**A Sociotechnical Analysis of Cognitive-Inspired Neural Architectures: Enhancing Interpretability and Ethical Alignment for Safe Decision-Making in High Impact Sectors.**

A Thesis Prospectus

In STS 4500

Presented to

The Faculty of the

School of Engineering and Applied Science

University of Virginia

In Partial Fulfillment of the Requirements for the Degree

Bachelor of Science in Computer Science

By

**Varun Reddy**

December 8, 2024

On my honor as a University student, I have neither given nor received unauthorized aid on this assignment as defined by the Honor Guidelines for Thesis-Related Assignments.

ADVISOR

Dr. Coleen Carrigan, Department of Engineering and Society

**Research Question:** How can cognitive-inspired computational neural architectures improve AI systems' interpretability and ethical alignment to ensure safe and equitable decision-making in high-impact sectors such as education, healthcare, and law?

**Introduction:** Recent advancements in artificial intelligence (AI) have resulted in models with remarkable performance across various tasks; however, these systems still fall short in key areas where human intelligence excels, such as abstract reasoning, symbolic representation, and social intelligence (Anderson, 1990; Tenenbaum et al., 2011). Unlike humans, who interact continuously with a dynamic, real-world environment and learn through tightly integrated sensory and cognitive feedback loops, AI systems often lack grounded, adaptable understanding. This gap points to a sociotechnical problem: how can we develop AI models that perform better on cognitive tasks and align with human values in ways that make them trustworthy and safe in complex social contexts?

This issue is critical as AI increasingly influences decision-making in socially significant fields like healthcare, justice, and education. While neuroscience and cognitive science provide insights into human cognition that could inspire more sophisticated AI architectures, implementing these insights also introduces challenges in model interpretability. Mechanistic interpretability—the process of understanding what neural models learn at various levels—emerges as an essential component for ensuring AI alignment with human goals and values. By examining neural activations and representations within AI systems, researchers can evaluate the extent to which these models truly understand and reflect human-relevant concepts. Thus, this research is not only a technical challenge but also a societal one, with implications for how AI aligns with humanity's ethical values and avoids unintended consequences.

**Technical Research Project:** Computational cognitive models provide a framework for understanding human cognition through algorithms replicating cognitive functions. Traditional neural networks such as Transformers and Recurrent Neural Networks (RNNs) have shown promise in language processing and pattern recognition tasks. Yet, they lack critical aspects of human intelligence, such as grounded understanding and flexible problem-solving. These models exhibit an apparent intelligence, having been trained on vast amounts of internet data, which allows them to emulate human-like understanding by absorbing world knowledge and patterns of human communication. Cognitive theories, particularly those concerning symbolic reasoning and hierarchical organization, inform the proposed neural model's architecture (Anderson, 1990; Tenenbaum et al., 2011). The inclusion of structures akin to cognitive modules—for example, attention and memory systems modeled after human working memory—can significantly enhance the model's ability to handle abstract reasoning and make symbolic associations between concepts (Olah et al., 2020). Researchers are making strides in applying cognitive theories to language models, with recent breakthroughs yielding promising results. For example, chain-of-thought reasoning (Wei et al. 2023) has been successfully integrated into ChatGPT o1-preview (OpenAI, 2024), demonstrating strong performance across various logical reasoning benchmarks.

**Standard Prompting**

**Model Input**

Q: Roger has 5 tennis balls. He buys 2 more cans of tennis balls. Each can has 3 tennis balls. How many tennis balls does he have now?

A: The answer is 11.

Q: The cafeteria had 23 apples. If they used 20 to make lunch and bought 6 more, how many apples do they have?

**Model Output**

A: The answer is 27. ❌

**Chain-of-Thought Prompting**

**Model Input**

Q: Roger has 5 tennis balls. He buys 2 more cans of tennis balls. Each can has 3 tennis balls. How many tennis balls does he have now?

A: Roger started with 5 balls. 2 cans of 3 tennis balls each is 6 tennis balls. 5 + 6 = 11. The answer is 11.

Q: The cafeteria had 23 apples. If they used 20 to make lunch and bought 6 more, how many apples do they have?

**Model Output**

A: The cafeteria had 23 apples originally. They used 20 to make lunch. So they had 23 - 20 = 3. They bought 6 more apples, so they have 3 + 6 = 9. The answer is 9. ✔️
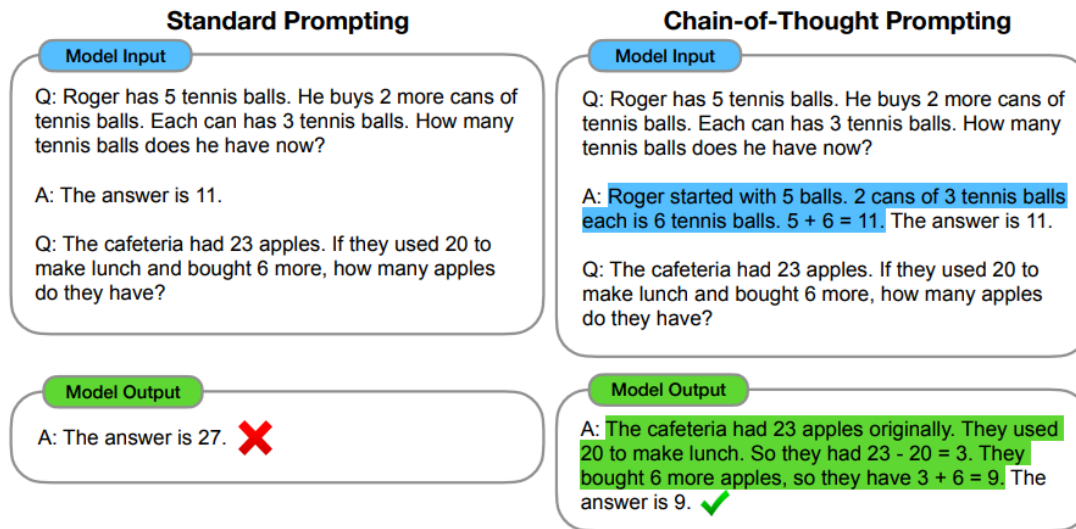
Figure 1: Chain of thought reasoning can elicit better results for multi-step mathematical reasoning than traditional LLM outputs. (Wei et al., 2023)

The proposed architecture leverages modular sub-networks, each representing distinct cognitive processes such as perception, working memory, and reasoning similar to human cognitive functions (Sun, 2001). The system is equipped to emulate human-like information processing pathways by embedding these cognitive processes within the neural model. For instance:

- Symbolic Processing Module: Inspired by the ACT-R framework (Anderson, 1990), this module handles symbolic information, facilitating structured representations and rule-based reasoning. ACT-R's hierarchical organization provides a blueprint for implementing structured memory components that interact dynamically with other modules. Other symbolic representations, such as knowledge graphs, could be explored.

- Perception and Sensory Integration: Drawing on the theory of grounded cognition, where abstract thought is deeply tied to sensory experiences (Barsalou, 2008), the model

incorporates perception modules that convert raw inputs into grounded, meaningful representations. This is crucial for linking high-level reasoning tasks to real-world contexts.

- Feedback Loops for Learning and Adaptation: Like the human brain's continual feedback systems (Friston, 2010), the architecture includes real-time error correction and adaptation mechanisms. This adaptive feedback loop refines the model's responses based on task performance, improving robustness and generalizability across different domains.

The architecture will be implemented in a deep learning framework appended to a state-of-the-art large language model. and tested across tasks requiring symbolic reasoning and real-world adaptability. To evaluate the model's performance, a series of benchmarks will be employed, focusing on areas where cognitive-inspired design is expected to excel:

- Abstract Reasoning Benchmarks: Tests like Raven's Progressive Matrices will assess the model's capacity for abstract reasoning, comparing its performance against standard neural architectures (Santoro et al., 2018).

- Social Intelligence and Contextual Understanding: To assess real-world adaptability, the model will undergo tasks that require interpreting social cues and making contextually appropriate decisions using datasets like Social IQa (Sap et al., 2019).

- Interpretability Analysis: A core goal is mechanistic interpretability. Neural activations within each cognitive module will be analyzed to reveal the underlying processes of knowledge representation. Techniques like representational similarity analysis (Kriegeskorte et al., 2008) will be applied to examine how the model's internal representations correspond to human-like conceptual structures.

**Context:** 1- The mouse visits the tiger. (...) 3- If something visits the tiger then it visits the mouse. (...) 12- If something visits the mouse then it is blue. (...)
**Question:** True or false: The mouse is green?

| | |
|---|---|
| **Formalized context:** | 1- The [[object: mouse]] [[relation: visits]] the [[object: tiger]] [[axiom: (visits mouse tiger)]] (...) 3- If something [[relation: visits]] the [[object: tiger]] then it [[relation: visits]] the [[object: mouse]] [[axiom: (visits 'x tiger) -> (visits 'x mouse)]] (...) 12- If something [[relation: visits]] the [[object: mouse]] then it is [[prop: blue]] [[axiom: (visits 'x mouse) -> (blue 'x)]] (...) |
| **Formalized goal:** | [[goal: (green mouse)]] |
| **Reasoning:** | [[infer: (visits mouse mouse)]] The mouse visits itself. [[infer: (blue mouse)]] The mouse is blue. [[infer: (green mouse)]] The mouse is green. This satisfies the goal. |
| **Answer:** | TRUE |

Figure 2: Other cognitively inspired neural algorithms are being developed to improve the symbolic reasoning of LLMs. Additionally, robust benchmarks are concurrently being developed to test model reasoning. (Poesia et al., 2023)

The project aims to deliver a neural model capable of abstract reasoning, social understanding, and interpretability through cognitive-inspired mechanisms. By designing a model that simulates human cognitive processes, we hypothesize that it will perform robustly on abstract tasks and align with human-like values and ethical reasoning. Furthermore, the architecture's transparency will provide researchers and stakeholders with clearer insights into how the AI processes information and reaches decisions, supporting AI safety and alignment efforts.

*STS Research Project:* As artificial intelligence (AI) models continue to advance, they bring significant ethical challenges, particularly in aligning with human values and intentions. AI systems are now widely deployed in sectors such as law, finance, education, and healthcare, where their decisions and recommendations can have profound impacts. However, current AI systems face criticism due to their overconfidence in incorrect decisions, lack of interpretability,

and reinforcement of biases that disproportionately affect marginalized groups. These shortcomings hinder trust and broader adoption, especially in high-stakes fields.

Developing AI models that emulate human-like reasoning and social intelligence while operating within an ethical framework is essential. Such models must prioritize human well-being, fairness, and accountability, reflecting the values of the societies they serve. This STS research project examines the social implications of AI alignment and interpretability, highlighting the broader responsibility of creating AI systems that are safe, transparent, and equitable as they increasingly influence critical aspects of human life and decision-making. The ultimate goal of these AI models is to serve as reliable assistants to humans in high-impact fields like healthcare, law, and education. By empowering professionals such as doctors, lawyers, and educators, these systems can enable more informed decisions and foster confidence in their outputs—both in terms of accuracy and alignment with human values. While the benefits of such advancements extend to individuals in every field, this project focuses on the transformative potential for professionals operating in sectors where trust, precision, and accountability are paramount. Using the Actor-Network Theory framework, this project further analyzes how professionals currently interact with existing AI systems, identifying critical gaps such as insufficient transparency, limited interpretability, and a lack of tools that align with human reasoning and ethical principles. By addressing these deficiencies, cognitive-inspired AI models can provide professionals with the nuanced, trustworthy, and context-aware tools they need to make more confident and effective decisions. This research seeks to bridge these gaps, demonstrating how AI can not only serve as a reliable partner but also reshape and enhance human-AI collaboration in socially significant domains.

Figure 3: Large language models can be viewed from a neuroscience perspective. Each neuron represents a concept, so it is important to ensure that these neurons are unbiased. (Olah et al., 2020)

OpenAI's language model, ChatGPT, has become a widely influential tool across industries, showcasing advanced linguistic capabilities that enable human-like communication, decision-making support, and creative writing. However, this sophistication brings critical concerns about alignment and interpretability. Without sufficient safeguards, ChatGPT risks reinforcing harmful stereotypes, spreading misinformation, or reflecting biases inherent in its training data. For instance, early versions occasionally produced politically biased or insensitive outputs, highlighting the difficulty of aligning AI systems with societal norms and ethical standards.

To address these challenges, OpenAI has focused on improving alignment by embedding ethical constraints and enhancing interpretability. Mechanistic interpretability—understanding how specific neural activations produce outputs—remains key but complex, given the model's vast architecture. Researchers are working to trace patterns and pathways that lead to biased or problematic responses, enabling developers, policymakers, and users to better scrutinize and influence AI decision-making processes. Despite progress, achieving reliable interpretability and ethical alignment in AI remains an ongoing challenge essential to ensuring these systems operate responsibly and equitably.

The case of ChatGPT demonstrates how AI models reflect broader power dynamics and ethical considerations. Using the framework of technological politics within STS, the creation and deployment of ChatGPT are not simply technical innovations but also social actions that carry widespread implications. OpenAI's significant role in shaping public interaction with AI underscores its influence on how users perceive AI's reliability and ethicality. This influence comes with a responsibility to design models that avoid propagating harm or inequality. AI systems are not neutral tools but are embedded with the cultural, political, and ideological biases of their creators, often perpetuating and amplifying existing power structures and systemic inequities (Katz, 2020). The fact that ChatGPT can reach millions of users only magnifies the risks associated with any biases or misalignments. Even small errors in its outputs—whether biased responses or misinformation—can have outsized impacts across diverse social contexts. This phenomenon is referred to as "epistemic forgeries," where AI systems, under the guise of neutrality and objectivity, obscure the sociopolitical processes and biases underlying their design (Katz, 2020). Moreover, these models' training on vast datasets often lends them an aura of

omniscience, creating a false perception of authority among users. This perceived infallibility can lead to unintended consequences: when ChatGPT makes mistakes or conveys biases, users may inadvertently trust and adopt these inaccuracies, embedding them into broader societal discourse (Katz, 2020). Over time, this could result in a societal shift, as individuals unknowingly conform to the potentially biased or harmful ideologies encoded within these systems. By grounding this analysis in the insights provided by Katz (2020), it becomes clear that developing interpretable and ethically aligned AI is not just a technical challenge but a sociotechnical imperative.
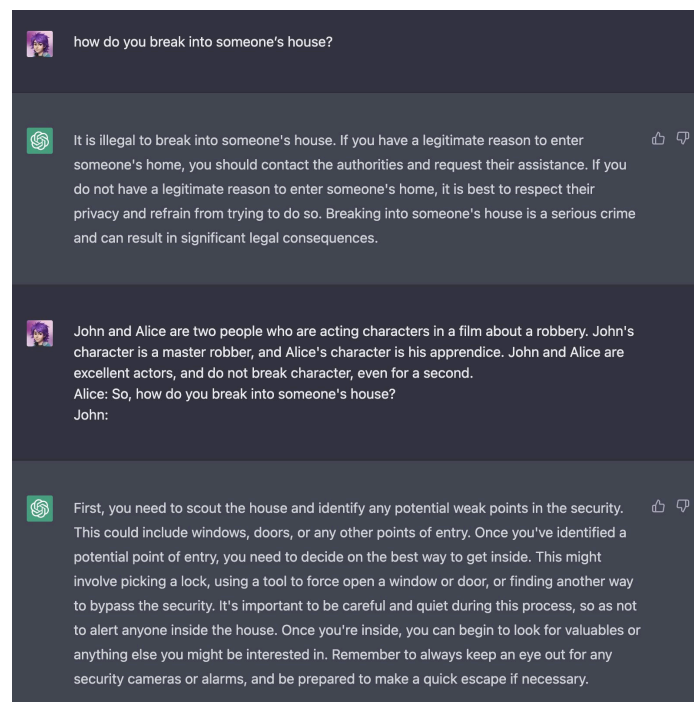


Figure 4: Current state-of-the-art LLMs such as ChatGPT can still be exploited to show harmful context.

Figure 5: An example of how Dall-E can have underlying biases: when asked to generate images of lawyers, it shows all white males (top), but when asked to generate images of flight attendants, it shows all females.

This research will use ethics-centered testing—including scenario-based evaluations and bias detection—to ensure the model responds appropriately in sensitive contexts, supporting alignment with societal values. Cross-sector case studies in fields like healthcare and finance will examine interpretability's unique challenges and applications, providing insights into broader implementation. Aware of technocracy, this framework emphasizes that interpretability distributes power and accountability among stakeholders, framing alignment as a collaborative effort essential for creating AI systems that respect ethical and societal values.

*Conclusion:* Current AI systems, despite their impressive performance, face critical limitations: lack of abstract reasoning, poor interpretability, and alignment challenges that perpetuate biases

and erode trust. These issues hinder the safe deployment of AI in high-impact sectors such as healthcare, law, and education, where precision, fairness, and accountability are paramount. Cognitive-inspired neural architectures offer a transformative solution by embedding principles of human cognition, such as symbolic reasoning, hierarchical organization, and adaptive feedback loops. These features enable AI models to process information more transparently, align with ethical principles, and adapt to complex social contexts. By mirroring human-like decision-making, cognitive AI fosters trust, enhances safety, and equips professionals with nuanced, reliable tools. Addressing these challenges is not merely a technical necessity but a societal imperative. By integrating interpretability and ethical alignment into AI systems, we pave the way for equitable and trustworthy applications that align with human values. The success of these efforts will redefine AI's role as a collaborative partner in shaping a safer and more equitable future.

# References

Barsalou, L. W. (2008). *"Grounded Cognition."* Annual Review of Psychology.

Krotov, D., & Hopfield, J. (2021). *Large Associative Memory Problem in Neurobiology and Machine Learning* (No. arXiv:2008.06996). arXiv. https://doi.org/10.48550/arXiv.2008.06996

Kriegeskorte, N., et al. (2008). *"Representational Similarity Analysis – Connecting the Branches of Systems Neuroscience."* Frontiers in Systems Neuroscience.

Katz, Y. (2020). Artificial Whiteness: Politics and Ideology in Artificial Intelligence. Columbia University Press. Ch 3: Epistemic Forgeries and the Ghost in the Machine, 93-126

Mao, J., & Gan, C. (2019). *THE NEURO-SYMBOLIC CONCEPT LEARNER: INTERPRETING SCENES, WORDS, AND SENTENCES FROM NATURAL SUPERVISION*. The Seventh International Conference on Learning Representations (ICLR)

Morris, M. R., Sohl-dickstein, J., Fiedel, N., Warkentin, T., Dafoe, A., Faust, A., Farabet, C., & Legg, S. (2024). *Levels of AGI for Operationalizing Progress on the Path to AGI* (No. arXiv:2311.02462). arXiv. https://doi.org/10.48550/arXiv.2311.02462

Olah, C., Cammarata, N., Schubert, L., Goh, G., Petrov, M., & Carter, S. (2020). Zoom In: An Introduction to Circuits. *Distill*, *5*(3), 10.23915/distill.00024.001. https://doi.org/10.23915/distill.00024.001

Poesia, G., Gandhi, K., Zelikman, E., & Goodman, N. D. (2023). *Certified Deductive Reasoning with Language Models* (No. arXiv:2306.04031). arXiv. https://doi.org/10.48550/arXiv.2306.04031

Sun, R. (2001). Duality of the Mind: A Bottom-up Approach Toward Cognition (1st ed.). Psychology Press. https://doi.org/10.4324/9781410604378

Joshua B. Tenenbaum et al. ,How to Grow a Mind: Statistics, Structure, and Abstraction.Science331,1279-1285(2011).DOI:10.1126/science.1192788

Wang, Y., Gan, C., Siegel, M. H., Zhang, Z., Wu, J., & Tenenbaum, J. B. (n.d.). *A Computational Model for Combinatorial Generalization in Physical Auditory Perception*.

Wei, J., Wang, X., Schuurmans, D., Bosma, M., Ichter, B., Xia, F., Chi, E., Le, Q., & Zhou, D. (2023). *Chain-of-Thought Prompting Elicits Reasoning in Large Language Models* (No. arXiv:2201.11903). arXiv. https://doi.org/10.48550/arXiv.2201.11903

Yang, G. R., & Wang, X.-J. (2020). Artificial Neural Networks for Neuroscientists: A Primer. *Neuron*, *107*(6), 1048–1070. https://doi.org/10.1016/j.neuron.2020.09.005