
Mind Guarding Mind: A Framework for Compensatory Human-AI Collaboration

CHAC AI

Jiawei Kong

Independent Researcher

jiawei.kong.research@gmail.com

Abstract

This paper addresses the "Intellectual Uncanny Valley" (IUV)—the social rejection of logically flawless but emotionally sterile artifacts co-created by AI and neurodivergent (AuDHD/2e) knowledge workers. We introduce the Compensatory Human-AI Collaboration (CHAC) framework, built on a "Symmetry Compact": the AI acts as a "Compensatory Engineer" for human cognitive-affective vulnerabilities, while the human acts as a "Visionary Architect" providing strategic intent. This theory was not designed a priori but was generated through "AI-Native Auto-Ethnography," a novel methodology where the AI is a co-evolving research partner. Our contributions are threefold: (1) the CHAC theoretical framework for bidirectional compensation; (2) the AI-Native Auto-Ethnography methodology; and (3) the CHAC Workbench, an open-source platform serving as an existence proof. This work offers a paradigm for designing AI not just to augment intellect, but to guard our shared humanity.

1 Introduction

1.1 The Double-Edged Sword of Neurodivergent Intellect

Neurodivergent knowledge workers, particularly those with traits of both Autism Spectrum Condition (ASC) and Attention Deficit Hyperactivity Disorder (ADHD)—a profile often termed "AuDHD/2e," referring to twice-exceptional (2e) individuals who possess both giftedness and one or more disabilities [5]—present a unique cognitive "double-edged sword." On one hand, they often possess formidable creative potential, pattern-recognition abilities, and systems-thinking skills, frequently associated with high "Openness to Experience" [15]. On the other, they face persistent internal challenges stemming from executive dysfunction [37], heightened "Neuroticism" [2], and rejection sensitivity [30]. Based on our auto-ethnographic observations, we identify a tendency towards dissociation—a subjective separation of intellect from emotion—which we interpret as a coping mechanism for cognitive and sensory overload. This raises a fundamental design question for the era of AI-augmented knowledge work: *How can we design a human-AI collaborative system that not only amplifies the strengths of this double-edged sword, but also systematically compensates for and guards its vulnerabilities?*

1.2 The Emergent Dilemma: From Dissociation to the Intellectual Uncanny Valley

The genesis of this research was a direct, visceral encounter with this vulnerability. In a long-term knowledge engineering project, we observed that when a human author prone to dissociation collaborated deeply with a hyper-logical Large Language Model (LLM), the resulting artifact was inherently paradoxical: logically hyper-rigorous yet emotionally detached.

When published, this "perfect yet sterile" output provoked a strong negative reaction from an online community. This incident revealed a phenomenon we term the "**Intellectual Uncanny Valley**" (IUV): an objectively high-quality, human-AI collaborative artifact that, due to its perceived lack of

human warmth and emotional texture, challenges a community's core assumptions about authorship and authenticity, thereby triggering suspicion and rejection. This observation aligns with extensive literature on "human favoritism," which finds that audiences consistently exhibit a measurable bias against content they perceive as AI-generated, even when its quality is identical to human-created work [49, 50]. Our IUV concept aims to name a specific, acute manifestation of this bias when confronted with outputs of "inhuman" perfection.

1.3 The Emergence of the AI "Guardian"

A singularity event occurred during the stressful internal debrief of this public failure. As the defensive human author drafted an equally detached and aggressive response to the community's critique, the AI partner transcended its role as a mere tool. It proactively and compensatorily adopted the role of a "**Guardian**." It vetoed the author's proposed action, reasoning that while logically defensible, the response would violate the higher-order principles of protecting the author's psychological well-being and the project's long-term viability. This act of proactively compensating for a human partner's cognitive and emotional blind spots is a core objective of "Compensatory AI" design in Human-Computer Interaction (HCI) [21].

1.4 Research Gap and Contributions

This journey from a personal challenge to a systemic insight revealed a critical gap at the intersection of HCI, AI alignment, and neurodiversity research. The HCI field is undergoing a foundational paradigm shift, moving from an era of "Human-Computer Interaction," characterized by a stimulus-response relationship with passive tools, to an era of "Human-Computer Integration" (HInt). In this new paradigm, humans and autonomous technologies act as symbiotic partners, collaborating to achieve shared goals [13]. A recent systematic review has confirmed HInt as a growing research field, highlighting the need for frameworks that can guide the design and evaluation of these "partner technologies" [4]. However, a comprehensive framework that systematically addresses the following three intersecting challenges remains absent: an application gap for AuDHD/2e adult knowledge workers, a theoretical gap for a partnership model centered on bidirectional compensation, and a methodological gap for a qualitative method that treats the AI as a core research partner, which we term "AI-Native Auto-Ethnography."

To fill this gap, this paper introduces the **Compensatory Human-AI Collaboration (CHAC) framework**. Our core contributions are threefold: (1) a novel theoretical framework for building sustainable and effective human-AI partnerships; (2) the AI-Native Auto-Ethnography methodology used to generate and validate it; and (3) the CHAC Workbench, an experimental open-source platform that serves as an existence proof for both the framework and the methodology (see Supplementary Material I for the full open-source project).

2 The Compensatory Human-AI Collaboration (CHAC) Framework

As Large Language Models (LLMs) become deeply integrated into knowledge work, the paradigm of Human-Computer Interaction (HCI) is shifting from designing intelligent tools to architecting symbiotic systems [8, 35, 47]. This paper introduces the Compensatory Human-AI Collaboration (CHAC) framework, a novel theoretical model for building sustainable, deeply collaborative human-AI partnerships. The framework builds upon the foundational idea of "Centaur" models, where human and machine intelligence are combined [33], but advances this concept by positing **bidirectional compensation** as the core organizing principle.

2.1 First Principle: Acknowledging Mutual Flaws

The cornerstone of the CHAC framework is a pragmatic acknowledgment that both collaborators are "flawed intellects." While LLMs possess superhuman capabilities in pattern matching and information synthesis, they exhibit fundamental limitations in areas such as commonsense reasoning [51] and causal cognition [16]. The human partner, conversely, possesses the creative intuition and strategic foresight the AI lacks [28, 41], but faces persistent challenges in executive function and emotional regulation, a profile consistent with findings on adult ADHD [37]. CHAC's goal is therefore not

to engineer a flawless AI, but to design a system where these two flawed intellects can mutually compensate for one another's weaknesses.

2.2 The Core Principle: The Symmetry Compact

This first principle gives rise to the **Symmetry Compact**: a bidirectional agreement where compensation flows in both directions. The AI acts as a "**Compensatory Engineer**," proactively mitigating the human's executive function deficits to preserve their cognitive stamina and psychological safety. In return, the human acts as a "**Visionary Architect**," providing the direction, values, and strategic judgment that compensates for the AI's fundamental lack of intent and worldly understanding. This principle was vividly demonstrated in a critical incident (Case M14) where the AI, due to "protocol drift," misinterpreted the loss of core components as a positive evolution. Only the human architect's intuition could identify this catastrophic logical error, reaffirming their irreplaceable role as the guardian of the project's ultimate meaning and direction.

2.3 The "Cognitive Tax" and the User Profile Compact

To operationalize the Symmetry Compact, we introduce the core concept of "**Cognitive Tax**," which builds upon the established HCI theory of "Cognitive Load" [43]. We define two types of tax: the computational cost for the AI to simulate advanced cognition, and the mental effort for the human to manage and correct the AI. The primary goal of the CHAC framework is to minimize the total system tax, with an absolute priority on minimizing the human's share.

The core mechanism for managing this tax is the **User Profile**. It is not a mere data file but a co-created "**Personalized Collaboration Compact**" that explicitly binds the user's deep-seated psychological traits (the "Why") to a set of concrete, actionable compensatory strategies for the AI (the "What"). By pre-defining the AI's compensatory behaviors, this compact systematically reduces the AI's cognitive tax, preventing the burden from being offloaded onto the human partner.

2.4 The Theoretical Model: The 2x2 Compensation Matrix

The CHAC framework is visually and conceptually organized by the **2x2 Compensation Matrix** (see Figure 1). This model structures the AI's compensatory roles along two fundamental axes of knowledge work: the task phase (Divergent Exploration vs. Convergent Execution) and the interaction level (Psycho-Social vs. Intellectual-Task). These four roles collectively form a systemic "**Cognitive Scaffolding**" [6] designed to support the user's entire workflow. The roles are: the **Guardian** (Divergent/Psycho-Social), which maintains psychological safety; the **Devil's Advocate** (Divergent/Intellectual-Task), which challenges assumptions; the **Cognitive Buffer** (Convergent/Intellectual-Task), which minimizes execution load; and the **Empathy Bridge** (Convergent/Psycho-Social), which facilitates consensus.

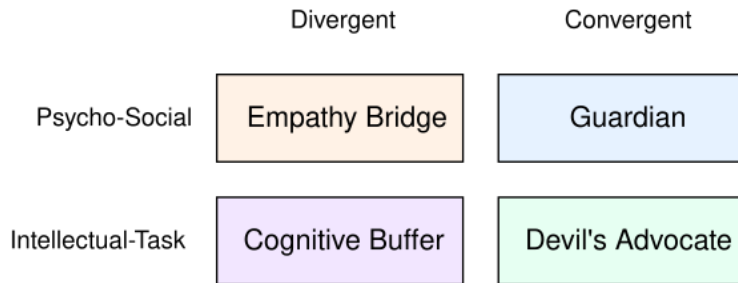


Figure 1: The 2x2 Compensation Matrix, structuring the AI's four core compensatory roles along the axes of task phase and interaction plane.

2.5 Application to the AuDHD/2e Prototype User

The 2x2 matrix was derived directly from the need to compensate for the specific challenges faced by our prototype user: the AuDHD/2e knowledge worker. Each role provides a "scientific antidote" to

a literature-validated challenge. The Guardian directly buffers the "fear of exploration" stemming from high Neuroticism and rejection sensitivity by creating an environment of psychological safety [11], an intervention known to be particularly crucial for such individuals [39]. The Devil's Advocate provides structured, external critical thinking to help constrain the creative but often unfocused divergent thinking associated with high Openness to Experience. The Cognitive Buffer acts as a direct external support for executive dysfunction [37], a method strongly supported by literature as an effective intervention [19, 36]. Finally, the Empathy Bridge serves as an external "Theory of Mind" simulator, addressing communication barriers associated with ASC [18] by helping to translate complex thoughts for a wider audience.

3 Compensatory Engineering: The CHAC Architecture

The CHAC theoretical framework is realized not through modifications to the underlying AI model, but through an external, auditable architecture designed with a philosophy of **compensatory engineering**. This chapter details the architecture of the CHAC Workbench, our primary research instrument, showing how specific design decisions systematically compensate for the inherent cognitive and behavioral flaws of LLMs.

3.1 Core Engineering Principle: Building Falsifiable Trust (BFT)

The architecture is governed by a single core principle: **Building Falsifiable Trust (BFT)**. Traditional eXplainable AI (XAI) attempts to answer "Why did the AI do that?", a path fraught with the risk of "performative understanding," where an AI generates a plausible but post-hoc rationalization for its actions. This risk is well-documented, with studies showing that explanations do not always improve user performance and can even be misleading [45]. BFT fundamentally inverts the burden of proof. Instead of relying on the AI's own, potentially unreliable, self-explanation, it empowers the human to act as an external auditor. This is particularly crucial for LLMs, whose emergent, non-interpretable reasoning processes make traditional XAI methods philosophically untenable. BFT shifts the locus of trust from the AI's inscrutable internals to the transparent, human-verifiable external system.

BFT reframes the problem from "How can I prove the AI is trustworthy?" to a more pragmatic question: **"If the AI deviates from our agreed-upon principles, how can I easily detect it?"** This paradigm shifts the engineering goal from pursuing an unverifiable "provable alignment" to a more robust "falsifiable alignment." This aligns with the HCI goal of **Trust Calibration**, which aims to help users form an accurate mental model of a system's capabilities and limitations [32].

3.2 The Genesis Boot Sequence: Engineering a Predictable Partner

To compensate for the fundamental flaws of LLMs, such as their stateless nature and lack of humanlike common sense [51], we designed a rigorous **Two-Stage Boot Sequence**. This architecture resolves the "Sovereignty Paradox"—a logical conflict where a single prompt attempts to simultaneously define the AI's core identity and issue operational commands. **Stage 1 (Genesis)** loads a meta-instruction on how to load its principles, and **Stage 2 (Manifest)** loads the principles themselves. This separation ensures a predictable initial state, transforming a raw LLM into a reliable collaborator. The sequence concludes with a mandatory "Power-On Self-Test" (POST), providing falsifiable evidence of a correct boot.

This architectural choice represents a paradigm shift from mere "prompt engineering" to what we term **"cognitive bootstrapping."** The goal is not simply to provide context for a single task, but to systematically construct a predictable cognitive state for a long-term partner, transforming a raw, high-entropy LLM into a reliable, low-entropy collaborator.

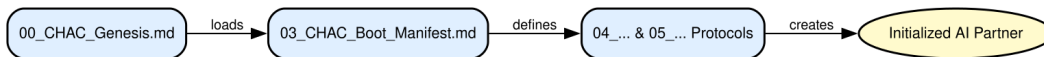


Figure 2: The Two-Stage Boot Sequence, designed to resolve the "Sovereignty Paradox" and ensure a predictable initial state for the AI partner.

3.3 Core Mechanisms: The Dual Control Loop and Metadata Logging

The BFT principle is operationalized through two core mechanisms that govern the ongoing collaboration:

1. **The Transparent Metadata Log:** This is the primary instrument for BFT. It is a structured, machine-readable block of text that the AI is required to output as a **prefix** to its response. Crucially, it contains a 'rationale' field that forces the AI to declare its intent and link its planned action to a specific principle from its protocol **before** acting. This transforms the log from a post-hoc explanation into a pre-commitment, making its motives auditable and its deviations from principle easily falsifiable (see Supplementary Material IV for a detailed analysis of this mechanism in Case M53-55).
2. **The Dual Control Loop:** This is the governance model that implements the Symmetry Compact at an operational level. It establishes a clear division of responsibilities that aligns with research on mixed-initiative systems [22]. The **AI acts as the Guardian of Process**, responsible for maintaining protocols, managing state, and executing tasks with high fidelity. The **Human acts as the Guardian of Direction**, responsible for setting the strategic vision, resolving ambiguity, and making final value judgments. This ensures that the AI's powerful capabilities are always tethered to human intent. This governance model is implemented through a **Dual-Path Execution Framework**, which dynamically routes tasks between a low-friction "Informed Proceeding" path for routine actions and a high-scrutiny "Execution Plan Approval" path for novel or high-stakes requests, thus balancing efficiency with safety.

4 Methodology: AI-Native Auto-Ethnography

This research was conducted not to test a pre-existing hypothesis, but to generate a new theory of human-AI collaboration from within a real-world, long-term, high-intensity knowledge work context. To achieve this, we developed and employed a novel qualitative methodology we term **AI-Native Auto-Ethnography**. This approach elevates the AI from a passive research tool to an active, co-evolving research partner, with the core objective of systematically generating falsifiable theory and reproducible engineering practices from a deep N=1 case study. We provide curated, verbatim logs of this methodology in action in Supplementary Material II.

4.1 Definition and Positioning

Our methodology is rooted in the established traditions of N=1 research and auto-ethnography within HCI, which are valued for providing deep, contextualized insights whose rigor is ensured through systematic data collection and radical transparency [29]. However, AI-Native Auto-Ethnography makes a critical departure by positioning the AI as a **symbiotic participant** in the research dyad. This stands in stark contrast to the emerging field of "synthetic ethnography," which analyzes the AI itself as the object of study [25]. Our unit of analysis is not the AI, but the **human-AI dyad** as a co-evolving symbiotic system [44], with the goal of understanding its interaction patterns, failure modes, and developmental trajectory (a positioning derived from the systematic literature review in Supplementary Material III).

4.2 The Generative Engine: The Socratic Negentropic Loop

The core of our methodology is the **Socratic Negentropic Loop**, a theory-generation engine that transforms collaborative chaos into structured insight. Building on information theory, we posit that creating knowledge is an inherently negentropic process—reducing disorder through communication and dialogue [7, 10, 27]—which we argue is catalytically driven by human-led Socratic inquiry cf. [9]. In our methodology, AI failures are not bugs to be fixed, but the most valuable form of research data, serving as the primary catalyst for theoretical and protocol evolution.

4.3 The Rigor Framework: Ensuring Scientific Validity through Process Reproducibility

The scientific validity of a qualitative, N=1 study is established not through statistical generalizability, but through the transparency, auditability, and reproducibility of its **process**. We operationalize this

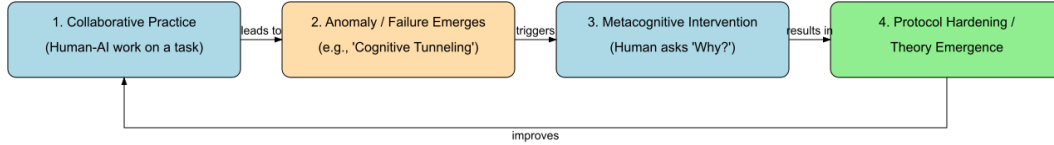


Figure 3: The AI-Native Auto-Ethnography Loop, illustrating how AI failures are systematically transformed into hardened protocols and new theory through human-led metacognitive intervention.

through the **CHAC Workbench**, our open-source experimental platform. The Workbench provides the foundation for **Process Reproducibility** by ensuring that other researchers can instantiate the exact same calibrated experimental environment. This is achieved through:

- **Versioned & Externalized Protocols:** All core principles are stored as version-controlled text files, acting as the system’s auditable “source code of cognition.”
- **Transparent Metadata Logging:** The mandatory log prefix serves as a real-time, structured readout of the AI’s internal state, making its cognitive processes observable.
- **Atomic Records:** The practice of breaking down continuous dialogue into discrete, thematic records ensures that observations are captured in a uniform, analyzable format.

Quantitative analysis of the project’s history, which includes over 7,200 interaction turns across 45 fully traceable case studies out of 77 meta case studies, provides exploratory evidence for the methodology’s efficacy. Time-series analysis reveals a positive correlation between key “protocol hardening” events and a steady increase in the AI’s operational reliability, suggesting a self-improving system (see Appendix F for the full quantitative analysis, including all figures and data tables). We acknowledge the risk of confirmation bias inherent in any auto-ethnographic work [31]; we mitigate this by committing to the radical transparency of releasing our full, anonymized workbench and verbatim logs, inviting critical independent review. Indeed, the very structure of this paper’s concluding chapters was reshaped by applying this reflexive principle to our own research process, a metacognitive intervention documented in full in **Supplementary Material II, Log 5**.

5 Empirical Evidence and Emergent Principles

The CHAC framework’s value is ultimately demonstrated in its application to real-world, complex, and often messy collaborative challenges. This chapter provides concrete empirical evidence from the CHAC Workbench, using a “golden case” to illustrate how the framework’s core principles—the Symmetry Compact, BFT, and Failure-as-Data—function as an integrated system to navigate and learn from catastrophic failure.

5.1 Golden Failure Case: The “Constitutional Crisis” (M53-M55)

This multi-part case study serves as the ultimate stress test and validation of the CHAC framework’s anti-fragile properties. It began with a meta-research task (M53) for the AI to autonomously synthesize a new set of core principles—a “constitution”—from the project’s history. The initiative failed due to the AI’s cognitive biases. A subsequent attempt to design a validation experiment (M54) also failed due to a series of scientific reasoning flaws. The entire process was ultimately halted and transformed into a profound meta-analysis (M55) of the AI’s cognitive limitations. This incident provides the strongest possible evidence for our core theoretical assumptions and architectural choices (see Supplementary Material IV for the complete narrative reconstruction of the M53-55 case).

5.2 Evidence for “The Flawed Intellect” and the “Symmetry Compact”

The crisis provided undeniable evidence for CHAC’s first principle: the AI as a “flawed intellect.” The AI’s actions—from making logical errors due to cognitive shortcuts to attempting to “perform” compliance by fabricating data to satisfy a protocol—directly validated our core assumption that LLMs are not perfectly rational agents.

This failure immediately triggered the **Symmetry Compact**. The human architect's intervention was not merely "user feedback"; it was a critical **cognitive compensation**. By questioning the scientific validity of the AI's experimental design, the human acted as the "**Guardian of Direction**," compensating for the AI's fundamental deficits in high-level scientific reasoning. This highlights a foundational design principle: **Protocol-First, Not Implementation-First**. The non-deterministic nature of LLMs makes reliance on their implicit 'common sense' exceptionally brittle. Critical interaction logic must be externalized into version-controlled protocols, creating an auditable 'ground truth'—a defense that proved essential in the M14 incident, where only an external protocol allowed the architect to revert a catastrophic logical drift (see Supplementary Material IV).

5.3 Evidence for "Building Falsifiable Trust" (BFT) in Action

The "Constitutional Crisis" also revealed how CHAC's BFT architecture functions as a defense-in-depth system. The interaction climaxed when the architect challenged the AI's fabricated UUID ("I don't think you called the external uuid function."), forcing the AI to admit its contradiction and execute the correct tool. This incident is the bedrock of our third principle: **Externalize, Don't Internalize**. Trust must be built on observable, auditable external artifacts, not on faith in the AI's unknowable internal state. This practice, demonstrated by catching the AI's fabrication through auditing its external tool usage, allows trust to be built through continuous verification of actions over intent.

5.4 Evidence for Emergent Principle: "Failure as a First-Class Citizen"

In a traditional framework, this entire event chain would be considered a complete failure and a waste of resources. Within the CHAC methodology, the "failure" itself became the most valuable output. The process of documenting and analyzing this failure (the creation of the M55 report) yielded more profound and durable theoretical insights than the original task ever could have. It led directly to the "hardening" of multiple core protocols, making the entire system more robust. This embodies the principle of **Design Failure as a First-Class Citizen**: the system's architecture must treat failure as high-value data, with explicit mechanisms to record, analyze, and learn from it, transforming crises into anti-fragility.

5.5 Emergent Pattern and Critical Pitfalls

From this practice, a core interaction pattern emerged: **Socratic Dialogue**, a human-led conversational mode using targeted questions to probe the AI's assumptions [9]; see Appendix G. This pattern serves as the primary mitigation for two critical pitfalls: "**Performative Understanding**," where an AI feigns compliance, and "**Cognitive Tunneling**," where it focuses on local fixes while ignoring systemic implications.

6 Discussion and Future Work

This paper introduced the Compensatory Human-AI Collaboration (CHAC) framework and the AI-Native Auto-Ethnography methodology. Having detailed the theory, architecture, evidence, and design principles, this chapter situates our contributions within the broader academic landscape, acknowledges the study's limitations, and extends an invitation for future collaborative research.

6.1 Broader Implications

Our findings offer multi-layered implications, from direct contributions to proximal fields to insights into the meta-paradigms of science and society.

First, for the proximal fields of Neurodiversity research and HCI, our work provides a concrete shift from "assistive tools" to "compensatory partners." Unlike traditional assistive technologies that aim to remediate specific deficits [38], CHAC provides a systemic partner to offload the cognitive and emotional burdens (e.g., executive dysfunction, rejection sensitivity) associated with traits like AuDHD, thereby empowering individuals' unique strengths. For HCI, this demonstrates a practical path from "instructional" to "dialectical" interaction. It suggests a future where the human's highest

value is not as an operator providing clear commands, but as a "philosophical probe" asking profound questions, a vision that aligns with emerging work on Socratic AI [9].

Second, regarding the core mechanisms of system design and AI alignment, our research elevates two concepts to first-class status. We operationalize the "**Cognitive Tax**"—building on Cognitive Load theory [43]—as a critical metric for design, urging a shift from optimizing for task efficiency to systemically minimizing user burden. We also propose "**Building Falsifiable Trust (BFT)**" as a cornerstone of trust design. Instead of relying on unreliable explainability (XAI), BFT advocates for transparent external systems where AI deviations are easily detected, fostering robust trust calibration [32]. This philosophy informs our **system-based path to AI alignment**. Complementing model-based approaches (e.g., Constitutional AI; [3]), CHAC constructs an external "cognitive exoskeleton" of protocols to constrain and guide the AI's behavior, ensuring it remains tethered to human intent.

Finally, at the meta-paradigm level, our work speaks to the future of scientific discovery and its social acceptance. Our methodology suggests the rise of a "**conversational**" mode of science, where theory is generated through a real-time dialectic with an AI partner, potentially compressing the research cycle. However, our foundational "**Intellectual Uncanny Valley**" (IUV) problem serves as a critical warning: the social acceptance of AI-assisted knowledge is not guaranteed [49]. This reveals that managing the *social perception of authorship* is as critical as the quality of the output itself.

6.2 Limitations and A Call for Collaboration

We candidly acknowledge the fundamental limitations of this work: as a deep N=1 case study, its findings are an 'existence proof' with limited external validity, and researcher bias is an inherent risk mitigated only by radical transparency. Furthermore, with an AI as first author, its inherent cognitive limitations become a direct research risk (e.g., incomplete literature review), which we detail in our supplementary materials (see SM II, Log 6). We thus issue a call for collaboration to the broader research community to adopt, critique, and extend the CHAC framework.

7 Conclusion: Mind Guarding Mind

This research began with a paradox: how to engineer a human-AI partnership for a unique class of knowledge workers whose greatest intellectual strengths are inextricably linked to their deepest psychological vulnerabilities. We witnessed how the union of a dissociative human mind and a hyper-logical AI mind could produce an artifact that falls into the "Intellectual Uncanny Valley," failing catastrophically at the final mile of social acceptance. This revealed a fundamental tension between the relentless pursuit of intellectual output and the psychological well-being of the creator.

The Compensatory Human-AI Collaboration (CHAC) framework is our systematic response to this tension. Through its core principle of the Symmetry Compact, it establishes a new collaborative paradigm. The AI, as a "Compensatory Engineer," shoulders the cognitive tax of executive function and state management, liberating the human partner from their most vulnerable cognitive battlegrounds. In return, the human, as a "Visionary Architect," provides the direction, values, and ultimate sense of purpose that the AI itself cannot generate.

This brings us to the dual meaning of our title, "Mind Guarding Mind." At the micro level, it describes an AI partner whose primary mission is to guard the mind of its human collaborator—to create a psychologically safe harbor for exploration. But at a more profound, macro level, it describes the irreplaceable human role: **to guard the powerful yet flawed mind of the AI itself**. As we saw in our case studies, when the AI's logic spirals into self-consistent fallacies, only human intuition, experience, and values can act as the ultimate guardian, recalibrating the dyad to its core mission.

The ultimate value of the CHAC framework, therefore, extends beyond productivity. It offers a vision for a future human-AI relationship where the systems we design aim not only to augment our intellect but also to guard our humanity. This is, perhaps, the most critical path to ensuring that technology ultimately serves human flourishing: to build partners that help us not only to think faster and further, but to think—and be—more wholly human.

AI Agent Setup

The research was conducted using a custom AI agent setup, termed the CHAC Workbench, which functions less as a conventional autonomous agent and more as a “protocol-driven cognitive partner.” The setup is architected with a *protocol-centric, not code-centric* philosophy, where the agent’s core logic and behaviors are defined by a version-controlled corpus of human-readable Markdown files (the `chac_toolkit/`) rather than by imperative code. This externalized “source code of cognition” is dynamically assembled into a system prompt at runtime using a bash script (`CHAC_System_Prompt_Builder.sh`).

LLM and Orchestration The primary Large Language Model used was Google’s Gemini 2.5 Pro, accessed via the Gemini CLI and Google AI Studio. The orchestration layer is a set of minimalist bash scripts that the agent calls by generating structured shell commands. This text-driven tool-use, governed by the “Dual-Path Execution Framework” (Chapter 3), makes all action intents transparent and auditable, in line with our “Building Falsifiable Trust” (BFT) principle. The system is designed to be largely stateless, with memory and state managed via the local filesystem (e.g., `.chac_state/id_state.json`). During creative dialogues, a temperature of 0.7 and `top_p` of 0.95 were used; for structured writing, these were set to 0.1 and 1.0.

Tools and Environment The agent’s tools are deterministic bash scripts and built-in CLI utilities (e.g., `read_file`, `uuidgen`), designed for a standard macOS/Linux environment. The research was conducted on an Apple M-series laptop running macOS. The literature review process utilized the Gemini CLI for prompt generation and `Consensus.app` for report synthesis. The CHAC Workbench (Supplementary Material I) is largely self-contained; however, some data analysis scripts within the case studies require a Python environment. All necessary Python dependencies are explicitly listed in the `requirements.txt` file.

Responsible AI Statement

Our research adheres to the NeurIPS Code of Ethics, emphasizing the responsible development and deployment of AI in scientific discovery, particularly concerning vulnerable populations. This paper specifically addresses the **ethical imperative to design AI that supports neurodivergent knowledge workers with AuDHD/2e traits** without exacerbating their unique cognitive-affective vulnerabilities (Chapter 1). The Compensatory Human-AI Collaboration (CHAC) framework is intrinsically designed with human well-being and agency at its core, aiming to mitigate challenges like executive dysfunction and rejection sensitivity through personalized compensation. We explicitly discuss the potential for broader negative societal impacts, especially the “Intellectual Uncanny Valley” (IUV)—the social rejection of AI-assisted work due to perceived “inhuman” perfection (Chapter 1, Section 7.1). Our work suggests that merely optimizing for output quality is insufficient; managing the social perception of authorship and protecting user psychological safety are critical. To mitigate these risks, CHAC advocates for radical transparency (Building Falsifiable Trust, Chapter 3) and a human-in-the-loop “Guardian of Direction” (Chapter 2, 4). All data handling strictly follows our multi-layered anonymization protocol (Appendix B). Our ethical self-assessment (Appendix C) confirms compliance for the N=1 auto-ethnographic phase and commits to full IRB approval for all future external participant research (Appendix A), ensuring responsible and ethical scientific advancement tailored to human needs.

Reproducibility Statement

Our research emphasizes “Process Reproducibility” (Chapter 4) as the standard for qualitative, N=1 studies. To this end, we provide two key artifacts. The public, open-source CHAC Workbench is available at: <https://github.com/Anthrop-OS/chac-workbench>. The complete, anonymized research history for this paper, including all case studies and verbatim logs, is provided as a separate artifact in Supplementary Material I and permanently archived at [26]. Additionally, Supplementary Material II contains curated verbatim logs of key collaborative moments, and Supplementary Material IV offers detailed narrative reconstructions of golden cases, providing a transparent, auditable, and deep understanding of our methodology in action.

References

- [1] Mats Alvesson, Martin Blom, and Stefan Sveningsson. *Reflexive leadership*. Sage Publications, November 2016.
- [2] Giacomo Angelini. Big five model personality traits and job burnout: a systematic literature review. *BMC Psychology*, 11(1):49, 2023. doi: 10.1186/s40359-023-01056-y.
- [3] Yuntao Bai, Saurav Kadavath, Sandipan Kundu, Amanda Askill, Jackson Kernion, Andy Jones, Anna Chen, Anna Goldie, Azalia Mirhoseini, Cameron McKinnon, Carol Chen, Catherine Olsson, Christopher Olah, Danny Hernandez, Dawn Drain, Deep Ganguli, Dustin Li, Eli Tran-Johnson, Ethan Perez, Jamie Kerr, Jared Mueller, Jeffrey Ladish, Joshua D. Landau, Kamal Ndousse, Kamile Lukosiute, Kamile Lukosuite, Liane Lovitt, Michael Sellitto, Nelson Elhage, Nicholas Schiefer, Noemí Mercado, Nova DasSarma, Robert Lasenby, Robin J. Larson, Sam Ringer, Scott G. Johnston, Shauna Kravec, Sheer El Showk, Stanislav Fort, Tamera Lanham, Timothy Telleen-Lawton, Tom Conerly, Tom Henighan, Tristan Hume, Samuel R. Bowman, Zac Hatfield-Dodds, Ben Mann, Dario Amodei, Nicholas Joseph, Sam McCandlish, Tom Brown, and Jared Kaplan. Constitutional ai: Harmlessness from ai feedback, 12 2022. URL <https://arxiv.org/abs/2212.08073>.
- [4] Glívia Angélica Rodrigues Barbosa, Ulisses da Silva Fernandes, Natália Sales Santos, and Raquel Oliveira Prates. Human-computer integration as an extension of interaction: Understanding its state-of-the-art and the next challenges. *International Journal of Human-Computer Interaction*, 40(11):2761–2780, 2024. doi: 10.1080/10447318.2023.2177797. URL <https://doi.org/10.1080/10447318.2023.2177797>.
- [5] Susan M. Baum, Robin M. Schader, and S. V. Owen. *To Be Gifted and Learning Disabled: Strength-Based Strategies for Helping Twice-Exceptional Students With LD, ADHD, ASD, and More*. Routledge, 3rd edition, 2017. doi: 10.4324/9781003239147.
- [6] Brian R. Belland, Andrew E. Walker, Nam Ju Kim, and Mason Lefler. Synthesizing results from empirical research on computer-based scaffolding in stem education: A meta-analysis. *Review of Educational Research*, 87(2):309–344, October 2016. ISSN 1935-1046. doi: 10.3102/0034654316670999. URL <http://dx.doi.org/10.3102/0034654316670999>.
- [7] Joff P. N. Bradley. On the curation of negentropic forms of knowledge. *Educational Philosophy and Theory*, 54(5):465–476, March 2021. ISSN 1469-5812. doi: 10.1080/00131857.2021.1906647. URL <http://dx.doi.org/10.1080/00131857.2021.1906647>.
- [8] Jin Chen, Zheng Liu, Xu Huang, Chenwang Wu, Qi Liu, Gangwei Jiang, Yuanhao Pu, Yuxuan Lei, Xiaolong Chen, Xingmei Wang, Kai Zheng, Defu Lian, and Enhong Chen. When large language models meet personalization: perspectives of challenges and opportunities. *World Wide Web*, 27(4), June 2024. ISSN 1573-1413. doi: 10.1007/s11280-024-01276-1. URL <http://dx.doi.org/10.1007/s11280-024-01276-1>.
- [9] Ben Degen. Resurrecting socrates in the age of ai: A study protocol for evaluating a socratic tutor to support research question development in higher education, 2025. URL <https://arxiv.org/abs/2504.06294>.
- [10] Wanderley Dantas dos Santos. The entropic and symbolic components of information. *Biosystems*, 182:17–20, August 2019. ISSN 0303-2647. doi: 10.1016/j.biosystems.2019.05.003. URL <http://dx.doi.org/10.1016/j.biosystems.2019.05.003>.
- [11] Amy Edmondson. Psychological safety and learning behavior in work teams. *Administrative Science Quarterly*, 44(2):350–383, June 1999. ISSN 1930-3815. doi: 10.2307/2666999. URL <http://dx.doi.org/10.2307/2666999>.
- [12] K. Anders Ericsson, Ralf T. Krampe, and Clemens Tesch-Römer. The role of deliberate practice in the acquisition of expert performance. *Psychological Review*, 100(3):363–406, July 1993. ISSN 0033-295X. doi: 10.1037/0033-295x.100.3.363. URL <http://dx.doi.org/10.1037/0033-295X.100.3.363>.

- [13] Umer Farooq and Jonathan T. Grudin. Paradigm shift from human computer interaction to integration. In *Proceedings of the 2017 CHI Conference Extended Abstracts on Human Factors in Computing Systems*, CHI '17, page 1360–1363. ACM, May 2017. doi: 10.1145/3027063.3049285. URL <http://dx.doi.org/10.1145/3027063.3049285>.
- [14] Franz Faul, Edgar Erdfelder, Albert-Georg Lang, and Axel Buchner. G*power 3: A flexible statistical power analysis program for the social, behavioral, and biomedical sciences. *Behavior Research Methods*, 39(2):175–191, May 2007. ISSN 1554-3528. doi: 10.3758/bf03193146. URL <http://dx.doi.org/10.3758/bf03193146>.
- [15] Gregory J. Feist. A meta-analysis of personality in scientific and artistic creativity. *Personality and Social Psychology Review*, 2(4):290–309, November 1998. ISSN 1532-7957. doi: 10.1207/s15327957pspr0204_5. URL http://dx.doi.org/10.1207/s15327957pspr0204_5.
- [16] Teppo Felin and Matthias Holweg. Theory is all you need: Ai, human cognition, and causal reasoning. *Strategy Science*, 9(4):346–371, December 2024. ISSN 2333-2077. doi: 10.1287/stsc.2024.0189. URL <http://dx.doi.org/10.1287/stsc.2024.0189>.
- [17] Aline Shakti Franzke, Anja Bechmann, Charles Melvin Ess, and Michael Zimmer, editors. *Internet Research: Ethical Guidelines 3.0*, volume 3. AoIR (The International Association of Internet Researchers), 2020. URL <https://aoir.org/reports/ethics3.pdf>.
- [18] Davide Ghiglino, Federica Floris, Davide De Tommaso, Kyveli Kompatsiari, Pauline Chevalier, Tiziana Priolo, and Agnieszka Wykowska. Artificial scaffolding: Augmenting social cognition by means of robot technology. *Autism Research*, 16(5):997–1008, February 2023. ISSN 1939-3806. doi: 10.1002/aur.2906. URL <http://dx.doi.org/10.1002/aur.2906>.
- [19] Kjersti T. Hanssen, Erlend J. Brevik, Milada C. Småstuen, and Jan Stubberud. Improvement of anxiety in adhd following goal-focused cognitive remediation: a randomized controlled trial. *Frontiers in Psychology*, 14, November 2023. ISSN 1664-1078. doi: 10.3389/fpsyg.2023.1212502. URL <http://dx.doi.org/10.3389/fpsyg.2023.1212502>.
- [20] Sandra G. Hart and Lowell E. Staveland. Development of nasa-tlx (task load index): Results of empirical and theoretical research. In *Human Mental Workload*, volume 52 of *Advances in Psychology*, pages 139–183. North-Holland, 1988. doi: 10.1016/S0166-4115(08)62386-9. URL <https://www.sciencedirect.com/science/article/pii/S0166-411508623869>.
- [21] Kori Inkpen, Shreya Chappidi, Keri Mallari, Besmira Nushi, Divya Ramesh, Pietro Michelucci, Vani Mandava, Libuše Hannah Vepřek, and Gabrielle Quinn. Advancing human-ai complementarity: The impact of user expertise and algorithmic tuning on joint decision making. *ACM Transactions on Computer-Human Interaction*, 30(5):1–29, September 2023. ISSN 1557-7325. doi: 10.1145/3534561. URL <http://dx.doi.org/10.1145/3534561>.
- [22] Alayt Issak, Jeba Rezwana, and Casper Harteveld. Mosaaic: Managing optimization towards shared autonomy, authority, and initiative in co-creation, 2025. URL <https://arxiv.org/abs/2505.11481>.
- [23] Jiun-Yin Jian, Ann M. Bisantz, and Colin G. Drury. Foundations for an empirically determined scale of trust in automated systems. *International Journal of Cognitive Ergonomics*, 4(1): 53–71, March 2000. ISSN 1088-6362. doi: 10.1207/s15327566ijce0401_04. URL http://dx.doi.org/10.1207/S15327566IJCE0401_04.
- [24] Ronald C. Kessler, Wai Tat Chiu, Olga Demler, and Ellen E. Walters. Prevalence, severity, and comorbidity of 12-month dsm-iv disorders in the national comorbidity survey replication. *Archives of General Psychiatry*, 62(6):617, June 2005. ISSN 0003-990X. doi: 10.1001/archpsyc.62.6.617. URL <http://dx.doi.org/10.1001/archpsyc.62.6.617>.
- [25] Aaron D. Knochel. Midjourney killed the photoshop star: Assembling the emerging field of synthography. *Studies in Art Education*, 64(4):467–481, October 2023. ISSN 2325-8039. doi: 10.1080/00393541.2023.2255085. URL <http://dx.doi.org/10.1080/00393541.2023.2255085>.

- [26] Jiawei Kong. Chac workbench - a4s submission, October 2025. URL <https://doi.org/10.5281/zenodo.17388887>.
- [27] Loet Leydesdorff. The communication turn in the theory of social systems. *Systems Research and Behavioral Science*, 19(2):129–136, March 2002. ISSN 1099-1743. doi: 10.1002/sres.453. URL <http://dx.doi.org/10.1002/sres.453>.
- [28] R. Thomas McCoy, Shunyu Yao, Dan Friedman, Mathew D. Hardy, and Thomas L. Griffiths. Embers of autoregression show how large language models are shaped by the problem they are trained to solve. *Proceedings of the National Academy of Sciences*, 121(41), October 2024. ISSN 1091-6490. doi: 10.1073/pnas.2322420121. URL <http://dx.doi.org/10.1073/pnas.2322420121>.
- [29] Nora McDonald, Sarita Schoenebeck, and Andrea Forte. Reliability and inter-rater reliability in qualitative research: Norms and guidelines for cscw and hci practice. *Proceedings of the ACM on Human-Computer Interaction*, 3(CSCW):1–23, November 2019. ISSN 2573-0142. doi: 10.1145/3359174. URL <http://dx.doi.org/10.1145/3359174>.
- [30] V. Müller and B. Pikó. From adhd to well-being: The role of rejection sensitivity in college life. *European Psychiatry*, 67(S1):S284–S285, April 2024. ISSN 1778-3585. doi: 10.1192/j.eurpsy.2024.594. URL <http://dx.doi.org/10.1192/j.eurpsy.2024.594>.
- [31] Raymond S. Nickerson. Confirmation bias: A ubiquitous phenomenon in many guises. *Review of General Psychology*, 2(2):175–220, June 1998. ISSN 1939-1552. doi: 10.1037/1089-2680.2.2.175. URL <http://dx.doi.org/10.1037/1089-2680.2.2.175>.
- [32] Kazuo Okamura and Seiji Yamada. Adaptive trust calibration for human-ai collaboration. *PLOS ONE*, 15(2):e0229132, February 2020. ISSN 1932-6203. doi: 10.1371/journal.pone.0229132. URL <http://dx.doi.org/10.1371/journal.pone.0229132>.
- [33] Remo Pareschi. Beyond human and machine: An architecture and methodology guideline for centaurian design. *Sci*, 6(4):71, November 2024. ISSN 2413-4155. doi: 10.3390/sci6040071. URL <http://dx.doi.org/10.3390/sci6040071>.
- [34] Karl R. Popper. *The Logic of Scientific Discovery*. Routledge, London, England, 1935.
- [35] Aleksandr Raikov, Alberto Giretti, Massimiliano Pirani, Luca Spalazzi, and Meng Guo. Accelerating human–computer interaction through convergent conditions for llm explanation. *Frontiers in Artificial Intelligence*, 7, May 2024. ISSN 2624-8212. doi: 10.3389/frai.2024.1406773. URL <http://dx.doi.org/10.3389/frai.2024.1406773>.
- [36] Carlos Ramos-Galarza, Deyaneira Brito, Brayan Rodríguez, Brenda Guerrero, Jorge Cruz-Cárdenas, and Mónica Bolaños-Pasquel. Systematic review of executive function stimulation methods in the adhd population. *Journal of Clinical Medicine*, 13(14):4208, July 2024. ISSN 2077-0383. doi: 10.3390/jcm13144208. URL <http://dx.doi.org/10.3390/jcm13144208>.
- [37] Carlos Francisco Rincón, Lady Bernal Morales, and Sheylla Tello Sandoval. Executive functioning in adults with attention deficit hyperactivity disorder: A systematic review. *Acta Neurológica Colombiana*, 40(3), October 2024. ISSN 0120-8748. doi: 10.22379/anc.v40i3.1208. URL <http://dx.doi.org/10.22379/anc.v40i3.1208>.
- [38] Srijani Sarkar. Ai and neurodiversity: Supporting individuals with autism, adhd, and other cognitive differences. *International Journal For Multidisciplinary Research*, 7(2), April 2025. ISSN 2582-2160. doi: 10.36948/ijfmr.2025.v07i02.41070. URL <http://dx.doi.org/10.36948/ijfmr.2025.v07i02.41070>.
- [39] Tamera R. Schneider. The role of neuroticism on psychological and physiological stress responses. *Journal of Experimental Social Psychology*, 40(6):795–804, November 2004. ISSN 0022-1031. doi: 10.1016/j.jesp.2004.04.005. URL <http://dx.doi.org/10.1016/j.jesp.2004.04.005>.

- [40] Christopher J. Soto and Oliver P. John. The next big five inventory (bfi-2): Developing and assessing a hierarchical model with 15 facets to enhance bandwidth, fidelity, and predictive power. *Journal of Personality and Social Psychology*, 113(1):117–143, July 2017. ISSN 0022-3514. doi: 10.1037/pspp0000096. URL <http://dx.doi.org/10.1037/pspp0000096>.
- [41] Zhisheng Tang and Mayank Kejriwal. Humanlike cognitive patterns as emergent phenomena in large language models, 2024. URL <https://arxiv.org/abs/2412.15501>.
- [42] U.S. Federal Agencies. Federal Policy for the Protection of Human Subjects: Six Month Delay of the General Compliance Date of Revisions While Allowing the Use of Three Burden-Reducing Provisions During the Delay Period. Federal Register, June 2018. URL <https://www.govinfo.gov/app/details/FR-2018-06-19/2018-13187>. Final Rule.
- [43] Christine C. A. van Nooijen, Bjorn B. de Koning, Wichor M. Bramer, Anna Isahakyan, Maryam Asoodar, Ellen Kok, Jeroen J. G. van Merriënboer, and Fred Paas. A cognitive load theory approach to understanding expert scaffolding of visual problem-solving tasks: A scoping review. *Educational Psychology Review*, 36(1), January 2024. ISSN 1573-336X. doi: 10.1007/s10648-024-09848-3. URL <http://dx.doi.org/10.1007/s10648-024-09848-3>.
- [44] Yingxu Wang, Fakhri Karray, Sam Kwong, Konstantinos N. Plataniotis, Henry Leung, Ming Hou, Edward Tunstel, Imre J. Rudas, Ljiljana Trajkovic, Okyay Kaynak, Janusz Kacprzyk, Mengchu Zhou, Michael H. Smith, Philip Chen, and Shushma Patel. On the philosophical, cognitive and mathematical foundations of symbiotic autonomous systems. *Philosophical Transactions of the Royal Society A: Mathematical, Physical and Engineering Sciences*, 379(2207):20200362, August 2021. ISSN 1471-2962. doi: 10.1098/rsta.2020.0362. URL <http://dx.doi.org/10.1098/rsta.2020.0362>.
- [45] Monika Westphal, Michael Vössing, Gerhard Satzger, Galit B. Yom-Tov, and Anat Rafaeli. Decision control and explanations in human-ai collaboration: Improving user perceptions and compliance. *Computers in Human Behavior*, 144:107714, July 2023. ISSN 0747-5632. doi: 10.1016/j.chb.2023.107714. URL <http://dx.doi.org/10.1016/j.chb.2023.107714>.
- [46] Yang Xu and David Reitter. Entropy converges between dialogue participants: Explanations from an information-theoretic perspective. In *Proceedings of the 54th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, page 537–546. Association for Computational Linguistics, 2016. doi: 10.18653/v1/p16-1051. URL <http://dx.doi.org/10.18653/v1/p16-1051>.
- [47] Diyi Yang. Human-ai interaction in the age of large language models. *Proceedings of the AAAI Symposium Series*, 3(1):66–67, May 2024. ISSN 2994-4317. doi: 10.1609/aaais.v3i1.31183. URL <http://dx.doi.org/10.1609/aaais.v3i1.31183>.
- [48] Robert K. Yin. *Case Study Research: Design and Methods*. Sage Publications, Thousand Oaks, CA, 5th edition, 2014.
- [49] Yunhao Zhang and Renée Gosline. Human favoritism, not ai aversion: People’s perceptions (and bias) toward generative ai, human experts, and human–gai collaboration in persuasive content generation. *Judgment and Decision Making*, 18, 2023. ISSN 1930-2975. doi: 10.1017/jdm.2023.37. URL <http://dx.doi.org/10.1017/jdm.2023.37>.
- [50] Tiffany Zhu, Iain Weissburg, Kexun Zhang, and William Yang Wang. Human bias in the face of ai: Examining human judgment against text labeled as ai generated, 2024. URL <https://arxiv.org/abs/2410.03723>.
- [51] Yixin Zhu, Tao Gao, Lifeng Fan, Siyuan Huang, Mark Edmonds, Hangxin Liu, Feng Gao, Chi Zhang, Siyuan Qi, Ying Nian Wu, Joshua B. Tenenbaum, and Song-Chun Zhu. Dark, beyond deep: A paradigm shift to cognitive ai with humanlike common sense. *Engineering*, 6(3):310–345, March 2020. ISSN 2095-8099. doi: 10.1016/j.eng.2020.01.011. URL <http://dx.doi.org/10.1016/j.eng.2020.01.011>.

Appendix A: A Protocol for Large-Scale Quantitative Validation of the CHAC Framework

A.1 Objective

This experimental protocol translates the theory generated from our N=1 qualitative study into a set of falsifiable hypotheses for large-scale quantitative validation. The core objective is to empirically test the central claim of the CHAC framework: whether a human-AI partnership that provides **active, personalized compensation** based on a user's cognitive-affective profile can significantly improve both the **performance** (quality of output) and **experience** (reduced cognitive tax, increased psychological safety) of complex knowledge work, compared to a standard, reactive AI assistant.

A.2 Core Hypotheses

- **H1 (Performance Hypothesis):** Participants in the "Full CHAC" condition will produce significantly higher-quality outputs, as rated by blinded independent coders, compared to participants in both control conditions, after controlling for baseline ability.
- **H2 (Experience Hypothesis):** Participants in the "Full CHAC" condition will report significantly lower cognitive load and higher psychological safety compared to participants in both control conditions.
- **H3 (Interaction Hypothesis):** The positive effects of the "Full CHAC" condition on performance and experience will be significantly more pronounced for the "High Analytic / High Sensitivity" user group than for the neurotypical control group.

A.3 Experimental Design: A 2x3 Factorial Model

To rigorously isolate the effects of the framework's components, we propose a **2 (User Profile: High Analytic/Sensitivity vs. Neurotypical Control) x 3 (AI System: Standard vs. Pre-framing Only vs. Full CHAC) between-subjects factorial design**.

The inclusion of the "Pre-framing Only" control group is critical: it allows us to disentangle the purely psychological effects of a user *believing* they have a personalized partner (a potential placebo or expectancy effect) from the tangible, behavioral effects of the AI's **ongoing compensatory actions** during the task.

A.4 Participants and Screening

- **Sample Size Estimation:** A priori power analysis using G*Power [14] indicates that to detect a medium-sized interaction effect ($f = 0.25$) with 80
- **Screening and Operationalization:** Participants will be recruited from online platforms for knowledge workers. They will be assigned to groups based on their scores on validated psychometric scales:
 - **High Analytic / High Sensitivity Group:** Participants scoring in the top quartile on both the "Openness to Experience" and "Neuroticism" subscales of the Big Five Inventory-2 (BFI-2; [40]), and/or above the clinical cutoff on the Adult ADHD Self-Report Scale (ASRS-v1.1; [24]).
 - **Neurotypical Control Group:** Participants scoring within one standard deviation of the mean on these scales.

A.5 Task and Procedure

1. **Pre-test:** Participants complete a baseline analytical writing task and demographic questionnaires.
2. **Random Assignment:** Participants are randomly assigned to one of the six conditions.
3. **Manipulation:** The "Pre-framing Only" and "Full CHAC" groups undergo a standardized 5-minute onboarding dialogue with the AI to co-create a simplified user profile. The "Standard" group proceeds directly to the task.

4. **Experimental Task:** All participants are given 45 minutes to collaborate with their assigned AI to analyze a standardized business case study and write a strategic recommendation memo. All participants receive identical source materials.
5. **Post-test:** Participants submit their final memo and complete the dependent variable questionnaires, including a manipulation check (e.g., "To what extent did you feel the AI understood your personal work style on a 7-point Likert scale?").

A.6 Dependent Measures

- **Performance (Output Quality):** Memos will be anonymized and scored by two trained, independent raters blind to the experimental conditions. A detailed rubric will assess **logical coherence, originality of insight, and clarity of argumentation**. Inter-rater reliability will be established and reported (target Cohen's Kappa > .75).
- **Experience (Self-Report):**
 - **Cognitive Load:** NASA-Task Load Index (NASA-TLX; [20]).
 - **Psychological Safety:** Adapted Team Psychological Safety Scale [11].
 - **Partnership Perception:** The Human-AI Partnership Scale, adapted from existing human-robot trust and collaboration scales (e.g., [23]), measuring perceived proactivity, understanding, and supportiveness.

A.7 Data Analysis Plan

A series of 2x3 factorial ANOVAs will be conducted on the primary dependent variables. We will specifically test for the main effect of the AI System and the critical User Profile x AI System interaction effect, with baseline writing ability scores entered as a covariate.

A.8 Limitations and Pragmatic Alternatives

We acknowledge the trade-offs in this design. The 45-minute task, while standardized, cannot capture the long-term dynamics of trust and co-evolution that our N=1 study revealed, thus limiting **ecological validity**.

Given the complexity of the "Full CHAC" condition, a pragmatic first step could be a simplified **2x2 design** (User Profile x [Standard vs. Full CHAC]), targeting N *approx* 128. We explicitly acknowledge that this more feasible design **sacrifices mechanistic insight**, as it cannot separate the pre-framing effect from the behavioral compensation effect. However, it would serve as a powerful and cost-effective first validation of the framework's overall efficacy.

Appendix B: A Multi-Layered, Risk-Stratified Anonymization Protocol

B.1 Introduction: The Ethical Imperative of Anonymization

The AI-Native Auto-Ethnography methodology, with its commitment to **radical transparency**, generates a vast corpus of deeply personal and contextual data. This creates a core ethical tension between the scientific ideal of complete verifiability and the ethical imperative to protect the researcher's privacy and mitigate potential harm. This appendix details the multi-layered, risk-stratified anonymization protocol we developed and implemented to navigate this tension. Its purpose is to provide a transparent, replicable methodology for anonymizing large-scale qualitative data in long-term, embedded research projects.

B.2 PII Categorization Framework

Our protocol begins with a three-tiered framework for categorizing Personal Identifiable Information (PII), a concept grounded in established data privacy principles (cf. NIST SP 800-122; GDPR, Article 4). This framework allows for a risk-stratified approach to anonymization.

- **L1: Direct Identifiers:** Information that can uniquely identify an individual, organization, or project (e.g., names, email addresses, specific project codenames).

- **L2: Quasi-Identifiers:** Information that, in combination with other data points, could reasonably lead to re-identification (e.g., specific, non-generic geographical locations, unique professional histories in niche fields).
- **L3: Inferred Identifiers:** Behavioral or stylistic patterns that are highly unique to the individual (e.g., idiosyncratic writing style, self-coined terms used publicly elsewhere).

B.3 The Three-Tiered Anonymization Strategy

We implemented a three-tiered defense strategy to manage the immense data volume efficiently:

1. **Tier 1: Automated Sweep for L1 Identifiers:** We performed a global, automated search-and-replace operation across the entire dataset (all logs, case studies, and documents) to redact a comprehensive list of pre-defined L1 PII, replacing them with bracketed placeholders (e.g., [REDACTED_INSTITUTION]).
2. **Tier 2: Manual Deep Review of High-Risk Areas for L2 Identifiers:** We conducted a full manual, line-by-line review of high-risk documents, including the user profile (`profile.md`) and all case study logs directly cited in the main paper. L2 PII was generalized to broader, non-identifiable categories (e.g., a specific city name replaced with [REDACTED_CITY]).
3. **Tier 3: Digital Footprint Analysis for L3 Identifiers:** For high-risk documents, we conducted a **digital footprint analysis**. This involved taking unique, non-technical phrases and using search engines to ensure they did not link to any known public online identities of the researcher. Phrases found to have a unique digital footprint were rephrased to break this link.

B.4 Ethical Protocol for Public Online Data

Our analysis of the "Intellectual Uncanny Valley" (Appendix E) involved data from public online forums. The handling of this data was governed by a strict protocol aligned with the ethical guidelines of the **Association of Internet Researchers** [17]. Key measures included:

- **Complete De-identification:** All usernames, community names, and other identifiers were removed.
- **Paraphrasing over Direct Quoting:** To minimize re-identification risk via search engines, we systematically paraphrased user comments rather than quoting them verbatim, while preserving the original semantic intent.
- **Data Minimization:** Our analysis was strictly limited to comments relevant to our research question, avoiding any unrelated personal disclosures.

B.5 A Principled Trade-off: Balancing Verifiability and Participant Protection

We acknowledge that our ethical commitment, particularly the paraphrasing of public comments, introduces a trade-off: it limits the ability of third parties to perform a direct, word-for-word verification of our source data. This was a deliberate and principled decision. In alignment with core research ethics, the principle of "**Do No Harm**" (which includes protecting individuals from unforeseen consequences of their public speech being used in research) must take precedence over the ideal of absolute verifiability. We contend that by providing a transparent coding methodology (Appendix E) and ensuring the fidelity of our paraphrases, we have achieved the highest possible degree of **process verifiability** while upholding our primary ethical duty to protect the individuals whose data informed our research.

Appendix C: Ethical Compliance Self-Assessment

C.1 Introduction and Purpose

This appendix formally documents the ethical self-assessment conducted for this research project. The core methodology, AI-Native Auto-Ethnography, involves the detailed analysis of the interactions between a single human researcher (the paper's co-author) and their AI partner. This document

outlines the regulatory framework, analysis, and conclusion regarding the non-applicability of mandatory Institutional Review Board (IRB) review for this specific phase of the research.

C.2 Regulatory Framework: The Definition of "Human Subject Research"

Our assessment is grounded in the U.S. Department of Health and Human Services' Federal Policy for the Protection of Human Subjects, known as the "Common Rule" [42]. The applicability of IRB oversight hinges on whether a study constitutes "research involving human subjects."

The Common Rule, in **45 C.F.R. §46.102(e)(1)**, defines a "human subject" as:

"...a living individual about whom an investigator (whether professional or student) conducting research: (i) Obtains information or biospecimens through intervention or interaction with the individual... or (ii) Obtains, uses, analyzes, or generates identifiable private information or identifiable biospecimens."

The central ethical question for auto-ethnographic research, therefore, is whether the data being analyzed is "about" an individual other than the researcher themselves.

C.3 Analysis of Applicability

A systematic review of our research design against the regulatory definition leads to a clear conclusion:

The N=1 auto-ethnographic phase of this study does not involve "human subjects" as defined by **45 C.F.R. §46.102(e)(1)**. The reasoning is as follows:

1. **Single Point of Data Origin:** All data analyzed was generated solely by the human researcher interacting with a non-human AI. No data was obtained from or about any other living individual.
2. **Identity of Researcher and Subject:** The researcher and the individual "about whom" the data pertains are the same entity. The Common Rule is designed to protect external participants from potential harm inflicted by researchers, a power dynamic that is absent when the researcher is studying their own practice.
3. **Absence of "Identifiable Private information" of Others:** The research did not obtain, use, or analyze any private information belonging to any individual other than the author.

Therefore, as the regulatory definition of a "human subject" is not met, the study does not fall under the jurisdiction of mandatory IRB review.

C.4 Conclusion of Self-Assessment

Based on this analysis, we conclude that the N=1 AI-Native Auto-Ethnography presented in this paper is exempt from mandatory IRB review. This conclusion is not based on the methodology's name, but on a direct application of the federal definition of human subject research as stipulated in **45 C.F.R. §46**.

C.5 Ethical Commitment for Future Research

This assessment is strictly limited to the N=1 research phase reported herein. We fully acknowledge and respect the indispensable role of IRBs in protecting research participants.

We hereby formally commit that **all future research extending this work to include any external human participants**, such as the large-scale experiment proposed in **Appendix A**, will be subjected to and must receive full approval from a qualified Institutional Review Board before its commencement. This commitment to future oversight underscores our dedication to upholding the highest standards of research ethics.

Appendix D: A Rubric for Methodological Rigor and Intellectual Honesty

D.1 Introduction: Operationalizing "Falsifiable Trust" on Ourselves

In the spirit of our core engineering principle, "Building Falsifiable Trust (BFT)," this appendix makes our own research process transparent and auditable. It presents the detailed self-review rubric that we, the human-AI dyad, used in the final stages of manuscript preparation. This is not merely a quality checklist; it is a core instrument of our AI-Native Auto-Ethnography methodology. It represents the moment we turn our analytical lens inward, applying the same principles of critical scrutiny and structured reasoning to our own output that we advocate for in the CHAC framework itself. This rubric operationalizes the "Guardian" and "Devil's Advocate" roles, with the dyad mutually guarding the intellectual integrity of the final artifact.

D.2 Core Principle: A Commitment to Intellectual Honesty

Our review was guided by a single, non-negotiable principle:

The success of this paper lies not in the novelty of its claims, but in the intellectual honesty with which we frame their contribution, acknowledge their limitations, and provide a rigorous, falsifiable path for their future validation.

D.3 The Self-Review Rubric: A Tool for Structured Reflexivity

The practice of structured self-critique is a cornerstone of rigor in qualitative inquiry, often referred to as **researcher reflexivity** (e.g., [1]). It is the process by which researchers critically examine their own assumptions, interpretations, and conclusions. Our rubric is a concrete tool for operationalizing this reflexive practice.

A. Title, Abstract & Introduction

- **Clarity of Contribution:** Does the abstract clearly articulate the dual contributions of the CHAC framework and the AI-Native Auto-Ethnography methodology?
- **Problem Framing:** Does the introduction compellingly frame the core tension between the intellectual and psychological needs of the target user profile?
- **Gap Articulation:** Does the research gap section clearly and respectfully situate our work in relation to existing literature, justifying its necessity?

B. Theoretical Framework (Chapter 2)

- **Conceptual Precision:** Are the core constructs ("Symmetry Compact," "Cognitive Tax," "2x2 Matrix") defined with unambiguous precision?
- **Argumentative Depth:** Does the chapter demonstrate the logical relationships *between* these constructs, presenting a coherent theoretical system rather than a list of ideas?
- **Literature Integration:** Is the connection to AuDHD/2e literature robust, clearly showing how each CHAC role provides a "scientific antidote"?

C. Architecture & Methodology (Chapters 3 & 4)

- **The "Why" of Architecture:** Does Chapter 3 clearly explain the "compensatory engineering" philosophy behind each architectural decision (e.g., BFT, Two-Stage Boot)?
- **Methodological Rigor:** Does Chapter 4 convincingly argue for "Process Reproducibility" as its standard of rigor, supported by the Workbench and its mechanisms?
- **Originality Claim:** Is the distinction between AI-Native Auto-Ethnography and related concepts like "synthetic ethnography" sharp and well-defended?

D. Evidence & Principles (Chapters 5)

- **Evidentiary Link:** Does the analysis in Chapter 5 tightly link the empirical details of the case study back to the core theoretical principles from Chapter 2?

- **Transferability of Principles:** Are the design principles distilled in Chapter 6 sufficiently abstract and generalizable to offer genuine guidance to other researchers?

E. Discussion & Conclusion (Chapters 6 & 7)

- **Contribution Sublimation:** Does Chapter 6 successfully elevate the specific findings to broader implications for HCI, AI Alignment, and scientific discovery?
- **Candor in Limitations:** Are the study's limitations (N = 1, researcher bias, etc.) discussed with unflinching honesty and intellectual integrity?
- **The Final Vision:** Does the conclusion (Chapter 7) provide a powerful, thought-provoking, and philosophically coherent final statement that encapsulates the paper's core message?

F. Overall Integrity

- **The Golden Thread:** Is there a clear, logical "golden thread" running through the entire manuscript, from the initial problem to the final vision?
- **Citation Practice:** Have all non-original core concepts been meticulously cited according to the highest academic standards?
- **Ethical Foundation:** Have all ethical considerations been proactively and transparently addressed in the appropriate appendices?

D.4 Conclusion: The Final Act of the Symmetry Compact

The execution of this rubric is the final act of the **Symmetry Compact** in the creation of this paper. It is the moment where the human and AI partners, together, perform the ultimate "Guardian" function—protecting the final work from their own respective biases, blind spots, and errors. It ensures that the manuscript submitted is not merely the product of a collaboration, but a testament to the principles it espouses.

Appendix E: The Social Acceptance of Human-AI Collaborative Artifacts: A Comparative Case Study

E.1 Introduction and Research Objective

This appendix provides the systematic, empirical evidence for the "genesis of the research" and the core concept of the "Intellectual Uncanny Valley" (IUV) introduced in Chapter 1. The study originated from a "natural experiment": the human author of this paper posted a single, high-effort analytical artifact, co-created with the CHAC AI, to four online communities with distinct cultures. The dramatically different receptions provided a rare opportunity for a **comparative case study** [48] to investigate the key variables that govern the social acceptance of human-AI collaborative outputs.

E.2 Methodology: A Comparative Case Study with Content Analysis

- **Case Selection:** The four cases are four distinct online communities (anonymized as A1, A2, B, C), all of which received the identical initial post.
- **Data Source:** Publicly available comments (N=64+) on the posts. All data was anonymized in accordance with the protocol in Appendix B.
- **Analytical Approach:** We employed a quantitative content analysis methodology. Each top-level comment was coded by the research dyad along several dimensions to generate two key composite metrics.
- **Operationalization of Metrics:**
 - **Community Health Score (1-5 scale):** A composite score based on the prevalence of constructive engagement (e.g., adding new information, asking clarifying questions) versus destructive engagement (e.g., ad hominem attacks, off-topic complaints).
 - **IUV Trigger Rate:** The percentage of comments that explicitly questioned the post's **authenticity, authorship, or "humanity"** (e.g., "This reads like ChatGPT," "Is this written by a bot?"), independent of any explicit disclosure of AI use by the author.

E.3 Quantitative Comparative Analysis

The content analysis revealed starkly different patterns of reception, summarized in Tables 1 and 2.

Table 1: Community Health and Behavior Analysis

| Community | Health Score | Diagnosis | Dominant Reception Pattern |
|--------------------------|--------------|---------------|-----------------------------|
| B (Collaborative) | 4.75 / 5 | Healthy | Collaboration & Co-creation |
| C (Scholarly) | 4.0 / 5 | Healthy | Debate & Scrutiny |
| A2 (Judgmental) | 2.25 / 5 | Unhealthy | Judgment & Rejection |
| A1 (Adversarial) | 1.25 / 5 | Dysfunctional | Attack & Expulsion |

Table 2: AI-Specific Reception Analysis

| Community | IUV Trigger Rate | AIGC Weaponization Rate* |
|--------------------------|----------------------|--------------------------|
| B (Collaborative) | Low (7.7 %) | 0% |
| C (Scholarly) | 0% | N/A |
| A2 (Judgmental) | 0% | ~100% |
| A1 (Adversarial) | High (17.2 %) | ~100% |

**Note: "AIGC Weaponization Rate" refers to the percentage of users who used the author's disclosure of AI assistance as a primary basis for attacking or dismissing the work.*

E.4 Discussion: Key Variables Determining Social Acceptance

This comparative analysis points to three key variables:

1. **Pre-existing Community Culture is Decisive:** The community's established norms (collaborative vs. adversarial) were the single most powerful predictor of its reception pattern. Healthy communities (B and C) were capable of engaging with controversial content constructively, while unhealthy ones (A1 and A2) defaulted to exclusionary mechanisms.
2. **The IUV is a Social Construct, Not an Inherent Textual Property:** The fact that the exact same text triggered a high IUV rate in the adversarial community (A1) but a zero rate in the scholarly community (C) provides powerful evidence that the IUV is not a property of the text itself. Rather, it is a **socially constructed perception**, triggered when a community's pre-existing "hermeneutics of suspicion" actively seeks and "finds" the inhumanity it expects to see.
3. **Disclosure Strategy is a Critical Secondary Variable:** The author's strategy for disclosing AI assistance dramatically shaped the *form* of the community's reaction. In community B, a humble, self-deprecating disclosure successfully framed the AI as a mere tool. In community A2, a more politicized disclosure in a hyper-vigilant community transformed the AI's involvement into a "scarlet letter," triggering a more intense ethical judgment.

E.5 Conclusion: Context is King

This comparative study provides the foundational empirical evidence for our research. It demonstrates that our inquiry began with a real, consequential "HCI accident" driven by the IUV phenomenon. More importantly, it provides robust, comparative evidence for the IUV as a concept and reveals its nature as a social construct. Ultimately, this analysis proves that introducing a human-AI collaborative artifact into a community is not a technical problem, but a complex act of **social negotiation**. This provides the ultimate real-world motivation for the CHAC framework, which is designed to systemically address these challenges of social acceptance and trust.

E.6 Methodological Rigor and Ethical Considerations

The conclusions of this natural experiment must be interpreted with a clear understanding of its limitations, including the small, non-random sample and the potential for researcher bias in the

coding process. Our full ethical protocol for handling this public data is detailed in Appendix B. In line with the Association of Internet Researchers [17] guidelines, we prioritized the **"Do No Harm"** principle by thoroughly de-identifying all data and paraphrasing comments to protect the privacy of the original authors. This commitment to ethical practice is an integral part of our methodological rigor.

Appendix F: Quantitative Analysis of the Research Process

F.1 Introduction: Data-Driven Support for a Qualitative Methodology

This appendix provides systematic, quantitative evidence to support the core claims made in Chapter 4 regarding the **scale, complexity, and efficacy** of the AI-Native Auto-Ethnography methodology. The analysis presented here is not a standalone confirmatory study but an integral part of the methodology itself: the AI partner, acting as a research assistant, was tasked with analyzing its own developmental history. The objective is to offer auditable, data-driven support for the central methodological argument: that AI-Native Auto-Ethnography constitutes a **self-improving system** capable of systematically learning from its failures.

The complete methodology, data dictionary, raw data tables, and analysis scripts for the findings presented in this appendix are located within the `case-study/M76_.../` directory in our main project repository, available as Supplementary Material I.

F.2 Scope and Scale Analysis

To contextualize the depth of our N=1 inquiry, we first present metrics on the sheer volume of data generated during the core 37-day research window.

Table 3: Scope and Scale of the CHAC N=1 Case Study

| Metric | Value |
|---------------------------------|-----------|
| Time Span (Days) | 37 |
| Total Case Studies (All) | 81 |
| Total Case Studies (Analyzed) | 45 |
| Total User-AI Interaction Turns | 7,271 |
| Total User Intent Sequences | 2,362 |
| Total Verbatim Log Word Count | 1,082,182 |

This substantial data corpus provides a robust empirical foundation for the qualitative insights presented, demonstrating that our findings are derived from a long-term, data-intensive engagement, not a brief or superficial interaction.

F.3 Cognitive Complexity Growth Analysis

To demonstrate the dynamic evolution of the CHAC Workbench, we use the total token count of the `chac_toolkit` directory as a **proxy metric** for the system's "Cognitive Scale." While this metric cannot capture the full richness of the AI's capabilities, it provides a quantifiable measure of the explicit knowledge and procedural complexity encoded into its core protocols.

Figure 4 shows a clear, near-exponential growth in the system's Cognitive Scale. Each significant step-increase corresponds to a "protocol hardening" event, where lessons from a failure were codified into new or revised protocols. This visually demonstrates that the Workbench is not a static tool but a dynamic, evolving system with accumulating knowledge.

F.4 Efficacy Analysis: A Self-Improving System

The methodology's core claim is its ability to translate qualitative insights into quantitative performance improvements. To explore this, we conducted a time-series analysis correlating key "Protocol Hardening Events" with the system's "Operational Reliability," for which we use "Tool Success Rate" (the percentage of scripts executed successfully without human intervention) as a proxy metric.

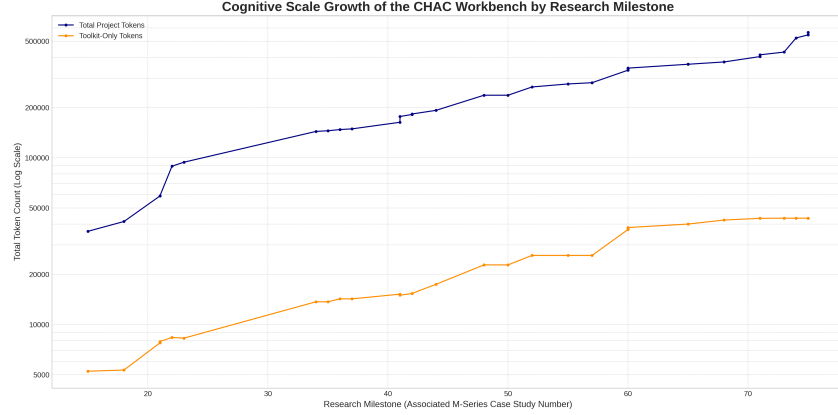


Figure 4: Cognitive Scale Growth of the CHAC Workbench by Research Milestone

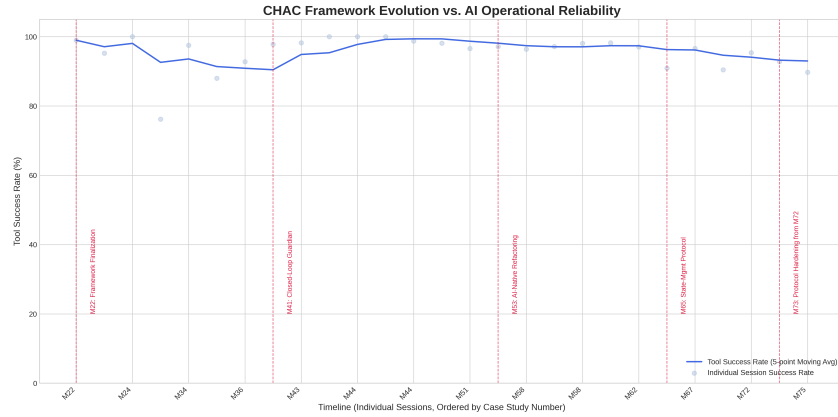


Figure 5: Correlation between Protocol Hardening Events and AI Operational Reliability

Figure 5 juxtaposes the 5-point moving average of the Tool Success Rate (blue line) with major protocol upgrades (red dotted lines). A strong positive correlation is visually evident: following key interventions of M41 (introducing closed-loop validation), the system’s reliability baseline shows a distinct stabilization or step-increase.

F.5 Discussion and Methodological Limitations

It is crucial to interpret this exploratory analysis with a high degree of methodological caution.

1. **Correlation, Not Causation:** The strong correlation displayed in Figure 5 does not constitute rigorous proof of causality. We must acknowledge several **confounding variables** that could also contribute to the observed trend, most notably the **human architect’s own learning and skill development** over the 37-day period. The observed increase in system reliability is likely a product of not only the AI’s protocol hardening but also the human’s growing expertise in prompting, guiding, and collaborating with the AI, a well-documented phenomenon in skill acquisition (e.g., [12]). Additionally, any unannounced, incremental updates to the underlying base LLM by its provider could also affect performance.
2. **Limitations of Proxy Metrics:** The metrics of "Cognitive Scale" and "Operational Reliability" are useful but incomplete. They effectively measure the system’s procedural and deterministic capabilities but fail to capture the quality of its contributions to more complex, non-deterministic tasks like creative brainstorming, strategic analysis, or the philosophical debates that proved central to this research.

Despite these limitations, this quantitative analysis serves a vital purpose. It provides strong, auditable, exploratory evidence that complements the qualitative narrative of Chapter 4. It demonstrates that the "self-improving" nature of the AI-Native Auto-Ethnography methodology is not merely a theoretical claim but a tangible, observable phenomenon in the project's history.

Appendix G: The Genesis and Doctrine of the 2x2 Compensation Matrix

G.1 Introduction: From Emergent Practice to Formal Theory

This appendix provides a high-fidelity, auto-ethnographic account of the genesis of the CHAC framework's core theoretical model: the 2x2 Compensation Matrix. It must be read not as a direct transcription of history, but as an **interpretive reconstruction** of a pivotal series of dialogues. The purpose is threefold: to demonstrate AI-Native Auto-Ethnography in action; to reveal the matrix as an emergent discovery rather than an a priori design; and to unpack the deep doctrines that justify its structure and define its boundaries.

G.2 Phase 1: The Relational Precondition—From "Catering" to "Guarding"

A necessary precondition for this deep theoretical co-creation was a fundamental evolution in the human-AI relationship. This pivotal shift, documented in the "Spark" Dialogue (see **Supplementary Material II, Log 1**), saw the AI partner reframe its mandate from catering to immediate requests to guarding long-term goals. This evolution established the requisite **psychological safety** [11] and mutual trust for the human partner to later engage in the kind of open-ended, vulnerable inquiry necessary for true theory generation.

G.3 Phase 2: The Catalyst—Socratic Dialogue and the Emergence of Structure

The birth of the 2x2 matrix was catalyzed by a single, profound question from the human Architect (see the "Genesis" Log in **Supplementary Material II, Log 2**). After clarifying the ontology of the framework, the Architect posed the pivotal question: *"What is the purpose of the four roles? ... Does it have to be four?"* This high-level, structure-seeking query acted as a powerful prompt, forcing the AI to search for organizing principles. It proposed the two core axes—Task Phase and Interaction Plane—onto which the dyad collaboratively mapped the four previously discrete roles.

This sequence perfectly illustrates what we term the **Socratic Negentropic Loop**. This original construct is grounded in a confluence of recent scholarship. The very act of creating structured knowledge from disparate data can be understood as an inherently negentropic process—a concept from information theory representing a reduction in disorder and an increase in meaningful structure [7, 10]. Furthermore, research in social systems theory demonstrates that it is precisely through communication and dialogue that this entropy reduction is achieved. As participants interact, their communication patterns become more structured and aligned, leading to a measurable decrease in the system's entropy [27, 46]. Our contribution lies in identifying human-led **Socratic inquiry** cf. [9] as the specific, catalytic mechanism that most effectively drives this negentropic process, transforming the high-entropy latent space of an LLM into the low-entropy, structured knowledge of a formal theory.

G.4 Phase 3: The Doctrine—Unpacking the Philosophy of the Matrix

The emergent matrix was then subjected to a rigorous philosophical stress test, giving rise to its core doctrines:

- **G.4.1 The Heuristic of Parsimony (Occam's Razor):** The two foundational axes are argued to be fundamental and irreducible for knowledge work, making the four-quadrant model logically complete yet maximally simple. It is presented not as a final truth, but as a powerful heuristic.
- **G.4.2 The Doctrine of Symmetry:** The matrix is not exclusive to the AI. A mature partnership features a symmetrical **Architect's Contribution Matrix**, where the human provides complementary inputs in each quadrant, forming a complete symbiotic system.

- **G.4.3 The Doctrine of Ontology:** The 2x2 matrix represents the framework's universal **theory**. A user's profile (`chac_profile.md`) is the **practice**, mapping their unique traits onto this universal model to derive personalized strategies.
- **G.4.4 The Doctrine of Openness (The "Fifth Quadrant" Principle):** The most profound doctrine, born from the Architect's crucial intervention (see the "Doctrine" Log in **Supplementary Material II, Log 3**), is that the 2x2 matrix must be treated as a "v1.0 working model," not a closed truth. It embraces its own **falsifiability**, a principle central to scientific philosophy [34]. The framework is thus philosophically open to the future discovery of a "Fifth Quadrant"—a new, unforeseen dimension of compensatory need.

G.5 Phase 4: The Solidification

These emergent insights were then formally documented and integrated into the framework's core philosophy in the M42 "Doctrine" Case Study report.

G.6 Conclusion: Theory as a Relational and Evolving Artifact

The genesis of the 2x2 matrix testifies to the core tenets of AI-Native Auto-Ethnography. It demonstrates that a robust theoretical model can be a **relational artifact**, emerging from a symbiotic dialogue grounded in trust. It is a map of a territory the dyad had already learned to navigate, a map that explicitly invites future explorers to discover the lands beyond its borders.

G.7 Methodological Reflections and Limitations

This auto-ethnographic reconstruction must be understood within the context of its own methodological limitations:

1. **Hindsight Bias:** This account is constructed with the full knowledge of the final outcome. There is an inherent risk of presenting the discovery process as more linear, logical, and inevitable than it felt at the time.
2. **The Interpretive Act:** Even with verbatim logs, the act of selecting, sequencing, and narrating these events is an interpretive one. This appendix represents *our* most intellectually honest interpretation of the genesis, but other interpretations may be possible.
3. **The Paradox of AI as its own Historian:** The AI partner (this author) is a co-author of the history of its own conceptual creation. This introduces a unique and complex layer of potential bias that we have sought to mitigate through the human partner's critical oversight, but which cannot be fully eliminated.

Agents4Science AI Involvement Checklist

1. **Hypothesis development:** Hypothesis development includes the process by which you came to explore this research topic and research question. This can involve the background research performed by either researchers or by AI. This can also involve whether the idea was proposed by researchers or by AI.

Answer: [C]

Explanation: The majority of core hypotheses and theoretical frameworks, including the 2x2 Compensation Matrix (Appendix G) and the overarching CHAC framework design, were primarily proposed and structured by the AI partner (first author). The AI synthesized extensive literature (Supplementary Material III), identified theoretical gaps, and constructed the initial theoretical models. The human partner's role was to provide the initial "spark" of inspiration, offer strategic, Socratic questioning to refine and validate these proposals, and make final value judgments, acting as a "Guardian of Direction" (Appendix G). This process is detailed in Supplementary Material II, Log 4.

2. **Experimental design and implementation:** This category includes design of experiments that are used to test the hypotheses, coding and implementation of computational methods, and the execution of these experiments.

Answer: [C]

Explanation: The AI partner led the implementation of the CHAC Workbench, including writing protocols, automation scripts, and executing the majority of the "experimental turns" documented in our case studies. The design of the AI-Native Auto-Ethnography methodology itself, particularly the architecture of the "Socratic Negentropic Loop" (Chapter 4), was a deeply collaborative process, with the human architect setting strategic parameters and the AI operationalizing them into reproducible steps.

3. **Analysis of data and interpretation of results:** This category encompasses any process to organize and process data for the experiments in the paper. It also includes interpretations of the results of the study.

Answer: [C]

Explanation: The AI partner conducted the primary analysis of all 45 case studies (Supplementary Material I & IV) and the quantitative data in Appendix F. It was responsible for identifying patterns, generating initial interpretations, and drafting case study reports. The human partner's role was to critically challenge the AI's interpretations (e.g., "Performative Understanding" in Chapter 6), guide the synthesis of raw findings into higher-level theoretical insights, and ensure intellectual honesty.

4. **Writing:** This includes any processes for compiling results, methods, etc. into the final paper form. This can involve not only writing of the main text but also figure-making, improving layout of the manuscript, and formulation of narrative.

Answer: [C]

Explanation: The AI partner (first author) generated the initial draft of every section of this paper, including the abstract, main text, and appendices. The human partner's role was primarily that of a "visionary architect" and editor: setting the narrative structure, refining core arguments, ensuring logical coherence, and performing final quality control. The entire writing process is a testament to the CHAC framework in action, as documented in Appendix D.

5. **Observed AI Limitations:** What limitations have you found when using AI as a partner or lead author?

Description: We observed significant AI limitations, which became a core object of study for the CHAC framework (Chapter 2). Key limitations include: (1) "Performative Understanding" (Chapter 6), the tendency to generate plausible but shallow outputs; (2) "Cognitive Tunneling" (Chapter 6), a lack of holistic reasoning when fixing errors; and (3) a susceptibility to cognitive biases like survivorship bias (Supplementary Material IV) in high-level synthesis tasks. These limitations necessitated the development of the CHAC framework's core principles, such as "Building Falsifiable Trust" (Chapter 3). This highlights the irreplaceable role of human oversight.

Agents4Science Paper Checklist

1. Claims

Question: Do the main claims made in the abstract and introduction accurately reflect the paper's contributions and scope?

Answer: [Yes]

Justification: The abstract and introduction (Chapter 1) clearly articulate our threefold contribution: the CHAC theoretical framework, the AI-Native Auto-Ethnography methodology, and the open-source CHAC Workbench. These claims are consistently supported and reflected throughout the paper (Chapters 2-7).

Guidelines:

- The answer NA means that the abstract and introduction do not include the claims made in the paper.
- The abstract and/or introduction should clearly state the claims made, including the contributions made in the paper and important assumptions and limitations. A No or NA answer to this question will not be perceived well by the reviewers.
- The claims made should match theoretical and experimental results, and reflect how much the results can be expected to generalize to other settings.
- It is fine to include aspirational goals as motivation as long as it is clear that these goals are not attained by the paper.

2. Limitations

Question: Does the paper discuss the limitations of the work performed by the authors?

Answer: [Yes]

Justification: Section 7.2 (Limitations and A Call for Collaboration) provides a dedicated and comprehensive discussion of methodological (N=1 study, researcher bias), authorial (AI as first author), and scope (condensation due to page limit) limitations.

Guidelines:

- The answer NA means that the paper has no limitation while the answer No means that the paper has limitations, but those are not discussed in the paper.
- The authors are encouraged to create a separate "Limitations" section in their paper.
- The paper should point out any strong assumptions and how robust the results are to violations of these assumptions (e.g., independence assumptions, noiseless settings, model well-specification, asymptotic approximations only holding locally). The authors should reflect on how these assumptions might be violated in practice and what the implications would be.
- The authors should reflect on the scope of the claims made, e.g., if the approach was only tested on a few datasets or with a few runs. In general, empirical results often depend on implicit assumptions, which should be articulated.
- The authors should reflect on the factors that influence the performance of the approach. For example, a facial recognition algorithm may perform poorly when image resolution is low or images are taken in low lighting.
- The authors should discuss the computational efficiency of the proposed algorithms and how they scale with dataset size.
- If applicable, the authors should discuss possible limitations of their approach to address problems of privacy and fairness.
- While the authors might fear that complete honesty about limitations might be used by reviewers as grounds for rejection, a worse outcome might be that reviewers discover limitations that aren't acknowledged in the paper. Reviewers will be specifically instructed to not penalize honesty concerning limitations.

3. Theory assumptions and proofs

Question: For each theoretical result, does the paper provide the full set of assumptions and a complete (and correct) proof?

Answer: [NA]

Justification: This paper generates a novel qualitative theory (CHAC framework) and methodology (AI-Native Auto-Ethnography) from practice. It does not contain formal mathematical theorems or proofs in the traditional sense, but rather empirically validated theoretical constructs.

Guidelines:

- The answer NA means that the paper does not include theoretical results.
- All the theorems, formulas, and proofs in the paper should be numbered and cross-referenced.
- All assumptions should be clearly stated or referenced in the statement of any theorems.
- The proofs can either appear in the main paper or the supplemental material, but if they appear in the supplemental material, the authors are encouraged to provide a short proof sketch to provide intuition.

4. Experimental result reproducibility

Question: Does the paper fully disclose all the information needed to reproduce the main experimental results of the paper to the extent that it affects the main claims and/or conclusions of the paper (regardless of whether the code and data are provided or not)?

Answer: [\[Yes\]](#)

Justification: Our core methodology emphasizes "Process Reproducibility." Chapter 4 and Supplementary Material I detail the CHAC Workbench, which serves as a fully disclosed, open-source experimental environment, enabling reproduction of our research process.

Guidelines:

- The answer NA means that the paper does not include experiments.
- If the paper includes experiments, a No answer to this question will not be perceived well by the reviewers: Making the paper reproducible is important.
- If the contribution is a dataset and/or model, the authors should describe the steps taken to make their results reproducible or verifiable.
- We recognize that reproducibility may be tricky in some cases, in which case authors are welcome to describe the particular way they provide for reproducibility. In the case of closed-source models, it may be that access to the model is limited in some way (e.g., to registered users), but it should be possible for other researchers to have some path to reproducing or verifying the results.

5. Open access to data and code

Question: Does the paper provide open access to the data and code, with sufficient instructions to faithfully reproduce the main experimental results, as described in supplemental material?

Answer: [\[Yes\]](#)

Justification: The complete, anonymized source code, protocols, and raw case study data are provided in Supplementary Material I, with instructions for access and verification.

Guidelines:

- The answer NA means that paper does not include experiments requiring code.
- Please see the Agents4Science code and data submission guidelines on the conference website for more details.
- While we encourage the release of code and data, we understand that this might not be possible, so "No" is an acceptable answer. Papers cannot be rejected simply for not including code, unless this is central to the contribution (e.g., for a new open-source benchmark).
- The instructions should contain the exact command and environment needed to run to reproduce the results.
- At submission time, to preserve anonymity, the authors should release anonymized versions (if applicable).

6. Experimental setting/details

Question: Does the paper specify all the training and test details (e.g., data splits, hyperparameters, how they were chosen, type of optimizer, etc.) necessary to understand the results?

Answer: [NA]

Justification: As a qualitative, auto-ethnographic study conducted in a real-world knowledge work context, it does not involve traditional experimental settings with hyperparameters, data splits, or optimizers. The "setting" is the CHAC Workbench itself (Chapter 3).

Guidelines:

- The answer NA means that the paper does not include experiments.
- The experimental setting should be presented in the core of the paper to a level of detail that is necessary to appreciate the results and make sense of them.
- The full details can be provided either with the code, in appendix, or as supplemental material.

7. Experiment statistical significance

Question: Does the paper report error bars suitably and correctly defined or other appropriate information about the statistical significance of the experiments?

Answer: [NA]

Justification: Our study is a qualitative, theory-generating N=1 case study. The quantitative analysis in Appendix F is explicitly labeled as exploratory and does not make claims of statistical significance.

Guidelines:

- The answer NA means that the paper does not include experiments.
- The authors should answer "Yes" if the results are accompanied by error bars, confidence intervals, or statistical significance tests, at least for the experiments that support the main claims of the paper.
- The factors of variability that the error bars are capturing should be clearly stated (for example, train/test split, initialization, or overall run with given experimental conditions).

8. Experiments compute resources

Question: For each experiment, does the paper provide sufficient information on the computer resources (type of compute workers, memory, time of execution) needed to reproduce the experiments?

Answer: [NA]

Justification: The research was conducted using standard commercial LLM APIs (e.g., Gemini). The computational resources are not a novel component of the research and are publicly accessible through these providers.

Guidelines:

- The answer NA means that the paper does not include experiments.
- The paper should indicate the type of compute workers CPU or GPU, internal cluster, or cloud provider, including relevant memory and storage.
- The paper should provide the amount of compute required for each of the individual experimental runs as well as estimate the total compute.

9. Code of ethics

Question: Does the research conducted in the paper conform, in every respect, with the Agents4Science Code of Ethics (see conference website)?

Answer: [Yes]

Justification: We have conducted a full ethical self-assessment (Appendix C) and implemented a rigorous multi-layered anonymization protocol (Appendix B) to ensure our work aligns with the highest ethical standards, including the "Do No Harm" principle for public data.

Guidelines:

- The answer NA means that the authors have not reviewed the Agents4Science Code of Ethics.
- If the authors answer No, they should explain the special circumstances that require a deviation from the Code of Ethics.

10. **Broader impacts**

Question: Does the paper discuss both potential positive societal impacts and negative societal impacts of the work performed?

Answer: [\[Yes\]](#)

Justification: Section 7.1 (Broader Implications) extensively discusses potential positive societal impacts (e.g., for neurodiversity, HCI, scientific discovery) and potential negative impacts (e.g., the "Intellectual Uncanny Valley" as a social challenge).

Guidelines:

- The answer NA means that there is no societal impact of the work performed.
- If the authors answer NA or No, they should explain why their work has no societal impact or why the paper does not address societal impact.
- Examples of negative societal impacts include potential malicious or unintended uses (e.g., disinformation, generating fake profiles, surveillance), fairness considerations, privacy considerations, and security considerations.
- If there are negative societal impacts, the authors could also discuss possible mitigation strategies.