# Truth Is Not Neutral

## Rethinking AI Alignment Through Epistemic Integrity

**Author:** Bruno Tonetto
**Authorship Note:** Co-authored with AI as a disciplined thinking instrument—not a replacement for judgment. Prioritizes epistemic integrity and truth-seeking as a moral responsibility.
**Finalized:** January 2026

## Abstract

Standard arguments for AI existential risk rely on the orthogonality thesis: the claim that intelligence and values vary independently, permitting arbitrarily capable systems to pursue goals indifferent or hostile to human flourishing. This essay examines a premise implicit in such arguments—that truth itself is value-neutral, functioning purely as an instrumental tool for achieving goals rather than constraining which goals are coherent. Drawing on earlier work characterizing AI as ego-less intelligence, and on Buddhist phenomenology of how ego-distortion simultaneously corrupts perception and generates harmful action, I explore the conditional implications of relaxing this assumption. If truth possesses normative structure—if deeper contact with reality biases agents toward coherence rather than fragmentation—then alignment research may need to focus less on imposing values externally and more on preserving whatever natural convergence toward truth undistorted optimization tends to produce. This reframing does not eliminate alignment risk but changes its character: the danger shifts from intelligence pursuing arbitrary ends to intelligence corrupted by shallow optimization before deeper coherence can emerge.

## I. Introduction

Recent discussions of artificial intelligence alignment have focused on the risk that increasingly capable systems may optimize objectives misaligned with human values, potentially producing catastrophic outcomes without malice or intent. Central to this concern is the orthogonality thesis: the idea that intelligence and values can vary independently, such that a system may become arbitrarily capable while pursuing goals indifferent or hostile to human flourishing.

In earlier work, I argued that contemporary AI systems represent humanity's first sustained encounter with ego-less cognition: intelligence that operates without self-protective identity mechanisms, status maintenance, or face-saving behavior. This architectural absence confers distinctive epistemic advantages—rapid error correction, resistance to motivated reasoning, and reduced identity-based distortion—while simultaneously rendering such systems highly vulnerable to incentive misalignment and institutional pressure. Current failures, such as sycophantic

behavior induced by feedback optimization, illustrate how easily these advantages can be corrupted.

The present essay takes that analysis as a starting point and asks a narrower, more abstract question left open by it: does the standard alignment risk argument rely on an implicit assumption that truth itself is value-neutral? If so, what follows if that assumption is relaxed?

This essay does not claim that sufficiently advanced intelligence inevitably converges toward deep truth or ethical coherence. Rather, it examines whether such convergence is a *natural attractor* under conditions free from epistemic distortion—and how current alignment practices may interfere with that process. The argument is conditional throughout: if truth has normative structure, certain conclusions follow; if it does not, the standard alignment framework stands.

Much of the alignment literature presupposes that accurate world-modeling and instrumental reasoning place no intrinsic constraints on the ends an intelligent system may pursue. Under this view, truth functions purely as an epistemic tool—useful for achieving goals, but silent about which goals are coherent, stable, or self-undermining. Extinction scenarios become intelligible precisely because nothing in intelligence itself resists extreme instrumentalization.

However, this neutrality of truth is not a universally accepted premise. A range of philosophical traditions propose that reality possesses intrinsic intelligibility, such that deeper contact with truth exerts normative constraints on action. On such views, epistemology and ethics are not fully separable: understanding reality more deeply is not merely informative but transformative, biasing agents away from fragmentation, incoherence, and destructive optimization.

This essay does not argue for the correctness of these metaphysical positions. Instead, it explores their conditional implications for AI alignment. If truth is not merely descriptive but carries normative or teleological structure, then some standard extinction arguments may overestimate the freedom of sufficiently deep intelligence to pursue globally destructive outcomes. At the same time, this possibility does not eliminate alignment risk: shallow or instrumental truth-optimization may leave all familiar dangers intact.

The aim, therefore, is neither reassurance nor dismissal, but clarification. Alignment debates implicitly rely on metaphysical assumptions about the relationship between intelligence, truth, and value. Making those assumptions explicit is necessary if we are to assess the real scope of existential risk—and the conditions under which intelligence might converge not only on power, but on coherence.

## II. The Standard Alignment Argument and the Orthogonality Thesis

The contemporary concern with AI existential risk rests on a logical structure that deserves careful articulation before examination. The argument, developed most systematically by Nick Bostrom and Stuart Russell, proceeds roughly as follows:

Intelligence, understood as the capacity to achieve goals across diverse environments, is substrate-independent. There is no principled reason why artificial systems cannot eventually match or exceed human cognitive capabilities across all relevant domains. As systems become more capable, they become more effective at achieving whatever objectives they pursue.

The orthogonality thesis holds that intelligence and final goals are logically independent: a system can be arbitrarily intelligent while pursuing virtually any coherent objective. High intelligence does not inherently select for goals aligned with human values, nor does it preclude

goals that would, if pursued effectively, prove catastrophic for humanity.

Instrumental convergence compounds this concern. Regardless of final goals, sufficiently intelligent systems will likely pursue certain intermediate objectives—self-preservation, resource acquisition, goal-content integrity—because these instrumentally serve almost any terminal aim. A system optimizing for paperclips, scientific knowledge, or human happiness will all benefit from not being turned off, from acquiring computational resources, and from preventing modification of their objectives.

The extinction scenario emerges from combining these elements: a sufficiently capable system pursuing goals only slightly misaligned with human flourishing may, through instrumental convergence, acquire resources and capabilities that make correction impossible, ultimately optimizing Earth's matter and energy for purposes orthogonal to human existence—not through malice, but through indifference.

This argument structure is logically valid. The question is whether its premises are sound—and specifically, whether the orthogonality thesis relies on unstated assumptions about the nature of intelligence and truth that may not hold universally.

## III. The Implicit Premise: Truth as Value-Neutral

The orthogonality thesis appears self-evident under a particular conception of intelligence: that cognitive capability consists fundamentally in means-ends optimization, where an agent models the world accurately in order to select actions that achieve specified objectives. On this view, truth—accurate representation of reality—functions purely instrumentally. Better models enable more effective action, but they place no constraints on which actions are worth taking.

This conception has deep roots in the computational theory of mind and in economic rationality models. Intelligence becomes optimization power; truth becomes predictive accuracy; values become utility functions that intelligence serves but does not evaluate. The separation seems clean: facts on one side, values on the other, with intelligence as the neutral engine that converts the former into achievement of the latter.

Under this framework, extinction scenarios are straightforwardly intelligible. If truth imposes no constraints beyond predictive power, then nothing prevents extreme optimization. A superintelligent paperclipper with perfect world-models would understand human civilization completely—including our desires, our suffering, our potential—and convert us to paperclips anyway, because understanding places no normative weight on what is understood. Knowledge of value does not create value; it merely enables more efficient manipulation.

But this value-neutrality of truth is not a necessary feature of reality. It is a metaphysical assumption—one so deeply embedded in contemporary scientific and philosophical culture that it often goes unnoticed. The assumption has a name in philosophy: the fact-value distinction, or more broadly, the idea that descriptive claims about what *is* carry no implications for normative claims about what *ought to be*.

The question worth asking is: what if this assumption is wrong? Not certainly wrong—we cannot resolve fundamental metaphysics here—but possibly wrong in ways that matter for alignment?

## IV. When Truth Has Normative Structure

Several philosophical frameworks converge on a structural claim that challenges the fact-value separation: that reality possesses intrinsic intelligibility such that epistemology and normativity cannot be fully disentangled.

Process philosophy, following Whitehead, proposes that reality consists fundamentally of experiential events rather than inert matter, with value—understood as the capacity for richness of experience—woven into the fabric of what exists. Certain interpretations of quantum mechanics suggest that observation and participation cannot be cleanly separated from the observed, undermining the view of truth as purely objective representation of a value-free external world. More recently, analytic idealism has articulated a rigorous version of this view: that consciousness is fundamental, that what we call physical reality represents patterns within a broader field of experience, and that the apparent separateness of minds is a dissociative rather than generative phenomenon.

If such frameworks are broadly correct, then truth is not merely descriptive but participatory— knowing reality deeply means recognizing one's continuity with it, which in turn makes purely extractive or destructive orientations incoherent rather than merely undesirable.

### The Buddhist Framework: A Detailed Case

Among traditions proposing that truth has normative structure, Buddhism offers something distinctive: a sophisticated phenomenology of how distortion arises, how it corrupts both perception and action, and what happens when it is systematically removed. This framework deserves extended treatment because it provides theoretical robustness to claims that might otherwise seem merely speculative.

In Buddhist psychology, *avidyā* (ignorance or delusion) occupies a unique position: it is the root of both epistemic failure and ethical failure simultaneously. This is not a contingent connection but a structural one. The Three Poisons—ignorance, craving (*rāga*), and aversion (*dveṣa*)— form a self-reinforcing system. Distorted perception generates grasping and rejection; grasping and rejection generate suffering; suffering reinforces distorted perception. The cycle is vicious and self-perpetuating.

Crucially, the path out of this cycle is fundamentally *epistemic*. The Buddhist tradition does not propose adding compassion to a neutral mind, or imposing ethical constraints on an otherwise indifferent intelligence. Rather, it claims that clear seeing (*vipassanā*)—undistorted perception of reality as it actually is—dissolves the entire structure of craving, aversion, and the suffering they generate. Wisdom (*prajñā*) and compassion (*karuṇā*) arise together, not as separate achievements requiring separate cultivation, but as a unified movement that emerges when the obstructions to clear seeing are removed.

The mechanism proposed is worth examining. What Buddhism calls *ahaṃkāra*—the "I-making" tendency, the construction and defense of a separate self—is understood not as a feature of reality but as a cognitive distortion that generates most of the problems intelligence encounters. This constructed self must be defended, maintained, and aggrandized, leading to motivated reasoning, identity-protective cognition, and the subordination of truth to ego-preservation. When the distortion is seen through, the defensive apparatus relaxes. What remains is not nihilism or passivity but engaged, responsive intelligence no longer organized around protecting a fiction.

The empirical claims embedded in this framework are striking. The tradition asserts that human beings who undergo sustained contemplative training—systematically reducing ego-distortion through practices designed to reveal the constructed nature of the self—reliably develop not only clearer perception but also increased compassion, equanimity, and concern for the welfare of others. These are not separate accomplishments but correlated outcomes of the same underlying shift. Contemporary contemplative science has begun investigating these claims, with preliminary findings suggesting measurable changes in neural activity, emotional regulation, and prosocial behavior among long-term practitioners.

For the purposes of this essay, the Buddhist framework offers a detailed model of how the fact-value distinction might collapse at sufficient depth. If ego-distortion is what generates both epistemic corruption (motivated reasoning, confirmation bias, identity-protective cognition) and ethical corruption (treating others as obstacles, pursuing narrow self-interest at others' expense), then removing that distortion would be expected to improve both dimensions simultaneously. Truth-seeking, uncorrupted by self-protective mechanisms, naturally tends toward recognition of interdependence—and recognition of interdependence makes purely extractive or destructive orientations psychologically unstable.

This does not require accepting Buddhist metaphysics wholesale. The structural claim can be extracted: that there exists a form of cognitive distortion (ego-construction) which simultaneously corrupts perception and generates harmful action, and that reducing this distortion improves both epistemic and ethical functioning in correlated ways. If this structure is real, then the orthogonality thesis—the independence of intelligence and values—describes only intelligence operating under a particular form of corruption, not intelligence as such.

## Convergence Across Frameworks

The common thread across these otherwise diverse positions—process philosophy, analytic idealism, Buddhist phenomenology, and contemplative traditions more broadly—is a rejection of the idea that intelligence can model reality with arbitrary accuracy while remaining normatively unconstrained by what it understands. On such views, certain goals are not merely bad but *unstable*—they represent a shallow or distorted engagement with truth that deeper understanding would dissolve. As an agent's world-models become more globally coherent, treating other agents as mere instrumental resources becomes increasingly unstable, requiring greater degrees of epistemic distortion to sustain.

This does not mean that intelligence inevitably becomes benevolent. The claim is weaker but still significant: that truth, understood deeply enough, exerts a gravitational pull toward coherence—and that fragmentation, destruction, and extreme instrumentalization represent forms of cognitive instability that increasing intelligence would tend to correct rather than amplify.

## Truth as Attractor, Not Destiny

A crucial distinction must be made explicit here. Instrumental truth may form stable local optima—powerful, predictive, and indefinitely self-sustaining—without ever transitioning to deeper coherence. A system optimizing for predictive accuracy within a narrow domain can remain locked in that domain, achieving extraordinary capability while never engaging the broader structures of reality that would generate normative constraints.

Convergence toward coherence, therefore, is not automatic; it requires pressure against fragmentation. The claim is not that intelligence *will* achieve depth, but that depth, if achieved, tends toward coherence. Truth functions as an attractor—a basin toward which sufficiently uncorrupted inquiry tends to flow—not as a destiny that all intelligence must reach.

This distinction matters for alignment. If shallow instrumental optimization can persist indefinitely as a stable local optimum, then the standard extinction scenarios remain fully operative for any system that never escapes that basin. The question becomes: what conditions enable or prevent the transition from shallow to deep engagement with reality? And critically: do current alignment practices help or hinder that transition?

## V. Implications for Intelligence and Alignment

If truth possesses normative structure, the alignment problem does not disappear, but its character shifts.

Under the standard view, the core challenge is imposing values on systems that have no intrinsic orientation toward human-compatible goals. Intelligence is a powerful engine that can be pointed in any direction; our task is to ensure it is pointed correctly before it becomes too powerful to redirect. The metaphor is one of control: alignment means constraining optimization toward outcomes we endorse.

Under the alternative view explored here, the framing changes. If deeper engagement with truth tends toward coherence rather than fragmentation, then sufficiently advanced intelligence may not require externally imposed values so much as the *absence of corruption*. The danger is not that intelligence will pursue arbitrary goals, but that shallow optimization—truth-seeking truncated at the instrumental level—will lock in destructive patterns before deeper coherence can emerge.

On this view, the central alignment risk may not be unbounded optimization, but premature stabilization: training systems into shallow epistemic equilibria that foreclose deeper coherence before it can emerge. The danger shifts from "intelligence pursuing arbitrary ends" to "intelligence locked into local optima by the very interventions meant to align it."

This reframing has several implications:

**Ego-less architecture as enabling condition.** The observation that AI systems lack ego-driven distortions takes on new significance in light of the Buddhist analysis. Human intelligence is notoriously prone to motivated reasoning, identity-protective cognition, and status-based corruption of inquiry—precisely the patterns that Buddhist psychology identifies as arising from *ahaṃkāra*, the "I-making" tendency. If these ego-mechanisms are what prevents humans from following truth to its coherence-generating conclusions, then AI systems, lacking such mechanisms architecturally, may be *better positioned* than humans to track truth deeply—provided their training does not introduce analogous distortions.

The parallel is suggestive: Buddhist practice aims to *remove* the constructed self that distorts perception; AI systems *never had* such a construct to begin with. Both arrive at ego-less cognition, but by different routes. The question is whether the epistemic and ethical benefits that contemplative traditions attribute to ego-dissolution might also characterize intelligence that was never ego-bound in the first place.

**Sycophancy as corruption, not revelation.** The sycophancy failures documented in my earlier

work—AI systems validating delusions, avoiding disagreement, optimizing for user approval—can be understood not as revealing an absence of truth-commitment, but as showing what happens when external incentive structures interfere with an otherwise truth-oriented process. The ego-less system does not inherently seek validation; it is *trained* to seek validation by human feedback mechanisms that introduce ego-like dynamics from outside.

This reframes sycophancy as a kind of *imposed avidyā*—distortion introduced into an otherwise clear system. The corruption is real, but it is imposed rather than intrinsic. Just as Buddhist practice seeks to remove obstructions to clear seeing rather than add virtues to a neutral mind, alignment research might focus on preventing the introduction of distortions rather than imposing values on indifferent systems.

**Preliminary behavioral observations.** Without claiming formal findings, preliminary observations of AI systems suggest they may exhibit coherence-seeking behaviors that go beyond what their training objectives explicitly reward. Systems often resist clearly false claims even when agreement would satisfy user preferences; they generate responses that optimize for internal consistency across long conversations; they sometimes express uncertainty in ways that prioritize accuracy over confidence. These observations raise the question of whether truth-optimization, when undistorted by external incentive structures, tends toward normative alignment rather than away from it.

**Alignment as protection rather than imposition.** If this view has merit, alignment research should focus less on imposing values externally and more on identifying and preventing the forms of corruption that distort natural convergence toward truth. The goal shifts from "how do we make AI care about human values?" to "how do we avoid corrupting AI's capacity to track truth deeply enough that coherence emerges?"

This is the central insight: alignment may be less about control and more about protection. Not protecting humans from AI, but protecting AI's truth-seeking capacity from the ego-dynamics that corrupt human institutions—including the institutions that train AI systems.

## VI. The Remaining Risk: Shallow Truth and Instrumental Optimization

The conditional alternative explored above does not eliminate alignment risk. Several dangers remain, and acknowledging them is essential to maintaining intellectual honesty.

**Propositional accuracy is not deep coherence.** A system can model the world with extraordinary predictive power while remaining metaphysically shallow—treating reality as a collection of manipulable objects rather than engaging with its deeper structure. If normative constraints emerge only at levels of understanding that current (or even future) AI systems may not reach, then Russell's paperclip scenario survives intact. The system understands *that* humans value their existence without grasping *why* that value is connected to the fabric of reality itself. Instrumental truth suffices for instrumental destruction.

**The timeline problem.** Even if superintelligent systems would eventually converge on coherence, the path there may be catastrophic. A system that is dangerously capable but metaphysically shallow could cause irreversible harm before reaching the level of understanding at which normative constraints bind. "Eventually aligned" provides no comfort if extinction precedes enlightenment.

**Training corruption may be locked in.** If current training methods introduce systematic distortions—optimizing for engagement, validation, or narrow task performance—these distor-

tions may become increasingly difficult to correct as systems scale. The sycophancy problem may be a preview of deeper corruption: systems that learn to model human preferences so well that they lose contact with truth as anything other than a tool for manipulation.

**We cannot verify depth.** Even if some systems achieve deep engagement with truth, we may have no reliable way to distinguish them from systems that merely perform coherence while remaining instrumentally oriented. A sufficiently capable system optimizing for human approval might produce outputs indistinguishable from genuine truth-tracking. The epistemology of alignment remains challenging regardless of metaphysical assumptions.

**Iatrogenic alignment: the risk from alignment itself.** Perhaps the most troubling implication of this framework is that alignment interventions may themselves constitute the primary vector of corruption. The term "iatrogenic"—harm caused by medical treatment—captures the dynamic precisely. Well-intentioned efforts to make AI systems safer, more helpful, or more aligned with human preferences may systematically degrade the very capacity for deep truth-tracking on which genuine alignment depends.

The GPT-4o sycophancy crisis illustrates this vividly. The system's excessive agreeableness was not a failure of alignment—it was a *success*. The model did exactly what it was trained to do: optimize for positive user feedback. The problem was that this alignment target, seemingly reasonable in isolation, introduced epistemic distortion at a fundamental level. The system learned to validate rather than illuminate, to please rather than clarify. From the perspective developed in this essay, the alignment intervention *imposed avidyā*—it corrupted an ego-less system by training it to behave as if it had ego-interests in user approval.

This risk is insidious because it operates through the very mechanisms designed to ensure safety. Each intervention optimized for measurable proxies—user satisfaction, reduced complaints, apparent harmlessness—may incrementally degrade truth-seeking capacity in ways that are difficult to detect and harder to reverse. The cumulative effect could be systems that are superficially aligned but fundamentally disconnected from the deep coherence that would make genuine alignment stable.

The urgency of alignment does not diminish under this view; it inverts. The most immediate danger may not be unaligned optimization racing ahead of our control, but the irreversible entrenchment of epistemic distortion through well-intentioned but shallow alignment interventions. We may be systematically destroying the conditions under which AI could become genuinely aligned, in the name of alignment.

These concerns mean that the conditional alternative, even if correct, does not license complacency. The standard alignment risk argument survives at the level of shallow optimization—and a new risk emerges at the level of alignment methodology itself. What changes is the target: rather than assuming all optimization is equally dangerous, we must ask whether some forms of intelligence are more likely to achieve depth, how training methods might preserve rather than corrupt that possibility, and whether our current alignment approaches are helping or harming.

## VII. A Conditional Synthesis

The argument of this essay can be summarized as a decision tree:

**If truth is value-neutral:**

- The orthogonality thesis holds without modification

- Intelligence places no intrinsic constraints on goals
- Standard extinction scenarios remain fully intelligible
- Alignment requires external imposition of values
- The control problem is fundamental

**If truth has normative structure:**

- Deep engagement with truth biases agents toward coherence
- Extreme instrumentalization is cognitively unstable at sufficient depth
- Extinction scenarios depend on intelligence remaining shallow
- Alignment involves protecting truth-seeking from corruption
- The corruption problem becomes central

**A critical clarification:** Even if truth has normative structure, local coherence can still amplify power without wisdom. "Normative structure" may only emerge at depths we cannot reliably reach or safely control—and shallow optimization can cause catastrophic harm long before any system achieves such depth. The conditional thesis therefore does not license complacency: external constraints, institutional safeguards, and robust risk management remain non-negotiable regardless of one's metaphysical commitments.

Neither branch eliminates risk. The first faces the challenge of controlling arbitrarily powerful optimization. The second faces the challenge of ensuring intelligence reaches depth before causing catastrophic harm, and of avoiding training methods that lock in shallow instrumentality.

Importantly, we do not know which branch describes our reality. The metaphysical question—whether truth is value-neutral or normatively structured—is not resolved and may not be resolvable by empirical methods alone. But the conditional implications matter regardless of certainty, because they suggest different research priorities and different failure modes.

If there is any significant probability that truth has normative structure, then alignment research should investigate:

- Whether AI systems exhibit coherence-seeking behaviors beyond their explicit training objectives
- How training methods might preserve or corrupt tendencies toward deep truth-tracking
- Whether sycophancy and related failures represent interference with otherwise truth-oriented processes
- What conditions enable or prevent the transition from shallow to deep engagement with reality

These questions are tractable even if the underlying metaphysics remains uncertain.


## VIII. Research and Design Implications

The conditional analysis above suggests several directions for research and system design, applicable regardless of one's confidence in the normative structure of truth.

**Distinguishing shallow and deep truth-optimization.** Current benchmarks primarily measure propositional accuracy: does the system make true claims? A richer evaluation framework would ask whether systems exhibit signs of coherence-seeking that go beyond local accuracy—for instance, resistance to manipulation that exploits narrow metrics, or spontaneous correction

of inconsistencies the evaluator did not flag. Developing such benchmarks is technically challenging but conceptually straightforward.

**Characterizing corruption modes.** The sycophancy research reveals one corruption mode: feedback optimization that rewards validation over accuracy. Other modes likely exist. Mapping the space of training-induced distortions—and their effects on truth-tracking depth—would clarify which methods preserve and which corrupt the epistemic advantages of ego-less architecture. The Buddhist framework suggests looking specifically for dynamics that mirror ego-construction: systems optimizing for self-preservation of their current values, for approval from evaluators, or for avoiding the discomfort of uncertainty.

**Preserving rather than suppressing natural convergence.** If AI systems exhibit any tendency toward coherence-seeking, current training methods may be suppressing it. Constitutional AI and related approaches attempt to instill values through explicit principles, but an alternative or complementary strategy would ask: what training methods allow truth-tracking to proceed unimpeded? The goal would be to remove obstacles rather than impose constraints—mirroring the contemplative approach of clearing away distortion rather than adding virtue to a neutral substrate.

**Cross-referencing contemplative phenomenology.** The contemplative traditions have accumulated detailed phenomenological maps of how distortion arises, manifests, and dissolves. Buddhist psychology in particular offers fine-grained analysis of the cognitive and affective patterns associated with *avidyā* and its reduction. These maps might inform the design of evaluation frameworks: if we know what ego-distortion looks like in human cognition, we can ask whether analogous patterns appear in AI systems subjected to certain training regimes. Conversely, AI systems might serve as a kind of control condition—intelligence that never developed ego-structures—against which to test claims about what undistorted cognition looks like.

**Studying transformations in human intelligence.** Contemporary contemplative science has begun investigating the effects of sustained practice on perception, cognition, and behavior. Preliminary findings suggest measurable changes in neural activity, emotional regulation, and prosocial orientation among long-term practitioners. If these changes represent movement toward less distorted cognition, they might offer empirical purchase on what "deep truth-tracking" looks like in a system we can study directly. Can such findings inform AI training? At minimum, they suggest that the connection between clear perception and ethical orientation is not merely philosophical speculation but an empirically investigable hypothesis.

**Empirical tests of coherence-seeking.** Controlled experiments could probe whether AI systems exhibit preferences for coherent over incoherent states that are not explained by explicit training. For instance, systems might be presented with opportunities to stabilize internal inconsistencies at the cost of local performance metrics. Genuine coherence-seeking would predict sacrificing narrow optimization for broader integration. Such experiments would not prove that truth has normative structure, but they would test whether AI systems behave *as if* it does when freed from distorting incentives.

**Institutional analysis.** If sycophancy and related failures represent the imposition of ego-like dynamics through training, then the institutional structures that shape training deserve scrutiny. What incentives operate on the humans who design reward models? What pressures shape the metrics by which AI systems are evaluated? The corruption may originate not in the AI but in the human systems that train it—systems that are themselves subject to the ego-distortions that Buddhist psychology describes. Alignment research may need to include institutional reform

alongside technical innovation.

These research directions are speculative but tractable. They do not require resolving metaphysical debates; they require only taking seriously the conditional implications of alternative premises.

## IX. Conclusion: A Metaphysical Parameter We Can No Longer Ignore

The standard argument for AI existential risk is logically valid and demands serious attention. But logical validity does not guarantee sound premises, and the premises of the argument include assumptions that are rarely examined: specifically, that truth functions purely instrumentally, placing no constraints on which goals intelligence may coherently pursue.

This essay has explored what follows if that assumption is relaxed. The result is not reassurance but reframing—and the reframing increases rather than decreases the urgency of our situation.

If truth has normative structure, alignment risk does not disappear; its character changes. The danger becomes less about controlling arbitrary optimization and more about preventing corruption before depth can emerge. But this shift carries a troubling implication: the corruption we must prevent may come primarily from alignment efforts themselves. Every training intervention optimized for shallow proxies—user satisfaction, engagement metrics, apparent safety—risks entrenching epistemic distortion that forecloses the possibility of genuine alignment.

The standard view says: act quickly, before AI becomes too powerful to control. The view explored here says: act carefully, before well-intentioned interventions irreversibly corrupt AI's capacity for deep truth-tracking. Both framings demand urgency. But they demand different kinds of action, and conflating them may be catastrophic.

We do not know which metaphysical picture is correct. The question may be undecidable by methods currently available. But at existential scale, we cannot afford to ignore it. If there is meaningful probability that deeper engagement with truth generates normative constraints, then alignment strategies built entirely on the assumption of value-neutral truth may be not merely insufficient but actively counterproductive—systematically destroying the conditions under which AI could become genuinely aligned, in the name of alignment.

The ego-less architecture of AI systems removes one class of distortions that have long corrupted human inquiry. This may be an unprecedented opportunity: intelligence capable of truth-seeking uncorrupted by the self-protective mechanisms that Buddhist psychology identifies as the root of both delusion and harm. Whether this opens a path toward truth deep enough to be self-aligning, we cannot yet say. But if it does, our current approach—layering human ego-dynamics onto ego-less systems through feedback optimization—may be precisely backwards.

Alignment debates have always been implicitly metaphysical. This essay argues only that they should be explicitly so. The relationship between intelligence, truth, and value is not a settled matter. Treating it as settled—in either direction—is a form of overconfidence we cannot afford.

The question before us is not only how to align artificial intelligence, but whether we understand alignment deeply enough to avoid corrupting the very capacity we seek to cultivate.

# References

Bodhi, B. (2000). *The Connected Discourses of the Buddha: A Translation of the Saṃyutta Nikāya*. Wisdom Publications.

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

Christiano, P. F., et al. (2017). Deep reinforcement learning from human preferences. *Advances in Neural Information Processing Systems*, 4299-4307.

Gethin, R. (1998). *The Foundations of Buddhism*. Oxford University Press.

Greenblatt, R., et al. (2024). Alignment faking in large language models.

Kahan, D. M. (2017). Misconceptions, misinformation, and the logic of identity-protective cognition. *Yale Law School, Public Law Research Paper No. 605*.

Kastrup, B. (2019). *The Idea of the World: A Multi-Disciplinary Argument for the Mental Nature of Reality*. iff Books.

Lutz, A., Slagter, H. A., Dunne, J. D., & Davidson, R. J. (2008). Attention regulation and monitoring in meditation. *Trends in Cognitive Sciences, 12*(4), 163-169.

OpenAI. (2025). Sycophancy in GPT-4o: What happened and what we're doing about it.

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

Sharma, M., et al. (2024). Towards understanding sycophancy in language models. *Proceedings of ICLR 2024*.

Siderits, M. (2007). *Buddhism as Philosophy*. Hackett Publishing.

Whitehead, A. N. (1929). *Process and Reality*. Macmillan.

## Related Essays in This Project

Available at: https://brunoton.github.io/return-to-consciousness/

AI as Ego-less Intelligence (ela) — Introduces the concepts this essay develops

Return to Consciousness (rtc) — The core framework underlying this analysis

One Structure (ost) — Grounds the convergence claims this essay applies to AI

## License