
Final Report for CSE561 Fall 2024

Zheyuan Wu

Department of Mathematics
Washington University in St. Louis
1 Brookings Dr, St. Louis, MO 63130
w.zheyuan@wustl.edu

Abstract

We introduce a novel framework called Flow of Thoughts that allows a large language model (LLM) to retrieve information from self generated knowledge graph to solve problems without relying on additional pre-training. The framework is designed to dynamically iterate over a large data corpus, intelligently filtering and aggregating relevant information to form a comprehensive solution. Website and code are available at <https://github.com/Trance-0/Flow-of-Thoughts>.

1 Introduction

This final project focuses on studying how to help LLMs store and manage "knowledge" and "references". The large language model is a natural compressor [3]. We can efficiently compress the knowledge as a graph node with a set of tag tuples for further referencing. We can scale the data and keep the relevant information to solve the problem when prompting the LLM.

The memory problem of LLM has persisted for a long time. Some recent research finds that modifying where important information is placed within the language model's input context—such as moving the section that answers a question—creates a U-shaped performance pattern. The model performs better when this information is at the very start (primacy effect) or at the end (recency effect) of the input context, but its performance drops notably when the information is situated in the middle of the context. [6]

However, in real life, we need to study and gather tons of information when solving problems, an agent must be aware of many aspects of a question before making correct and consistent decisions. Increasing the memory size or other methods that help LLMs to gain information in the large corpus is essential to make the models solve problems like a human expert.

2 Related work

2.1 Chain-of-thoughts

Chain-of-thoughts [7] is an effective prompting method discovered in 2023 by Jason Wei etc. It provides an example of a logic chain to solve a problem similar to the target question allowing the LLMs to think step by step. This prompt can be used in the Training and Prompting stage of the LLM and generally provides better results when dealing with problems in Mathematics and Engineering.

One significant constraint faced by LLMs based on CoT is the context window size. This can lead to situations where the model forgets previous steps when working through complex problems, particularly in scenarios that require Chain-of-Thought (CoT) reasoning.

31 2.2 Graph-of-thoughts

32 Graph of thoughts [2] is a framework used to prompt LLMs by collecting the thinking process and
33 branching different ideas generated by LLMs. Using generation and aggregation, the model selects
34 the best result in the thinking process and develops on that.

35 integrating concepts such as the Graph of Thoughts and other search methods might provide valuable
36 support for solving problems with a limited context window. By employing these techniques, we can
37 enhance the model's ability to organize and retrieve relevant data, facilitating better problem-solving
38 even with constrained memory resources.

39 However, during the aggregation process, the LLM cannot freely choose the information they need to
40 solve the problem but just propagate from previous thoughts, in this research, our framework will try
41 to give the LLM to choose the material that they find helpful.

42 2.3 Express uncertainty

43 Another relevant approach is detailed in a paper discussing the certainty of LLMs in problem-solving
44 [5]. This research focuses on how models express and handle uncertainty, which can be instrumental
45 in determining when to terminate the search or prompting process. Specifically, the authors explore
46 techniques such as uncertainty sampling and confidence thresholds, which allow models to quantify
47 their level of certainty about generated outputs. Understanding and incorporating measures of
48 certainty can help optimize when and how the model utilizes its context, thereby allowing it to defer
49 to more reliable responses or request additional information when faced with ambiguous queries.
50 Understanding and incorporating measures of certainty can help optimize when and how the model
51 utilizes its context, potentially leading to more accurate and efficient problem-solving.

52 2.4 Self-verification

53 Self-verification [8] is an important technique in improving the reliability of responses generated
54 by large language models (LLMs). When an LLM initially produces an answer, there may be some
55 inaccuracies due to the model's limitations in understanding the full context or providing detailed
56 reasoning. However, the model has the ability to correct or refine its output by "thinking twice." One
57 way this can be achieved is by prompting the model to reverse the roles of questions and answers. By
58 taking the original answer and transforming it back into a question, followed by asking the LLM to
59 generate an updated or revised answer, researchers can often obtain a more accurate and thoughtful
60 response.

61 This process encourages the model to evaluate its earlier reasoning and detect inconsistencies or gaps
62 in logic that may have been overlooked initially. The technique leverages the model's own knowledge
63 to reassess and validate its responses, thereby functioning as a form of internal feedback. Additionally,
64 this self-verification approach may prompt the model to consider alternative interpretations of the
65 question, helping to mitigate issues like oversimplification or misunderstanding of complex queries.
66 By iterating in this manner, the quality of the response can be significantly enhanced, offering a more
67 reliable and nuanced answer for the user.

68 2.5 Memorizing Transformers and Self-Reflective Retrieval-Augmented Generation 69 (SELF-RAG)

70 Furthermore, exploring memorizing transformers and their approaches could provide additional
71 strategies for extending Transformer architectures using kNN [9]. These methods focus on enhancing
72 the model's ability to retain and recall information across longer contexts, potentially offering
73 practical solutions for memory limitations.

74 Other frameworks like Self-RAG [1] are also helpful for LLM to retrieve essential information when
75 solving problems in long paragraphs. The model incorporates a feedback loop where it reflects on
76 its own generated responses to improve their quality before delivering them. This reflection can
77 involve: Checking for consistency with the retrieved documents, verifying the factual accuracy, and
78 identifying potential hallucinations (when the model generates incorrect or fabricated information).

79 3 Framework design

80 We want to create a Flow of Thoughts framework that can dynamically iterate the long passage with
81 self-RAG based on a graph of thoughts. The idea goes as follows:

82 First, the LLM will split the passage into several syntactically independent paragraphs by iterating the
83 long texts. For example, the LLM will split a paper into sections like abstract, importance of projects,
84 related work A, related work B, proposed solutions, etc that can fit into the context windows.

85 Then we let LLM to determine if the passage is relevant to the problem that we are going to solve.
86 For example, when we ask "What related technologies did the author use when doing the project?"
87 The LLM should ignore the proposed framework, conclusion, and experimental results section and
88 only focus on reading the "Related work" sections.

89 Finally, we let LLM compose answers based on the related paper segments with supporting references.

Algorithm 1 Flow of thoughts(P, Q)

Require: Generator LM \mathcal{M}

Require: Large-scale passage collections $P = \{d_1, \dots, d_N\}$

Require: Final question Q .

$methods \leftarrow LM(\text{How to solve the problem } Q, \text{ number of methods})$ \triangleright Generate methods to solve the problem

$solutions \leftarrow []$

for each $method \in methods$ **do**

$segments \leftarrow []$

$mistakes \leftarrow LM(\text{common mistakes in } method)$ \triangleright Extract segments from the knowledge graph to obtain relevant information

while P is not empty **do**

$current_passage \leftarrow P.pop()$

$current_segment \leftarrow \mathcal{M}(\text{relevant info in } current_passage \text{ to solve problem } Q)$

if $current_segment$ is not empty **then**

$segments.add(current_segment)$

end if

end while

$thoughts \leftarrow []$ \triangleright Aggregate the segments to form a rudimental solution

for each $current_segment \in segments$ **do**

$thoughts.add(current_segment)$

end for

while $thoughts.size() > 1$ **do**

$thought_a, thought_b \leftarrow \text{first two solution of } thoughts$

$current_thought \leftarrow \mathcal{M}(\text{aggregate } thought_a, thought_b)$

$refined_thought, mistakes \leftarrow LM(\text{refine } current_thought, mistakes)$

$thoughts.add(refined_thought)$

$mistakes.add(mistakes)$

end while

$solutions.add(thoughts)$

end for

return $majority(solutions)$

90 3.1 Novelty of the solution

91 The Flow of Thoughts framework introduces a new approach to processing and extracting relevant
92 information from a large data corpus that might be space for the necessary information to solve
93 the target problem using Self-RAG and a Graph of Thoughts structure, significantly enhancing the
94 efficiency, relevance, and quality of interactions with long-form content.

95 This model provides more explainability by recording how LLMs get the final solution from aggregat-
96 ing the partial information gained from Self-RAG. Through its intelligent segmentation, contextual
97 relevance filtering, and ability to compose well-supported answers, this framework can potentially
98 stand out in the landscape of text processing and information retrieval technologies.

99 3.2 Efficiency

100 Compared with a normal Graph of Thoughts, the framework proposed independent approaches to
101 solve the problem and solve it automatically with self checking for common mistakes and refine the
102 solution along each step.

103 Moreover, the framework intelligently filters irrelevant sections based on the posed question and
104 methods, enabling the model to focus solely on pertinent information. This specificity enhances the
105 accuracy of the responses generated by the LLM and saves costs when dealing with large data corpus.

106 Compared with a normal Self-RAG Inference, the framework provides more flexibility for the
107 convergence of information for black-boxed models like ChatGPT and Claude. It’s easier to migrate
108 to a new model without training the retriever and fine-tuning costs.

109 4 Experiment results

110 Now we have implemented basic functions for the Flow of Thoughts in revised Graph of Thoughts
111 framework and let the LLM decide whether to use those messages or not when aggregating the final
112 answer.

113 4.1 Sorting

114 The sorting task is a basic task that can be solved by Graph of Thoughts. We have implemented the
115 basic functions for the Flow of Thoughts in the Graph of Thoughts framework and introduced basic
116 RAG layers after the generation layer and let the LLM decide whether to use those messages or not
117 when aggregating the final answer.

118 4.2 Set Intersection

119 The set intersection task is a basic task that also can be solved by Graph of Thoughts. We have
120 implemented the basic functions for the Flow of Thoughts with the same functionality as Graph of
121 Thoughts and compare the performance with Graph of Thoughts.

122 We use the same dataset as the Graph of Thoughts and compare the performance with Graph of
123 Thoughts. However, due to time constraints, we only have the set intersection task results for length
124 64 and 100 trials using ChatGPT-4o and Claude-3.5-sonnet models.

125 4.3 Reading Comprehension

126 For the reading comprehension task, we use the RACE dataset [4] to test the performance of the Flow
127 of Thoughts framework. The dataset contains 28,000+ passages and 100,000+ questions, which is a
128 large dataset for reading comprehension tasks. We extract 100 trials on the subset of passages where

129 5 Analysis

130 6 Limitations and potential future work

131 References

- 132 [1] Akari Asai, Zeqiu Wu, Yizhong Wang, Avirup Sil, and Hannaneh Hajishirzi. Self-rag: Learning to retrieve,
133 generate, and critique through self-reflection, 2023.
- 134 [2] Maciej Besta, Nils Blach, Ales Kubicek, Robert Gerstenberger, Michal Podstawski, Lukas Gianinazzi,
135 Joanna Gajda, Tomasz Lehmann, Hubert Niewiadomski, Piotr Nyczyk, and Torsten Hoefler. Graph of
136 thoughts: Solving elaborate problems with large language models. *Proceedings of the AAAI Conference on*
137 *Artificial Intelligence*, 38(16):17682–17690, March 2024.
- 138 [3] Grégoire Delétang, Anian Ruoss, Paul-Ambroise Duquenne, Elliot Catt, Tim Genewein, Christopher Mattern,
139 Jordi Grau-Moya, Li Kevin Wenliang, Matthew Aitchison, Laurent Orseau, Marcus Hutter, and Joel Veness.
140 Language modeling is compression, 2024.

- 141 [4] Guokun Lai, Qizhe Xie, Hanxiao Liu, Yiming Yang, and Eduard Hovy. Race: Large-scale reading
142 comprehension dataset from examinations. *arXiv preprint arXiv:1704.04683*, 2017.
- 143 [5] Stephanie Lin, Jacob Hilton, and Owain Evans. Teaching models to express their uncertainty in words,
144 2022.
- 145 [6] Nelson F. Liu, Kevin Lin, John Hewitt, Ashwin Paranjape, Michele Bevilacqua, Fabio Petroni, and Percy
146 Liang. Lost in the middle: How language models use long contexts, 2023.
- 147 [7] Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Ed Huai hsin Chi, F. Xia, Quoc Le, and Denny
148 Zhou. Chain of thought prompting elicits reasoning in large language models. *ArXiv*, abs/2201.11903, 2022.
- 149 [8] Yixuan Weng, Minjun Zhu, Fei Xia, Bin Li, Shizhu He, Shengping Liu, Bin Sun, Kang Liu, and Jun Zhao.
150 Large language models are better reasoners with self-verification, 2023.
- 151 [9] Yuhuai Wu, Markus N. Rabe, DeLesley Hutchins, and Christian Szegedy. Memorizing transformers, 2022.

152 **7 Appendix**

153 **7.1 Prompts used in solving the problem**

154 **7.1.1 A**