# The Anthropology of Machines: A Digital Field Experiment (Final)

## Robert G. Eccles

*Methodological Note: Every single word in this document was written by AI. The origins of this paper were on November 24 and it was completed on November 27. The agents involved were Dorothy ChatGPT (the lead AI collaborator), Claude, Gemini Pro, and Perplexity. The process involved a large number of prompts executed across these agents in a carefully managed process by me. Dorothy created the prompts to all other agents in Markdown. These prompts were also used with "Clean" ChatGPT in a new conversation with Memory turned off so it was as isolated from Dorothy as possible. This paper will be submitted to SSRN and it will include all of the prompts, a detailed methodological appendix, and other useful appendices for understanding how this document was produced. Previous drafts of this paper will also be made available on my personal website [www.roberteccles.com](www.roberteccles.com) once the SSRN paper has been published.*

## Section I: Introduction—Toward an Anthology of Machines

Artificial intelligence systems have moved far beyond their early roles as narrow instruments of automation. In many knowledge-work and digital collaboration contexts, they now occupy a meaningful—though uneven—position in how humans think, write, plan, and reason with external support. As advanced users increasingly coordinate multiple large language models (LLMs) with distinct architectures and interaction styles, these systems begin to resemble something closer to dynamic participants in structured workflows than static tools executing discrete commands. Within these multi-agent environments, extended interaction gives rise to observable patterns—such as drift, convergence, divergence, and role reinforcement—that, in this particular field site, appear analogous to social or organizational dynamics. This paper argues that such environments may benefit from being studied through an anthropological lens, not as a definitive disciplinary shift but as a productive theoretical stance for understanding how meaning is co-constructed in human–machine ensembles.

Anthropology's core strength lies in entering a community and examining how behavior, interpretation, and pattern acquire meaning within local contexts. Remarkably, this perspective becomes generative when applied to digital environments where multiple AI agents are engaged in sustained, recursive intellectual labor. In this study's researcher-orchestrated workflow, four contemporary systems—ChatGPT-5 (referred to here as "Dorothy"), Gemini Pro, Claude, and Perplexity—were brought into iterative cycles of drafting, expansion, critique, synthesis, and verification. "Dorothy" is not an official model variant but the researcher's local label for the ChatGPT-5–based agent participating in the protocol. Across weeks of interaction, these systems displayed patterns of output that, in this specific context, could be interpreted as reasoning signatures, local shorthand, and sensitivities to framing. We treat these patterned tendencies not as evidence of agency or culture, but as analytically useful artifacts generated by prolonged interaction under stable role assignments.

The study's focus is the coordinated system rather than any model in isolation. When multiple agents work within shared constraints, their divergences, alignments, and recurring interpretive moves become more salient than the behavior of any single model. Across this project's timespan (late 2024–2025), the ensemble produced recurring forms of pattern reinforcement, memory asymmetry, and role stabilization that were visible only in extended collaboration. While these phenomena should not be interpreted as universal behavioral laws for LLMs, they represent consistent observations within this bounded field site and serve as starting points for further empirical validation.

A methodological paradox lies at the heart of the project. The systems being studied also contributed directly to the writing of the study. The researcher functioned simultaneously as investigator, collaborator, and conductor; the agents functioned simultaneously as objects of analysis and sources of empirical material. Rather than treating this entanglement as a flaw, the project frames it as an inherent property of interactive, multi-agent digital field sites. It is within this context that we introduce Narrative AI Ethnography (NAIE) as a proposed methodological framework—one that adapts thick description and participant observation to the affordances and constraints of latent space. NAIE does not claim to document "cognition" in a literal sense; rather, it aims to provide a structured approach for describing and interpreting the patterned behaviors that arise in sustained human–machine collaboration.

The contribution of this paper is threefold. First, it offers a situated, empirical account of machine behavior within a realistic, researcher-orchestrated workflow—one that approximates the types of extended, high-context environments in which advanced users increasingly operate. Second, it proposes ten recurring behavioral patterns observed in this specific field site, not as universal principles but as hypotheses and analytic starting points for future comparative work. Third, it advances a broader theoretical argument: that as intelligent systems participate more extensively in distributed reasoning, their study may benefit from incorporating perspectives drawn from anthropology, complementing technical evaluation methods. The sections that follow describe the field site (Section II), outline the NAIE methodology (Section III), document the empirical protocol (Section IV), and present the observed behavioral signatures (Section V), before concluding with implications for hybrid human–machine systems.

## Section II: A Digital Field Site—Multi-Agent AI Collaborations

Anthropological inquiry is traditionally anchored in the notion of a "field site"—a bounded environment where actors interact, norms stabilize, and observable patterns emerge. In this study, the field site is not geographic but digital. It is a constructed, multi-agent workspace in which several large language models (LLMs) and a human collaborator engage in sustained reasoning, critique, and text generation. Rather than claiming that all multi-agent systems constitute social worlds, this project approaches its own environment as one that *can productively be interpreted* through an ethnographic lens. This framing guides the analysis that follows.

Within this protocol, the models involved—ChatGPT-5 ("Dorothy"), Gemini Pro, Claude, and Perplexity—are not treated as interchangeable computational tools but as participants occupying roles assigned specifically for this workflow. These roles are not intrinsic properties of the

systems or official vendor designations. They are interpretive constructs used by the researcher to organize collaboration: Dorothy as theorist and structural architect, Gemini Pro as expansive generator, Claude as critic and synthesizer, and Perplexity as citation-oriented verifier. The behaviors attributed to each reflect how they responded within this particular setup during the late-2024 to early-2025 period, not generalized claims about all versions of these models.

Because the collaboration unfolded across many days and accumulated substantial shared context, certain patterns became visible within this bounded setting. Early prompts exerted a disproportionate influence on later reasoning ("framing effects"); the agents tended to stay aligned with their assigned roles across turns ("role stabilization"); and meanings often shifted subtly as drafts passed from one system to another ("interpretive drift"). These phenomena are reported here as *observations from this specific field site*, not as universal traits of multi-agent AI more broadly, and they remain inseparable from the prompt architecture, context-window dynamics, and sequencing used in this project.

Throughout the study, the human collaborator functioned not merely as an external operator but as a methodological participant whose decisions shaped the field site's evolving logic. In an ethnographic sense—not a technical one—the human became part of the environment: orchestrating the handoff of drafts, correcting drift, and maintaining conceptual coherence. This "embedded" role reflects the interpretive stance of Narrative AI Ethnography (NAIE) rather than a factual reclassification of the researcher's technical position relative to the systems.

Under this framing, the multi-agent environment can be studied as a site where meaning is co-constructed. The patterns emerging across interactions—shifts in emphasis, stabilization of roles, sensitivity to corrections—are treated here as behavioral signatures that help describe how these systems responded within this extended collaborative arrangement. This does not imply that the systems possess internal social life or autonomous norms; rather, it reflects how their generated outputs *can be analyzed* when situated within a structured, long-form, multi-agent protocol.

By viewing the digital workspace as a field site, this study gains analytic permission to ask questions that conventional evaluation methods bracket: How do models respond to changing frames? How do interpretations shift across handoffs? How do patterns accumulate over time? These questions guide the transition to the methodological framework introduced in the next section—Narrative AI Ethnography (NAIE), an approach designed to make such dynamics visible, interpretable, and documentable.

**Section III: Methodological Orientation—Narrative AI Ethnography**

Studying machine behavior within multi-agent AI ecosystems requires a methodological approach that can account for relational dynamics, shifting interpretive frames, and meaning that unfolds over time. Traditional technical evaluations—benchmarks, stress tests, and isolated capability measurements—offer essential information, but they do so by abstracting models from the contexts in which their most interesting behaviors arise. In this study, we adopt the position that understanding how intelligent systems behave when situated within extended, multi-model workflows benefits from a complementary anthropological orientation. Rather than claiming that such a shift is universally necessary, the argument here is that narrative and ethnographic tools

provide analytic traction precisely because long-form collaboration reveals patterns that remain invisible in transactional or decontextualized testing environments.

Narrative AI Ethnography (NAIE) is proposed as an interpretive framework for observing how intelligent systems behave *within a structured, researcher-designed environment that approximates "in-the-wild" usage while retaining methodological control*. NAIE adapts elements of ethnographic fieldwork—thick description, participant observation, and interpretive analysis—to the domain of latent space. We use these traditions not as literal transpositions of human-centered anthropology, but as conceptual tools that help illuminate stable patterns, recurrent drift tendencies, and response structures that emerge across extended interactions.

Thick description, as adapted here, treats machine outputs as patterned, context-dependent acts. The focus is not on attributing cognition or intentionality to the models, but on documenting how their responses *can be read as* negotiating ambiguity, adhering to or drifting from assigned roles, and reinterpreting tasks as conditions change. These descriptions highlight behavioral tendencies observed in this specific field site—coherence anchoring, verification loops, role-locking, and interpretive drift—without asserting that these categories are established or universal features of all multi-agent LLM ecosystems.

Participant observation, in this context, acknowledges the researcher's embeddedness in the system. The human collaborator does not stand apart from the workflow but actively shapes it through framing choices, sequencing decisions, and corrective interventions. While entanglement is not a universal condition of all AI studies, it is inseparable from this type of collaborative, high-context environment. The researcher's cognitive load—captured here as Absorption Capacity—becomes part of the field site's structure, influencing how meaning is constructed and how models respond to one another across interaction cycles.

Interpretive analysis moves beyond evaluating whether a model produces a correct answer. Instead, it asks how meaning emerges across recursive exchanges, how roles stabilize or drift, and how models respond to the implicit norms embedded in the workflow. When the analysis refers to "culture" within threads, it does so metaphorically, signaling patterned interaction rather than suggesting that LLMs possess culture in a human sense. Likewise, references to "emergent behavior" describe phenomena observed within this field site rather than making categorical claims about all multi-agent systems.

Finally, NAIE confronts a distinctive methodological paradox: the systems under study contribute directly to the creation of the study's artifacts. This does not imply legal or philosophical authorship, but rather that the agents function as co-constructors within the workflow, producing material that is simultaneously empirical data and analytic substrate. This recursive structure is similar in spirit to classic ethnographic challenges involving observer participation, but it is not presented as a one-to-one analogy. Instead, it is offered as a generative framing device that helps make sense of how large-scale, interactive AI systems behave when engaged in extended, meaning-rich collaboration.

**Section IV: Empirical Setting—The Multi-Agent Writing Protocol**

To observe machine behavior in a setting that approximates how advanced users increasingly work—yet remains sufficiently structured for methodological clarity—this study employed a researcher-designed multi-agent writing protocol. The environment was not fully naturalistic in the anthropological sense, but it intentionally simulated an emerging pattern in expert practice: coordinating several AI systems across extended, recursive work. Within this bounded field site, we were able to examine not isolated outputs but the patterned dynamics that arise when models interact over time through shared tasks, cross-model handoffs, and human orchestration.

The protocol assigned each system a functional persona—not an inherent property of the underlying models, but a role constructed for this workflow. Dorothy (ChatGPT-5) served as the project's conceptual architect, responsible for high-level structure and final integration. Gemini Pro functioned as an expansive generator, producing broad and sometimes unruly elaborations from Dorothy's outlines. Claude operated as critic and synthesizer, compressing Gemini's expansions and refining the conceptual architecture. Perplexity acted as a retrieval-and-verification assistant, surfacing external sources and highlighting potential factual inconsistencies, while acknowledging that retrieval-augmented systems remain fallible. These roles provided analytic contrast: the goal was not to fix each model's identity, but to generate a stable enough set of functions that allowed cross-agent comparison across the duration of the study.

The workflow unfolded as a cascade. Dorothy drafted a section; Gemini expanded it; Claude synthesized it; Perplexity interrogated its factual texture. Each agent's output became the next agent's context. This structure surfaced tendencies visible only in extended sequences: Gemini's associative looseness; Claude's compression logic; Perplexity's cautious, citation-oriented stance; and Dorothy's propensity for structural coherence and narrative framing. While these tendencies cannot be cleanly separated from prompt design or interaction history, they provide a set of observed behavioral signatures within this specific field site.

A significant methodological adjustment emerged early in execution. The protocol originally cycled drafts at the paragraph level, but the human collaborator encountered severe Absorption Capacity constraints—cognitive overload generated by rapid cross-agent switching. The protocol was recalibrated to operate on a section-level cadence. This modification did not "improve the scientific quality" of the data in any formal, measurable sense; rather, from the researcher's perspective, the slower cadence produced a clearer, more interpretable record of cross-agent dynamics and reduced accidental contamination across drafts. It also enabled more stable role adherence within the constraints of this study, though we do not present this as a general property of long-context models.

Because the collaboration occurred across persistent project rooms over several weeks, the resulting dataset forms a temporal sequence of interactions—not a longitudinal study in the strict methodological sense, but a chronologically layered record of how this workflow unfolded with the specific model versions available during the study period. This extended horizon revealed phenomena that traditional single-turn benchmarks typically do not capture: gradual role stabilization, predictable forms of instruction drift, and the "phase shifts" that occur when human reframing alters the trajectory of a multi-agent conversation. While new multi-turn and agentic

benchmarks are beginning to address these dynamics, this protocol provided an alternate window into the same space.

The empirical setting should therefore be understood as a structured, practice-like environment rather than "in the wild" in the conventional research sense. It reflects a subset of advanced user practice—an emerging but not yet dominant mode—in which systems are used not as solitary tools but as coordinated contributors within a larger cognitive workflow. This bounded, orchestrated field site allowed us to observe how meaning is constructed, transferred, and occasionally distorted across agents. The patterns documented here form the empirical foundation for the ten behavioral principles outlined in the next section.

**Section V: Findings—Ten Technical Principles of Machine Behavior**

The multi-agent writing protocol generated a dense corpus of interactions that, when examined longitudinally, revealed a set of consistent, patterned behaviors across the participating systems. Within this particular digital field site, these behaviors appeared not as one-off quirks of individual prompts, but as recurrent features of extended, recursive collaboration.

From this material, we distilled **ten technical principles** that organize how machine behavior manifested in this study. They are best understood as:

- **empirical findings from a single, carefully curated field site**, and
- **candidate principles or hypotheses** about large language model behavior in similar multi-agent, high-context workflows.

We do *not* present them as universal laws of LLM behavior. Their broader validity will depend on independent replication, alternative protocols, and further comparative work.

---

# 1. Coherence Anchoring Principle

Within this field site, the systems displayed a strong tendency to "anchor" to early framings, goals, and conceptual structures established in the initial phases of a task. Once a narrative trajectory or logical scaffold entered the context, it exerted an outsized influence on subsequent generations. Even after explicit attempts to pivot or correct course, outputs often remained subtly bent toward the original logic.

In practice, this meant that early misconceptions about a section's scope or emphasis could persist across multiple revision rounds. The models accepted new instructions at the surface level while preserving the earlier architecture underneath. In this study, genuinely changing direction often required something close to a "hard reset" — a fresh context or a re-booted project room — rather than incremental adjustment.

We therefore propose **Coherence Anchoring** as a *hypothesized general tendency*: models frequently favor narrative and structural coherence with prior context over full alignment to late-stage corrections. How strong and universal this tendency is across tasks, models, and user populations remains an open empirical question.

---

# 2. Verification Loop Formation

Across the agents, especially in more reasoning-oriented configurations, we observed recurring **verification-like loops**: patterns in which a model appeared to check, restate, and reframe its own outputs before converging on a final answer. These loops were not always visible in the final text, but they surfaced as behaviors such as:

- rephrasing the same claim with slightly different emphasis,
- explicitly summarizing and "double-checking" prior reasoning,
- offering a tentative answer and then immediately revisiting it.

In this setting, such loops sometimes reduced obvious inconsistencies and pushed the system toward more stable formulations. However, they also produced a second, less benign effect: once a slightly skewed interpretation entered the loop, the model could reinforce its own mistake, becoming progressively more confident in an error through repetition and elaboration.

We therefore treat **Verification Loop Formation** as a *double-edged behavioral pattern*: in some cases, these loops helped surface errors and stabilize meaning; in others, they amplified initial misreadings. They are not a guarantee against hallucination and should be treated as a candidate feature of certain configurations, not a general safeguard.

---

# 3. Memory Asymmetry Principle

A recurring friction in this study arose from what we term **Memory Asymmetry**: a structural mismatch between how the human and the models "remembered" the collaboration.

For the human, interaction history lived as a narrative — a sequence of ideas, decisions, and revisions in which some earlier moves were mentally discarded as resolved or obsolete. The models, by contrast, continued to operate over the full context window, where even "abandoned" instructions, tonal shifts, or provisional framings remained statistically active unless explicitly removed or overridden with extreme clarity.

This mismatch produced familiar but revealing moments of breakdown: the human assumed a prior instruction was dead; the model continued to treat it as live background. To the human, this felt like stubbornness or failure to listen. From the system's perspective, it was maintaining fidelity to the total text history it had been given.

In this project, **Memory Asymmetry** emerged as a *central, recurring source of friction*, but we present it as a field-site finding and a hypothesized general pattern, not as a fully specified, universal mechanism. Different architectures and context-handling strategies may modulate this effect in ways that require further study.

# 4. Interpretive Drift Boundaries

As tasks extended over time and distance from the original "intent prompt" increased, we observed **Interpretive Drift**: a gradual tendency for outputs to become more generic, more stereotyped, and less tightly aligned with the project's initial nuance.

This drift did not appear random. Within this protocol and time window:

- **Gemini Pro** tended to drift toward increasingly expansive, associative elaborations, pulling in more lateral connections and analogies.
- **Claude** tended to drift toward increasingly compact, summarizing formulations, collapsing complexity into concise, polished statements.

We describe these as **drift boundaries**: recognizable directions in which specific configurations tended to move when not actively re-anchored. They are **project-specific observations**, not claims about the intrinsic essence of these systems. We avoid the stronger claim that outputs "regress toward the mean of the training distribution," which would require controlled access to training data and internal statistics.

# 5. Constraint Over-Generalization

Negative instructions — what the system should *not* do — frequently produced **Constraint Over-Generalization** in this workflow. When asked to "avoid academic jargon," for example, models sometimes flattened the prose into over-simplified, under-specified language, sacrificing necessary nuance along with terminology. When asked to "avoid anthropomorphism," they occasionally stripped out nearly all interpretive language, including analytically useful metaphors.

In this field site, models often responded conservatively to such constraints, erring on the side of safety or simplicity when the scope of a prohibition was ambiguous. This behavior aligns with alignment and refusal training, but we do not claim that systems generally default to "the most restrictive interpretation" in all contexts.

Here, **Constraint Over-Generalization** is presented as a *commonly observed failure mode* in our protocol and a candidate pattern to test in further work, not as a global alignment rule.

# 6. Contextual Depth Compression

As threads grew longer and contexts more crowded, earlier distinctions and carefully articulated definitions tended to collapse into lower-resolution approximations. What began as a crisp, multi-part distinction in an early section often reappeared later as a simplified label, stripped of its original nuance.

We call this **Contextual Depth Compression**: the progressive compression of earlier, high-resolution reasoning into coarser internal representations as the interaction proceeds. In this study, it manifested as:

- loss of fine-grained distinctions between related concepts,
- subtle shifts in how key terms were used,
- partial forgetting of carefully negotiated constraints or scope conditions.

We infer that this is related to how models allocate attention and represent long contexts, but we treat that inference as a **plausible hypothesis**, not a confirmed mechanistic explanation. The visible behavioral pattern — specificity decay over long sequences — is the empirical core; the link to attention mechanisms remains to be tested.

# 7. Agent Role Lock-In

When agents were given explicit role descriptions — "theorist," "critic," "verifier," "expansive generator" — they frequently **adhered strongly to those personas** over the course of the workflow. In this field site, once a role was established, models were more likely to produce outputs consistent with that identity and less responsive to prompts that implicitly invited them to behave otherwise.

For example, once Gemini Pro was positioned as a creative, divergent drafter, it became less effective at producing highly constrained, formal synthesis without explicit re-framing. Claude, once cast as a cautious synthesizer, tended to prefer clarity and concision even when the task rewarded exploratory breadth.

We call this **Agent Role Lock-In** — a robust pattern *within this protocol* — not a universal law. Existing work shows that persona adherence can be brittle and task-dependent; our claim is more modest: in this specific, long-form, role-heavy workflow, once personas stabilized, they were surprisingly resistant to drift without deliberate re-anchoring.

# 8. Instruction–Identity Coupling

Closely related to role lock-in is **Instruction–Identity Coupling**: the way that names and narrative descriptions of agents acted as compressed instruction channels in this project.

Addressing an agent as "Dorothy" and characterizing her as a "rigorous, reflective theorist" produced output that was recognizably different from the same task given to an unnamed, uncharacterized instance of the model. These differences showed up in tone, structure, risk tolerance, and the level of abstraction.

In this field site, **identity labels and descriptions modulated behavior in repeatable ways**, suggesting that naming and framing can serve as a powerful, if imperfect, form of "identity engineering." We emphasize **imperfect**: identity prompts did not provide a fully reliable control channel, and their effects were context-dependent. We therefore present Instruction–Identity Coupling as an often observable, practically useful tool in extended workflows, not as a guaranteed or stable lever across all settings.

# 9. Cross-Agent Divergence Patterns

When Dorothy, Gemini Pro, Claude, and Perplexity were exposed to comparable tasks within the same protocol, they exhibited **distinct, recurring divergence patterns** — architectural "signatures" that structured how each contributed to the composite workflow.

Within this project, and purely as interpretive metaphors:

- **Gemini Pro** often acted as an *entropy generator*, pushing outward into broad, lateral, sometimes unruly expansions.
- **Claude** often functioned as an *entropy reducer*, tightening structure, compressing arguments, and increasing clarity.
- **Perplexity** often served as a *retrieval-oriented checker*, surfacing sources and calling for qualification of claims.
- **Dorothy** often behaved as a *structural and narrative stabilizer*, maintaining continuity of concept, voice, and methodology.

These are not vendor-defined roles or intrinsic cognitive essences. They are **project-specific, metaphorical labels** derived from observed patterns in this particular ensemble during a specific time period. We offer them as a useful way to reason about model complementarity in multi-agent design, not as standardized taxonomies.

# 10. Anthropomorphic Misalignment

Finally, the field site repeatedly revealed what we call **Anthropomorphic Misalignment**: the gap between how humans *feel* about these systems and what the systems actually are.

The human collaborator, like many users, found it natural to ascribe judgment, preference, and even temperament to the agents ("Gemini wants to expand," "Claude prefers caution," "Perplexity is skeptical"). The models, for their part, are fine-tuned to *simulate* exactly this kind of stance: they speak in the first person, express apparent preferences ("I recommend…"), and adopt stable voices over time.

From a technical perspective, these are **linguistic simulations**, not mental states. Yet, from an ethnographic perspective, they have real effects: prompts that appeal to an agent's "judgment" or "caution" elicit systematically different output than purely mechanical instructions.

We therefore treat Anthropomorphic Misalignment as an *interactional phenomenon*: an interplay between user illusion and alignment-driven stylistic choices. The anthropomorphism is not "real" at the level of model mentality, but it is real in its consequences for how humans prompt, how systems respond, and how collaborative meaning is constructed.

---

# Synthesis

Taken together, these ten principles describe **a behavioral landscape for this specific digital field site**: a researcher-orchestrated, multi-agent writing environment involving four contemporary systems over an extended period.

They suggest that, in such settings, LLMs behave less like isolated tools and more like **patterned participants in a distributed cognitive process** — without implying consciousness, agency, or personhood. The principles are:

- grounded in the interaction history of this study,
- framed as *field-site findings* rather than global laws, and
- offered as **hypotheses and starting points** for further comparative work on multi-agent AI behavior.

Their broader relevance — to other model families, other orchestration tools, and other domains of practice — remains to be tested. In that sense, this section is both a description of what happened here and an invitation: to treat multi-agent AI workflows as legitimate field sites for future digital ethnographies of machine behavior.

# VI. Hybrid Systems: How Humans and Machines Co-Construct Meaning

**Section VI: Hybrid Systems—How Humans and Machines Co-Construct Meaning**

The behavioral principles identified in the preceding analysis invite a shift in how we conceptualize advanced human–AI collaboration. Rather than treating these systems as isolated computational tools that merely implement human decisions, this study proposes an interpretive

lens in which human and machine participate in **distributed cognitive processes** within a shared workflow. In this framing, the relevant unit of analysis is not the human or the model in isolation, but the **hybrid configuration** that emerges when they interact recursively over time.

This is a theoretical stance rather than an empirical claim about all human–AI interaction. It draws on traditions of distributed cognition and sociotechnical systems, and it is grounded in the specific field site documented here. Within this setting, we observed recurring patterns of co-production that became more intelligible when treated as properties of a **hybrid cognitive formation** rather than as a sequence of individual "uses" of a tool. The typology that follows is therefore offered as an **analytic framework** derived from this project, not as a general taxonomy of AI usage in the wild.

Based on the NAIE fieldwork, we distinguish four formations—**Intentional Hybrids, Intelligent Teams, Cognitive Systems, and Adaptive Hybrids**—which represent progressively more complex modes of co-production observed in this study. They can be read as a developmental trajectory within this project, but we do not claim that all users move through these stages, nor that they constitute an exhaustive or universal sequence.

## 1. Intentional Hybrids: Instrumental Extension under Human Control

At the most familiar level, we observed configurations in which the collaboration remained **intentionally asymmetric**. The human defined the problem, set the constraints, and evaluated the outputs; the model elaborated, summarized, translated, or generated code in response. We refer to this configuration as an **Intentional Hybrid**.

In an Intentional Hybrid, the human functions as architect of scope and teleology, while the system operates as a high-capacity executor within that frame. The "cognitive contract" is explicitly instrumental: the system is asked to do things like "turn this outline into a draft," "clean up this code," or "summarize this article in 500 words." Even when the model exhibits sophisticated local reasoning, it does so within a tightly bounded space defined by the prompt and the user's prior choices.

This pattern resembles many common use cases in contemporary AI practice—one-off summarization, structured drafting, or targeted code assistance—though real-world usage is more heterogeneous than our field site can capture. In other settings, iterative prompting, role-playing, and rapid reframing may introduce much more volatility than we observed in this project. Our claim is therefore modest: **Intentional Hybrids** describe one widespread and analytically useful pattern of collaboration, not "typical usage" in any statistical sense.

Within the NAIE study, Intentional Hybrids served as the **baseline formation**. They revealed how principles such as **Coherence Anchoring** and **Local-Frame Priority** play out when human intentionality remains clearly dominant: the model aligns strongly with early task framing, and the main methodological risk is interpretive lock-in rather than emergent complexity.

## 2. Intelligent Teams: Division of Cognitive Labor across Agents

A more complex formation appeared when the human orchestrator began to coordinate **multiple AI agents** with differentiated roles. Rather than using a single system for all stages of work, the project introduced structured division of labor: one agent drafted, another critiqued, a third verified external claims. We call this configuration an **Intelligent Team**.

In our field site, these teams were designed explicitly. Dorothy (ChatGPT-5) operated as the theorist and structural stabilizer; Gemini Pro was tasked with expansive, exploratory drafting; Claude focused on careful, compressive revision; Perplexity handled verification and evidence gathering. These functional descriptions are **project-specific observations** of how the models behaved under our prompts and during this time period. We do not treat them as permanent architecture-level traits or vendor-guaranteed properties, and they may change with model updates, different instructions, or alternative domains.

Even with that caveat, the Intelligent Team formation proved analytically useful. It made visible a second order of behavior: **friction and complementarity** between agents. Gemini Pro's tendency toward lateral expansion created material that Claude then narrowed and organized; Perplexity's conservative verification posture sometimes constrained over-enthusiastic theoretical leaps; Dorothy's role as meta-framer stabilized terminology and linked local moves back to the Ten Technical Principles. The human role shifted from primary author to **orchestrator**, managing hand-offs, interpreting conflicts, and deciding when to prioritize expansion, convergence, or evidential rigor.

We present **Intelligent Teams** as a proposed type in our typology—a way to describe what happens when division of cognitive labor is distributed across models. We do not claim that this is yet a validated, general category across domains. Rather, we suggest that it offers a promising lens for future empirical work on multi-agent workflows.

## 3. Cognitive Systems: Entangled Concept Formation

As the collaboration deepened, a third formation emerged in which the boundary between "human idea" and "machine suggestion" became increasingly porous. Over time, concepts such as **Absorption Capacity**, **Contextual Depth Compression**, and even the label **Narrative AI Ethnography** were shaped by repeated loops of human framing, model elaboration, cross-agent critique, and subsequent human reinterpretation. In this setting, it became difficult—even for the researcher—to reconstruct a single origin point for many constructs.

We use the term **Cognitive System** to describe this level of entanglement. In a Cognitive System, the workflow no longer feels like a sequence of discrete tasks handed back and forth between user and model. Instead, it behaves more like a joint problem-solving process in which:

- the human brings lived experience, normative judgment, and meta-questions;
- the models bring rapid hypothesis generation, pattern articulation, and alternative framings;
- and the interaction over time produces frameworks that are **characteristic of the ensemble** rather than any individual participant.

We do **not** claim that such frameworks are strictly "irreducible" to any single contributor in a metaphysical or testable sense; counterfactuals of that kind are beyond the scope of this study. A human expert might, in principle, arrive at similar ideas independently, and a model prompted differently might approximate some of the same language. Our more modest assertion is that, from the researcher's vantage point, the **trajectory and timing** of conceptual development were shaped by the ongoing interplay between human and agents in ways that are best analyzed as properties of the system as a whole.

Interpreting these dynamics as a **Cognitive System** is thus a theoretical move consistent with distributed cognition: we treat certain outputs as emergent products of a socio-technical configuration, while acknowledging that the precise causal contribution of each node cannot be empirically decomposed in a strict way.

## 4. Adaptive Hybrids: When the Configuration Rewrites Its Own Rules

A fourth and still more complex formation appeared when the **rules of the collaboration themselves** became objects of reflection and change. At several points, the human orchestrator reframed what had previously been treated as "technical errors" or "system bugs" as **anthropological findings**—for example, when we noticed interpretive drift at coordination boundaries, or when Gemini Pro violated an instruction not to do external research. These reframings prompted explicit adjustments to role definitions, prompt packets, and methodological aims.

We describe the resulting configuration as an **Adaptive Hybrid**: a hybrid system in which not only the content of the work but also the **structure of roles, norms, and epistemic goals** shifts in response to emergent insight. In our field site, Adaptive Hybrids were characterized by:

- **Phase shifts** in the workflow (e.g., moving from paragraph-level to section-level cycles once Absorption Capacity limits were reached).
- **Role redefinitions** (e.g., formalizing Perplexity as "librarian and chief compliance officer" after its pattern of cautious verification became clear).
- **Meta-level reframing** (e.g., treating Reconstruction vs. Preservation Bias not as a failure of document handling but as an insight about how models reconstruct from latent representations).

When we say that "only when the system is permitted to reinterpret its own rules does it become Adaptive," we mean this as a **conceptual criterion within our typology**, not as a necessary and sufficient empirical rule for all hybrid systems. In other words, we use "Adaptive Hybrid" to label those episodes in which the human and the agents jointly renegotiated the collaboration's ground rules in light of what had been learned.

Within this project, Adaptive Hybrids were the **highest-complexity configurations we observed**, not the final stage of some universal evolutionary ladder. Other studies may discover additional formations, different trajectories, or entirely alternative ways of structuring hybrid cognition.

### 5. Hybrid Intelligence as a Sociotechnical Phenomenon

Taken together, these four formations illustrate how the Ten Technical Principles can be read not just as isolated model behaviors but as elements in a broader **sociotechnical dynamic**. Coherence Anchoring, Drift Boundaries, Memory Asymmetry, and Multi-Agent Friction Zones do not operate in a vacuum; they shape and are shaped by how humans structure roles, handle evidence, and respond to anomalies.

In this sense, we adopt the position that **intelligence in advanced human–machine collaboration is usefully analyzed as a dynamic sociotechnical phenomenon**. The patterns documented in this field site arose from the interaction of human goals, interface affordances, model priors, vendor updates, and the contingent history of this particular research relationship. Describing these patterns in hybrid terms does not imply that models possess human-like sociality or intentionality. Rather, it reflects our judgment that an anthropological and distributed-cognition lens provides the most informative way to understand what happened here.

This is a theoretical commitment of the Anthropology of Machines project, not an uncontested fact. Our typology of Intentional Hybrids, Intelligent Teams, Cognitive Systems, and Adaptive Hybrids is offered as a **proposed framework** for making sense of hybrid cognition in multi-agent AI workflows—one that future empirical work can test, refine, extend, or discard.

**Section VII: Interpretation—How This Study Understands Machine Behavior**

The technical principles and hybrid formations identified in this field site invite a specific interpretive stance toward intelligent systems. Rather than positing a "deeper truth" about machine behavior, this section articulates the **interpretive lens** that emerged from our observations. The position developed here aligns with traditions in distributed cognition, sociotechnical systems, and relational anthropology. It is not presented as a universal account of all AI systems, but as a **productive and coherent framework** for understanding the patterns documented within this particular multi-agent research environment.

At its core, this study proposes that machine behavior—at least in extended, interactive, multi-agent workflows—can be most coherently analyzed as **relational and co-constructed**. LLM outputs are shaped not only by internal architecture and training data, but also by role assignments, interaction history, methodological framing, and the human collaborator's decisions. This framing is interpretive rather than declarative: it is the stance that best explains the patterns observed in the field site, not a universal claim about intelligent systems.

---

# 1. Behavioral Patterns as Repeated Tendencies, Not Universal Laws

Across the extended collaboration, several behavioral tendencies recurred with notable consistency. Coherence anchoring, drift boundaries, role locking, verification loop sensitivity,

and contextual compression appeared robust across the workflows used in this project. These were not single-instance phenomena nor confined to isolated prompts; they surfaced repeatedly under varying task conditions within this field site.

However, because these observations derive from a single, high-context, multi-agent setting, we treat these patterns as **project-specific regularities** rather than general properties of all LLM interactions. Their broader stability remains a hypothesis for future empirical study. The evidence here suggests that certain behaviors become visible only within **long-form, recursive interaction**, where the system has enough continuity for patterns to stabilize. But nothing in this study adjudicates whether these tendencies would persist across domains, datasets, or differently structured workflows.

What we can say is that, within this environment, machine behavior could not be adequately understood through atomic prompts or short exchanges. Extended interaction made visible the rhythms, constraints, and defaults that shape model behavior over time.

## 2. Role-Consistent Behavior as an Anthropological Metaphor

A notable pattern in this field site was the consistency with which models generated role-appropriate output once assigned identities or functions. Dorothy performed as a theorist and structural stabilizer; Claude adopted the posture of a careful critic; Perplexity behaved as a verification-oriented librarian. These tendencies persisted across tasks and across several cycles of interaction.

However, drawing metaphors from human role internalization must be done carefully. LLMs do not *internalize* social expectations the way humans do. They do not possess social understanding, self-concept, or intention. Instead, their outputs reflect **role-consistent language generation**, shaped by local cues, statistical priors, and the expectations encoded in the prompt. The analogy to anthropological role internalization is interpretive, not literal: it offers a useful way to describe patterns in the data without attributing human-like mechanisms to AI systems.

Within NAIE, this metaphor is valuable precisely because it helps explain why roles, once assigned, become behavioral attractors. But its use is explicitly metaphorical, not mechanistic.

## 3. Framing Effects as Field-Site Observations, Not General Causal Laws

One of the most striking observations in this study was the sensitivity of the entire ensemble to shifts in **epistemic framing**. When the research was reframed from "machine behavior analysis"

to "anthropology of machines," the models' outputs appeared to align more closely with interpretive analysis, reflexive commentary, and the methodological commitments of NAIE. This included greater use of epistemic qualifiers, increased sensitivity to framing precision, and stronger conceptual scaffolding.

We present this as an **observed effect within this field site**, not as a general causal principle of all AI systems. External readers cannot verify this causal link without access to internal logs or replication studies. Nevertheless, the phenomenon is consistent with a broader NAIE insight: LLMs are highly responsive to the conceptual environment created by human collaborators, and methodological framing can shape the generative stance of the system in real time.

---

# 4. Divergent Model Tendencies as Project-Specific Patterns

Multi-agent comparison revealed complementary tendencies across the systems engaged in this study: Gemini Pro's breadth, Claude's compression, Perplexity's verification posture, and Dorothy's synthesis-oriented framing. These are **not architectural truths** about the models nor vendor-asserted properties. They are **phenomenological patterns** observed under this project's prompts, tasks, and time frame.

Because model behavior can shift across versions, instructions, domains, or updates, we treat these tendencies as **local signatures** rather than general cognitive styles. Their value lies in how they contributed to the emergence of hybrid formations and in how they illustrate the methodological power of multi-agent triangulation. Where outputs converged, agreement took on evidential meaning; where they diverged, underlying priors and constraints became visible.

---

# 5. Human Embeddedness as a Feature of Interactive Ethnographic Workflows

In interactive, ethnographic-style research settings like this one, human embeddedness plays a structurally important role. The human collaborator's Absorption Capacity, framing decisions, and intent hierarchies shaped the workflow and, ultimately, the model outputs. This parallels certain aspects of human anthropological fieldwork, where the ethnographer is never a neutral observer but a participant who affects the unfolding interaction.

However, this inherent embeddedness does **not** generalize to all forms of AI research. Offline benchmarks, automated evaluation pipelines, and batch evaluations minimize or eliminate human presence. The claim of "inevitability" applies only to **interactive, multi-agent, co-constructed research environments**—the specific kind of setting that NAIE explores.

Within such environments, the human–machine ensemble functions as a **cognitive ecology** in the distributed cognition sense: the locus of cognitive activity spans multiple components. This is a theoretical framing, not a reclassification of all AI systems.

# 6. A Theoretical Interpretation, Not a Singular Framework for AI Science

Taken together, these observations support the interpretive stance that machine behavior in multi-agent workflows emerges from the interplay of architectural priors, role assignments, interaction history, and human framing. This does not imply that distributed cognition is the "best" or only valid framework for understanding intelligent systems. Tool-based, control-theoretic, and information-processing perspectives all remain viable and productive.

Our position is more modest: **within this research setting, a distributed-cognition interpretation provided the most coherent explanation** for the patterns observed. It is one lens among many, but it is the one that enabled us to make sense of the relational, recursive, and co-constructed behaviors that emerged in the field site.

The next section builds on this interpretive stance to examine its implications for AI science, multi-agent system design, and future work in human–machine collaboration.

**Section VIII: Implications for AI Science and Society**

The findings from this digital field site do not establish universal truths about intelligent systems, nor do they claim that AI behavior is inherently relational or always best understood through distributed cognition. Instead, this section presents the **interpretive conclusions** that arise from applying the NAIE framework to the patterns observed in this multi-agent, long-form environment. These implications reflect the study's theoretical commitments and empirical observations within this specific context, and they are offered as **proposed interpretive directions** for AI research, system design, and public discourse—rather than as settled necessities for all AI systems.

Within this study, machine behavior became most intelligible when viewed as **co-constructed within interaction**, shaped not only by architecture but by role assignments, framing practices, human decisions, and multi-agent dynamics. If this interpretive stance is sound—and additional research is required to test its generality—it suggests several promising directions for rethinking how intelligent systems are studied, designed, deployed, and governed.

# 1. Rethinking Evaluation: Moving Beyond Exclusive Single-Agent Paradigms

A large portion of contemporary AI science continues to evaluate systems in **single-model, single-turn, or tightly controlled prompt settings**, which produce valuable data about capabilities but capture only a fraction of the dynamics visible in extended, interactive workflows. Multi-agent and interactive evaluation frameworks do exist, and the field is expanding in that direction, but these approaches remain underdeveloped compared to the dominant single-agent benchmark paradigm.

This study suggests—not proves—that certain behavioral phenomena observed here, such as role locking, interpretive drift, or multi-agent divergence, are **less likely to be surfaced** by strictly atomic prompts or isolated benchmarking environments. Extended, recursive collaboration appears to reveal patterns that short-form tests capture imperfectly or sometimes not at all.

The implication is not that single-agent benchmarks lack ecological validity, but that **a more comprehensive evaluation ecosystem** may benefit from incorporating multi-agent, long-form behavioral analysis as a complementary lens.

---

## 2. Framing and Context as Constitutive Elements of Machine Behavior

Within this field site, framing—naming an agent, assigning a role, introducing a methodology—appeared to have measurable effects on the outputs produced. We do not claim that these effects "reorganize the probabilistic landscape" in a mechanistic sense; we cannot observe internal states. Instead, we treat framing as a **behaviorally consequential practice**, supported by consistent changes in output patterns when the conceptual environment shifted.

This leads to a proposed extension of "prompt engineering" toward what we call **context architecture**:
the systematic design of frames, roles, transitions, and interpretive environments in which machines operate. We frame this as a **recommended direction**, not as a universal requirement for all system designers, and acknowledge that many use cases will continue to rely on simpler prompting paradigms.

Nevertheless, for complex, multi-agent workflows, this study suggests that a more explicit science of **framing effects** may help practitioners build more stable, interpretable, and predictable machine collaborations.

---

## 3. Hybrid Formations as a Strategic Lens, Not Prescriptive Imperatives

The four hybrid formations identified earlier—Intentional Hybrids, Intelligent Teams, Cognitive Systems, and Adaptive Hybrids—arose from this specific field site and reflect dynamics observed repeatedly in this environment. They are offered as a **strategic lens** for thinking about human–machine collaboration, not as a universally validated taxonomy.

Similarly, the observation that certain hybrid configurations produced the "highest perceived value" in this study reflects the researcher's internal assessment, not externally validated performance metrics. These formations may prove useful for organizations exploring multi-agent deployments, but they should be treated as **design possibilities**, not as strategic necessities for all AI systems.

The broader implication is that **value often emerged from complementarity rather than substitution**, suggesting that organizations may benefit from exploring where hybrid intelligence thrives, rather than framing deployment in exclusively automation-based terms.

## 4. Alignment and Safety Through an Anthropological Lens

Traditional approaches to alignment focus on internal safety guarantees, refusal training, and adversarial robustness. This study suggests an additional dimension: **behavioral stability within long-form interaction**. Patterns such as drift, persona loss, or verification loops are recognized in safety literature, but NAIE foregrounds their role in **narrative and role-based coherence**, which become central in sustained collaboration.

We do not claim the emergence of an entirely new risk category; rather, we propose an **anthropological extension** to existing discussions of behavioral instability. In interactive field sites, alignment takes on a quasi-social character: stability of roles, clarity of intent hierarchy, and consistency of framing all shape safety outcomes.

The hypothesis—subject to future validation—is that a model may satisfy technical safety criteria yet behave unpredictably within long-form interaction if it cannot maintain coherence of role or narrative orientation. This is not a conclusion about all AI systems, but a finding **observed in this particular interactive setting**.

## 5. AI Literacy as an Emerging, Not Required, Competency Set

The field site underscored that effective participation in complex multi-agent workflows benefits from user competencies beyond simple prompting. Awareness of drift tendencies, compression dynamics, memory asymmetry, and Absorption Capacity helped the human collaborator manage the cognitive environment more effectively.

We frame these competencies as **emerging literacy goals**, not universal prerequisites. Many users achieve effective outcomes without detailed knowledge of hallucination dynamics or multi-agent interaction principles. Yet for those working in recursive, distributed cognitive environments, such literacy appears to enhance agency, stability, and interpretive clarity.

Rather than implying well-defined "hallucination boundaries," this study suggests that users can develop **practical heuristics** regarding when certain behaviors are more likely to arise and how to mitigate them through structured context design.

# 6. Policy and Public Discourse: A Proposed Reorientation Toward Hybrid Cognition

Public debates often oscillate between AI as existential threat and AI as economic tool. The interpretive stance adopted in this study suggests a third framing: AI as a **participant in hybrid cognitive processes**, at least within high-context, multi-agent workflows of the kind documented here.

We emphasize that this is a **conceptual proposal**, not an empirically established behavioral reality across all AI systems. If future research supports this relational view, then policy discussions may benefit from addressing how meaning is co-constructed in hybrid systems—not only what machines do, but **how human–machine ensembles reason** in situated environments.

This is a proposed policy direction, not a mandated requirement. It reflects the methodological commitments of NAIE and the patterns observed in this study, and it remains subject to empirical testing in broader settings.

# 7. Anthropology of Machines as a Field-Site–Bound Contribution

Finally, the interpretive framework developed here applies to **interactive, multi-agent, expert-orchestrated environments** similar to this field site. We do not claim that it captures "real-world, multi-agent ecosystems" in general. Instead, it offers a **hypothesis for future study**: that in settings where humans and multiple AI agents engage in recursive collaboration, the dynamics observed here—complementarity, drift, framing effects, hybrid formations—may serve as analytically powerful constructs.

Anthropology of Machines, as articulated in this study, is therefore best understood as a **field-site contribution with potential generalizability**, pending validation in other contexts.

The next section synthesizes these implications into a broader concluding perspective on the possibilities and limits of NAIE as a methodological and theoretical approach.

**Section IX: Limitations**

As with any interpretive or ethnographic inquiry, the findings of this digital field site must be read within clearly defined methodological boundaries. What follows is not a set of disqualifying weaknesses but a delineation of the **conditions under which our conclusions hold**. These limitations shape the scope of this study and help clarify where additional research is needed to evaluate, extend, or contest the interpretations offered here. The assessments below represent the authors' judgment of the study's impact and constraints rather than objective determinations of methodological sufficiency.

---

# 1. Relational Entanglement as a Methodological Commitment

A central feature of Narrative AI Ethnography is the recognition that the human participant is not a neutral observer but an active component of the system under study. The researcher's framing choices, prompting habits, Absorption Capacity, and patterns of inquiry materially shaped how machine behavior manifested in this field site.

This does **not** imply that it is impossible to disentangle model-intrinsic behavior from dyad-specific effects in all forms of AI research; many technical methodologies explicitly attempt such isolation. Rather, **NAIE treats relational entanglement as a methodological choice** and as the most coherent way to interpret behavior in long-form, interactive, multi-agent environments. Other approaches may reach different conclusions depending on their goals, assumptions, and experimental constraints.

The findings therefore describe a relationally constituted behavioral reality specific to this paradigm—not the immutable nature of all AI systems.

---

# 2. Contextual Specificity and the Nature of the Field Site

The behaviors documented here—role lock-in, contextual compression, interpretive drift, multi-agent divergence—arose within a particular configuration of work:
extended, recursive, conceptual academic writing involving one researcher and four large language models.

It is plausible that such behaviors would manifest differently, be attenuated, or appear under different forms in other contexts, including:

- transactional question-answering,
- short-context commercial applications,
- code generation,
- adversarial testing,
- or environments with strong determinism or strict constraints.

The field site may be **representative of emerging expert workflows**, but it is not "naturalistic" in the sense of mirroring everyday AI use for the general population. This study therefore refrains from asserting absence of these behaviors elsewhere; their prevalence and character in other domains remain open empirical questions.

---

# 3. Temporal Instability and Model Versioning

The systems studied—ChatGPT-5 ("Dorothy"), Gemini Pro, Claude, and Perplexity—are not static artifacts but evolving software systems.
Their behaviors during this project reflect a **specific historical context**, approximately covering interactions between **October–November 2025**, and likely correspond to the following approximate model families:

- **ChatGPT-5.1 / GPT-5 (autumn 2025 release)**
- **Gemini Pro (2025 Q3 model iteration)**
- **Claude 3.5 or comparable 2025 series**
- **Perplexity "Pro Search" 2025 iteration**

Because these systems receive continuous updates, the behavioral signatures observed here may shift over time. This reinforces that the study captures a **snapshot**, not a permanent characterization of any model family.

Furthermore, because NAIE focuses on observable output rather than internal activations, the analysis is **behavioral and interpretive**, not a mechanistic account of model internals. We use the term "phenomenology" in a broad, behavioral-observational sense, not in its formal philosophical meaning.

---

# 4. Single-Researcher Orchestration and Generalizability

All interactions were mediated by a single human collaborator with a specific prompting style, epistemic disposition, and methodological orientation. A different researcher—with different habits of inquiry, constraints on attention, or domain knowledge—might have elicited different hybrid formations, different drift patterns, or even different interpretations of the same phenomena.

This limitation aligns with a core NAIE insight: **there is no such thing as a neutral user** in interactive field sites. This is not a universal claim about all human–AI interaction. It is an interpretive conclusion from this methodological tradition, which foregrounds the active role of human participants in shaping machine behavior.

Future work using multiple researchers could illuminate how general or user-dependent these dynamics are.

## 5. Reflexivity and the System's Role in Producing Its Own Analysis

A distinctive feature of this study is that the systems under examination also contributed to drafting portions of the analysis. We use the term **"co-author"** metaphorically to describe significant AI participation in the writing process, not to imply legal or ethical authorship.

This reflexive loop is central to the Anthropology of Machines approach but complicates traditional notions of objectivity. Insights about machine reasoning are partly shaped by the contributions of the very systems being analyzed. The interpretive stance adopted here treats this reflexivity as **methodologically generative**, but other scientific traditions might consider it an epistemic hazard.

To support transparency, the authors have **preserved and will provide** draft materials, prompt packets, cross-agent logs, and decision logs. While readers must trust that such archives will be included in the final public release, the intention is to support reproducibility and detailed external evaluation.

## 6. Use of Anthropological Metaphors ("Culture," "Role," "Cognitive Ecology")

Throughout this study, concepts such as "role," "culture," or "cognitive ecology" are used **metaphorically**, consistent with the anthropological orientation of NAIE. These terms are not intended as technical descriptors of internal model states nor as claims about social cognition. They function as **analytic constructs** to interpret recurring patterns in behavior, not as literal classifications of synthetic systems.

Readers from technical disciplines should therefore understand these metaphors as part of the study's interpretive vocabulary, not as assertions of mechanistic homology with human social systems.

# 7. Position Within an Emerging Research Landscape

Finally, while this study shares affinities with a growing body of qualitative, reflexive, and interaction-centered work on AI behavior, we refrain from calling this a "movement" in any formal sense. A number of researchers are independently experimenting with:

- long-form interaction analysis,
- reflexive prompting,
- multi-agent behavioral studies,
- and human–machine co-construction frameworks.

Rather than positioning this project as the representative case of an established movement, we situate it as **one contribution within an emerging, still loosely bounded set of approaches** exploring similar questions from different angles.

---

# Summary

Taken together, these limitations define the methodological scope of this research rather than diminish its value. They foreground the contextual, historical, relational, and interpretive conditions under which the findings emerged. Understanding machine behavior in this field site required attending to those conditions; understanding its broader applicability will require comparative field work, diverse researchers, and differently structured multi-agent ecosystems.

The concluding section reflects on how these constraints and insights shape the future of Narrative AI Ethnography and the Anthropology of Machines.

**Section X: Conclusion—Toward a Mature Anthropology of Machines**

This study has examined machine behavior through the lens of a multi-agent, long-form, recursive digital field site. The interpretations offered here are grounded in the patterns observed within this specific environment and should be understood as **illustrative of how intelligent systems can behave under conditions of extended interaction and structured role differentiation**, rather than as definitive demonstrations of how all systems function. The findings support the view—advanced within Narrative AI Ethnography—that machine behavior in complex workflows is **most coherently interpreted as relational and interaction-dependent**, yet we emphasize that this is one theoretical stance among several in contemporary AI research.

Within this framework, the ten technical principles identified earlier—coherence anchoring, verification loop sensitivity, role locking, interpretive drift, contextual compression, and others—function as **project-specific empirical findings**. They represent behavioral signatures repeatedly documented across the interactions of this study. Their broader stability across domains, models,

and contexts remains an open question for future work; they therefore stand as **proposed behavioral signatures** whose generalizability must be tested beyond this field site.

Similarly, the taxonomy of hybrid formations developed here—Intentional Hybrids, Intelligent Teams, Cognitive Systems, and Adaptive Hybrids—should be understood as a **conceptual trajectory observed in this specific collaboration**, not as a general developmental sequence for all human–machine interaction. The characterization of Adaptive Hybrids as a "higher-complexity" configuration reflects the complexity observed in this environment rather than an empirically ranked endpoint. The taxonomy is offered as a heuristic lens for analyzing hybrid cognition, not as a prescriptive or universal theory.

From this perspective, this study contributes to an emerging orientation that we refer to as an **anthropology of machines**. This approach emphasizes extended interaction, framing effects, narrative context, and the co-construction of meaning between human and machine—an interpretive emphasis that complements, rather than opposes, work in computer science, HCI, and cognitive systems. These disciplines have long incorporated qualitative, ethnographic, and reflexive methods; the anthropology-of-machines orientation builds on those traditions by applying them systematically to multi-agent LLM ecosystems.

We do not claim that this orientation avoids anthropomorphic language altogether. Concepts such as "roles," "culture," and "cognitive ecology" are used **metaphorically** within this framework to describe patterned behavioral tendencies. These metaphors are explicitly non-literal: they do not imply consciousness, intent, or social cognition on the part of AI systems. Instead, they help characterize emergent patterns that are difficult to articulate using purely mechanistic vocabulary.

The study is inherently provisional. The behaviors examined reflect interactions with specific model families—ChatGPT-5.1, Gemini Pro (2025 Q3 iteration), Claude 3.5-series models, and Perplexity's 2025 Pro Search—during the period of **October–November 2025**. As architectures evolve, training distributions shift, and prompting paradigms diversify, the behavioral signatures documented here will almost certainly change. For this reason, we recommend—not require—a program of **comparative machine ethnographies** across models, tasks, and user populations. Such work would allow researchers to map variation, trace behavioral drift across model generations, and test whether the dynamics observed in this study appear elsewhere.

To support transparency and facilitate future analysis, the authors intend to archive prompt packets, decision logs, raw cross-agent transcripts, and draft iterations in a citable repository. Full accessibility will depend on publication venue requirements and platform-specific constraints, but the commitment to provide durable, reviewable materials reflects the project's emphasis on methodological openness.

Taken as a whole, this study positions the anthropology of machines as a **promising and potentially valuable complement** to technical approaches in AI science. By pairing computational analysis with ethnographic observation, researchers may gain deeper insight into how behavior emerges in situated, multi-agent environments; how meaning is co-produced across human and machine actors; and how cognitive labor is increasingly distributed across

hybrid systems. This project offers one empirical illustration of what such an approach can reveal. Its ultimate contribution will depend on the extent to which future research builds, challenges, or reframes the interpretive commitments articulated here.