digital-duck / **agreeability**    Public

<> Code     ⊙ Issues     ⅄ Pull requests     ▷ Actions     ⊞ Projects     ⊘ Security     ⟋ Insights

**agreeability** / **paper-v0.1.md**    ⧉                                                                    ⋯

wgong  upd                                                                    daa6974 · now    ⟲

119 lines (60 loc) · 18.3 KB

# Beyond Agreeability: The Critical Thinking Gap Between Large Language Models and Human Cognition

Wen Gong ([Digital-Duck](#))

Abstract: Large Language Models (LLMs) have demonstrated impressive capabilities in various natural language processing tasks, generating human-like text, translating languages, and even writing different kinds of creative content. However, a crucial limitation remains: a lack of independent critical thinking. This paper argues that current LLMs primarily operate using System 1 thinking (fast, intuitive, pattern recognition) but lack the robust System 2 capabilities (slow, deliberate, logical reasoning, meta-cognition) necessary for true understanding and general intelligence. This deficiency is illustrated through the "Emperor's Court" analogy, where a lack of critical feedback leads to flawed reasoning and a "chamber echo effect." We focus on the specific challenge of commonsense reasoning as a key example of this System 2 gap and discuss the implications for AI development and the importance of maintaining human agency. We further propose that developing a "model of self," involving meta-cognition, is a necessary step towards more balanced and robust AI systems.

# Introduction:

The rapid advancement of Artificial Intelligence (AI), particularly in the realm of Large Language Models (LLMs), has sparked widespread excitement and speculation about the future of technology and the potential arrival of Artificial General Intelligence (AGI). LLMs have demonstrated remarkable abilities in generating human-like text, translating languages, writing various creative content, and even producing precise working codes and software applications. These impressive feats have led to a surge of interest and investment in AI, with many predicting a transformative impact on various industries and aspects of human life. However, amidst this enthusiasm, it is crucial to maintain a critical perspective and carefully examine the fundamental limitations of current AI systems. This paper argues that while LLMs excel at certain cognitive tasks, primarily those involving pattern recognition and prediction—which align with what psychologists call System 1 thinking (fast, intuitive, and associative)—they fundamentally lack crucial aspects of human cognition, most notably the capacity for independent critical thinking, abstract reasoning, and a robust "model of self"—capabilities associated with System 2 thinking (slow, deliberate, and analytical). This deficiency has significant implications for AI development, the responsible deployment of AI systems, and the broader societal impact of this rapidly evolving technology. To illustrate this critical gap, we draw a parallel with historical power structures, specifically the "Emperor's Court" analogy, demonstrating how a lack of critical feedback and independent thought can lead to flawed reasoning and a "chamber echo effect," mirroring the limitations observed in current LLMs. We will then delve into a technical analysis of this deficiency, focusing on the specific challenge of commonsense reasoning as a key example of the gap between current AI and human-level general intelligence. Finally, we will propose research directions that could lead to more balanced, robust and ethical AI systems capable of more sophisticated forms of cognition.

# The Emperor's Court Analogy (System 1 in Action):

Throughout history, rulers have often surrounded themselves with advisors and courtiers. While wise rulers sought out advisors who provided honest counsel and critical feedback, others found themselves surrounded by "yes-men"—individuals who primarily agreed with the ruler's opinions, reinforcing existing beliefs and suppressing dissent. This dynamic, often referred to as a "chamber echo effect," can have disastrous consequences, leading to flawed decisions, a lack of innovation, and ultimately, harm to the ruler and their realm. A compelling example is the reign of Zhu Yuanzhang, the founding emperor of the Ming Dynasty in China (1368-1644). Zhu Yuanzhang, after rising from humble beginnings to become emperor, was known for his suspicious nature, autocratic rule, and tendency to distrust his advisors. He established a highly centralized government and often relied on advisors who catered to his whims and avoided contradicting him, fearing his wrath. This created an environment where critical feedback was suppressed, and dissenting voices were silenced. The lack of independent evaluation and critical thinking within the court created a "chamber echo" where the emperor's own thoughts, biases, and sometimes even paranoia, were constantly amplified and reflected back to him, without any external correction. This lack of critical input led to numerous misjudgments, purges of officials, and policies that ultimately hampered the long-term stability and prosperity of the Ming Dynasty. Note: The selection of Zhu Yuanzhang as one historic data point was inspired by a 50 episode TV series titled "传奇皇帝朱元璋 Legendary emperor Zhu Yuanzhang". Similar examples exist in history, and even in current organizations and affairs.

This historical example provides a powerful analogy for understanding the limitations of current LLMs. Like an emperor surrounded by "yes-men," LLMs are trained on vast amounts of data and are designed to generate outputs that align with the statistical patterns they have learned from that data. They are essentially trained to "agree" with the data, reinforcing existing biases and lacking the ability to independently evaluate the validity, truthfulness, or even the coherence of the information they process. This is a clear manifestation of System 1 thinking: fast, automatic responses based on learned associations, without the deliberate, reflective analysis characteristic of System 2. Just as the Emperor's court lacked the critical thinking necessary to challenge the Emperor's views, LLMs lack the internal mechanisms to critically evaluate the information they are given as input or the outputs they generate. This "chamber echo effect" in LLMs can manifest in various ways, such as the amplification of biases present in the training data, the generation of factually incorrect or nonsensical outputs that are nonetheless grammatically correct and contextually plausible within the limited scope of the training data, and the inability to adapt to novel situations that require a deeper understanding of context, meaning, and real-world knowledge.

# Technical Analysis: The System 1/2 Divide in AI

### System 1: Pattern Recognition Without Understanding (Current AI)

Current AI systems, particularly those based on deep learning architectures like Large Language Models (LLMs), excel at identifying statistical patterns within vast datasets. They operate through a process of pattern matching, predicting the most probable next word or sequence of words based on the statistical relationships they have learned. This allows them to perform remarkably well on tasks like text generation, translation, and summarization. However, this proficiency is achieved through a form of "surface-level understanding." While LLMs can manipulate symbols effectively, generating grammatically correct and contextually relevant text, they lack a true comprehension of the underlying meaning or the real-world implications of their outputs. This is analogous to System 1 thinking in humans: fast, automatic, intuitive, and associative. It's a system that excels at rapid responses and pattern recognition but lacks the depth of analysis and critical evaluation. Furthermore, these models are susceptible to biases present in their training data and often struggle to generalize to unseen situations or contexts.

### The Challenge of Commonsense Reasoning (A System 2 Deficiency)

Commonsense reasoning requires understanding and applying everyday knowledge about the world. It involves abilities like understanding physical laws, social norms, and causal relationships. Current AI struggles significantly with this, indicating a deficiency in System 2 thinking. Examples like the Visual Commonsense Reasoning (VCR) dataset (Zellers et al., 2018), SWAG (Zellers et al., 2018), and research on physical and social commonsense (Bisk et al., 2020; Sap et al., 2019) demonstrate this. These failures highlight the lack of key System 2 components: causal reasoning, counterfactual thinking, and mental models of the world. For example, understanding that a "suit of armor cannot conduct an orchestra" (Bisk et al., 2020) requires knowledge of physical properties, human capabilities, and social contexts – knowledge that current AI lacks. This lack of grounding in the physical world is a significant obstacle for current AI systems.

### The Missing Ingredient: Meta-cognition and the "Model of Self" (Essential for System 2)

Meta-cognition, "thinking about thinking," is crucial for critical thinking and a "model of self." It enables reflection on cognitive processes, bias identification, and reasoning evaluation. Current AI lacks this, operating in a "black box" fashion. This prevents them from understanding why they make mistakes, correcting biases, planning learning, or assessing output certainty. This lack directly impairs commonsense reasoning. Without reflecting on their understanding, they cannot identify inconsistencies or flawed inferences.

## The "Aha!" Moment: Insight and the Role of Parallel Processing (A System 2 Phenomenon)

A distinctive characteristic of human cognition is the experience of sudden insight, often referred to as an "aha!" moment. These moments of clarity often arise after a period of focused thought, incubation, or even seemingly unrelated experiences. They represent a shift from System 1's associative thinking to System 2's deliberate and analytical processing. The human brain, being a massively parallel neural network, can process information simultaneously across multiple interconnected regions. This parallel processing allows for the integration of diverse information and the emergence of novel connections and insights. These insights often feel spontaneous or unexpected, even though they are the result of underlying cognitive processes.

A classic example of this is the story of Isaac Newton and the falling apple. While observing an apple falling from a tree, Newton is said to have had a sudden insight into the nature of gravity, realizing that the same force that pulled the apple to the ground also kept the moon in orbit around the Earth. This "aha!" moment, triggered by a seemingly simple observation, led to the formulation of his groundbreaking theory of universal gravitation. This exemplifies how the human brain can connect seemingly disparate pieces of information to generate novel and profound insights.

Current AI systems, while employing parallel processing within their neural network architectures, lack the global integration and flexible connectivity of the human brain. They operate primarily in a feed-forward manner, processing information sequentially through layers of neurons. This architecture makes it difficult for them to generate the kind of sudden insights that characterize human creativity and problem-solving. The lack of a "model of self" and meta-cognition further hinders their ability to engage in the kind of reflective processing that often precedes "aha!" moments.

## Beyond Scaling: The Need for New Approaches

The "scaling hypothesis"—that simply increasing model size and data will lead to AGI—is challenged by these limitations. While scaling improves performance on some tasks, it does not address the fundamental lack of System 2 capabilities. New approaches are needed, including:

- Integrating symbolic reasoning and knowledge graphs (Ahn et al., 2021).

- Developing methods for causal inference.

- Exploring architectures that support meta-cognition.

- Embodied AI and grounded cognition.

## Implications:

The limitations of current AI have significant real-world implications:

- Misleading Marketing and Hype: The marketing of AI often overstates its capabilities, creating unrealistic expectations and obscuring limitations. This can lead to misinformed decisions and a lack of preparedness for the actual capabilities and limitations of deployed systems.

- Ethical Concerns: Deploying AI systems lacking critical thinking in high-stakes domains (healthcare, law, finance) raises serious ethical concerns. Flawed reasoning, bias amplification, and a lack of accountability can have severe consequences.

- Impact on Human Work and Society: Overreliance on AI could diminish human critical thinking skills and agency. If we become accustomed to accepting AI outputs without question, we risk losing our ability to make informed judgments.

- The Importance of Human-in-the-Loop Systems: Human oversight is crucial, especially in situations requiring critical thinking or commonsense judgment. Human-in-the-loop systems, where humans provide feedback and intervention, are essential for responsible AI deployment.

## Recommendations:

To address these limitations and move towards more robust and beneficial AI, we recommend:

- Focus on System 2 Capabilities: Research should prioritize developing AI systems with System 2 functions like causal reasoning, meta-cognition, and abstract thinking.

**agreeability** / **paper-v0.1.md**                                    ↑ Top

| Preview |   Code      Blame |                        Raw  ⧉  ⤓      ☰

experience necessary for developing commonsense understanding.

- Developing Better Evaluation Metrics: New metrics are needed to assess System 2 capabilities, going beyond simple accuracy to measure reasoning, causal understanding, and meta-cognition.

- Promoting Public Understanding of AI Limitations: Public education is crucial to prevent unrealistic expectations and ensure responsible AI development and deployment.

## Conclusion:

This paper has argued that while Large Language Models (LLMs) have demonstrated impressive capabilities in various natural language processing tasks, they fundamentally lack the capacity for independent critical thinking, a key characteristic of human cognition. By applying the System 1/2 framework, we have shown that current LLMs operate primarily through System 1 thinking—fast, intuitive, and associative—excelling at pattern recognition and prediction. However, they struggle with tasks that require System 2 thinking—slow, deliberate, analytical reasoning, and meta-cognition. The "Emperor's Court" analogy effectively illustrated how a lack of critical feedback and independent thought, mirroring the limitations of current LLMs, can lead to flawed reasoning and a "chamber echo effect." We focused on the specific challenge of commonsense reasoning as a prime example of this System 2 deficiency, highlighting the lack of causal reasoning, counterfactual thinking, and robust mental models of the world in current AI systems. Furthermore, we argued that meta-cognition, the ability to "think about thinking" and develop a "model of self," is a crucial missing ingredient for achieving true general intelligence. Simply scaling up current models is unlikely to bridge this gap; instead, new approaches are needed that explicitly target the development of System 2 capabilities. These include integrating symbolic reasoning and knowledge graphs, developing methods for causal inference, exploring new architectures that support meta-cognition, and pursuing research in embodied AI and grounded cognition. Ultimately, fostering public understanding of AI's current limitations and promoting responsible AI development are essential to ensure that this powerful technology is used for the benefit of humanity.

## Acknowledgments:

The author also acknowledges Albert W. Gong for suggesting System 1/2 framework by psychologist Daniel Kahneman in the context of this paper.

# Appendix:

### Data-Copilot: A Generative AI Application

Source code: https://github.com/digital-duck/data-copilot

Demo videos:

- Data Copilot demo 2024-12-10
- [Movie Explorer demo](to be made)

This appendix will provide a detailed description of the data-copilot generative AI system, including its architecture, functionality, and potential applications. It will showcase how the system can be used to gather empirical data on human cognition, such as attention span, task switching, and thought processes, and how this data can inform the development of more advanced AI systems. Further details, including a link to the GitHub code repository and accompanying YouTube video demonstrations, will be provided in the final version of the paper.

# References:

- Ahn, D., Shin, S., & Choi, H. (2021). Improving Commonsense Reasoning in Language Models via Knowledge Graph Transformations. arXiv preprint arXiv:2106.04419.

- Ba, J., & Caruana, R. (2014). Do Deep Nets Really Need to be Deep?. Advances in Neural Information Processing Systems, 27.

- Bisk, Y., Zellers, R., Lin, X., & Choi, Y. (2020). Can a Suit of Armor Conduct an Orchestra? Common Sense Reasoning for Physical Interactions. arXiv preprint arXiv:2005.00693.

- Hendrycks, D., Burns, C., Basart, S., Zemel, R., Misra, S., Parikh, J., ... & Jurafsky, D. (2020). Measuring massive multitask language understanding. arXiv preprint arXiv:2009.03300.

- Sap, M., Rashkin, H., Chen, D., Le Bras, R., & Choi, Y. (2019). Social IQA: Commonsense Reasoning about Social Interactions. arXiv preprint arXiv:1904.09728.

- Zellers, R., Bisk, Y., Schwartz, R., & Choi, Y. (2018). SWAG: A Large-Scale Adversarial Dataset for Commonsense Reasoning. arXiv preprint arXiv:1808.05326.

- Zellers, R., Yatskar, M., Yih, S. W., & Choi, Y. (2018). From Recognition to Cognition: Visual Commonsense Reasoning. Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, 6720-6729.

- Daniel Kahneman. 2011 book: Thinking, Fast and Slow. https://www.wikiwand.com/en/articles/Thinking,_Fast_and_Slow

- Legendary emperor Zhu Yuanzhang 传奇皇帝朱元璋