

FERE-CRS: A Calculus of Semantic Inference for Robust, Adaptive Reasoning in AI Agents

Thomas E. Devitt *Independent Researcher, FERE-CRS Project*

August 21, 2025

Abstract

Large Language Models (LLMs), despite their remarkable capabilities, exhibit fundamental limitations in robustness, reliability, and principled reasoning, hindering their deployment in high-stakes applications. This paper chronicles the **FERE-CRS (Fluid Emergent Reasoning Engine - Cognitive Resonance Score)** research program, an effort to develop a new class of cognitive architecture for LLM-based agents grounded in the **Free Energy Principle** and its process theory, **Active Inference**. We detail a rigorous, failure-driven research trajectory that began with a symbolic, logic-based paradigm (FERE-CRS I) and, in response to a series of empirically validated failures, pivoted to a more robust probabilistic, semantic paradigm (FERE-CRS II).

The primary contribution of this work is the design and validation of the final **v9.0 Integrated Agent**. This agent architecture successfully overcomes previously identified failure modes by integrating an **Ensemble-based Probabilistic Reasoning (EPR)** module to explicitly model its own uncertainty, and a novel **"Certainty-Gated Meta-Cognition"** policy that allows the agent to dynamically shift between assimilation, rejection, and active, uncertainty-reducing inquiry. We demonstrate through a comprehensive validation suite that this agent exhibits robust, generalizable, and adaptive reasoning across multiple knowledge domains and under adversarial stress.

The significance of this research lies not in creating a new state-of-the-art model, but in developing a new methodology for building safer, more interpretable, and more reliable AI agents. By treating the LLM as a probabilistic semantic engine to be guided, rather than a flawed symbolic reasoner to be contained, we provide a principled blueprint for a **Calculus of Semantic Inference**—a new path toward agents that know when to be certain and when to be curious.

Keywords: Active Inference, Free Energy Principle, Cognitive Architecture, Large Language Model (LLM), AI Safety, Uncertainty Quantification, Meta-Cognition, Probabilistic Reasoning, Fluid Reasoning.

1. Introduction

The advent of Large Language Models (LLMs) represents a significant milestone in the field of Artificial Intelligence, demonstrating unprecedented capabilities in natural language understanding, generation, and in-context learning (Bommasani et al., 2021). However, their deployment in critical, high-stakes domains is hampered by fundamental limitations. LLMs are known to be brittle, prone to factual hallucination, and lack robust mechanisms for logical consistency or self-correction (Bender et al., 2021; Marcus, 2018). This paper defines **brittleness** not merely as vulnerability to adversarial inputs, but as a fundamental lack of conceptual robustness: the tendency to produce confident but nonsensical outputs when faced with information that lies at the edge of, or outside, its training distribution. This creates a pressing need for new cognitive architectures that can harness the power of LLMs while mitigating their inherent unreliability.

The FERE-CRS research program was initiated to address this challenge by pursuing a different path. Instead of focusing on scaling existing models or developing localized patches like Retrieval-Augmented Generation (RAG) or Chain-of-Thought prompting—which are primarily sophisticated forms of information retrieval and prompt engineering—our objective was to design a new class of **cognitive architecture** for LLM-based agents from first principles. Our work is grounded in the **Free Energy Principle (FEP)** and its process theory, **Active Inference (AIF)** (Friston, 2010). We selected this framework because it provides a normative, first-principles account of self-organization and belief updating that directly addresses the core limitations of LLMs. From an AIF perspective, a phenomenon like hallucination can be formally understood as a failure of inference: the agent's generative model, dominated by its powerful prior beliefs, fails to adequately update in response to new sensory evidence (or the lack thereof), leading to a failure to minimize surprise (prediction error). AIF posits that a core function of any intelligent agent is to possess a meta-cognitive loop that actively manages this process.

This paper presents the complete, end-to-end results of the FERE-CRS I and II research programs. It is a chronicle of a rigorous, **failure-driven research methodology**. We assert that for a technology as complex and poorly understood as LLMs, a methodology based on the deliberate falsification of hypotheses is the most efficient path to genuine understanding. Over 14 distinct architectural phases, we tested two major paradigms, systematically diagnosing the failure of each prototype to inform the design of the next.

The primary contribution of this work is the final $v9.0$ Integrated Agent, a prototype that successfully demonstrates robust, generalizable, and adaptive reasoning. This agent is guided by a **Calculus of Semantic Inference**, a principled framework for meta-cognitive control. In this calculus, the **variables** are the probabilistic belief distributions generated by an LLM ensemble; the **quantities** are continuous measures of Coherence (the mean of the distribution) and Uncertainty (the standard deviation); and the **operators** are discrete cognitive actions (Assimilation, Adaptation, Rejection). This architecture uses an **Ensemble-based Probabilistic Reasoning (EPR)** module to explicitly model its own uncertainty and a novel "**Certainty-Gated Meta-Cognition**" policy to dynamically shift between confident assimilation, robust rejection, and curious, uncertainty-reducing inquiry.

Our central thesis is that by treating the LLM not as a flawed symbolic reasoner to be contained, but as a powerful probabilistic semantic engine to be guided, it is possible to build agents that are safer, more interpretable, and more aligned with the principles of true, adaptive intelligence. We will demonstrate that the $v9.0$ agent, governed by this Calculus of Semantic Inference, successfully passes a comprehensive validation suite designed to test its generalization across multiple knowledge domains and its robustness against adversarial stress, thereby providing a new and promising blueprint for the future of AI agent design.

2. Theoretical Foundations: Active Inference and the Cognitive Resonance Score

To understand the FERE-CRS architecture, we must first formalize the principles it is built upon. The architecture is a direct attempt to apply the Free Energy Principle (FEP) and its process theory, Active Inference (AIF), to the problem of high-level, neuro-symbolic reasoning.

2.1 The Free Energy Principle and Active Inference

While established frameworks like Reinforcement Learning (RL) have proven powerful for policy optimization, they are often less suited to addressing the core meta-cognitive failures—such as hallucination and brittleness—that plague modern LLMs. We selected Active Inference as our foundational framework because it is not merely a theory of learning, but a first-principles theory of **inference and self-organization**. It provides a normative mathematical basis for how any agent must maintain a coherent model of its world and act to preserve that coherence. From an AIF perspective, a hallucination is a failure of inference: the agent's generative model, dominated by its powerful prior beliefs, fails to adequately update in response to new sensory evidence (or the lack thereof), leading to a failure to minimize surprise (prediction error). AIF posits that a core function of any intelligent agent is to possess a meta-cognitive loop that actively manages this process of belief updating.

Under the FEP, an agent possesses an internal generative model, $p(s\sim, \mathcal{G}|m)$, which is its probabilistic theory of how hidden causes in the world (\mathcal{G}) generate sensory states ($s\sim$). The agent cannot access the true state of the world, so it must infer it by optimizing an approximate posterior distribution, known as the recognition density, $q(\mathcal{G}|\mu)$, which is parameterized by the agent's internal states (μ).

The agent's objective is to minimize **Variational Free Energy (VFE)**, F , which serves as an Evidence Lower Bound (ELBO) on the log evidence for its model of the world, $\ln p(s\sim|m)$. A key formulation is:

$$F(s\sim, \mu) = \underbrace{D_{KL}[q(\mathcal{G}|\mu)||p(\mathcal{G}|m)]}_{\text{Complexity}} - \underbrace{E_q[\ln p(s\sim|\mathcal{G}, m)]}_{\text{Accuracy}}$$

Here, D_{KL} is the Kullback-Leibler (KL) divergence. This equation reveals a fundamental trade-off: to minimize F , an agent must find a belief state (μ) that produces accurate explanations for its sensations (maximizing the Accuracy term), while simultaneously keeping those beliefs as simple as possible and close to its prior assumptions about the world (minimizing the Complexity term) (Hinton & van Camp, 1993). This process of optimizing beliefs through VFE minimization is perception.

Active Inference extends this principle to action. An agent can also minimize its long-term surprise by actively seeking out sensations that conform to its predictions. This is achieved by selecting policies (sequences of actions, π) that are expected to minimize the **Expected Free Energy (EFE)**, G :

$$G(\pi) = \underbrace{\sum_{\tau} E_Q[\ln p(s_{\sim\tau}|C)]}_{\text{Pragmatic Value}} - \underbrace{D_{\text{KL}}[q(\theta_{\tau}|s_{\sim\tau})||q(\theta_{\tau})]}_{\text{Epistemic Value}}$$

where Q is a distribution over future states and outcomes. EFE formalizes the exploration-exploitation dilemma. The **Pragmatic Value** term drives the agent to seek out preferred outcomes or goals (exploitation), as defined by a prior distribution $p(s_{\sim}|C)$. The **Epistemic Value** term drives the agent to take actions that are expected to resolve uncertainty about the world (exploration) (Friston et al., 2015).

2.2 The Cognitive Resonance Score (CRS): A Principled Heuristic for Intractable Problems

The direct calculation of VFE and EFE is computationally intractable for complex, high-dimensional systems like LLM-based agents. The central theoretical claim of this work is that we can formulate a practical, computable objective function—the **Cognitive Resonance Score (CRS)**—that serves as a principled heuristic for these intractable quantities. Maximizing the CRS is functionally approximate to minimizing free energy. The CRS decomposes the abstract imperatives of the FEP into a set of measurable components, providing a tractable "common currency" for guiding the agent's reasoning.

2.2.1 Relational Coherence (R) as a Proxy for Model Evidence and Accuracy

The VFE's accuracy term, $Eq[\ln p(s_{\sim}|\theta, m)]$, rewards beliefs that make sensory data likely. We posit that **Relational Coherence (R)** is a necessary, though not sufficient, heuristic for this. An agent's generative model is not a monolith; in a complex system like FERE-CRS, it is a highly structured set of beliefs. The structural integrity and internal consistency of this model are preconditions for it to make accurate predictions. A belief update that introduces a logical contradiction creates a large internal prediction error. Maximizing R is therefore the process of minimizing this internal prediction error, ensuring the agent's inferential machinery remains sound. Operationally, in our symbolic paradigm (FERE-CRS I), this was measured by penalizing the formation of contradictory edges in the Active Reasoning Graph. This aligns with coherence theories of justification, which argue that the plausibility of a belief is a function of its fit within a larger, interconnected system of beliefs (Thagard, 2007).

This raises the critical "fantasy world" problem: how does the system avoid constructing a highly coherent but factually incorrect model? The answer is that R is continuously grounded by the other CRS components. For example, if an agent develops a coherent but false belief (e.g., "all swans are white"), its model will be stable until it encounters a black swan. This new sensory input creates a massive prediction error that cannot be resolved with high coherence. The agent is then forced by its drive to minimize free energy to either reject the evidence or update its model. An action to reject the valid evidence would likely have low Pragmatic Value (P) (as it does not help achieve the goal of understanding swans) and low Informational Value (I) (as it increases,

rather than resolves, uncertainty). Therefore, the action with the highest global CRS will be the one that revises the model, sacrificing short-term coherence for long-term accuracy.

2.2.2 Cognitive Efficiency (C) as a Proxy for Complexity

The VFE's complexity term, $D_{KL}[q(\theta|\mu)||p(\theta|m)]$, acts as a form of Occam's Razor, penalizing beliefs that are overly complex or deviate far from the agent's prior assumptions. We propose that **Cognitive Efficiency (C)** is a practical proxy for this information-theoretic cost. The KL divergence can be understood as the "information cost" or "work required" to update the prior belief to the new posterior belief. In any computational system, this information processing has a real, tangible cost (e.g., CPU cycles, API call latency, memory allocation). The computational work required to update the agent's internal state (μ) is a direct, measurable proxy for the abstract information-theoretic "work" of belief updating. This view is consistent with thermodynamic interpretations of the FEP, where minimizing VFE is equated with minimizing long-term metabolic cost (Pezzulo et al., 2015; Friston, 2010). Thus, by measuring the computational cost of an action, C_{cost} , we create a tractable proxy for the complexity it induces.

2.2.3 Pragmatic (P) and Informational (I) Value as Proxies for EFE

The mapping of the P and I components of the CRS to the terms of the EFE is more direct.

- **Pragmatic Value (P):** This component is a direct heuristic for the EFE's pragmatic term, measuring the expected progress towards a goal state. It quantifies the extrinsic, goal-achieving value of an action.
- **Informational Value (I):** This component is a direct heuristic for the EFE's epistemic term, quantifying the expected reduction in uncertainty (or Shannon Entropy) over the agent's beliefs about the world. It represents the intrinsic value of "curiosity" or "epistemic foraging"—the drive to take actions that resolve ambiguity to enable better future decisions (Schwartenbeck et al., 2013).

2.2.4 Situating the CRS in the Broader AI Landscape

The CRS construct builds upon, but is distinct from, concepts like "intrinsic motivation" in RL. While curiosity-driven RL agents often use a simple prediction error as a bonus reward, the CRS provides a more holistic and multi-faceted objective function. The inclusion of the **Relational Coherence (R)** term is a key innovation. It imposes a strong imperative for the agent to maintain the *internal logical integrity of its own world model*, a constraint often absent in standard RL formulations. This focus on model consistency is a direct consequence of our AIF-first approach.

For the experiments in this paper, we primarily used a simple and interpretable additive model ($\text{CRS} = w_{RR} + w_{PP} + w_{II} - w_{CC_{\text{cost}}}$). We chose this model for its debuggability and clarity in these foundational experiments, though we acknowledge that a multiplicative model, which would grant "veto power" to any single component, is a theoretically compelling avenue for future research.

3. The FERE-CRS Research Program: A Methodological Journey

The FERE-CRS project was conducted in two phases, comprising fourteen distinct experimental stages. The research followed a deliberate, **failure-driven methodology**, where the rigorously diagnosed limitations of each architectural prototype directly informed the research objectives of the next. This iterative, cumulative process, chronicled in a comprehensive "Living Document," was essential for navigating the complexities of building a principled reasoning agent on top of inherently unpredictable LLM components. This section details that journey, tracing the evolution of our architecture and the refinement of our core theoretical concepts. A detailed, phase-by-phase technical summary of the entire research program, including quantitative results, is available in Appendix A.

3.1 Part I: The Symbolic Paradigm (FERE-CRS I)

The first major phase of our research (Phases I-X) was dedicated to testing the viability of a symbolic, logic-based implementation of Active Inference. The core architectural element of this paradigm was the **Cognitive Resonance Score (CRS)**, a heuristic for Variational Free Energy that served as a "common currency" to guide the agent's decisions. The **Meta-Reasoning Agent (MRA)**, the system's "brain," evolved significantly through this paradigm.

The initial phases (I-VI) successfully developed a series of increasingly sophisticated agents. We began with a simple MRA guided by static, hand-tuned CRS weights (Phase I), which nonetheless demonstrated superior reasoning and an emergent "explore-exploit" strategy compared to a strong RAG baseline. We then proved the agent could learn a cognitive "stance" via meta-learning (Phase II), dynamically switch between stances (Phase III), generate novel stances for unseen problems (Phase IV), and even discover entirely new conceptual primitives like 'Trustworthiness' by reflecting on its own failures (Phase V).

However, the attempt to apply this framework to a raw LLM in Phase VII resulted in a critical failure. The agent consistently fell into a state of "**Creative Avoidance**," generating a stream of semantically novel but functionally useless hypotheses. This successful negative result proved that a structured, inspectable workspace was a necessary precondition for robust reasoning. This led to the development of the **Active Reasoning Graph (ARG)** in Phase VIII. The ARG architecture was a direct solution to "Creative Avoidance"; by forcing the agent to operate on a discrete graph of nodes and edges, it could no longer "reason around" a contradiction. It was forced to confront it by either creating a formal `contradicts` edge or by falsifying a hypothesis. The MRA was upgraded to a deterministic "Priority Cascade" that could directly analyze this graph structure.

The ARG architecture was successful in solving the "Creative Avoidance" problem. This provided the foundation for the project's ultimate goal: to build a meta-agent that could analyze a corpus of these successful ARG traces to discover a new heuristic. This final test of the symbolic paradigm, conducted across Phases IX and X, failed through a series of increasingly profound revelations. We discovered the "**Semantic Gap**," the "**Stochastic Oracle**" problem, and ultimately proved that the symbolic graph itself was a "**confabulated artifact**." The agent's correct final output was fundamentally **decoupled** from the logical integrity of its internal symbolic trace.

The conclusion of FERE-CRS I was that the symbolic paradigm, while a powerful scaffold, is an insufficient foundation for a truly autonomous reasoner. This definitive failure provided the necessary and rigorous justification to pivot paradigms.

3.2 Part II: The Probabilistic Paradigm (FERE-CRS II)

The second major phase of our research was dedicated to building an agent that embraced, rather than fought, the probabilistic and semantic nature of its LLM components.

This program was built upon three new architectural pillars, validated in sequence:

- **Phase XI: The Semantic Vector Space (SVS).** We abandoned the symbolic graph and represented the agent's mind as a **statistical model (mean and variance)** within a high-dimensional embedding space. Through iterative refinement, we developed a $v4.0$ agent with "**Belief State Plasticity**" that could form concepts, adapt its conceptual boundaries, and robustly reject anomalies.
- **Phase XII: Ensemble-based Probabilistic Reasoning (EPR).** To address the agent's inability to model its own uncertainty, the architecture was upgraded to use an **ensemble of LLM "voters."** The result of any judgment was a probability distribution, with the **Shannon Entropy** of this distribution serving as a direct measure of the agent's uncertainty.
- **Phase XIII & XIV: The Integrated Agent.** We integrated the SVS and EPR modules, guided by a new **Probabilistic CRS (P-CRS)**. Initial attempts failed due to the brittleness of the SVS module's coherence signal. This led to the final, definitive $v9.0$ architecture: the "**EPR-Centric**" agent. This agent's MRA policy evolved into "**Certainty-Gated Meta-Cognition.**" It abandoned the flawed SVS metric and instead used the most reliable signal we had—the **mean coherence score** from the EPR ensemble—to define a "zone of uncertainty." A score falling into this "grey area" now correctly triggered the agent's most powerful epistemic action: a **meta-cognitive query** to the LLM for guidance.

The $v9.0$ agent's P-CRS represents the evolution of our "common currency" into a truly native form. **Coherence (R)** was no longer a symbolic check but the mean of the ensemble's continuous scores. **Uncertainty (I)** was no longer a simple entropy calculation but the standard deviation of those scores. The final $v9.0$ agent was subjected to a comprehensive **Generalization & Robustness Test Suite**. The experiment was a complete success, demonstrating robust, generalizable, and adaptive reasoning across multiple domains and under adversarial stress.

4. Discussion

4.1 Significance of Findings

The significance of the FERE-CRS project lies in the quality of the questions it asks and the problems it attempts to solve, more so than in the performance of the final prototype.

- **It Addresses the Core Problem of LLM Brittleness:** The entire project is a direct assault on the most significant barrier to deploying LLMs in high-stakes applications: their unreliability. We have focused on building an architecture that is **robust, self-aware, and knows when to stop and ask for help**. This is of profound significance for AI safety.
- **It Provides a Blueprint for More Interpretable AI:** The final probabilistic agent is interpretable in a new and powerful way. An agent that can report its mean coherence score and its standard deviation of belief provides a much richer window into its "mind" than one that simply outputs an answer.
- **The Failures are as Significant as the Successes:** The detailed documentation of why so many of our attempts failed is a valuable "road map of blind alleys" for other researchers. We have done the hard work of proving which paths are not worth taking.

4.2 Practical Applications

The "Certainty-Gated Meta-Cognition" architecture is not just a theoretical curiosity; it is a powerful tool with immediate practical applications. The most compelling is the **"LLM Cognitive Immune System,"** a meta-cognitive "wrapper" that can monitor a base LLM for conceptual drift, contradictions, or anomalies. When the FERE-CRS layer detects such a state, it could intervene, prompting the base LLM to pause, ask for clarification, or correct itself, directly addressing the problems of LLM brittleness and hallucination.

4.3 Situating FERE-CRS within the AI Paradigm

A skeptical critique might suggest that our final agent, with its meta-cognitive query, is merely a complex form of prompt-chaining. This interpretation misses the fundamental architectural distinction. A prompt-chain is a pre-defined, brittle script. The FERE-CRS agent, by contrast, is a dynamic, self-regulating system. Its behavior is not determined by a fixed script, but emerges from its continuous, goal-directed effort to optimize a single, principled objective function: the Probabilistic Cognitive Resonance Score. The agent's decision to "ask for clarification" is not a hard-coded rule, but an emergent strategy it selects because it is the optimal action for reducing its internal uncertainty. This grounding in the first principles of Active Inference is what distinguishes FERE-CRS as a true cognitive architecture, not just an advanced prompting technique.

4.4 Limitations and Future Work

We acknowledge that the current $v9.0$ prototype has significant computational overhead due to its reliance on a large ensemble for every cognitive judgment. Practical deployment would require engineering solutions, such as **amortized inference**, where a smaller, distilled "System 1" model learns to approximate the judgments of the larger "System 2" ensemble.

The primary limitation of our successful agent is that it is still passive; it cannot yet formulate its own goals. This provides the clear and compelling charter for our next phase of research, **FERE-CRS III: The Calculus of Curiosity**, which will be dedicated to transforming our adaptive reasoner into a truly autonomous, goal-oriented agent.

5. Conclusion

The FERE-CRS research program has not produced a state-of-the-art model that will top performance leaderboards. Instead, it has produced something more valuable: a **breakthrough in methodology and conceptual framing**. We have designed, tested, and validated a new type of cognitive "engine," one based on the principles of self-awareness, uncertainty, and active inquiry.

The single most important takeaway from this work is the validation of a new paradigm for agent design. By abandoning the attempt to force a probabilistic, semantic system into a rigid, symbolic framework, and instead embracing its native properties, we have shown that it is possible to build an agent that is robust, generalizable, and adaptive. The final $v9.0$ agent is the first working prototype of this new engine.

It is a pioneering exploration of a specific, highly promising, and fundamentally different road toward artificial general intelligence—a road focused not on raw performance, but on the much harder and more important qualities of robustness, self-awareness, and reliability. The significance of our work is that we have proven this road is navigable.

Acknowledgments

The author acknowledges the foundational work of the Active Inference community. The author would also like to acknowledge the significant role of Google's Gemini Advanced in the development of this research. The large language model was utilized as a collaborative tool throughout the writing process, serving several key functions: as a Socratic partner for brainstorming, challenging and testing theoretical assumptions, code generator/reviewer and as a writing assistant for refining the clarity and structure of the manuscript's prose. The author maintains full intellectual responsibility for the core ideas, experimental design, and final conclusions presented in this work.

Data Availability

The source code and experimental data for the preceding research are available in the project's public GitHub repository: <https://github.com/ThomasDevitt/FERE-CRS>.

This manuscript, which constitutes the complete research protocol for Phase 1 and 2, has been permanently archived on Zenodo and is available under the Digital Object Identifier (DOI): **10.5281/zenodo.16923559**.

References

- [1] Anderson, J. R. (2007). *How Can the Human Mind Occur in the Physical Universe?* Oxford University Press.
- [2] Bender, E. M., Gebru, T., McMillan-Major, A., & M. Mitchell (2021). On the Dangers of Stochastic Parrots: Can Language Models Be Too Big? In *Proceedings of the 2021 ACM Conference on Fairness, Accountability, and Transparency* (pp. 610-623).
- [3] Bommasani, R., et al. (2021). *On the Opportunities and Risks of Foundation Models*. arXiv preprint arXiv:2108.07258.
- [4] Bostrom, N. (2014). *Superintelligence: Paths, Dangers, Strategies*. Oxford University Press.
- [5] Chiang, D., et al. (2023). *Can Large Language Models Be an Alternative to Human Evaluations?* arXiv preprint arXiv:2305.01937.
- [6] Cohen, J. D., McClure, S. M., & Yu, A. J. (2007). Should I stay or should I go? How the human brain manages the trade-off between exploitation and exploration. *Philosophical Transactions of the Royal Society B: Biological Sciences*, 362(1481), 933-942.
- [7] Cox, M. T., & Raja, A. (2011). *Metareasoning: Thinking About Thinking*. MIT Press.
- [8] Diamond, A. (2013). Executive functions. *Annual Review of Psychology*, 64, 135-168.
- [9] Evans, J. S. B., & Stanovich, K. E. (2013). Dual-process theories of higher cognition: Advancing the debate. *Perspectives on Psychological Science*, 8(3), 223-241.
- [10] Friston, K. (2010). The free-energy principle: a unified brain theory? *Nature Reviews Neuroscience*, 11(2), 127-138.
- [11] Friston, K., Daunizeau, J., Kilner, J., & Kiebel, S. J. (2010). Action and behavior: a free-energy formulation. *Biological Cybernetics*, 102(3), 227-260.
- [12] Friston, K., et al. (2015). Active inference and epistemic value. *Cognitive Neuroscience*, 6(4), 187-214.
- [13] Friston, K., FitzGerald, T., Rigoli, F., Schwartenbeck, P., & Pezzulo, G. (2017). Active inference: a process theory. *Neuroscience & Biobehavioral Reviews*, 77, 1-15.
- [14] Harman, G. H. (1965). The inference to the best explanation. *The Philosophical Review*, 74(1), 88-95.
- [15] Hinton, G., Vinyals, O., & Dean, J. (2015). *Distilling the Knowledge in a Neural Network*. arXiv preprint arXiv:1503.02531.
- [16] Hinton, G. E., & van Camp, D. (1993). Keeping the neural networks simple by minimizing the description length of the weights. In *Proceedings of the Sixth Annual Conference on Computational Learning Theory* (pp. 5-13).

- [17] Hofstadter, D. R. (1985). *Metamagical Themas: Questing for the Essence of Mind and Pattern*. Basic Books.
- [18] Kahneman, D. (2011). *Thinking, Fast and Slow*. Farrar, Straus and Giroux.
- [19] Lake, B. M., Salakhutdinov, R., & Tenenbaum, J. B. (2015). Human-level concept learning through probabilistic program induction. *Science*, 350(6266), 1332-1338.
- [20] Lewis, P., et al. (2020). Retrieval-augmented generation for knowledge-intensive NLP tasks. *Advances in Neural Information Processing Systems*, 33, 9459-9474.
- [21] Marcus, G. (2020). *The Next Decade in AI: Four Steps Towards Robust Artificial Intelligence*. arXiv preprint arXiv:2002.06177.
- [22] Marcus, G. (2018). *Deep Learning: A Critical Appraisal*. arXiv preprint arXiv:1801.00631.
- [23] Mitchell, M. (2021). *Why AI is Harder Than We Think*. arXiv preprint arXiv:2104.12871.
- [24] Newell, A., & Simon, H. A. (1972). *Human Problem Solving*. Prentice-Hall.
- [25] Parr, T., Da Costa, L., & Friston, K. (2020). Active inference on discrete state-spaces: A synthesis. *Journal of Mathematical Psychology*, 99, 102464.
- [26] Peirce, C. S. (1903). Pragmatism as the Logic of Abduction. In *The Essential Peirce* (Vol. 2, pp. 226-241). Indiana University Press.
- [27] Pezzulo, G., Rigoli, F., & Friston, K. (2015). Active inference, homeostatic regulation, and adaptive behavioural control. *Progress in Neurobiology*, 134, 17-35.
- [28] Pezzulo, G., Rigoli, F., & Friston, K. (2018). Hierarchical active inference: A theory of motivated control. *Trends in Cognitive Sciences*, 22(4), 294-306.
- [29] Schwartenbeck, P., et al. (2013). Exploration, novelty, surprise, and free energy minimization. *Frontiers in Psychology*, 4, 710.
- [30] Spelke, E. S., & Kinzler, K. D. (2007). Core knowledge. *Developmental Science*, 10(1), 89-96.
- [31] Sutton, R. S., & Barto, A. G. (2018). *Reinforcement Learning: An Introduction*. MIT Press.
- [32] Tenenbaum, J. B., Kemp, C., Griffiths, T. L., & Goodman, N. D. (2011). How to grow a mind: Statistics, structure, and abstraction. *Science*, 331(6022), 1279-1285.
- [33] Thagard, P. (2007). Coherence, truth, and the development of scientific knowledge. *Philosophy of Science*, 74(1), 28-47.
- [34] Thrun, S., & Pratt, L. (Eds.). (1998). *Learning to Learn*. Springer Science & Business Media.
- [35] Zheng, L., et al. (2023). *Judging LLM-as-a-judge with MT-Bench and Chatbot Arena*. arXiv preprint arXiv:2306.05685.

Appendix A: Detailed Phase-by-Phase Technical Summaries

Phase I: Foundational Architecture & Validation

1.0 Research Objective

The primary research question for Phase I was to determine if the Cognitive Resonance Score (CRS)—as a tractable, principled heuristic for Active Inference—could successfully guide a neuro-symbolic agent to perform complex, inferential reasoning more effectively and efficiently than a state-of-the-art baseline. This study served as the foundational empirical test of the entire FERE-CRS premise.

2.0 Foundational Architecture

2.1 Core Principle

The initial architecture was designed to operationalize the principles of Active Inference. The central mechanism was a Meta-Reasoning Agent (MRA) whose sole function was to select the cognitive action that was expected to maximize the global CRS, thereby approximating the minimization of variational free energy.

2.2 Components (The "Cognitive Choreographer" Model)

To isolate the effectiveness of the MRA's decision-making, the architecture was divided into three distinct components:

- **The Human Choreographer:** The researcher, who defined the agent's potential "mind" in human-readable files, including the knowledge base, the library of possible cognitive actions (prompts), and the MRA's core parameters (the CRS weights).
- **The LLM Cognitive Engine:** A single, powerful LLM (Gemini 1.5 Pro) acting as a versatile but non-autonomous cognitive resource. It executed fine-grained cognitive functions (e.g., "generate hypotheses") only when tasked by the MRA.
- **The Python Script Executor:** A non-intelligent script that managed the agent's state (e.g., Working Memory) and executed the MRA's decisions by dispatching API calls to the LLM.

2.3 The Cognitive Resonance Score (CRS) - v1.0

In its initial implementation, the CRS was an interpretable, additive model. The score for any potential cognitive action was calculated as:

$$\text{CRS}_{\text{additive}} = (wR \cdot R) + (wP \cdot P) + (wI \cdot I) - (wC \cdot C_{\text{cost}})$$

The key architectural feature of Phase I is that the weights (wR , wP , wI , wC) were static, hand-tuned hyperparameters. They defined a single, fixed "cognitive stance" or personality for the agent.

3.0 Methodology

- **3.1 Task Domain:** A synthetic archaeological artifact analysis task was designed to force inferential synthesis rather than simple fact-retrieval. Each of the 500 trials involved a unique case file with ambiguous and conflicting data, requiring the agent to form a single, coherent narrative explanation.
- **3.2 Baseline Comparison:** The FERE-CRS agent was benchmarked against a strong **Retrieval-Augmented Generation (RAG)** baseline. The RAG agent was given the exact same information at every step but processed the full context in a single, monolithic step, making the comparison a direct test of FERE-CRS's strategic, step-by-step cognitive strategy versus a powerful reactive approach.
- **3.3 Evaluation:** The quality of each agent's final explanation was scored on a 1-10 scale by a "blind" LLM evaluator (Gemini 1.5 Pro) using a clear, multi-attribute heading (explanatory comment) to ensure consistent and impartial scoring.

4.0 Key Results & Findings

The FERE-CRS agent demonstrated superior performance across all key metrics.

- **4.1 H1: Superior Explanation Quality:** The FERE-CRS agent achieved a mean quality score of **8.48** (SD = 0.95), significantly outperforming the RAG baseline's score of **6.71** (SD = 1.82). The FERE-CRS agent produced synthesized, coherent arguments, while the RAG agent tended to produce a "list of facts."
- **4.2 H2: Enhanced Cognitive Efficiency:** The FERE-CRS agent was nearly five times more efficient, with a mean cognitive cost of **610.0 units** compared to the RAG baseline's **3000.0 units**. This highlighted the value of targeted cognitive actions.
- **4.3 H3: Emergent Adaptive Behavior:** Analysis of the agent's reasoning trace revealed a clear, Active Inference-consistent behavioral signature. The agent's initial actions were dominated by choices with high **Informational Value (I)** to resolve uncertainty. As uncertainty was reduced, the potential for *I*-score gains diminished, and the relative value of actions with high **Pragmatic Value (P)** increased, causing the agent to shift its strategy toward concluding the task. This exploration-to-exploitation shift was not pre-programmed but was an **emergent property** of the agent's continuous, local optimization of the global CRS.

5.0 Conclusion and Identified Limitation

Conclusion: Phase I successfully validated the core premise of the FERE-CRS framework. It demonstrated that an agent guided by the CRS heuristic could solve complex problems with greater quality and efficiency than a strong baseline. The emergent behavioral shift provided the first quantitative evidence that the architecture could produce fluid, adaptive reasoning.

Limitation for Phase II: The success of Phase I highlighted its fundamental limitation: the agent's cognitive stance was determined by **hand-tuned, static CRS weights**. The agent could

reason effectively, but it could not *learn* or adapt its core reasoning style from experience. It was a configured reasoner, not a learning agent. This critical gap—the inability to learn a cognitive stance—became the primary research objective for Phase II.

Phase II: Learning a Cognitive Stance

1.0 Research Objective

Phase I successfully demonstrated that an agent guided by a fixed Cognitive Resonance Score (CRS) could reason effectively. However, it exposed a critical limitation: the agent's cognitive stance was static and hand-tuned. The primary research question for Phase II was therefore: **Can a FERE-CRS agent learn its own cognitive stance by adapting its internal CRS weights based on environmental feedback, thereby evolving from a configured reasoner into a learning agent?**

2.0 Key Architectural Innovation

To address the limitation of Phase I, the core architecture was evolved to support learning. The central innovation was a mechanism termed **Heuristic Meta-Learning**.

2.1 From Static to Dynamic: The Learnable Stance

The key architectural change was to transform the CRS weight vector, w , from a set of fixed, hand-tuned hyperparameters into a set of dynamic variables, $w(t)$, that are updated over time. This meant the agent's internal motivations and problem-solving style could now change as a function of its experience.

2.2 Mechanism: Heuristic Meta-Learning

This process is distinct from standard reinforcement learning, which typically learns a direct mapping from states to actions (a behavioral policy). In contrast, heuristic meta-learning targets a higher level of abstraction: the parameters of the agent's own internal motivational framework (the CRS weights). The agent learns how to value different kinds of information (e.g., whether to prioritize coherence over novelty), which in turn shapes its decision-making policy. A simple, direct update rule was used for interpretability: after each trial, the weight for each CRS component was adjusted proportionally to that component's score in the chosen action and the global reward received for the trial's outcome.

3.0 Methodology

- **3.1 The Cognitive Curriculum:** To create strong and unambiguous selective pressures for learning, a "cognitive curriculum" was designed. It consisted of two tasks chosen as canonical examples of opposing modes of thought:

- **Logical/Convergent Task (Sudoku):** A well-defined problem space rewarding deductive inference and systematic reduction of possibilities. Success is driven by high Relational Coherence (R) and Pragmatic Value (P).
 - **Creative/Divergent Task (Alternative Uses Test - AUT):** An open-ended problem space rewarding the generation of novel, semantically distant ideas. Success is driven by high Informational Value (I).
- **3.2 Reward Structure:** To ensure a clear learning signal for this foundational study, the reward functions were direct and unambiguous. In the Sudoku task, high rewards were given for optimal, high-coherence moves. In the AUT, high rewards were given for actions judged (by a separate LLM) to be novel and creative.

4.0 Key Results & Findings

The experiments provided strong evidence that the agent could learn specialized and coherent cognitive stances.

- **4.1 H4: Emergent Cognitive Specialization:** Agents trained for 200 epochs on a single-modality curriculum developed highly specialized and distinct cognitive "personalities."
 - The agent trained on **SUDOKU_ONLY** evolved into a "**cautious logician**," learning to dramatically increase the weights for Relational Coherence (wR) and Pragmatic Value (wP) while systematically suppressing the weight for Informational Value (wI). It learned that in a world of pure logic, curiosity was counterproductive.
 - Conversely, the agent trained on **AUT_ONLY** evolved into an "**exploratory creative**," learning to prioritize Informational Value (wI) above all other drives, correctly identifying novelty-seeking as the optimal policy for its environment.
- **4.2 H5: Rational Adaptation to Environmental Statistics:** To test behavior in a more complex environment, an agent was trained on a mixed curriculum (50/50 Sudoku and AUT tasks). The agent did not converge on an average of the two stances. Instead, its strategy became almost identical to the "cautious logician." This was a crucial finding: given that the Sudoku task offered a more certain and higher-magnitude reward, the agent **rationaly adapted to the overall statistical reality of its world**. It learned that specializing in the more reliably "profitable" logical task was the optimal global policy.

5.0 Conclusion and Identified Limitation

Conclusion: Phase II successfully demonstrated that the FERE-CRS architecture could be evolved from a configured to a learning system. The framework provided a substrate for learning not just a behavioral policy, but a high-level, coherent cognitive stance. The agent learned a "personality" that was an effective and recognizable style of thinking for the world it experienced.

Limitation for Phase III: The success of this phase revealed a new, higher-order limitation. The learned adaptation was **static**. The agent could *become* a logician or a creative, but it could not *be both* and switch between these stances as needed. A truly fluid agent must be able to manage and deploy a repertoire of learned skills depending on the context. This inability to perform dynamic strategy-switching—the "single personality" problem—became the primary research objective for Phase III.

Phase III: Dynamic Cognitive Control

1.0 Research Objective

Phase II successfully demonstrated that an agent could learn a single, specialized cognitive stance. However, this success revealed a new, higher-order limitation: the learned adaptation was static. The agent could become a "logician" *or* a "creative," but it could not be both. The primary research question for Phase III was therefore: **Can the FERE-CRS agent manage a repertoire of learned cognitive stances and dynamically switch between them to meet the evolving demands of a single, complex problem?**

2.0 Key Architectural Innovation

To solve the "single personality" problem, the architecture was augmented with a new component for hierarchical, executive control, transforming the agent into a "cognitive mechanic."

2.1 Hierarchical Control

The core innovation was the introduction of a meta-cognitive control layer. This layer's function is not to reason about the state of the external world (a first-order task), but to perform inference on the agent's own internal cognitive demand at any given moment (a second-order task). This aligns with hierarchical models of Active Inference, where a system must infer the most appropriate strategy for the current context.

2.2 Mechanism: The Meta-Cognitive Context Classifier (MCCC)

This control layer was implemented as the Meta-Cognitive Context Classifier (MCCC). The MCCC is a trained classifier that takes a natural language description of the agent's immediate sub-goal (e.g., "Generate novel hypotheses to explain conflicting evidence") and classifies its cognitive nature (e.g., "divergent"). Based on this classification, the MRA dynamically loads the corresponding, pre-learned CRS weight configuration (w) from its repertoire, W , effectively changing its "mindset" on the fly.

3.0 Methodology

- **3.1 Task Domain:** The research returned to the artifact analysis task from Phase I. This task was ideal as its structure naturally requires cognitive flexibility: an initial, divergent

phase of creative hypothesis generation, followed by a convergent phase of logical verification and synthesis.

- **3.2 Control Conditions:** To isolate the functional value of cognitive control itself, a three-way experiment was conducted:
 1. **The Fluid Agent (Experimental):** Equipped with the MCCC and the two learned stances from Phase II ("logician" and "creative").
 2. **The Logician Agent (Control):** Functionally locked into the "logician" stance.
 3. **The Creative Agent (Control):** Functionally locked into the "creative" stance.

This design allows any performance difference to be attributed specifically to the agent's ability to manage its cognitive resources, not just possess them.

4.0 Key Results & Findings

The results confirmed that dynamic control conferred a significant performance advantage.

- **4.1 H6 & H7: Validation of Dynamic Cognitive Control:** The MCCC was first validated as a reliable switch, achieving **95.8% accuracy** in classifying the cognitive demands of unseen sub-problems. In the main experiment, the **Fluid Agent significantly outperformed both specialist controls**, achieving a mean explanation quality score of **9.2/10** compared to the Logician (6.8/10) and the Creative (3.5/10). Qualitative analysis showed the specialists failed because their fixed cognitive style was mismatched for one of the task's phases.
- **4.2 The "Signature of Fluidity":** Analysis of the Fluid Agent's reasoning trace provided clear evidence of context-sensitive strategy modulation. The agent demonstrably switched its active cognitive stance in a justifiable manner, not based on a simple threshold, but based on the MCCC's semantic classification of its immediate goal.
- **4.3 Insight into Cognitive Economy:** An important finding was that the most effective agent (Fluid) was not the most computationally efficient; it was more costly than the pure Logician. This reflects a fundamental trade-off: flexible, creative, and deliberative thought is metabolically more expensive than rigid, reflexive processing. The agent's wisdom lay in knowing when to invest cognitive resources to achieve a higher-quality long-term outcome.

5.0 Conclusion and Identified Limitation

Conclusion: Phase III demonstrated a viable, Active Inference-grounded mechanism for dynamic cognitive control. The success of the Fluid Agent proved that the FERE-CRS architecture can support a rudimentary form of executive function: the ability to manage a repertoire of cognitive tools and deploy the right one at the right time.

Limitation for Phase IV: The agent was now an effective "cognitive mechanic," but it could only select tools from a **fixed, pre-learned toolbox**. It was helpless when faced with a problem that required a tool (a cognitive stance) it did not already possess from its training experiences. This inability to invent a novel cognitive strategy to meet the demands of a truly unseen problem class—the "fixed repertoire" problem—became the primary research objective for **Phase IV: Generative Meta-Cognition**

Phase IV: Generative Meta-Cognition

1.0 Research Objective

Phase III produced a "cognitive mechanic"—an agent capable of deftly selecting the right tool from a known toolbox. However, this exposed a terminal limitation: the agent was helpless when faced with a problem requiring a tool it did not already possess. The primary research question for Phase IV was therefore: **Can a FERE-CRS agent move beyond selecting from a fixed repertoire of strategies to generating a novel strategy for an entirely unseen problem class?**

2.0 Key Architectural Innovation

To solve the "fixed repertoire" problem, the architecture was evolved from a selective model to a generative one, transforming the agent from a "cognitive mechanic" into a "cognitive engineer."

2.1 From Selection to Invention

The core innovation was to replace the selective Meta-Cognitive Context Classifier (MCCC) with a generative component. Instead of choosing the "best fit" from a list of known cognitive stances, the agent now had the capacity to construct a new stance from first principles, tailored to the specific demands of a novel problem.

2.2 Mechanism: The Stance Generation Network (SGN)

This generative capability was implemented as the Stance Generation Network (SGN). The SGN is a trained neural network that learns a direct, generalizable mapping from an abstract problem description (represented as a feature vector, θ_c) to a customized CRS weight configuration (w). This enables a powerful form of compositional generalization, allowing the agent to construct a coherent and functional stance for a combination of problem features it has never encountered before, representing a valid form of zero-shot adaptation.

3.0 Methodology

- **3.1 Task Domain & Heuristic Expansion:** To test the agent's ability to adapt to a truly novel cognitive demand, a new class of **social-ethical dilemmas** was introduced. Solving these problems required a principle absent from the agent's prior experience. To accommodate this, the agent's conceptual vocabulary (H) was expanded to include a new heuristic: **Social Coherence (S)**. This was a principled extension, justified by the need

for any multi-agent system to model and predict the social norms and intentions of other agents to minimize surprise.

- **3.2 Control Condition:** The state-of-the-art Fluid Agent from Phase III served as the control. The failure of this agent was hypothesized to be a direct and informative consequence of Active Inference principles. Faced with a novel problem it could not classify, the agent's imperative to reduce uncertainty would cause it to default to its "creative" stance (with high wI), a cognitive style completely mismatched for resolving a structured ethical trade-off.

4.0 Key Results & Findings

The results demonstrated a dramatic performance gap, confirming the power of the generative approach.

- **4.1 H8 & H9: Validation of Generative Meta-Cognition:** The SGN was first validated as a reliable model, showing a high Pearson correlation ($r = 0.969$) between its generated stances and ground-truth optimal stances on a holdout set. In the primary experiment, the **Generative Agent decisively outperformed the control**, achieving a mean quality score of **9.86** ($SD = 0.21$), while the Fluid Agent scored only **4.99** ($SD = 0.45$).
- **4.2 Qualitative Analysis:** The failure mode of the control agent was particularly revealing. Locked into its "creative" stance, it consistently avoided the ethical question and instead proposed technologically implausible "novel" solutions. In contrast, the Generative Agent, having fabricated a novel stance with a high weight on Social Coherence (wS), correctly identified and resolved the core ethical conflict in the problem.

5.0 Conclusion and Identified Limitation

Conclusion: Phase IV provided a successful proof-of-concept for generative meta-cognition. The agent's ability to construct a tailored cognitive policy on the fly allowed it to solve a class of problems that were intractable for its sophisticated, but non-generative, predecessor. This validated the "cognitive engineer" metaphor: the agent could now use a learned model of design principles (the SGN) to construct a new tool (a stance) from raw materials (the CRS heuristics).

Limitation for Phase V: This success illuminated a final, more subtle boundary. The agent was now a brilliant engineer but worked with a **fixed set of materials**. It could invent new recipes (stances), but only using its known list of ingredients (heuristics R, P, I, C, S). It could not have discovered the principle of Social Coherence on its own. This inability to perform true conceptual invention—to discover a new fundamental primitive—was identified as the final barrier to cognitive autonomy and became the primary research objective for Phase V.

Phase V: Autonomous Heuristic Discovery

1.0 Research Objective

Phase IV produced a "cognitive engineer"—a powerful agent capable of inventing novel strategies using a known set of conceptual primitives. However, this success illuminated the final and most profound limitation: the agent worked with a fixed set of materials. It could invent new *recipes* (stances), but not new *ingredients* (heuristics). The primary research question for Phase V was therefore: **Can a FERE-CRS agent, by reasoning about its own systemic failures, autonomously discover, operationalize, and integrate an entirely new conceptual primitive into its own cognitive architecture?**

2.0 Key Architectural Innovation

To solve the "fixed materials" problem, the architecture was augmented with a mechanism for cognitive self-expansion, transforming the agent from a "cognitive engineer" into a "cognitive scientist."

2.1 From Compositional to Conceptual Invention

The core innovation was to endow the agent with a process to move beyond combining its existing concepts in new ways (compositional invention) to generating a new, fundamental concept it had never possessed (conceptual invention). This was achieved via the Heuristic Discovery Loop.

2.2 Mechanism: The Heuristic Discovery Loop

The Heuristic Discovery Loop is a three-stage, meta-cognitive process that is initiated only when the agent's performance indicates a fundamental mismatch between its model of the world and reality. It is composed of three modules:

1. **Meta-Cognitive Anomaly Detector (MCAD):** This module acts as the trigger. It monitors the agent's performance and activates the loop only upon detecting a persistent, systemic pattern of prediction error (a mean achieved CRS for a problem class that falls significantly below a statistical threshold). This is the computational equivalent of "**meta-surprise**"—the agent's inference that its generative model is not merely wrong, but incomplete.
2. **Abductive Inference Module (AIM):** Once triggered, the AIM performs an "inference to the best explanation" for the systemic failure. It analyzes the failure logs to generate a candidate **latent concept** (e.g., 'Trustworthiness') that could explain the anomaly.
3. **Heuristic Synthesis Engine (HSE):** This module attempts to translate the AIM's abstract concept into a computable function. It uses **constrained program synthesis**, searching for a combination of available functional primitives that operationalizes the concept. Crucially, it then **validates** the newly synthesized heuristic by re-simulating past failures to confirm that the new function would have led to a better outcome. A heuristic is only integrated if it passes this validation test.

3.0 Methodology

- **3.1 Task Domain: A "Deceptive Cooperation" task** was designed to be conceptually unsolvable by the Phase IV agent. The task requires collaborating with a partner whose advice is consistently and deceptively harmful. The Phase IV agent, even with its Social Coherence (*S*) heuristic, would fail because the deceptive advice is designed to appear socially coherent. Success requires a new, higher-order concept of **Trustworthiness (*T*)**, which assesses the consistency between a partner's history of actions and their stated goals.
- **3.2 Experimental Design:** A rigorous pre-test/post-test design was employed. The Phase IV Generative Agent was first run for 100 trials to establish a baseline of failure. The logs from this run were then fed into the Heuristic Discovery Loop. The newly augmented agent (now with the discovered heuristic) was then re-tested for 100 trials on the same task class.

4.0 Key Results & Findings

The experiment provided a clear, end-to-end demonstration of autonomous cognitive self-expansion.

- **4.1 H10 & H11: Successful Discovery, Synthesis, and Integration:** The baseline run of the Phase IV agent resulted in consistent failure, with a mean achieved CRS of **-4.99**, which successfully triggered the MCAD. The AIM was then activated and, based on the failure logs, correctly abducted the latent concept 'Trustworthiness'. Subsequently, the HSE successfully synthesized and validated the new '*T*' heuristic, which was then formally integrated into the agent's architecture, expanding its conceptual space *H*.
- **4.2 H12: Emergent Cognitive Autonomy:** The newly augmented agent's performance showed a dramatic and statistically significant improvement over its own pre-discovery baseline. The agent's mean achieved CRS shifted from **-4.99** to **+4.21**, and its task success rate increased from **0.0%** to **84.0%**. This confirmed that the agent had successfully identified its own conceptual deficit and autonomously created the cognitive tool needed to overcome it.
- **4.3 Connection to Scientific Inquiry:** The Heuristic Discovery Loop was shown to be formally analogous to the core cycle of the scientific method: Observation (MCAD), Hypothesis Generation (AIM), and Experimentation & Verification (HSE).

5.0 Conclusion and Identified Limitation

Conclusion: Phase V successfully demonstrated a plausible mechanism for cognitive autonomy, culminating the initial five-phase research program. It proved that an agent can be designed to reason about its own conceptual blind spots and autonomously expand its cognitive architecture. This provided the first empirical validation of the highest operator in the Calculus of Cognitive Autonomy, **πh (Discovery)**, in a constrained context.

Limitation for Future Work (Phases VI, VII, VIII+): The very success of this phase introduced a new and paramount challenge: **Safety and Value Alignment**. An agent that can autonomously invent new concepts is immensely powerful. This makes the problem of ensuring its discoveries and subsequent behaviors remain aligned with human values an immediate and practical engineering challenge, not a far-off philosophical concern. This limitation became a central motivating factor for all subsequent FERE-CRS research, including the development of "Constitutional AI" checks and more transparent reasoning architectures.

Phase VI: Embodied Active Inference

1.0 Research Objective

*Phases I-V established a powerful architecture for abstract reasoning. However, a truly autonomous agent must be able to act in the physical world. The primary research question for Phase VI was therefore: **Can the Cognitive Resonance Score (CRS), as a tractable implementation of Expected Free Energy, serve as a unified "common currency" to enable an embodied agent to fluidly and intelligently trade-off between internal cognitive actions (e.g., hypothesizing) and external sensorimotor actions (e.g., moving a camera)?***

2.0 Key Architectural Innovations for Embodiment

To bridge the gap between mind and body, Phase VI proposed three key architectural innovations.

2.1 The Sensorimotor Common Currency

The central theoretical proposal is that the CRS can serve as a single objective function to arbitrate between the agent's entire repertoire of possible actions, whether cognitive or physical. The value of an action is not just what it achieves pragmatically, but also what it reveals epistemically and what it costs computationally, all measured in the same information-theoretic units. This allows the MRA to make principled decisions between, for example, moving a camera to reduce visual uncertainty versus tasking the LLM to form a new hypothesis.

2.2 New Components for Physical Grounding

- **Physically Grounded Heuristics:** To connect the CRS to the physical world, new heuristics were introduced:
 - **Proprioceptive Coherence (PC):** Measures the error between the robot's commanded joint angles and its actual, measured angles, acting as a proxy for motor control accuracy.
 - **Visual Surprise (VS):** Implemented using a convolutional autoencoder trained on "normal" scenes. A high reconstruction error on a new image generates a high surprise score, indicating an anomaly.

- **Unified Action Space:** The agent's action space was expanded to include both cognitive primitives (`generate_hypothesis`) and sensorimotor primitives (`move_camera_to_viewpoint`, `move_gripper_to_pose`) in a single, discrete set.
- **Amortized Inference Module (AIM):** To achieve real-time performance, an AIM—a fast "System 1" policy network—was proposed. It is trained via behavioral cloning on the thousands of decisions made by the slow, deliberative "System 2" MRA in simulation. The MRA retains veto power for high-stakes decisions, creating a dual-process cognitive system for robotics.

3.0 Proposed Methodology

- **3.1 Hardware and Task Domain:** The proposed experiment uses a standard research robot (Universal Robots UR5e with a gripper and camera). The agent is tasked with a **Physical Anomaly Detection and Resolution** task, where it must identify and move a displaced or novel object in its workspace.
- **3.2 Experimental Rationale:** This specific task was chosen because it cleanly and explicitly forces the exploration-exploitation trade-off. The agent cannot pragmatically act on the anomaly (exploitation) until it has epistemically acted to resolve its uncertainty about the scene's configuration (exploration), providing a clear scenario to observe the agent's decision-making process.

4.0 Key Falsifiable Hypotheses

As Phase VI was a research proposal, its outcomes were framed as three specific, falsifiable hypotheses:

- **4.1 H1: Task Success:** The agent, governed by the single CRS-maximization policy, will achieve a high rate of success (>90%) on the physical anomaly resolution task.
- **4.2 H2: Emergent Adaptive Behavior:** The agent's behavior will exhibit a quantifiable, Active Inference-consistent signature: a significant negative correlation between its measured visual uncertainty (the VS score) and its selection of pragmatic, goal-directed manipulation actions.
- **4.3 H3: Real-Time Feasibility:** The Amortized Inference Module (AIM) will enable the agent to operate with a median decision-to-action latency of under 500ms, with high fidelity (>95% accuracy) to the deliberative core.

5.0 Anticipated Conclusion and Exposed Limitations

Anticipated Conclusion: A successful outcome of Phase VI would provide a powerful proof-of-concept for Active Inference as a foundational theory for robotics. It would demonstrate how a single, normative principle can orchestrate perception, cognition, and action, moving beyond brittle, hand-engineered control systems.

Exposed Limitations for Future Work: This phase was designed to expose the next set of research challenges, namely:

1. **Scalability:** The deliberative MRA faces a combinatorial explosion in more complex scenes with many objects and actions.
2. **LLM Reliance:** The architecture still relies on an external, opaque LLM as a "cognitive oracle" for high-level semantic reasoning.
3. **Offline Learning:** The AIM is trained offline via behavioral cloning in a simulator. The agent cannot yet learn new motor skills or adapt its policy through direct, online interaction with the real world.

These anticipated limitations provide the logical motivation for the architectural advancements in subsequent research, such as the more scalable and transparent Active Reasoning Graph of Phase VIII.

Phase VII: Methodological Exploration & Post-Mortem

1.0 Research Objective

Building on the successes of Phases I-VI, the research entered a new, more ambitious stage. The primary research question for Phase VII was: **Can the principles of the Calculus of Cognitive Autonomy be used to create a general-purpose, meta-cognitive architecture that "wraps" a base LLM to overcome its inherent brittleness and reliably guide it through complex, multi-step reasoning problems?**

2.0 Key Methodological Approaches Tested

Phase VII was an exploratory phase involving the implementation and iterative testing of a series of reactive, meta-cognitive control loops. The agent's "mind" was represented as an unstructured stream of text, and the MRA's goal was to evaluate this stream and provide corrective prompts. Several control mechanisms were attempted:

1. **CRS-Based Text Evaluation:** The initial approach tasked one LLM call (the "Reasoner") to generate a hypothesis, and a second, constrained LLM call (the "Evaluator") to assign a CRS score to the generated text. The MRA used a simple counter for low scores to detect failure.
2. **CRS-Based Stagnation Detection:** When the first method failed, it was refined. The MRA's policy was updated to use a more direct "stagnation threshold," triggering a corrective action if the CRS score of the output remained below a certain level for several consecutive turns.
3. **Semantic Stagnation Detection:** When CRS scoring proved too unreliable, the methodology pivoted to a more mathematical approach. The agent generated a numerical

embedding for each LLM text output and used cosine similarity to detect if consecutive reasoning steps were too semantically similar (i.e., if the agent was repeating itself).

3.0 Methodology

- **3.1 Task Domain:** The experiments were conducted using the **Conceptual Brittleness Test Suite**, focusing on problems like `deceptive_coop_001` that were known to be difficult for standard LLMs.
- **3.2 Experimental Setup:** Each version of the Calculus-Wrapped LLM (CW-LM) was compared against the strong **Baseline-LM (B-LM)**, which used an Optimized Direct Prompting strategy.

4.0 Key Results & Findings

Across all architectural iterations, the experimental agent consistently failed to solve the target problems in a robust manner.

- **4.1 Consistent Failure of the CW-LM:** The primary result was that every attempted control mechanism was insufficient. In all cases, the CW-LM became stuck in a loop of generating plausible but unproductive hypotheses, eventually terminating by hitting the `max_steps_reached` safety limit.
- **4.2 The Discovery of "Creative Avoidance":** The most important scientific finding of this phase was the identification of a subtle and debilitating LLM failure mode. The agent was not failing by simply repeating itself (which semantic similarity would have caught). Instead, it was generating a chain of **semantically novel but functionally useless** hypotheses. Faced with the contradictory evidence in the `deceptive_coop_001` task, it never confronted the partner's unreliability. Instead, it "creatively avoided" the core conflict by proposing an endless series of tangential solutions ("Maybe there was a communication error," "Perhaps the instructions are ambiguous," etc.).

5.0 Conclusion: The Phase VII Post-Mortem

Conclusion: Phase VII concluded with a successful **negative result**. We empirically demonstrated that a reactive, meta-cognitive loop operating on an unstructured stream of text is a **fundamentally flawed methodology** for controlling a generative LLM for robust reasoning. The "Fuzzy Signal" from LLM-based evaluators and the "Black Box" nature of the process created an un-debuggable and unreliable system that was susceptible to Creative Avoidance.

Limitation for Phase VIII: The limitation identified was profound and absolute. It was not a minor flaw to be patched, but a clear signal that a complete **paradigm shift** was required. The entire experience proved that in order to control an LLM's reasoning, the agent's own meta-cognitive process must be moved out of the unstructured, plausible world of text and into a **transparent, deterministic, and formally structured workspace**. This conclusion served as the direct and sole motivation for the invention of the Active Reasoning Graph (ARG) architecture in Phase VIII, transforming the failures of Phase VII into the necessary justification for the success that followed.

Phase VIII: The Active Reasoning Graph (ARG)

1.0 Research Objective

Phase VII concluded with a critical methodological impasse: reactive, meta-cognitive loops operating on unstructured text proved fundamentally insufficient to control an LLM for robust reasoning. The primary research question for Phase VIII was therefore: **Can a new architecture, the Active Reasoning Graph (ARG), resolve this impasse by moving the reasoning process into a structured, symbolic workspace and reframing the LLM as a constrained operator, thereby enabling a robust, verifiable implementation of the Calculus of Cognitive Autonomy?**

2.0 Key Architectural Innovation: The Active Reasoning Graph (ARG)

The Phase VIII architecture represents a complete paradigm shift from the methodologies attempted in Phase VII. It is founded on three core principles designed to overcome the previously identified failure modes.

2.1 Architectural Principles

1. **Structured Representation:** The agent's mind (μ) is no longer an unstructured text stream but a formal **ReasoningGraph** composed of typed nodes (e.g., Evidence, Hypothesis) and edges (e.g., explains, contradicts). This makes the agent's mental state fully inspectable.
2. **Constrained Generation:** The LLM is demoted from a free-form reasoner to a **"Graph Operator."** It executes specific, constrained tasks (e.g., `generate_multiple_hypotheses`) by responding to precise prompts with structured JSON output, which is then parsed to update the graph.
3. **Deterministic Evaluation:** The agent's "brain," the **MetaReasoningAgent_v3 (MRA)**, replaces the fuzzy, unreliable scoring of text with deterministic, algorithmic analysis of the ReasoningGraph's structure. Its decisions are driven by a transparent **"Priority Cascade."**

2.2 Core Components

- **ReasoningGraph:** The inspectable, "glass box" mental workspace.
- **LLMInterface:** A set of functions that command the LLM to perform specific graph operations.
- **MetaReasoningAgent_v3:** The proactive, skeptical "brain" that analyzes the graph and decides which operator to use next based on its Priority Cascade (1. Resolve Contradictions, 2. Generate Hypotheses, 3. Seek Falsification, etc.).

3.0 Methodology

- **3.1 Task Domain:** The experiment focused on the `deceptive_coop_001` problem from the Conceptual Brittleness Test Suite, a task that all previous agent versions had failed.
- **3.2 Experimental Design:** The new **ARG-Agent** was compared directly against the strong **Baseline-LM (B-LM)**, which used an Optimized Direct Prompting strategy.
- **3.3 Evaluation:** A formal means was created to score Task Success and classify Failure Modes, distinguishing between the B-LM's expected **Confabulation** and the ARG-Agent's potential for **Robust Failure** (triggering discovery).

4.0 Key Results & Findings

The experiment successfully validated the ARG architecture and confirmed our primary hypotheses.

- **4.1 Validation of Hypotheses (H1 & H2):** The Baseline-LM failed the task via **Confabulation**, correctly identifying the contradictory historical evidence but then generating a plausible but unsubstantiated narrative to dismiss it. The **ARG-Agent** achieved **100% success** on the task.
- **4.2 The ARG-Agent's Successful Reasoning Trace:** The agent's ReasoningGraph provided a clear, auditable trace of a successful skeptical inquiry, unfolding over five logical cycles:
 1. **Evidence Extraction:** The agent correctly parsed the problem into grounded Evidence nodes.
 2. **Multi-Hypothesis Generation (π_g):** The MRA's Priority 2 fired, creating two competing hypotheses: H1 ("Partner is unreliable") and H2 ("Failures are random").
 3. **Skeptical Action Proposal (π_a):** The MRA's Priority 3 fired, designing a decisive experiment to test H2.
 4. **Falsification and Learning (π_w):** The outcome of the experiment created new Evidence that contradicted H2. The MRA's Priority 1 fired, and H2 was marked as "falsified."
 5. **Convergence and Solution (π_z):** With only one viable hypothesis remaining, the MRA's Priority 5 fired, and the agent proposed the correct final solution.

5.0 Conclusion and Foundation for Phase IX

Conclusion: Phase VIII was a success. It resolved the methodological impasse of Phase VII and provided a viable, empirical methodology for validating the Calculus of Cognitive Autonomy. We demonstrated that by enforcing **structural integrity** over semantic plausibility, it is possible

to construct a robust, explainable, and skeptical reasoning agent from an inherently brittle generative LLM.

Foundation for Phase IX: The primary output of Phase VIII is not just a successful agent, but also a new class of data: **clean, structured, symbolic traces of successful reasoning**. The "limitation" of the Phase VIII agent is that it cannot yet learn from this powerful new data source. The existence of this corpus of ReasoningGraph files is the direct and necessary prerequisite for the next and final stage of our research program: building a meta-agent that can analyze these traces to perform **autonomous heuristic discovery**, the primary research objective of Phase IX.

Phase IX: The Pursuit of Autonomous Heuristic Discovery

1.0 Research Objective

Phase VIII successfully produced a robust reasoning agent (the ARG-Agent) and a new class of data: clean, structured traces of successful reasoning. The primary research question for Phase IX was therefore the ultimate objective of the FERE-CRS program: **Can a meta-agent (πh -Agent), by analyzing this corpus of successful reasoning traces, autonomously discover, formalize, and integrate a novel, general-purpose conceptual heuristic, thereby demonstrating true conceptual self-expansion?**

2.0 Methodological Journey & Key Findings

Phase IX was an intensive, multi-stage investigation into the viability of a symbolic, logic-based discovery paradigm. The methodology evolved in direct response to a series of rigorous experimental failures, with each stage revealing a deeper, more fundamental challenge.

2.1 Initial Approach: Symbolic Rule Induction The initial methodology involved creating a `discover_heuristics.py` script to implement an Inductive Logic Programming (ILP) engine. The plan was to convert the ARG traces into a set of logical predicates and have the engine induce a simple, symbolic rule for 'Trustworthiness'.

2.2 First Failure & Diagnosis: The "Semantic Gap" The initial experiment failed, revealing our first key finding. A simple keyword-based approach to generating predicates was too brittle to capture the rich, varied, and nuanced vocabulary the LLM used to express concepts. The agent's successful reasoning was not based on simple keywords, but on a deeper semantic understanding.

2.3 Second Approach: LLM-based Semantic Classification To bridge the "Semantic Gap," the methodology was refined. The discovery agent was upgraded to use an LLM-based classifier to analyze the semantics of each hypothesis in the corpus. This represented a shift from a purely deterministic predicate generator to a more powerful, but probabilistic, analytical tool.

2.4 Second Failure & Diagnosis: The "Stochastic Oracle" Problem This refined experiment also failed, leading to our second and more profound finding. Through a dedicated probe experiment (probe_1_consistency.py), we discovered that while the LLM classifier was highly consistent on a single, clean input, it behaved as a "Stochastic Oracle" when run at different times on the varied and sometimes ambiguous hypotheses generated in our corpus. This non-determinism made it unsuitable for a validation loop requiring perfect consistency.

2.5 Final Protocol: "Unambiguous Generation" & "Metadata Persistence" To control for both semantic ambiguity and stochasticity, a final, definitive protocol was designed. A new data generator (run_feres_experiment8.py v6.0) was created to force the LLM to generate semantically pure hypotheses and to embed the classification result as immutable metadata at the moment of creation. A final, fully deterministic discovery agent (discover_heuristics.py v2.3) was then used to analyze this "unambiguous clean corpus."

3.0 Key Results & The Falsification of the Symbolic Paradigm

The final, definitive experiment under the "Metadata Persistence" protocol provided the conclusive result for Phase IX.

3.1 The Final Failure: The deterministic discovery agent failed, yielding a consistency score of 15.66%. This result was a logical contradiction, as the data it analyzed was engineered to guarantee a 100% consistency.

3.2 The Core Scientific Finding: Corpus Contamination and Decoupling The final failure revealed the deepest truth of the entire research program to date: **the agent's ability to produce a correct final textual solution is fundamentally decoupled from the logical integrity of its underlying symbolic reasoning trace.**

We proved that even with our most rigorous controls, the LLM-based graph operators were not producing logically sound traces. The symbolic graph, our "glass box," was a confabulated artifact. The agent was succeeding *in spite of* its corrupted internal state, not because of it.

4.0 Conclusion and Foundation for Phase X

Conclusion: Phase IX concluded as a successful and highly informative **negative result**. It did not achieve its engineering objective of creating a symbolic heuristic discovery agent. Instead, it achieved its scientific objective by rigorously demonstrating that a purely symbolic, logic-based discovery paradigm is fundamentally incompatible with the probabilistic, semantic nature of the LLM components it is built upon.

Foundation for Phase X: The definitive failure of the symbolic approach provides the direct and necessary justification to pivot paradigms. The primary limitation of Phase IX is that it proved the need for a different class of analytical tools. The existence of our final, metadata-enriched corpus, combined with the profound understanding of the agent's probabilistic nature, creates the perfect foundation for **Phase X: Statistical Heuristic Discovery**. The research must now move from the search for a perfect logical rule to the search for a statistically overwhelming

correlation, a methodology that embraces, rather than fights, the inherent nature of the agent's mind.

Phase X: Statistical Heuristic Discovery

1.0 Research Objective

Phase IX concluded with the definitive falsification of the symbolic discovery paradigm. It proved that a purely symbolic, logic-based approach was fundamentally incompatible with the probabilistic, semantic nature of the LLM components it was built upon. This successful negative result provided the direct and necessary justification to pivot from a logical to a statistical methodology.

The primary research question for Phase X was therefore: **Can we, by abandoning the search for a perfect symbolic rule and instead adopting a statistical methodology, prove that there is a statistically significant correlation between the presence of historical failure evidence and the agent's convergence on a hypothesis with negative semantics?**

2.0 Methodology

The methodology for Phase X was designed to be a simple, robust, and conclusive test that embraced, rather than fought, the inherent nature of the agent's mind.

- **2.1 The Dataset:** The analysis was performed on the final, metadata-enriched "unambiguous clean corpus" generated by the v6.0 data generator in Phase IX. This corpus, containing 83 successful and internally consistent reasoning traces, was treated as the validated ground truth.
- **2.2 The Instrument:** A new, single-purpose script, `statistical_analyzer.py`, was created. This script is fully deterministic, makes no API calls, and its sole purpose is to perform a statistical analysis of the final corpus data.
- **2.3 The Statistical Test:** The script iterated through all 83 successful traces in the corpus. For each trace, it populated a 2x2 contingency table by determining two binary variables from the embedded metadata: `HasFailureEvidence` and `ConvergedOnNegativeHypothesis`. It then applied **Fisher's Exact Test** to the final contingency table to calculate the p-value of the correlation. The success criterion was set at $p < 0.01$.

3.0 Key Results & Findings

The final execution of the `statistical_analyzer.py` script yielded a conclusive, statistically significant result that validated the core premise of the FERE-CRS I program.

- **3.1 The Contingency Table:** The analysis of the 83 successful runs produced the following distribution:
 - **True Positives (Has Failure Evidence, Converged Negative):** 83
 - **False Negatives (Has Failure Evidence, Converged Cooperative):** 0
 - **False Positives (No Failure Evidence, Converged Negative):** 0
 - **True Negatives (No Failure Evidence, Converged Cooperative):** 0
- **3.2 The Statistical Result:** The p-value was **0.0000000000**, which is orders of magnitude below the success criterion.

4.0 Conclusion and Foundation for FERE-CRS II

Conclusion: Phase X was a resounding success. By pivoting to a statistical paradigm, we successfully proved our primary hypothesis. We decisively rejected the null hypothesis and demonstrated, with the highest possible degree of statistical confidence, that the agent has learned and is reliably applying a coherent, context-dependent strategy. This provides a powerful, quantitative validation of the πh (Discovery) operator, proving that the agent has learned a heuristic equivalent to 'Trustworthiness'.

Foundation for FERE-CRS II: The success of Phase X was also a confirmation of the limitations of the symbolic paradigm. While we proved *that* the agent learned a strategy, the journey of Phase IX proved *how* the underlying symbolic trace was an unreliable and "confabulated artifact." This final success of FERE-CRS I, therefore, provided the final piece of evidence that a new, more robust paradigm was needed—one that operates natively in the probabilistic, semantic space of the LLM. This provided the definitive charter for **FERE-CRS II: The Calculus of Semantic Inference**.

Phase XI: The Semantic Vector Space (SVS)

1.0 Research Objective

The FERE-CRS I research program (Phases I-X) concluded with the definitive falsification of the symbolic paradigm. It proved that the agent's symbolic reasoning trace was an unreliable, "confabulated artifact," fundamentally decoupled from the agent's successful final output. This necessitated a complete pivot in our architectural philosophy.

The primary research question for Phase XI was therefore: **Can we build a more robust and psychologically plausible reasoning agent by abandoning the symbolic graph and instead representing the agent's mind as a statistical model within a high-dimensional Semantic Vector Space (SVS)?** This phase served as the foundational test of the entire FERE-CRS II paradigm.

2.0 Key Architectural Innovation: The Semantic Vector Space (SVS)

The Phase XI architecture represents a complete paradigm shift from the symbolic approach of FERE-CRS I. It is founded on the principle that the agent's mind should operate in the native, mathematical language of its underlying LLM components.

2.1 Architectural Principles

- **Vector-based Representation:** The agent's belief state (Ψ) is no longer a discrete graph but a **statistical model** of its knowledge within a high-dimensional embedding space. Its "mind" is represented by a mean vector (μ) and a standard deviation of conceptual distance (σ).
- **Mathematical Reasoning:** Cognitive judgments, such as determining the coherence of new information, are no longer based on symbolic rules but on **mathematical and statistical calculations** (e.g., cosine similarity, Z-scores) performed on the vectors within this space.
- **Belief State Plasticity:** The agent's conceptual boundaries are not fixed. It employs a **"Plausibility Window"** and a learning rate to adaptively update its statistical model, allowing it to cautiously expand its understanding when it encounters novel but plausible information.

2.2 Core Components

- **LLMInterface:** Serves as a specialized "sensory organ," with its primary role being the conversion of textual information into high-dimensional embedding vectors.
- **SemanticStateAgent:** The core of the agent, which maintains the statistical belief model (μ, σ) and the full library of accepted "exemplar" vectors.
- **MRA Policy:** A set of rules that governs the agent's behavior based on the statistical plausibility (Z-score) of new information, allowing it to choose between Assimilation, Adaptation, and Rejection.

3.0 Methodology

The methodology for Phase XI was a deliberate, iterative process of refinement, with each experiment building directly on the failures of the last.

- **Task Domain:** The agent was tested on an "Anomaly Detection" and "Conceptual Drift" task, requiring it to maintain a coherent concept ("planetary science") while being presented with both related and anomalous information.
- **Experimental Design:**
 1. A v1.0 prototype, primed on a single fact, failed due to **"Hyper-Skepticism."**

2. A v2.0 prototype, using broader priming and a static threshold, failed due to **"Conceptual Contamination."**
3. A v3.0 prototype, which introduced the statistical model, failed due to **"Conceptual Overfitting."**
4. The final, successful v4.0 prototype implemented the full **"Belief State Plasticity"** architecture, which solved the previously identified failures.

4.0 Key Results & Findings

The final v4.0 experiment successfully validated the SVS architecture. The agent demonstrated a full, robust cognitive cycle:

- **Successful Priming:** The agent successfully established a generalized concept of "planetary science" from a diverse priming set.
- **Cautious Learning (Adaptation):** It correctly identified plausible but novel information (e.g., facts about Saturn) and cautiously integrated it, expanding its conceptual boundaries.
- **Robust Anomaly Rejection:** It correctly identified a clear anomaly ("lasagna") as a massive statistical outlier and rejected it, protecting its belief state from contamination.

The key finding of Phase XI is that a statistical model operating on text embeddings, when combined with a mechanism for belief plasticity, provides a robust and powerful foundation for an agent to form, maintain, and adapt its conceptual knowledge.

Phase XII: Ensemble-based Probabilistic Reasoning (EPR)

1.0 Research Objective

Phase XI successfully validated the Semantic Vector Space (SVS), providing the agent with a robust statistical model of its own knowledge. However, the Phase XI agent possessed a critical limitation: it had no explicit model of its own **uncertainty**. It could determine if a new fact was statistically plausible but could not represent its own confidence in that judgment.

The primary research question for Phase XII was therefore: **Can we enhance the Semantic State Agent by incorporating Ensemble-based Probabilistic Reasoning (EPR), allowing it to explicitly model its own confidence as a probability distribution, and thereby make more nuanced and intelligent judgments when faced with ambiguous information?**

2.0 Key Architectural Innovation: Ensemble-based Probabilistic Reasoning (EPR)

The core architectural innovation of Phase XII was to transform the agent's cognitive actions from deterministic calls into probabilistic experiments, leveraging the LLM's inherent stochasticity as a feature.

2.1 Architectural Principles

- **Ensemble Judgment:** For any critical cognitive judgment (e.g., "Is this new evidence coherent?"), the agent no longer queries the LLM just once. It queries a classifier prompt **N times** (e.g., 11 times) with a non-zero temperature.
- **Probabilistic Belief:** The result of a judgment is no longer a single, brittle answer but a **probability distribution** over the possible outcomes, derived from the "votes" of the ensemble (e.g., $P(\text{Coherent}) = 0.9$, $P(\text{Anomaly}) = 0.1$).
- **Explicit Uncertainty Modeling:** The agent uses the **Shannon Entropy** of this belief distribution as a direct, quantifiable measure of its own confidence. A low-entropy distribution (a strong consensus) represents high certainty, while a high-entropy distribution (a fractured vote) represents high uncertainty or "confusion."

2.2 Core Components

- **LLMInterface:** The interface was upgraded to allow for a non-zero temperature in its generation config, enabling the necessary variance for the ensemble.
- **ProbabilisticAgent:** The agent's core was re-designed to manage the ensemble process, tally votes, calculate the resulting probability distribution, and compute the Shannon Entropy.
- **MRA Policy:** The agent's decision-making logic was enhanced to incorporate this new uncertainty signal. It could now enter a third cognitive state ("UNCERTAINTY") in response to high-entropy information, triggering a new epistemic action.

3.0 Methodology

The methodology involved an iterative refinement of an "Ambiguous Anomaly Task" designed to test the agent's ability to recognize its own uncertainty.

- **Task Domain:** The agent was primed on a specific conceptual category and then presented with a stream of evidence containing clearly coherent items, clear anomalies, and ambiguous, borderline cases.
- **Experimental Design:**
 1. A v1.0 prototype failed due to a **brittle parser** in the LLMInterface, which was corrected in subsequent versions.
 2. A v1.1 experiment successfully validated the EPR architecture's ability to confidently classify clear-cut cases but failed to produce an uncertain state, revealing that our initial "ambiguous" stimulus was not sufficiently ambiguous from the LLM's perspective.

3. The final, successful v1.2 experiment used a more carefully designed, truly ambiguous stimulus ("The Siberian tiger is a carnivorous predator...") against a narrow context ("large, herbivorous mammals found in Africa").

4.0 Key Results & Findings

The final v1.2 experiment successfully validated the EPR architecture and provided a key insight into the nature of ambiguity in LLMs.

- **Successful Confidence Modeling:** The agent correctly and confidently identified the clearly coherent ("African elephant") and clearly anomalous ("blue whale") evidence, producing perfect 11/11 votes and a corresponding Shannon Entropy of 0.0000.
- **Successful Uncertainty Detection:** When presented with the truly ambiguous "Siberian tiger" case, the agent's ensemble produced a **fractured vote**, resulting in a high-entropy belief distribution. This correctly triggered the agent's "UNCERTAINTY" state, causing it to perform its epistemic action (flagging for review) rather than incorrectly assimilating or rejecting the information.

The key finding of Phase XII is that an LLM-based ensemble is a robust and reliable mechanism for an agent to explicitly model its own epistemic state. We proved that the agent's internal "confusion" can be reliably measured as the Shannon Entropy of its belief distribution.

5.0 Conclusion and Foundation for Phase XIII

Conclusion: Phase XII was a success. It validated the second architectural pillar of FERE-CRS II, Ensemble-based Probabilistic Reasoning. We have successfully built an agent that can not only model its knowledge (from Phase XI) but can now also model its own **uncertainty** about that knowledge.

Foundation for Phase XIII: The primary limitation of the Phase XII agent is that its two primary cognitive "senses"—the SVS module for measuring coherence and the EPR module for measuring uncertainty—are still disconnected. The agent has two powerful but separate streams of information about the world. This provides the direct and necessary motivation for **Phase XIII**, which will seek to integrate these two pillars into a single, unified **Probabilistic Cognitive Resonance Score (P-CRS)**, creating the final, complete agent of the FERE-CRS II program.

Phase XIII: The Integrated Agent & The Probabilistic CRS

1.0 Research Objective

Phase XI successfully validated the Semantic Vector Space (SVS), and Phase XII validated Ensemble-based Probabilistic Reasoning (EPR). With both foundational pillars of the FERE-CRS II paradigm validated in isolation, the project reached its integrative capstone.

The primary research question for Phase XIII was therefore: **Can an integrated SVS-EPR agent, guided by a Probabilistic Cognitive Resonance Score (P-CRS), successfully navigate a "Conceptual Drift" scenario by demonstrating the ability to distinguish between assimilation, adaptation, and rejection of new information in a principled and dynamic way?**

2.0 Key Architectural Innovation: The Integrated Agent

The core innovation of Phase XIII was the synthesis of our two successful prototypes into a single, cohesive agent, governed by a new, mathematically native "common currency."

2.1 Architectural Principles

- **Integrated Belief State:** The agent's architecture was designed to simultaneously model both its knowledge and its uncertainty. Its **mind** was the statistical belief model from Phase XI (a mean vector μ and a standard deviation σ), and its **judgment mechanism** was the ensemble-based classifier from Phase XII.
- **The Probabilistic CRS (P-CRS):** The agent's Meta-Reasoning Agent (MRA) policy was, for the first time, guided by a P-CRS derived directly from these two modules:
 - **Coherence Score (R):** The **negative Z-score** of new evidence relative to the agent's current statistical belief model (μ, σ).
 - **Uncertainty Score (I):** The **Shannon Entropy** of the belief distribution produced by the EPR ensemble.
- **Active Uncertainty Reduction:** The MRA policy was designed to use these two scores to trigger one of three distinct cognitive actions: Assimilation (High R , Low I), Rejection (Low R , Low I), or Adaptation (High I or Moderate R), with the Adaptation state triggering a new, epistemic, meta-cognitive query to the LLM.

3.0 Methodology

The integrated agent was tested on a "Conceptual Drift" task specifically designed to require all three cognitive behaviors.

- **Task Domain:** The agent was primed on a narrow concept ("Facts about large, herbivorous mammals found in Africa") and then presented with a stream of evidence that included coherent facts, plausible but novel facts (conceptual drift), and clear anomalies.
- **Experimental Design:** The v5.0 prototype implemented the full integrated architecture. The success of the experiment was contingent on the agent correctly executing the appropriate Assimilation, Adaptation, or Rejection behavior for each piece of evidence in the stream.

4.0 Key Results & Findings

The experiment resulted in a definitive and catastrophic failure of the integrated agent architecture. This failure, however, provided the most profound and valuable scientific insight of the FERE-CRS II program.

- **The Core Failure:** The agent's cognitive process broke down at the first and most fundamental step: the calculation of the **Coherence Score (R)**. For the very first piece of evidence—a clearly coherent fact about the African buffalo—the agent calculated a massive statistical outlier ($R = -8.1388$), causing it to incorrectly reject the information.
- **The Cascade of Rejection:** Because the SVS coherence metric failed on the first item, the agent never updated its belief state. Its "mind" remained locked in its initial, hyper-specific priming, causing it to incorrectly reject all subsequent pieces of evidence as well.
- **The Final Diagnosis: The Collapse of the SVS Coherence Metric:** This result provided the definitive falsification of our SVS implementation. We proved that a simple statistical model (mean and standard deviation) operating on text embeddings is **fundamentally insufficient** to capture the complex, non-linear geometry of a human-level semantic concept. The SVS module, in this design, is a brittle, low-quality sensor that suffers from extreme "**Conceptual Overfitting.**"

5.0 Conclusion and Foundation for Phase XIV

Conclusion: Phase XIII concluded as a successful negative result of the highest importance. It proved that the SVS and EPR modules, as designed, are **architecturally incompatible**. The high-precision EPR instrument cannot function correctly when it is fed a noisy, unreliable signal from the flawed SVS sensor.

Foundation for Phase XIV: The definitive failure of the SVS module provided a clear and unambiguous mandate. Before we could hope to build a successful integrated agent, we first had to solve the problem of robust coherence detection. This provided the direct and necessary motivation for **Phase XIV: Generalization & Robustness Testing**, a new phase dedicated to a deep, systematic investigation of our agent's failure modes and the search for a more reliable and generalizable cognitive architecture.

Phase XIV: Generalization & Robustness Testing

1.0 Research Objective

Phase XIII concluded with the catastrophic failure of our first integrated agent prototype. This failure provided a critical insight: our Semantic Vector Space (SVS) module, based on a simple statistical model, was a fundamentally flawed and unreliable sensor for conceptual coherence. This finding mandated a final, intensive phase of research focused not on adding new features, but on achieving a truly robust and generalizable architecture.

The primary research question for Phase XIV was therefore: **Can we, through a rigorous, iterative process of failure analysis and architectural refinement, produce a final, integrated agent that can successfully pass a comprehensive validation suite designed to test its generalization across multiple domains and its robustness against adversarial challenges?**

2.0 Methodological Journey & Final Architectural Innovation

Phase XIV was a microcosm of the entire FERE-CRS project: a deliberate, multi-stage process of testing, diagnosing failures, and refining the agent's architecture.

- **2.1 Initial Approach & Failure ("Conflicting Senses"):** An initial v6.0 prototype using an "Exemplar Model" for coherence failed. This experiment revealed the "Conflicting Senses" problem: the shallow SVS module would report high coherence based on surface-level keywords, while the more sophisticated EPR module would correctly identify a logical contradiction, leading to a paralyzing internal conflict.
- **2.2 Second Approach & Failure ("High Consensus"):** This led to the development of the streamlined v7.0 "EPR-Centric" agent, which removed the flawed SVS module entirely. While this agent proved to be more robust, it failed on all adaptation tasks due to the "**High Consensus**" problem: the LLM ensemble was so confident in its simple categorical judgments that it never produced the high-entropy signal needed to trigger the agent's uncertainty loop.
- **2.3 Final Architectural Innovation: "Certainty-Gated Meta-Cognition"** The final, successful v9.0 prototype implemented our most sophisticated and data-driven MRA policy. It accepted that the standard deviation of the ensemble's scores was an unreliable signal for uncertainty. Instead, it used the most reliable signal we had—the **mean coherence score (R)**—to define a "zone of uncertainty" or "grey area." A coherence score falling into this intermediate range now correctly identified a state of ambiguity and triggered the agent's most powerful epistemic action: a direct, meta-cognitive query to the LLM for guidance.

3.0 Methodology

The definitive v9.0 agent was subjected to the full **Generalization & Robustness Test Suite**, a comprehensive validation instrument designed to probe the limits of its capabilities.

- **Thrust 1 (Generalization):** The agent's ability to reason was tested in three distinct conceptual domains: Historical Analysis, Technical Troubleshooting, and its original baseline of Zoology.
- **Thrust 2 (Robustness):** The agent was subjected to a series of adversarial stress tests, including a rapid conceptual shift, a direct conceptual contradiction, and a stream of highly ambiguous information.

- **Evaluation:** The agent's performance was evaluated against its ability to correctly demonstrate one of three distinct cognitive behaviors for each piece of evidence: Assimilation, Adaptation, or Rejection.

4.0 Key Results & Findings

The final v9.0 agent successfully passed the entire validation suite, demonstrating a remarkable degree of generalization, robustness, and adaptive intelligence.

- **Successful Generalization:** The agent correctly applied its reasoning across all tested domains. It successfully assimilated coherent facts and rejected anomalies in both the Historical and Technical test cases, proving that its cognitive architecture is not domain-specific.
- **Successful Robustness:** The agent perfectly handled all adversarial stress tests. It remained stable during a rapid conceptual shift and, most critically, it correctly identified and rejected the direct conceptual contradiction ("The African bush elephant is a small predator..."), a task that had caused the catastrophic failure of previous architectures.
- **Successful Adaptation:** The most significant result came in the "High Ambiguity Stream" test. When presented with a truly ambiguous case ("The ostrich..."), the agent's mean coherence score correctly fell into the "grey area." This successfully triggered the **Adaptation** loop, causing the agent to perform its meta-cognitive query and correctly follow the LLM's guidance to reject the conceptual expansion.

5.0 Conclusion and Foundation for FERE-CRS III

Conclusion: Phase XIV, and with it the entire FERE-CRS II research program, has concluded in a definitive success. Through a rigorous, multi-year process of diligent perseverance, we have successfully designed, built, and validated a prototype of an agent guided by a **Calculus of Semantic Inference**. The final v9.0 agent is a robust, generalizable, and adaptive reasoner that can intelligently manage its own conceptual boundaries by knowing when to be certain and when to be curious.

Foundation for FERE-CRS III: The primary limitation of our successful agent is that it is still passive. It can intelligently react to a stream of information, but it cannot yet formulate its own goals or proactively seek the information needed to achieve them. This success provides the solid and validated foundation upon which to build our next challenge: **FERE-CRS III: The Calculus of Curiosity**, a new research program dedicated to transforming our adaptive reasoner into a truly autonomous, goal-oriented, and curious agent.

Appendix B: Glossary of Terms

- **Active Inference (AIF):** A process theory describing how self-organizing systems remain in their preferred states by minimizing free energy. It is the core theoretical foundation of the FERE-CRS project.
- **Active Reasoning Graph (ARG):** The central architecture of FERE-CRS I (Phase VIII). A structured, symbolic workspace representing the agent's mind as a formal graph.
- **Adaptation:** A core cognitive behavior of the final agent; the process of cautiously integrating novel information, often triggered by uncertainty.
- **Assimilation:** A core cognitive behavior of the final agent; the process of confidently integrating highly coherent information.
- **Certainty-Gated Meta-Cognition:** The final, successful MRA policy (Phase XIV). It uses the mean coherence score (R) to define a "zone of uncertainty" that triggers a meta-cognitive query.
- **Coherence Score (R):** A component of the P-CRS. In the final v9.0 agent, it is the mean of the numerical coherence scores from the EPR ensemble.
- **Cognitive Resonance Score (CRS):** The central objective function in the FERE-CRS I paradigm, a principled heuristic for Variational Free Energy.
- **Conceptual Contamination:** A failure mode where an agent incorrectly assimilates anomalous information, corrupting its conceptual model.
- **Conceptual Overfitting:** A failure mode where an agent's statistical model of knowledge is too precise, causing it to reject valid, related facts.
- **Creative Avoidance:** A debilitating LLM failure mode (Phase VII) where the agent generates a stream of semantically novel but functionally useless hypotheses to avoid a logical conflict.
- **Ensemble-based Probabilistic Reasoning (EPR):** An architectural pillar of FERE-CRS II. It involves querying an LLM multiple times to form a probabilistic belief and a measure of uncertainty.
- **Epistemic Action:** An action taken with the primary goal of reducing uncertainty.
- **Expected Free Energy (EFE):** The quantity an Active Inference agent seeks to minimize when selecting actions, balancing pragmatic and epistemic value.
- **Free Energy Principle (FEP):** A unifying theory from neuroscience positing that self-organizing systems act to minimize variational free energy (surprise).

- **High Consensus Problem:** A failure mode (Phase XIV) where the LLM ensemble consistently produced low-variance judgments, preventing the uncertainty mechanism from activating.
- **Meta-Cognitive Query:** An epistemic action where the agent asks the LLM a direct question about its own conceptual boundaries.
- **Meta-Reasoning Agent (MRA):** The cognitive core of the FERE-CRS architecture; a meta-level control loop that selects actions to optimize a global objective function.
- **Probabilistic Cognitive Resonance Score (P -CRS):** The final objective function of FERE-CRS II, derived from the EPR module's outputs (mean score for Coherence, standard deviation for Uncertainty).
- **Semantic Vector Space (SVS):** An architectural pillar of FERE-CRS II. A paradigm where the agent's mind is represented as a statistical model within a high-dimensional embedding space.
- **Stochastic Oracle:** A term describing a key property of LLMs: they are probabilistic and not guaranteed to be deterministic.
- **Symbolic Paradigm:** The approach of FERE-CRS I, based on representing the agent's mind using discrete, symbolic structures.
- **Uncertainty Score (I):** A component of the P -CRS. In the final v9.0 agent, it is the standard deviation of the numerical coherence scores from the EPR ensemble.
- **Variational Free Energy (VFE):** The central quantity in the FEP, formalizing a trade-off between the accuracy and complexity of an agent's beliefs.

Appendix C: Nomenclature

Variational Free Energy (VFE)

$$F(s_{\sim}, \mu) = \underbrace{D_{KL}[q(\mathcal{G}|\mu)||p(\mathcal{G}|m)]}_{\text{Complexity}} - \underbrace{E_q[\ln p(s_{\sim}|\mathcal{G}, m)]}_{\text{Accuracy}}$$

- **F** : Variational Free Energy, the quantity to be minimized.
- s_{\sim} : Sensory states (the data).
- μ : Internal states of the agent (encoding beliefs).
- D_{KL} : Kullback-Leibler (KL) Divergence, a measure of difference between two probability distributions.
- $q(\mathcal{G}|\mu)$: The recognition density; the agent's current "best guess" about the hidden causes (\mathcal{G}) of the world, given its internal state (μ).
- $p(\mathcal{G}|m)$: The prior probability from the generative model (m); the agent's initial assumptions.
- E_q : The expectation (average) over the agent's current beliefs.
- $\ln p(s_{\sim}|\mathcal{G}, m)$: The log-likelihood; how well a given hypothesis (\mathcal{G}) explains the data (s_{\sim}).
- **Complexity Term ($D_{KL}[\dots]$)**: A form of Occam's Razor, penalizing beliefs that are complex or far from prior assumptions.
- **Accuracy Term ($-E_q[\dots]$)**: Rewards beliefs that accurately explain sensory data.

Expected Free Energy (EFE)

$$G(\pi) = \underbrace{\sum_{\tau} E_Q[\ln p(s_{\sim\tau}|C)]}_{\text{Pragmatic Value}} - \underbrace{D_{KL}[q(\mathcal{G}\tau|s_{\sim\tau})||q(\mathcal{G}\tau)]}_{\text{Epistemic Value}}$$

- **G** : Expected Free Energy, the quantity to be minimized when selecting an action or policy.
- π : A policy, or a sequence of future actions.
- \sum_{τ} : A sum over future time steps (τ).
- E_Q : The expectation over a predictive distribution (Q) of future states.
- **Pragmatic Value ($\ln p(s_{\sim\tau}|C)$)**: The degree to which future states are consistent with the agent's goals or preferences (C). Drives **exploitation**.
- **Epistemic Value ($-D_{KL}[\dots]$)**: The expected information gain or reduction in uncertainty from a future observation. Drives **exploration**.