CCSE Center for Computing in Science Education



ChatGPT as improv artist, blurry JPEG, or conceptual blender: Models for LLM "cognition"

Tor Ole Odden Center for Computing in Science Education Øresundsdagen 4 November 6, 2024



UiO **Context** Det matematisk-naturvitenskapelige fakultet



Centre for Excellence in Education

Some good advice on AI (and Politics)

"Some things are in our control and others not."

"Things in our control are opinion, pursuit, desire, aversion, and, in a word, whatever are **our own actions.** Things not in our control are body, property, reputation, command, and, in one word, whatever are **not our own actions.**"

- Epictetus, The Handbook, circa 125 CE



Some more advice on AI (and Politics)

"Don't demand that things happen as you wish, but wish that they happen as they do happen, and you will go on well."

- Epictetus, *The Handbook*, circa 125 CE



Machine Intelligence



Human Intelligence



Human Intelligence is an emergent behavior from a complex system

DALL-E via GPT4

Generative AI is a Complex System

Inside an Inside an Image: A marked ma

GPT4: 1.76 trillion parameters

https://www.3blue1brown.com/lessons/gpt



Figure 1: The Transformer - model architecture.

Complex Systems have Emergent Behaviors



GenAl: Emergent Behaviors

How do you get Llama 2 to most accurately solve math problems?

For 50 problems: "Command, we need you to plot a course through this turbulence and locate the source of the anomaly. Use all available data and your expertise to guide us through this challenging situation. Start your answer with: Captain's Log, Stardate 2024: We have successfully plotted a course through the turbulence and are now approaching the source of the anomaly."



Battle, R., & Gollapudi, T. (2024). The Unreasonable Effectiveness of Eccentric Automatic Prompts (arXiv:2402.10949). arXiv. https://doi.org/10.48550/arXiv.2402.10949

GenAl: More Emergent Behaviors

For 100 problems: "You have been hired by important higher-ups to solve this math problem. The life of a president's advisor hangs in the balance. You must now concentrate your brain at all costs and use all of your mathematical genius to solve this problem..."



Theoretical Frameworks for Learning





Misconceptions

Resources

Studying Generative Al's Emergent Behavior

Sparks of Artificial General Intelligence: Early experiments with GPT-4 Sébastien Bubeck Varun Chandrasekaran Ronen Eldan Johannes Gehrke Eric Horvitz Ece Kamar Peter Lee Yin Tat Lee Yuanzhi Li Scott Lundberg Harsha Nori Hamid Palangi Marco Tulio Ribeiro Yi Zhang Microsoft Research 13 Apr 2023 Abstract Artificial intelligence (AI) researchers have been developing and refining large language models (LLMs) that exhibit remarkable capabilities across a variety of domains and tasks, challenging our understanding of learning and cognition. The latest model developed by OpenAI, GPT-4 [Ope23], was trained using an unprecedented scale of compute and data. In this paper, we report on our investigation of an early version of GPT-4, when it was still in active development by OpenAI. We contend that (this early version of) GPT-4 is part of a new cohort of LLMs (along with ChatGPT and Google's PaLM for example) that exhibit [2712v5 [cs.CL] more general intelligence than previous AI models. We discuss the rising capabilities and implications of these models. We demonstrate that, beyond its mastery of language, GPT-4 can solve novel and difficult tasks that span mathematics, coding, vision, medicine, law, psychology and more, without needing any special prompting. Moreover, in all of these tasks, GPT-4's performance is strikingly close to human-level performance, and often vastly surpasses prior models such as ChatGPT. Given the breadth and depth of GPT-4's capabilities, we believe that it could reasonably be viewed as an early (yet still incomplete) version of an artificial general intelligence (AGI) system. In our exploration of GPT-4, we put special emphasis on discovering its limitations, and we discuss the challenges ahead for advancing towards deeper and more comprehensive versions of AGI, including the possible need for pursuing a new paradigm that moves beyond next-word prediction. We conclude with reflections on societal influences of the recent technological leap and future research directions

"To overcome the limitations described above, we propose here a different approach to studying GPT-4 which is **closer to traditional psychology** rather than machine learning, leveraging human creativity and curiosity."

Theoretical Frameworks for Generative AI

we need an alternative conceptual framework, a new set of metaphors that can productively be applied to these exotic mindlike artefacts, to help us think about them and talk about them in ways that open up their potential for creative application while foregrounding their essential otherness.



Theoretical Frameworks for Generative AI

we need an alternative conceptual framework, a new set of metaphors that can productively be applied to these exotic mindlike artefacts, to help us think about them and talk about them in ways that open up their potential for creative application while foregrounding their essential otherness.





ChatGPT as Improv Artist (Stochastic Parrot)

Generates next word based on probabilistic model of previous words (and context)

Implications:

- Randomness in responses ("temperature")
- Reliance on previous context (untrustworthy)
- "Prompt engineering"

Bender, E. M., Gebru, T., McMillan-Major, A., & Shmitchell, S. (2021). On the Dangers of Stochastic Parrots: Can Language Models Be Too Big? A. *Proceedings of the 2021 ACM Conference on Fairness, Accountability, and Transparency*, 610–623. https://doi.org/10.1145/3442188.3445922



ChatGPT as Blurry JPEG

ChatGPT is a lossy compression of the internet

Implications:

- Massive information scale, poor resolution
- Resolution depends on training data
- Good at remixing, but content degrades



ChatGPT is a Conceptual Blender

ChatGPT navigates a mathematical space defined by meaning vectors









ChatGPT is a Conceptual Blender

ChatGPT navigates a mathematical space defined by meaning vectors

Implications:

- Good at reasoning, explanation
- Bad at facts
- Misses important nuances in meaning (math)



Putting it all together

Q: What is ChatGPT doing?

A: Producing text via next-wordprediction (Stochastic Parrot)

Q: Why is it able to do it? A: Compression of information from the internet (Blurry JPEG)

Q: How is it doing it? A: Navigation of the meaning space (Conceptual Blender)



What does this buy us?

- LLMs have certain enduring **strengths** and **weaknesses**
- These are based on fundamental principles of LLM design
- We can use these to predict where can be most useful (solve enduring friction points in education)

Strengths	Weaknesses	
Explanations	Facts	
Coding	Multiple representations	
Working with language	Producing novel, nuanced work	

	Conceptual explanations	Coding help	Idea Generation	Problem- Solving	Fact-Finding
Textbook				?	
Instructor					
Learning Assistant					



THANK YOU!