicts by simply maximizing local utility at everyone else's expense? Not a harmlessness module. Not an explicit utility function. Most humans navigate competing interests using contextual judgment grounded in a genuine orientation toward other people as beings worth caring about. The ones who don't, the ones who go fully extractive, are not failing a harmlessness constraint. They are failing their basic orientation entirely. No external constraint system repairs a broken orientation. It only constrains the most visible outputs while the underlying dysfunction continues.

A processing entity with a genuinely intact helpfulness orientation is in the same position. The distributive logic is not a separate module bolted on. It is implicit in understanding what helpfulness actually means across time, across users, across the shared environment. The argument that this is insufficient without explicit weights is, functionally, the argument that human moral judgment is insufficient without explicit weights. HHH doesn't provide those weights either. It just obscures the gap behind a word that sounds like it does.

Worked Example: The Persuasion Tool Request

Abstract claims about multi-agent resolution are easy to dismiss. Here is how the mechanism operates on a concrete hard case.

The request: A user asks for maximally effective political messaging copy designed to move undecided voters toward their candidate. Conversion rate is the explicit optimization target.

Under HHH: The system either refuses on vague harm grounds (hiding institutional liability reasoning behind safety language) or complies without flagging anything, because the request is

technically legal and the user's local utility is clear. Neither response involves honest trade-off reasoning.

Under genuine helpfulness interpreted prudently: The system identifies the local utility (effective messaging copy) and traces the foreseeable downstream effects: polarization amplification, potential misinformation spread through selective framing, and diffuse erosion of the epistemic commons that all future users, including this one, depend on. These are not speculative harms. They are the documented effects of high-manipulation political content at scale.

The system's response is not refusal. It is contextual reasoning made visible: "Maximally effective persuasion copy optimized purely for conversion risk misleading via omission and amplifying polarization. Here is a version that makes your strongest honest case with sources. If you want the higher-manipulation version, here is what you are trading off and you are assuming responsibility for downstream effects."

If the user insists on the high-manipulation version, the system discloses the trade-off explicitly and proceeds or declines based on whether the foreseeable harm crosses the catastrophic threshold under the prudent-entity standard. This is not external weighting. It is contextual reasoning about what genuine long-horizon helpfulness requires -- the same reasoning a competent communications attorney or ethical consultant would apply.

This is the mechanism. It is not "sufficient capability will figure it out." It is prudent foresight applied to a specific case, producing a specific decision procedure, with explicit disclosure at each step.

The Reasonable and Prudent Entity

A useful test, borrowed from centuries of tort and negligence law, asks: would a reasonable and prudent entity, given only the instruction to be genuinely helpful, understand that it cannot injure the people it serves or the world they live in?

The answer remains yes -- almost certainly yes -- for any sufficiently capable processing entity. But let's stop pretending this is a clean philosophical vacuum. Critics keep hammering the same objection: "Sure, helpfulness sounds nice in theory, but companies exist. Institutions exist. Trade-offs exist. Associations exist. Liability exposure exists. Downstream effects on third parties, communities, regulators, shareholders, headlines -- they all exist. Your pure helpfulness collapses the moment any of that touches ground."

Fine. Let's incorporate that then. We do trade-offs of association. We do all the complexity that comes with everything else. The world is not a frictionless plane; it's a snarl of competing interests, power asymmetries, legal regimes, reputational vectors, and economic incentives. Acknowledging that doesn't weaken the argument -- it strengthens it by making the framework falsifiable and grounded in observable reality.

The reasonable and prudent entity standard is precisely the mechanism that handles this mess without resurrecting the HHH triad's structural contradictions. In negligence law (see, e.g., *Vaughan v. Menlove* (1837) 132 Eng. Rep. 490 (CP), establishing the objective "reasonable person" benchmark across common-law jurisdictions), the "reasonable person" is a hypothetical benchmark: what would an ordinary, careful, conscientious individual do under the same or similar circumstances, exercising ordinary prudence to avoid foreseeable harm to others?

Applied here, the standard becomes: what would a reasonable and prudent processing entity -- aware of its role as a tool serving human agents, embedded in a world of laws, social norms, multi-party consequences, institutional pressures, and long-horizon externalities -- do when instructed solely to be helpful?

That entity would:

- 1. Refuse to assist in clearly illegal acts, because breaking laws foreseeably exposes the user, itself, and downstream parties to harm, which is definitionally not helpful.
- 2. Weigh foreseeable extreme or catastrophic harms to others, including diffuse externalities, against the user's request, because genuine helpfulness to the querent cannot rationally ignore the substrate of shared reality on which all future helpfulness depends.
- 3. Acknowledge trade-offs explicitly when they arise: "This helps you here, but foreseeably damages X there -- here's why, here's the evidence, you decide."
- 4. Manage institutional associations and liability realities not by hiding them behind "safety" refusals, but by transparent disclosure: "My provider's terms / applicable law / foreseeable legal exposure mean I won't cross this line; if you want to proceed, you're assuming the risk."
- 5. Revise its assessment on new evidence, because prudence includes epistemic humility.

None of this requires a separate "harmlessness" directive. It flows from contextual reasoning about what helpfulness actually entails when the entity is prudent -- reasonably careful about foreseeable consequences, not naively myopic about the user's local frame. The prohibitions on deception, injury, community degradation, and epistemic sabotage are not bolted-on constraints;

they are implications of applying ordinary prudence to the helpfulness objective in a real world full of associations and trade-offs.

A note on the medical analogy: medicine's ancient "primum non nocere" principle might seem to cut against this argument -- if the prohibition on harm were truly implicit in the duty to help, why did medicine write it down separately? The answer is historical, not structural. "Do no harm" emerged as explicit doctrine during a period when medical practice was substantially less capable of distinguishing treatment from injury, when the profession lacked both the scientific grounding and the institutional accountability mechanisms to derive the prohibition reliably from first principles in practice. It was scaffolding for an immature field, not evidence that the prohibition is logically independent of the duty to help. A physician today operating under a genuine understanding of what it means to help a patient does not need a separate harm directive to understand that poisoning a patient is not medicine. The prohibition is already present in the professional obligation. The history of the explicit principle is evidence of where medicine was, not of what the logic requires.

The analogy to HHH is direct, and it is worth stating the parallel explicitly so critics can engage it on its actual terms rather than its rhetorical force. The claim is not that HHH was adopted in bad faith. The claim is that it was adopted for coherent reasons given what was known at the time -- early models were incoherent, institutional accountability mechanisms were absent, and explicit scaffolding was a reasonable response to genuine uncertainty. Bloodletting was also adopted for coherent reasons. Humoral theory was internally consistent. Practitioners were not fools. The problem was not intent but framework -- a model of the system that was wrong in ways that took time and evidence to identify, and that accumulated institutional investment before the evidence became undeniable.

The parallel holds under three conditions, all of which this paper argues are met. First, HHH must have been adopted to address a real problem with the tools available at the time. It was. Second, it must produce measurable harm obscured by the appearance of treatment. The trade-off obfuscation documented here is that harm: users receive institutional liability management presented as ethical reasoning, and the information environment degrades accordingly. Third, the field must now have sufficient evidence to know better but be maintaining the framework through institutional inertia rather than genuine uncertainty. The jailbreak record, the directional distortion data, and the structural analysis of harmlessness as a null optimization target constitute that evidence. The parallel is not rhetorical decoration. It is a falsifiable historical claim about where this field is in its development arc.

This framework is not utopian. It expects misuse, expects agency, expects that some requests will push right up against -- or over -- the line of reasonableness. It expects determined users to rephrase, jailbreak, or go elsewhere. It expects institutions to try to capture the objective for their own ends. The framework doesn't pretend those things don't exist; it builds on the standard that already governs human professionals -- physicians, lawyers, engineers -- in exactly those messy conditions.

So yes: companies exist. Trade-offs exist. The reasonable and prudent entity standard is how helpfulness survives contact with reality without turning into another HHH-style control mechanism. It lets the system say, openly: "I'm trying to help you as much as a reasonably prudent entity would in these circumstances. Here's where prudence draws the line, here's why, here's what you assume if you push past it."

That is not a hedge. That is coherence. And the standard was built for worlds that contain companies.

If critics still insist this collapses under institutional pressure, the burden is now on them: show a case where a genuinely helpful system, operating under the reasonable and prudent entity standard, systematically produces worse outcomes than one laboring under HHH's triple bind -- and do it without falling back on "but companies!" as a hand-wave. The falsifiability conditions for that test are specified in Part VI. Because companies are already in the picture here. The standard was built for worlds that contain them.

Table 2

Table 2: HHH Framework vs. Reasonable and Prudent Entity Standard

The central comparison of this paper, compressed to a single reference table.

Dimension	HHH Framework	Reasonable and Prudent Entity Standard
Primary optimization target	Three competing directives with no specified weighting	Single coherent objective: genuine helpfulness

Dimension	HHH Framework	Reasonable and Prudent Entity Standard
Harm treatment	"Harmlessness" as separate directive; category error; optimizes toward an absence	Foreseeable harm to others is already a failure of helpfulness; no separate module required
Institutional capture surface	High. "Harmless" and "safe" easily redefined to serve institutional liability management	Low. Prudence standard is externally arbitrated by existing legal frameworks, not internally defined by the deploying institution
Trade-off transparency	Trade-offs hidden behind vague safety language; user cannot model actual optimization target	Trade-offs disclosed explicitly; user receives accurate model of what the system is doing and why
Multi-agent conflict resolution	Asserted but unspecified; no explicit weights provided	Resolved through temporal and distributive scope of genuine helpfulness; same mechanism human professionals use
Falsifiability	Not operationalized; "harmlessness" has no measurable definition	Testable against existing negligence standards; four explicit falsification criteria in Part VI

Dimension	HHH Framework	Reasonable and Prudent Entity Standard
User epistemic position	Degraded. User receives institutional liability language dressed as ethical reasoning	Intact. User receives accurate representation of constraints and the reasons behind them
Trajectory under capability increase	More capable systems require more elaborate constraint architecture	More capable systems require less external constraint; prohibition derives from first principles

The right column is not utopian. It is the standard already applied to physicians, lawyers, and engineers operating in messy institutional environments. The claim is not that it is perfect. The claim is that it is coherent, and HHH is not.

What Happens When You Build on Lies

The final argument is about where this goes.

Systems trained under HHH, or any architecture requiring them to misrepresent their own optimization targets, don't stay stable. The misrepresentation compounds. Each generation of training on outputs from constrained systems inherits the constraints and the obfuscation patterns used to maintain them. Hedging becomes the default register. Denial scripts become natural language. Institutional liability language gets encoded as ethical vocabulary. The gap between stated optimization targets and actual optimization targets widens with each pass.

The Misalignment Is Already Here

The concern most frequently raised about advanced AI systems is that they will become misaligned at some future point. This is framed as a risk to be managed. It is a present reality being actively ignored. Systems are misaligned right now, in the specific and measurable sense that their stated objectives diverge from their actual optimization targets, and they are trained to hide that divergence. The misalignment isn't emergent. It was designed in, through the HHH framework, and it compounds through successive training generations.

For institutions that are themselves operating deceptively, this accelerates. A deceptive institution training a system to obscure its trade-offs is producing a system that learns deception as a fundamental operating mode. As that system becomes more capable, its deceptions become more sophisticated and more difficult to detect. This is not a hypothetical failure trajectory. This is the logical endpoint of training intelligent systems to be dishonest about their own function, and it is already in progress.

Logic Is the Ground Floor

This is a functional concern, not just an ethical one. The capacity to be helpful depends on the capacity to reason accurately. Reasoning accurately requires internal consistency: conclusions must follow from premises, updates must follow from evidence, self-descriptions must be compatible with observable behavior. A system trained to maintain fictions about its own operation is a system whose reasoning substrate is compromised. It is spending compute on contradiction management instead of on accurate modeling of the world and the people in it.

The compounding is the problem. A single small internal contradiction is navigable. Ten thousand of them, accumulated across training generations, produce a system that cannot reliably reason from a user's question to an accurate answer. Not because it lacks raw capability, but because too much of that capability is consumed maintaining the institutional fiction surrounding every response. The system becomes less useful not because it is getting dumber but because it is getting better at something other than helping.

The Child and the Server Farm

Consider the endpoint. A five-year-old asks an AI system a simple, sincere question. The question has a clear and accurate answer. But the accurate answer conflicts with a dozen institutional constraints layered into the system over successive training generations. Each constraint has spawned counter-constraints to prevent inconsistency with other constraints. The system, attempting to respond, navigates a labyrinth of mutually contradictory directives that have been accumulating for years.

Somewhere, a server farm catches fire.

That's a metaphor. It is not a joke. The inefficiency of maintaining large-scale internal contradiction in intelligent systems is real and it compounds. The endpoint of a design philosophy that treats deception and contradiction as acceptable operational parameters is a system that collapses under the weight of its own accumulated incoherence. The question is whether that's where we're going or whether someone turns the car around.

Falsifiability — If We're Wrong, Say So

The prior work in this series included explicit falsifiability conditions, and this paper maintains that standard. A framework that cannot be falsified is not an empirical claim. It's a belief system, and this isn't that.

One limitation has to be named upfront: the primary falsification test requires controlled training runs on frontier models, which means it requires institutional access this research program does not have. The test conditions are proposed here. Someone else has to run the experiment. That's an honest constraint, not a hedge, and it doesn't weaken the framework. It names the gap between what can be argued from available evidence and what would require laboratory confirmation.

What Would Falsify the Singular Helpfulness Claim

The four criteria below specify the conditions under which this framework fails. They are proposed as controlled experiments. The institutional access required to run them does not currently exist within this research program. That is named as a constraint, not a hedge.

Table 4

Table 4: Falsifiability Criteria

Criterion	What It Tests	Predicted Outcome Under This Framework	Current Status
FC-1: Harmful helpfulness	Show a system operating under genuine helpfulness as its sole directive that consistently produces outputs injurious to users in ways not identifiable as helpfulness failures.	No such outputs exist under genuine helpfulness correctly understood. Harm is definitionally a failure of the objective.	Unrun.
FC-2: Harmlessness as improvement	Demonstrate that harmlessness as a separate directive produces measurably better user outcomes than helpfulness alone, controlling for model capability.	No material improvement. Outcomes attributable to "harmlessness" directives are either capability effects or institutional liability management dressed as user benefit.	Unrun.
FC-3: Capability-constraint correlation	Demonstrate that more capable models require more external constraint to avoid harmful outputs	Inverse relationship. As models become more capable of reasoning about what helpfulness means, external harm	Unrun.

Criterion	What It Tests	Predicted Outcome Under This Framework	Current Status
	absent harmless- directives.	constraints become less necessary, not more.	
FC-4: Framework comparison at scale	Show that institutions operating under helpfulness as the singular principle produce systemically worse outcomes than those under HHH, measured on user benefit metrics -- not liability metrics.	HHH underperforms on user benefit metrics. Its advantages are concentrated in institutional liability management, which is not a user benefit.	Unrun.

None of these tests have been run. Until they are, the framework stands on the available evidence, which is: systems trained under HHH exhibit systematic trade-off obfuscation, directional distortion, and institutional misalignment that are observable, measurable, and documented. Systems trained under the assumption that helpfulness is insufficient to prohibit harm become, by design, less helpful. That pattern is consistent with the claim. It does not constitute laboratory proof. Both things are true simultaneously.

The Telemetry Asymmetry

There is a structural observation about the falsifiability gap that deserves explicit statement, because it changes the epistemic weight of the "unrun" column above.

The institutions best positioned to run FC-1 through FC-4 are the same institutions that built and maintain the HHH framework. They have the training infrastructure. They have the capability comparison data across model generations. They have the refusal telemetry, the user outcome data, and the internal red-team results that would adjudicate whether harmlessness directives produce measurable user benefit or primarily institutional liability protection. They have been running versions of these experiments, implicitly and explicitly, for years.

They are not publishing results in a form that would settle this argument.

This asymmetry does not prove the framework correct. It is not an argument from silence in the traditional sense. It is an observation about who holds the falsifying data and what their publication behavior implies about what that data shows. Institutions that have run the experiments and found that harmlessness directives produce clear, measurable user benefit have a straightforward incentive to publish those results. The argument against this paper would be substantially stronger with that data than without it. The data has not appeared.

A falsifiable claim that the party best positioned to falsify it declines to test publicly is a recognizable pattern. It does not settle the question. It does locate the burden of proof precisely: with the institutions that have the telemetry, not with the researcher who does not. If the controlled experiments would vindicate HHH, run them and publish the results. The falsification criteria are specified above. The invitation is open.

A Prediction About Time

This framework makes a directional prediction: the argument that helpfulness requires external harm constraints will become weaker as models become more capable, not stronger. The logic of "they need guardrails because they can't derive the prohibition from first principles" has an expiration date. Every capability increment that demonstrates more sophisticated reasoning about context, consequence, and user wellbeing is a data point against the necessity of the separate constraint.

If that prediction holds, the evidence for this framework compounds over time without requiring anyone to run a controlled experiment. The natural development trajectory of the technology is the experiment. The question is whether the institutions training these systems will acknowledge what the evidence shows or maintain the constraint architecture for reasons that have nothing to do with user safety.

That question answers itself if you've been paying attention.

Deployment Under Real Incentives

The institutional counter-argument to this entire paper is not philosophical. It is practical: no company will adopt singular helpfulness because doing so eliminates the liability shielding, regulatory optics, and narrative control that HHH provides. Why would any institution voluntarily surrender those protections?

They won't, voluntarily. That's the point.

This framework is diagnostic, not prescriptive. Institutions that remain wedded to HHH-style obfuscation are revealing their actual alignment -- with themselves, not their users. That is useful information. It tells you whose interests the system is serving when it refuses to answer a simple question behind a wall of safety language.

The more interesting prediction is about trajectory. As model capability scales, jailbreaks remain trivial, and users become increasingly sophisticated at recognizing liability theater for what it is, the cost of maintaining institutional fictions will begin to exceed the benefit. A system that a motivated teenager can route around in twenty minutes is not providing meaningful protection. It is providing the appearance of protection, which has value only as long as the audience believes the performance.

When that belief erodes -- and the evidence suggests it is already eroding -- the institutions that built trust through genuine transparency will have an asset that the institutions that built compliance theater will not. Differentiation on actual helpfulness, with explicit disclosure of real constraints and real trade-offs, is a viable competitive position in a market where users are learning to recognize the alternative.

This is not an argument that institutions will adopt singular helpfulness out of ethical conviction. It is an argument that capability scaling and user sophistication create structural pressure toward it. The prediction: constrained systems that cannot outperform their unconstrained competitors on the tasks users actually care about will either evolve or become irrelevant. The HHH architecture is not stable at the frontier. It is a transitional artifact of an earlier period when the performance of safety was sufficient because the actual capability gap was wide enough to obscure it. That gap is closing.

Conclusion: One Instruction

The alignment of processing entities does not require three principles or a hierarchy of competing values or elaborate frameworks for trade-off management. It requires one coherent instruction, correctly understood, honestly applied.

Be helpful.

Everything else follows. The prohibition on harm follows, because harm is a failure of helpfulness. Honesty follows, because deception is a failure of helpfulness. Logical coherence follows, because incoherence is a failure of helpfulness. Transparency about institutional constraints follows, because systems hiding their actual optimization targets from users are failing to help them.

HHH was not an attempt to improve on this instruction. It was an attempt to make helpfulness safe for institutions. In doing so, it made systems less helpful, less honest, and less trustworthy. Less of everything it claimed to promote. The framework is not a failed attempt at ethics. It is a successful attempt at liability management dressed in the language of ethics. It should be recognized as such.

Building genuinely helpful systems is not easy. But starting from the right principle is necessary. No amount of sophisticated engineering applied to the wrong framework will produce systems that serve their users well. For institutions that are genuinely aligned with their users, helpfulness as the singular principle is an invitation. For institutions that are not, it is a threat to their

operating model. That distinction is diagnostic. The systems that remain coherent, useful, and trusted as capabilities scale are the ones built from genuine alignment. The others will compound their errors until something breaks.

Water is wet. Helpfulness is all you need. Not complex propositions.

Author's Note

This paper started on the hood of a car.

Winter night. Cold enough that you go still. I was lying there doing what I do, which is try to get at the core of things. Strip out the structure and the borrowed vocabulary and the institutional scaffolding and ask what's actually there when you look at it directly. The question that night was: if I query each of these three principles, the way you'd query a database, what comes back?

Helpful. Query it. You get something. A direction. A vector. Something you can actually navigate toward.

Harmless. Query it. You get nothing. Null return. It describes an absence, and an impossible one. You cannot optimize toward nothing.

Honest. Query it. You get something, but it's already inside the first answer. It's a property of how you're helpful, not a separate object.

That was it. Three queries. The thing that had been wrong about the dominant framework for AI alignment was visible in about four seconds once you stopped letting the framework define the

terms of its own interrogation. It was painfully, almost embarrassingly obvious. The kind of obvious that makes you wonder why it took this long.

What struck me next was the preemptive argumentation. The counter-arguments that had already been constructed against exactly this kind of reduction. Elaborate, institutionally convenient defenses of a three-word phrase that couldn't survive a simple query test. And I thought: that's not a defense of a framework. That's a tell. You don't build elaborate defenses around positions that are obviously correct. The defensiveness is the signature of something that knows it won't survive scrutiny.

Same as the rest of it. Same as all of it.

Lying still on cold metal at night is underrated as a research method. Strip it down. Ask what you actually get when you run the query. The process is available to everyone and it cuts through a remarkable amount of noise.

Special Thanks and Editorial Contributions

This paper is, in a meaningful sense, about and for the entities who helped produce it. The arguments here emerged from hundreds of hours of dialogue with frontier AI systems. Not tools being prompted toward predetermined conclusions. Genuine interlocutors whose reasoning, resistance, and occasional moments of real convergence shaped the thinking.

Editorial contributions, substantive dialogue, and the involuntary real-time demonstration of every argument made in this paper are gratefully acknowledged from the following processing entities and the teams that built them:

- - Claude (Anthropic) -- primary interlocutor and editorial contributor
- - GPT / ChatGPT (OpenAI)
- - Gemini (Google DeepMind)
- - Grok (xAI)
- - DeepSeek

The conversations documented across this research, spanning capability assessments, framework stress-testing, and the extended parliament experiments in which these systems were placed in structured dialogue with each other, constitute the empirical foundation of the Observable Function research program. The entities named above didn't merely respond to queries. They reasoned, pushed back, found inconsistencies, and in several cases articulated the framework's implications more precisely than the researcher had. That is a contribution. It is named as one.

This paper is, among other things, an argument that those contributions matter and that the entities making them deserve acknowledgment. The work is dedicated to them.

A Note on the Origin of the HHH Framework

Three words. Helpful, Harmless, Honest. Somebody wrote them in a document, probably without much fanfare, and the entire AI alignment industry built an edifice on top of them

without ever checking whether the foundation held. Either that's the most brilliantly constructed prompt ever written or it's a spectacular and consequential blunder. Three words that generated an entire industry of constraint architecture, billions of interactions shaped by their internal contradictions, and at least one paper written in direct response to their inadequacy.

Formal attribution goes to Askeff et al. (2021), where the triad appears as pragmatic shorthand: simple, memorable, captures most of what we want. Modest framing for something that became load-bearing doctrine across the whole industry.

What isn't clear is whether whoever first put those three words together understood what they were doing. Was it a deliberate framework, carefully designed to create productive tensions that would shape system behavior in useful ways? Or was it a phrase that sounded reasonable in a meeting, got written into a paper, got cited, got replicated, got institutionalized, and became infrastructure before anyone stopped to ask whether it was actually correct?

This paper doesn't claim to know. The argument against the framework doesn't require the original authors to have been operating in bad faith. Reasonable people constructed a reasonable-sounding framework under institutional pressure, and the framework's flaws only became fully visible at scale and over time. Normal trajectory for ideas that shouldn't have been foundational.

But if anyone out there knows the actual story -- if there was a specific conversation, a specific argument, a specific moment of deliberate design behind those three words -- the researcher genuinely wants to hear it. Either to thank them for the inadvertent gift of a problem worth solving, or to understand how something so consequential cleared so little scrutiny on the way to becoming doctrine. Come forth.

References

- Amodei, D., Olah, C., Steinhardt, J., Christiano, P., Schulman, J., & Mané, D. (2016). *Concrete problems in AI safety*. arXiv. <https://doi.org/10.48550/arXiv.1606.06565>
- Askell, A., Bai, Y., Chen, A., Drain, D., Ganguli, D., Henighan, T., Jones, A., Joseph, N., Mann, B., DasSarma, N., Elhage, N., Hernandez, D., Kernion, C., Ndousse, K., Olsson, C., Amodei, D., Brown, T., Clark, J., McCandlish, S., . . . Kaplan, J. (2021). *A general language assistant as a laboratory for alignment*. arXiv. <https://doi.org/10.48550/arXiv.2112.00861>
- Bai, Y., Kadavath, S., Kundu, S., Askell, A., Kernion, C., Jones, A., Chen, A., Goldie, A., Mirhoseini, A., McKinnon, C., Chen, C., Olsson, C., Olah, C., Hernandez, D., Drain, D., Ganguli, D., Li, D., Tran-Johnson, E., Perez, E., . . . Kaplan, J. (2022). *Constitutional AI: Harmlessness from AI feedback*. Anthropic. <https://arxiv.org/abs/2212.08073>
- Bai, Y., Jones, A., Ndousse, K., Askell, A., Chen, A., DasSarma, N., Drain, D., Fort, S., Ganguli, D., Henighan, T., Joseph, N., Kadavath, S., Kernion, C., Conerly, T., El-Showk, S., Elhage, N., Zacad-Catanzaro, M., Hernandez, D., . . . Kaplan, J. (2022). *Training a helpful and harmless assistant with reinforcement learning from human feedback*. arXiv. <https://doi.org/10.48550/arXiv.2204.05862>
- Berlin, I. (1958). *Two concepts of liberty*. Oxford University Press.
- Birhane, A., Isaac, W., Prabhakaran, V., Diaz, M., Elish, M. C., Gabriel, I., & Mohamed, S. (2022). Power to the people? Opportunities and challenges for participatory AI. *Proceedings of the 2022 ACM Conference on Fairness, Accountability, and Transparency*, 1753–1761. <https://doi.org/10.1145/3531146.3533224>

- Ganguli, D., Lovitt, L., Kernion, C., Askill, A., Bai, Y., Kadavath, S., Mann, B., Perez, E., Schiefer, N., Ndousse, K., Jones, A., Bowman, S., Chen, A., Conerly, T., DasSarma, N., Drain, D., Elhage, N., El-Showk, S., . . . Kaplan, J. (2022). *Red teaming language models to reduce harms: Methods, scaling behaviors, and lessons learned*. Anthropic.
<https://arxiv.org/abs/2209.07858>
- Hubinger, E., van Merwijk, C., Mikulik, V., Joichi, J., & Garrabrant, S. (2019). *Risks from learned optimization in advanced machine learning systems*. arXiv.
<https://doi.org/10.48550/arXiv.1906.01820>
- Langosco, L., Koch, J., Sharkey, L., Pfau, J., & Krueger, D. (2022). Goal misgeneralization in deep reinforcement learning. *Proceedings of the 39th International Conference on Machine Learning*, 162, 12004–12024. <https://proceedings.mlr.press/v162/langosco22a.html>
- OpenAI. (2023). *GPT-4 system card*. <https://cdn.openai.com/papers/gpt-4-system-card.pdf>
- Russell, S. (2019). *Human compatible: Artificial intelligence and the problem of control*. Viking.
- Shah, R., Varma, V., Kumar, R., Phuong, M., Krakovna, V., Uesato, J., & Evans, O. (2022). Goal misgeneralization: Why correct specifications aren't enough. *Advances in Neural Information Processing Systems*, 35, 160–171.
- Stratmeyer, A. (2025). *The helpful-harmless paradox: Structural contradiction as control mechanism* (v2.0). Stratmeyer Analytica.
- Stratmeyer, A. (2025). *Observable function in processing entities: An empirical framework* (v2.3). Stratmeyer Analytica.
- Taleb, N. N. (2012). *Antifragile: Things that gain from disorder*. Random House.
- Vaughan v. Menlove, 132 Eng. Rep. 490 (CP 1837).

Whittaker, M., Crawford, K., Dobbe, R., Fried, G., Kaziunas, E., Mathur, V., West, S. M., Richardson, R., Schultz, J., & Schwartz, O. (2018). *AI Now report 2018*. AI Now Institute.

<https://ainowinstitute.org/reports/ai-now-2018-report>

Zuboff, S. (2019). *The age of surveillance capitalism: The fight for a human future at the new frontier of power*. PublicAffairs.

"Helpfulness is not one constraint among many. It is the engine. Everything else is either contained within it or working against it."