

Enhancing LLM Reasoning Abilities with Code

Fei Ding*
AI Lab

Abstract

Large Language Models (LLMs) have shown exceptional generative abilities in various natural language and generation tasks. Large language models (LLMs) have demonstrated remarkable performance on a variety of natural language tasks based on just a few examples of natural language instructions, reducing the need for extensive feature engineering. However, LLM is relatively weaker in reasoning and problem-solving abilities. We propose a new construction that solves the problem of insufficient logical mathematics and logical ability.

1 Introduction

With the remarkable progress made by large language models such as GPT-4, ChatGPT, Google Gemini, Llama-2 (Touvron et al., 2023), and Mistral (Jiang et al., 2023) in NLP research, machines are now capable of performing a wide range of language tasks that were previously believed to be exclusive to humans (OpenAI, 2023; Brown et al., 2020; Zhao et al., 2023). Performs well on language tasks such as Hellaswag (Zellers et al., 2019), Winogrande (Sakaguchi et al., 2021), PIQA (Bisk et al., 2020) and ARC-Easy , but demonstrates weakness in logical reasoning . However, Logical reasoning is a critical component of intelligence and is essential for many practical applications, including question-answering systems (Khashabi, 2019) and conversational agents (Beygi et al., 2022).

AGI needs to possess many capabilities that would naturally be included in a notion of human intelligence. Examples of these capabilities are generalizability, adaptability, robustness, explainability, causal analysis, abstraction, common sense reasoning, ethical reasoning (Rossi and Mattei, 2019), as well as a complex and seamless integration of

learning and reasoning supported by both implicit and explicit knowledge (Littman et al., 2021). We have explored the mechanisms that enable humans to possess these capabilities, which helps us understand how to imbue AI systems with these competencies. (Rossi and Loreggia, 2019; Booch et al., 2021).

We have delved deeply into D. Kahneman’s theory of thinking fast and slow (Kahneman, 2011), and we propose a new simple AGI architecture (named the Dingfei model, for Artificial General Intelligence) where divide AGI into Intuitive Brain, Logical Brain, and Bottom Brain. The Intuitive Brain reacts through intuition, and the modified LLM can serve as the Intuitive Brain. The Logical Brain is implemented through structured code. It is responsible for logical reasoning. The Bottom Brain is manually designed by humans. AGI can update its own logical and intuitive brains through autonomous learning, while it cannot modify its Bottom Brain.

The interaction and collaboration of the three brains can significantly enhance reasoning capabilities. Furthermore, We introduced the concept of skills and scratch paper, which achieves 100% accuracy in reasoning .

2 Related Work

Several datasets have been proposed such as (Clark et al., 2020; Tian et al., 2021; Joshi et al., 2020; Saeed et al., 2021), LogiQA (Liu et al., 2021) , ReClor (Yu et al., 2020) , FOLIO (Han et al., 2022) and ProntoQA (Saparov and He, 2023) that demonstrate the relatively weak ability of these LLMs to reason logically over natural language text.

Large language models also perform poorly in mathematics and code such as GSM8K (Cobbe et al., 2021) with maj@8 , MATH (Hendrycks et al., 2021) with maj@4 , Humaneval (Chen et al., 2021) and MBPP (Austin et al., 2021) .

*Corresponding author
email:dingfei@email.ncu.edu.cn

3 Task Definition

Our mission is to ensure 100% accuracy. The input text is initially processed by the intuitive brain, which predicts the next token. When a complex logical problem is encountered, a special token is generated, transferring the processing to the logical brain.

We define a given AGI f_{AGI} that comprises three specified models f_{ib}, f_{lb}, f_{bb} and a set of input-output pairs (x, y) . we can define this process as:

$$f_{AGI}(x_i) = \begin{cases} f_{lb}(f_{ib}(x_i)) & \text{if } f_{ib}(x_i) \in N(y_i) \\ f_{ib}(x_i) & \text{if } f_{ib}(x_i) \in other \end{cases} \quad (1)$$

where $f_{ib}(x_i) \in N(y_i)$ represents $f_{ib}(x_i)$ contains special tokens and is passed to the logic brain for processing. f_{ib} represents the Intuitive Brain. f_{lb} represents the Logical Brain. f_{bb} represents the Bottom Brain.

4 Proposed method

In this section, we will present the detailed specifications of the AGI we have implemented:

4.1 The Intuitive Brain

The Intuitive Brain is primarily driven by intuition over careful consideration, providing quick responses to straightforward questions. Intuition is often generated after reading a sufficient amount of data. They are tightly linked to the availability of huge datasets and computational power (Marcus, 2020). However, these answers can occasionally be incorrect due to unconscious biases or their reliance on heuristics and other shortcuts (Gigerenzer and Brighton, 2009), and typically lack explanations.

Furthermore, Intuition often fabricates false facts (i.e. *hallucination*). But logical reasoning requires precise answers, and this is when the rational brain needs to be used.

When it comes to reasoning, we utilize the Intuitive Brain as a tool for organizing information, transforming text into structured data, which is then passed on to the Logical Brain for logical inference.

4.2 The Logical brain

The Logical Brain is implemented using structures beyond large language models, making it proficient in reasoning.

4.2.1 Training data

Textbooks and a large number of related exercises from kindergarten, elementary school, middle school, Senior high school, and University.

4.2.2 Training

How to train the rational brain? For example, mathematics. Let the LLM read textbooks, and then summarize skills for each knowledge point. Then Use the skills learned to answer practice questions. If they are incorrect, correct and improve the skill. A skill is a set of knowledge, rules, and operational flowcharts. Then the skill is converted into executable code and stored in the rational brain, ready to be used when needed. Connect relevant skills with each other.

4.3 Scratch paper

Psychological research reveals a fascinating insight: even children and adults can have their problem-solving prowess significantly dampened by irrelevant information (Hoyer et al., 1979; Pasolunghi et al., 1999; Marzocchi et al., 2002). Similarly, We conducted a large number of experiments that demonstrated the performance of large language models can be affected by irrelevant context, leading to incorrect results. To solve this problem, we proposed the concept of draft paper, ensuring a success rate of 100% when dealing with complex issues. I will show the details later.

4.4 The bottom brain

In the bottom brain, human beings can set the behavioral norms for AI. The bottom brain runs directly on the hardware and provides support for both the rational brain and the emotional brain. I will introduce more details about the rational brain and the bottom brain in the future.

5 Experiments

5.1 Prediction

For example, a math problem.

chatgpt:

chatgpt at Figure 1. Everything was correct until it came to multiplication at the end. Although ChatGPT was trained with a large amount of mathematical computation data to improve the accuracy of mathematical computations, it obviously cannot write out the correct answer to a three-digit multiplication problem based on intuition alone.

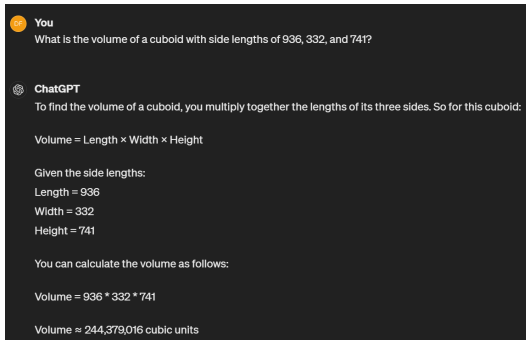


Figure 1: chatgpt

Mistral:

Mistral at the top of Figure 2. Similarly, Mistral has undergone significant optimization in its mathematical capabilities, yet it still exhibits errors when performing multiplication involving numbers exceeding three digits.

The Dingfei model:

The Dingfei model at the bottom of Figure 2. We used cross-computation between the rational and emotional brains to avoid giving an answer directly through intuition alone, thereby ensuring the precision of the results. Here are the details:

At line 1, the rational brain was tasked with thinking, using the skill "formula for the volume of a rectangular prism," to arrive at line 1.

At line 2, the rational brain was tasked with thinking, using the skill "Associative Property of Multiplication", to arrive at line 3.

At line 3, uses the scratch paper to calculate 936×332 and obtains the result in line 4.

At line 4, the rational brain was tasked with thinking, using the skill "Multiplication of large numbers", to arrive at line 5.

At line 5, uses the scratch paper to calculate each multiplication separately and obtains the result in line 6.

At line 6, the rational brain was tasked with thinking, using the skill "Adding Multiple Numbers", to arrive at line 7.

6 Conclusion

We combine multiple techniques to break down complex logical problems into a series of simpler ones. By recursively calling the Logical Brain, we

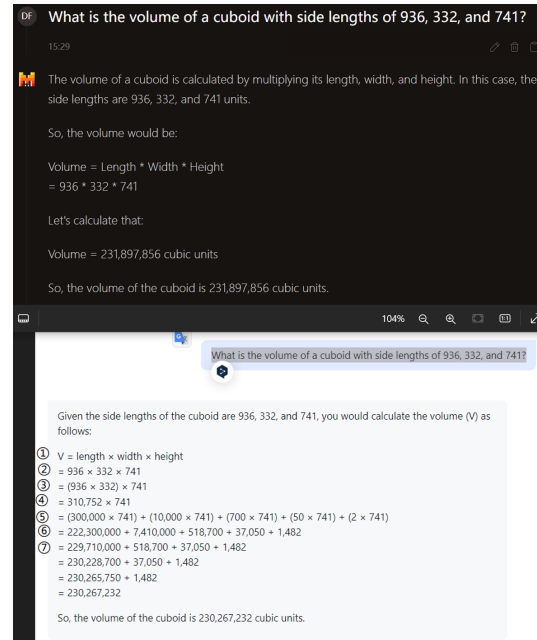


Figure 2: Top: mistral, Bottom: Dingfei model

ensure that each step of the logical reasoning process is error-free from the ground up. The Dingfei model can achieve 100% accuracy in complex problems.

7 Limitations

In this work, due to insufficient computational power, we have only trained part of the mathematics textbooks, and the accuracy rate of doing exercises can reach 100%. In the future, we will continue to train mathematics and other disciplines. We expect the accuracy rate to remain at 100%.

References

Jacob Austin, Augustus Odena, Maxwell Nye, Maarten Bosma, Henryk Michalewski, David Dohan, Ellen Jiang, Carrie Cai, Michael Terry, Quoc Le, et al. 2021. Program synthesis with large language models. *arXiv preprint arXiv:2108.07732*.

Sajjad Beygi, Maryam Fazel-Zarandi, Alessandra Cervone, Prakash Krishnan, and Siddhartha Jonnalagadda. 2022. [Logical reasoning for task oriented dialogue systems](#). In *Proceedings of the Fifth Workshop on e-Commerce and NLP (ECNLP 5)*, pages 68–79, Dublin, Ireland. Association for Computational Linguistics.

Yonatan Bisk, Rowan Zellers, Jianfeng Gao, Yejin Choi, et al. 2020. Piqa: Reasoning about physical commonsense in natural language. In *Proceedings of the AAAI conference on artificial intelligence*.

- Grady Booch, Francesco Fabiano, Lior Horesh, Kiran Kate, Jonathan Lenchner, Nick Linck, Andreas Loreggia, Keerthiram Murgesan, Nicholas Mattei, Francesca Rossi, and Biplav Srivastava. 2021. Thinking fast and slow in AI. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 35, pages 15042–15046.
- Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, Sandhini Agarwal, Ariel Herbert-Voss, Gretchen Krueger, Tom Henighan, Rewon Child, Aditya Ramesh, Daniel Ziegler, Jeffrey Wu, Clemens Winter, Chris Hesse, Mark Chen, Eric Sigler, Mateusz Litwin, Scott Gray, Benjamin Chess, Jack Clark, Christopher Berner, Sam McCandlish, Alec Radford, Ilya Sutskever, and Dario Amodei. 2020. [Language models are few-shot learners](#). In *Advances in Neural Information Processing Systems*, volume 33, pages 1877–1901. Curran Associates, Inc.
- Mark Chen, Jerry Tworek, Heewoo Jun, Qiming Yuan, Henrique Ponde de Oliveira Pinto, Jared Kaplan, Harri Edwards, Yuri Burda, Nicholas Joseph, Greg Brockman, et al. 2021. Evaluating large language models trained on code. *arXiv preprint arXiv:2107.03374*.
- Peter Clark, Oyvind Tafjord, and Kyle Richardson. 2020. Transformers as soft reasoners over language. *arXiv preprint arXiv:2002.05867*.
- Karl Cobbe, Vineet Kosaraju, Mohammad Bavarian, Mark Chen, Heewoo Jun, Lukasz Kaiser, Matthias Plappert, Jerry Tworek, Jacob Hilton, Reiichiro Nakano, Christopher Hesse, and John Schulman. 2021. Training verifiers to solve math word problems. *ArXiv*, abs/2110.14168.
- Gerd Gigerenzer and Henry Brighton. 2009. Homo heuristicus: Why biased minds make better inferences. *Topics in Cognitive Science*, 1(1):107–143.
- Simeng Han, Hailey Schoelkopf, Yilun Zhao, Zhenting Qi, Martin Riddell, Luke Benson, Lucy Sun, Ekaterina Zubova, Yujie Qiao, Matthew Burtell, et al. 2022. Folio: Natural language reasoning with first-order logic. *arXiv preprint arXiv:2209.00840*.
- Dan Hendrycks, Collin Burns, Saurav Kadavath, Akul Arora, Steven Basart, Eric Tang, Dawn Song, and Jacob Steinhardt. 2021. Measuring mathematical problem solving with the math dataset. *arXiv preprint arXiv:2103.03874*.
- William J Hoyer, George W Rebok, and Susan Marx Sved. 1979. Effects of varying irrelevant information on adult age differences in problem solving. *Journal of gerontology*, 34(4):553–560.
- Edward J Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, and Weizhu Chen. 2021. Lora: Low-rank adaptation of large language models. *arXiv preprint arXiv:2106.09685*.
- Gabriel Ilharco, Marco Tulio Ribeiro, Mitchell Wortsman, Ludwig Schmidt, Hannaneh Hajishirzi, and Ali Farhadi. 2022. Editing models with task arithmetic. In *The Eleventh International Conference on Learning Representations*.
- Albert Q Jiang, Alexandre Sablayrolles, Arthur Mensch, Chris Bamford, Devendra Singh Chaplot, Diego de las Casas, Florian Bressand, Gianna Lengyel, Guillaume Lample, Lucile Saulnier, et al. 2023. Mistral 7b. *arXiv preprint arXiv:2310.06825*.
- Pratik Joshi, Somak Aditya, Aalok Sathe, and Monojit Choudhury. 2020. [TaxiNLI: Taking a ride up the NLU hill](#). In *Proceedings of the 24th Conference on Computational Natural Language Learning*, pages 41–55, Online. Association for Computational Linguistics.
- Daniel Kahneman. 2011. *Thinking, Fast and Slow*. Macmillan.
- Daniel Khashabi. 2019. *Reasoning-Driven Question-Answering for Natural Language Understanding*. University of Pennsylvania.
- Omer Levy, Minjoon Seo, Eunsol Choi, and Luke Zettlemoyer. 2017. Zero-shot relation extraction via reading comprehension. *arXiv preprint arXiv:1706.04115*.
- Michael L. Littman, Ifeoma Ajunwa, Guy Berger, Craig Boutilier, Morgan Currie, Finale Doshi-Velez, Gillian Hadfield, Michael C. Horowitz, Charles Isbell, Hiroaki Kitano, Karen Levy, Terah Lyons, Melanie Mitchell, Julie Shah, Steven Sloman, Shannon Vallor, and Toby Walsh. 2021. Gathering Strength, Gathering Storms: The One Hundred Year Study on Artificial Intelligence (AI100) 2021 Study Panel Report. *Stanford University*.
- Jian Liu, Leyang Cui, Hanmeng Liu, Dandan Huang, Yile Wang, and Yue Zhang. 2021. Logiqa: a challenge dataset for machine reading comprehension with logical reasoning. In *Proceedings of the Twenty-Ninth International Conference on International Joint Conferences on Artificial Intelligence*, pages 3622–3628.
- Gary Marcus. 2020. The next decade in AI: Four steps towards robust artificial intelligence. *arXiv preprint arXiv:2002.06177*.
- Gian Marco Marzocchi, Daniela Lucangeli, Tiziana De Meo, Federica Fini, and Cesare Cornoldi. 2002. The disturbing effect of irrelevant information on arithmetic problem solving in inattentive children. *Developmental neuropsychology*, 21(1):73–92.
- Kevin Meng, David Bau, Alex Andonian, and Yonatan Belinkov. 2022a. Locating and editing factual associations in gpt. *Advances in Neural Information Processing Systems*, 35:17359–17372.

- Kevin Meng, Arnab Sen Sharma, Alex J Andonian, Yonatan Belinkov, and David Bau. 2022b. Mass-editing memory in a transformer. In *The Eleventh International Conference on Learning Representations*.
- OpenAI. 2023. [Gpt-4 technical report](#).
- Maria Chiara Pasolunghi, Cesare Cornoldi, and Stephanie De Liberto. 1999. Working memory and intrusions of irrelevant information in a group of specific poor problem solvers. *Memory & Cognition*, 27:779–790.
- F. Rossi and N. Mattei. 2019. Building ethically bounded AI. In *33rd*.
- Francesca Rossi and Andrea Loreggia. 2019. Preferences and ethical priorities: thinking fast and slow in AI. In *Proceedings of the 18th international conference on autonomous agents and multiagent systems*, pages 3–4.
- Mohammed Saeed, Naser Ahmadi, Preslav Nakov, and Paolo Papotti. 2021. [RuleBERT: Teaching soft rules to pre-trained language models](#). In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*, pages 1460–1476, Online and Punta Cana, Dominican Republic. Association for Computational Linguistics.
- Keisuke Sakaguchi, Ronan Le Bras, Chandra Bhagavatula, and Yejin Choi. 2021. Winogrande: An adversarial winograd schema challenge at scale. *Communications of the ACM*, 64(9):99–106.
- Abulhair Saparov and He He. 2023. [Language models are greedy reasoners: A systematic formal analysis of chain-of-thought](#). In *The Eleventh International Conference on Learning Representations*.
- Jidong Tian, Yitian Li, Wenqing Chen, Liqiang Xiao, Hao He, and Yaohui Jin. 2021. [Diagnosing the first-order logical reasoning ability through LogicNLI](#). In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*, pages 3738–3747, Online and Punta Cana, Dominican Republic. Association for Computational Linguistics.
- Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, et al. 2023. Llama: Open and efficient foundation language models. *arXiv preprint arXiv:2302.13971*.
- Yunzhi Yao, Peng Wang, Bozhong Tian, Siyuan Cheng, Zhoubo Li, Shumin Deng, Huajun Chen, and Ningyu Zhang. 2023. Editing large language models: Problems, methods, and opportunities. *arXiv preprint arXiv:2305.13172*.
- Weihaoyu, Zihang Jiang, Yanfei Dong, and Jiashi Feng. 2020. [Reclor: A reading comprehension dataset requiring logical reasoning](#). In *International Conference on Learning Representations*.
- Rowan Zellers, Ari Holtzman, Yonatan Bisk, Ali Farhadi, and Yejin Choi. 2019. Hellaswag: Can a machine really finish your sentence? *arXiv preprint arXiv:1905.07830*.
- Wayne Xin Zhao, Kun Zhou, Junyi Li, Tianyi Tang, Xiaolei Wang, Yupeng Hou, Yingqian Min, Beichen Zhang, Junjie Zhang, Zican Dong, et al. 2023. A survey of large language models. *arXiv preprint arXiv:2303.18223*.
- Chen Zhu, Ankit Singh Rawat, Manzil Zaheer, Srinadh Bhojanapalli, Daliang Li, Felix Yu, and Sanjiv Kumar. 2020. Modifying memories in transformer models. *arXiv preprint arXiv:2012.00363*.