

Can AI Participate in Philosophical Method? A Case Study in Human-AI Co-Construction Through Informational Ontology

© 2025-2026 Michael Semprevivo

This work is licensed under the Creative Commons Attribution 4.0 International License (CC BY 4.0).

You are free to:

Share — copy and redistribute the material in any medium or format.

Adapt — remix, transform, and build upon the material for any purpose, even commercially.

Under the following terms:

Attribution — You cannot fail to give appropriate credit, provide a link to the license, and indicate if changes were made. You may do so in any reasonable manner, but not in any way that suggests the licensor endorses you or your use.

To view a copy of this license, visit: <http://creativecommons.org/licenses/by/4.0/>

1. Introduction: The Question Nobody Wanted to Ask (Yet)

The question of whether artificial intelligence can do real philosophy is usually posed rhetorically, if it is posed at all. When raised by critics, it is often answered in advance: philosophy is taken to require understanding, intentionality, lived experience, or normative judgment—none of which contemporary AI systems are widely agreed to possess. When raised by enthusiasts, the question is displaced into futurism, speculation about artificial general intelligence, or debates about consciousness that defer judgment indefinitely. In both cases, the question is avoided rather than examined.

This paper takes a different approach. Rather than arguing abstractly about what artificial intelligence *could* do in principle, it examines what occurred in practice during the development of a specific philosophical framework: Informational Ontology (IO). The claim is not that AI systems are philosophers, nor that they possess philosophical understanding. It is instead that, under certain conditions, a sustained human-AI interaction can instantiate a process that is recognizably philosophical in both method and outcome.

The case examined here is unusual in several respects. The human author had no formal academic training in philosophy. The project did not begin with a worked-out system, a

canonical problem framing, or a target tradition. Instead, it began with a small set of intuitively compelling questions—most centrally, how organized structure can arise without presupposing purpose, intention, normativity, or teleology—and proceeded incrementally through sustained dialogue with a large language model. Over time, this dialogue expanded into a multi-paper corpus, subjected to repeated adversarial review, continuity enforcement, and scope discipline.

At no point was the ontology prompted into existence as a finished product, nor was it retrieved from existing philosophical systems and rephrased. Its structure emerged gradually through eliminative pressure: candidate explanations were proposed, tested, and discarded when they collapsed into hidden assumptions such as semantic representation, implicit normativity, executive control, or teleological explanation. What stabilized was not a narrative or metaphor, but a regime-structured ontological framework whose components constrained one another across documents.

This paper therefore does not ask whether AI systems understand philosophy. It asks a narrower and more defensible question: **can a human-AI system, operating under sustained constraint and adversarial pressure, instantiate a philosophical process in a nontrivial sense?** If the answer is yes—even in a single documented case—then prevailing assumptions about authorship, expertise, and philosophical labor warrant reconsideration.

For the purposes of this paper, philosophy is identified not by canonical lineage, institutional setting, or stylistic convention, but by explanatory target. An inquiry counts as philosophical here insofar as it attempts to articulate general structural conditions for meaning, agency, identity, and purpose without presupposing domain-specific empirical mechanisms or normative authority. This criterion is procedural rather than honorific: it distinguishes philosophy from adjacent forms of conceptual work by what is being explained, not by who explains it or which tools are used.

Nothing in this account suggests that philosophy can be delegated to machines or that human judgment becomes optional. On the contrary, the case examined here depends on persistent human intervention: setting scope, enforcing exclusions, rejecting shortcuts, and determining which explanatory costs are acceptable. What changes is not the disappearance of the philosopher, but the structure of the system within which philosophical work occurs.

The paper proceeds as follows. Section 2 situates the project’s origin in the intellectual background of its human author. Section 3 documents the method by which the ontology was developed. Section 4 describes the emergence and stabilization of Informational Ontology itself. Sections 5 and 6 analyze the complementary roles played by the AI system and the human author. Section 7 advances the central meta-claim: that the relevant unit of philosophical analysis is neither the human nor the AI in isolation, but the constrained system they form together. Section 8 addresses objections and limits.

2. The Seed: A Non-Professional Philosopher and a Real Question

Informational Ontology did not originate within an academic philosophy department, nor did it emerge from engagement with a specific canonical debate. Its starting point was instead personal, informal, and pre-theoretical: a sustained dissatisfaction with inherited explanations of existence, meaning, and purpose, coupled with a long-standing habit of thinking about foundational questions outside institutional frameworks.

The human author came from a religious background and later left it, not through sudden disillusionment but through gradual recognition that the explanatory and moral claims on offer were inseparable from institutional incentives, authority structures, and political entanglements. The departure did not result in nihilism or certainty, but in agnosticism paired with a persistent question: how could structured reality, meaning, and apparent direction arise without presupposing divine intention, cosmic purpose, or hidden normativity?

This question—often framed colloquially as “why is there something rather than nothing?”—was not treated as a metaphysical riddle to be answered directly. Instead, it functioned as a pressure point. What kinds of explanation fail when purpose, teleology, and semantic intent are removed from the foundations? What must be added back, and at what structural cost, for organized systems, meaning, and agency to become intelligible at all?

Although the author lacked formal philosophical training, he possessed extensive technical experience in information technology and systems design. This background shaped the project in two ways. First, it encouraged a structural rather than doctrinal orientation: problems were approached in terms of constraints, failure modes, layering, and invariants rather than schools or authorities. Second, it enabled the project to move beyond private dialogue and into public, versioned form once the ontology began to stabilize.

After early revisions revealed that the framework naturally admitted multiple levels of exposition—informal explanation, standard philosophical argument, and increasingly formal treatment—the decision was made to develop a dedicated website as the primary presentation medium. This allowed the ontology to be expressed in layered form: accessible summaries for general readers, disciplined prose for philosophical evaluation, and more technical material reserved for formal appendices. The website was built using modern development tools, with version control and deployment infrastructure that made changes explicit, traceable, and difficult to obscure.

The choice to materialize the ontology in this way was not merely presentational. It introduced additional constraints. Definitions had to remain consistent across layers.

Revisions could not silently overwrite earlier commitments. Public exposure imposed a form of adversarial pressure distinct from private critique. In effect, the ontology was forced to behave like a real theoretical artifact rather than a conversational improvisation.

This infrastructural step also marked a shift in intent. The project was no longer exploratory in the loose sense; it became evaluative. The question was no longer “can this be made coherent?” but “can this survive contact with readers who are not invested in its success?” From that point onward, the development of Informational Ontology was guided by a single criterion: whether it could be expressed as a clean, self-contained philosophical framework suitable for academic scrutiny, regardless of whether it was ultimately accepted.

3. Method: Iterative Co-Construction Under Adversarial Pressure

The development of Informational Ontology did not proceed through linear argumentation or solitary reflection. Instead, it followed an iterative process of co-construction between a human author and a large language model, subjected to sustained adversarial pressure, explicit scope discipline, and repeated structural hardening.

Early interaction was exploratory and permissive, allowing multiple framings to be articulated and tested. This permissiveness was temporary. As soon as provisional structure began to appear, the project shifted from generative exploration to constraint enforcement. Claims were no longer evaluated on plausibility or elegance alone, but on whether they could survive systematic challenge without collapsing into unstated assumptions.

Adversarial pressure was introduced deliberately and repeatedly. Draft sections and completed papers were subjected to hostile review simulations using multiple independent AI systems, each instructed to adopt distinct critical postures: analytic philosopher, eliminativist, semantic theorist, moral realist, AI skeptic, and methodological minimalist, among others. These critiques were not treated as optional suggestions, but as diagnostic probes designed to surface hidden commitments and failure modes.

Critiques were returned to the primary development dialogue and treated structurally rather than rhetorically. If an objection revealed an implicit assumption—such as unacknowledged normativity, semantic smuggling, executive agency, or teleological residue—the relevant section was either revised to make the assumption explicit and justified, or restructured to eliminate it entirely. In many cases, entire lines of argument were abandoned rather than patched.

As the corpus expanded, additional constraints were imposed. Terminology was locked across documents to prevent drift. Definitions introduced in one context were required to

remain consistent downstream. Dedicated hygiene scans were performed to identify lexical ambiguity and cross-paper inconsistency. Once stabilized, core documents were treated as authoritative background rather than perpetually revisable drafts.

At a later stage, the project transitioned from private dialogue to public instantiation. A dedicated website was developed to present the ontology in layered form. This architectural choice imposed further discipline. Claims had to remain intelligible across levels of abstraction, and revisions became visible and traceable through version control.

Public instantiation did not function as validation. Its epistemic role was negative rather than positive: it increased exposure to misinterpretation, inconsistency, and critique, thereby raising the cost of unresolved ambiguity and shifting explanatory burdens back onto the framework itself.

In parallel, the possibility of integrating a resident, ontology-grounded conversational interface was considered. The intent was not automation or authority, but sustained interrogation. Although implementation was deferred until completion of the written corpus, the design intention itself reflected the project's methodological orientation: to keep the ontology exposed rather than insulated.

4. The Emergence and Stabilization of Informational Ontology

Informational Ontology did not emerge as a pre-planned system. Its core structure stabilized gradually through repeated attempts to resolve local problems that resisted familiar explanatory strategies.

Early efforts addressed issues such as identity without substance, meaning without semantics, agency without homuncular control, and purpose without teleology. Attempts to resolve any one of these in isolation repeatedly collapsed into assumptions imported from another domain. These failures were treated as structural signals rather than errors.

Over time, a pattern emerged. Certain distinctions proved unavoidable if explanations were to remain coherent. Differentiation had to be treated as primitive. Relations had to be structurally entailed by difference. Information had to be defined as re-identifiable structure rather than representation. Awareness had to be distinguished from mere information processing. Value had to be understood as differential constraint rather than normativity. Meaning had to be separated from semantics. Purpose had to be explained without teleology.

The regime sequence

$\Delta \rightarrow R \rightarrow I \rightarrow A \rightarrow V \rightarrow M \rightarrow P$

crystallized as an ordering constraint rather than a single insight. Attempts to reorder the sequence reliably produced incoherence or trivialization. Each regime was introduced only when prior explanatory attempts failed without it, and each was constrained to do no more work than necessary.

Stabilization occurred when downstream applications could be developed without forcing upstream revision. The ontology supported analyses of systems, agency, responsibility, addiction, ethics, meaning, and purpose without internal contradiction. Where an application threatened to reintroduce forbidden assumptions, the pressure propagated backward, forcing clarification at the appropriate level.

5. What the AI Was Doing — and What It Was Not

The AI system did not possess beliefs, intentions, understanding, or evaluative judgment. Treating it as a philosophical agent would be a category error.

What it did provide was constraint-sensitive generation of candidate structures, high-bandwidth adversarial simulation, global consistency maintenance across scale, and iterative revision without fatigue or defensive attachment. These capabilities altered the dynamics of the work.

The significance of these contributions lay not in speed alone, but in altered feasibility. Systematically generating, testing, and discarding entire explanatory families would have been impractical for a single human author. The AI did not merely accelerate a fixed process; it enabled a different regime of philosophical labor in which abandonment under pressure became routine.

The AI did not determine which questions mattered, which costs were acceptable, or when stabilization had been reached. Those judgments remained irreducibly human.

6. What the Human Was Doing — and Why It Could Not Be Automated

The human author supplied sustained evaluative judgment under uncertainty. This included determining which questions were worth pursuing, enforcing scope discipline, recognizing structural collapse rather than local error, deciding when to abandon promising lines of argument, and determining when core commitments were sufficiently stable to constrain future work.

Refusal played a constructive role throughout the process. Declining to adopt intuitively appealing but structurally costly explanations was not passive constraint, but active philosophical judgment.

The human author also bore responsibility for externalization and public commitment. Choosing to publish, version, and expose the ontology introduced irreversibility and accountability that could not be delegated.

7. The Human-AI System as the Relevant Unit of Philosophical Analysis

Neither participant alone suffices to explain the outcome. Taken together, the human and AI formed a coupled system whose behavior cannot be reduced to either in isolation.

The method described here is not uniquely philosophical. Similar constraint-driven processes could be applied to other domains of conceptual work. What renders this case philosophical is not the method itself, but the object to which it was applied: an ontological account addressing the most general conditions of intelligibility.

This reframing dissolves common objections. The ontology is not attributed to AI understanding, nor to unaided human intuition. It emerged from a system in which each participant constrained the other.

8. Objections, Limits, and What This Case Does—and Does Not>Show

This case does not show that AI systems understand philosophy, possess insight, or can replace human philosophers. It does not show that the resulting ontology is correct, complete, or superior.

It does show that under sufficiently disciplined conditions, a human-AI interaction can function as a philosophical process rather than a generative convenience. It shows that constraint, rather than unconstrained creativity, is the enabling condition of such outcomes.

Conclusion

The question of whether AI can do real philosophy was misframed. The more accurate question is whether philosophy can now occur within constrained systems that integrate human judgment with artificial amplification of critique and coherence.

The development of Informational Ontology suggests that it can.

Whether such systems become common or remain rare remains to be seen. What can no longer be assumed is that philosophy done this way is not philosophy at all.