

---

# 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 Flow of Thoughts, a novel framework that enables large language models (LLMs) to effectively process and reason about long passages by combining self-RAG based generation with knowledge graph construction. The framework dynamically builds and traverses a knowledge graph while solving problems, allowing flexible information retrieval and aggregation. Through independent solution paths with built-in self-verification, Flow of Thoughts achieves competitive performance compared to existing approaches like Graph of Thoughts, particularly excelling at tasks requiring synthesis of information from large texts. Our experiments on sorting, set intersection, and reading comprehension demonstrate the framework's effectiveness across different problem domains. Website and code are available at <https://github.com/Trance-0/Flow-of-Thoughts>.

## 1 Introduction

This final project investigates novel approaches to enhance how Large Language Models (LLMs) store, manage, and retrieve knowledge and references. Recent research has shown that LLMs inherently act as natural compressors of information [3]. Building on this insight, we propose efficiently representing knowledge as graph nodes annotated with semantic tag tuples, enabling flexible and contextual retrieval. This graph-based approach allows us to scale to large amounts of data while maintaining the ability to surface relevant information during inference.

A fundamental challenge with LLMs is their limited ability to effectively utilize long-range context. Studies have revealed that LLMs exhibit a distinctive U-shaped performance pattern based on information positioning - they perform optimally when critical information appears either at the beginning (primacy effect) or end (recency effect) of the input context, but struggle significantly when important details are embedded in the middle [7]. This "lost in the middle" phenomenon poses a significant barrier to processing longer documents.

Also, LLMs are prone to generating hallucinations when solving problems. They may generate incorrect or fabricated statements that snowball into more errors [12], making the model overconfident in its own answers, especially when the input is long. The credibility of the model's responses will be lost along long induction chains and retrieval processes.

This limitation becomes particularly problematic when considering real-world applications, where problem-solving often requires synthesizing information from multiple sources and considering various aspects simultaneously. Just as human experts draw upon extensive knowledge bases to make informed decisions, LLMs need mechanisms to efficiently access and reason over large bodies of information. While simply increasing context window sizes offers one potential solution, we argue that more sophisticated approaches to knowledge management and retrieval are essential for enabling LLMs to achieve human-like reasoning capabilities across complex, information-rich domains.

## 36 2 Related work

### 37 2.1 Chain-of-thoughts

38 Chain-of-thoughts [9] is an effective prompting method discovered in 2023 by Jason Wei et al. It  
39 provides an example of a logic chain to solve a problem similar to the target question, allowing the  
40 LLMs to think step by step. This prompt can be used in the Training and Prompting stage of the LLM  
41 and generally provides better results when dealing with problems in Mathematics and Engineering.

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

### 45 2.2 Graph-of-thoughts

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

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

53 However, during the aggregation process, the LLM cannot freely choose the information they need to  
54 solve the problem but must propagate from previous thoughts. In this research, our framework will  
55 try to give the LLM the ability to choose the material that they find helpful.

### 56 2.3 Express uncertainty

57 Another relevant approach is detailed in a paper discussing the certainty of LLMs in problem-solving  
58 [6]. This research focuses on how models express and handle uncertainty, which can be instrumental  
59 in determining when to terminate the search or prompting process. Specifically, the authors explore  
60 techniques such as uncertainty sampling and confidence thresholds, which allow models to quantify  
61 their level of certainty about generated outputs. Understanding and incorporating measures of  
62 certainty can help optimize when and how the model utilizes its context, thereby allowing it to defer  
63 to more reliable responses or request additional information when faced with ambiguous queries.  
64 Understanding and incorporating measures of certainty can help optimize when and how the model  
65 utilizes its context, potentially leading to more accurate and efficient problem-solving.

### 66 2.4 Self-verification

67 Self-verification [10] is an important technique in improving the reliability of responses generated  
68 by large language models (LLMs). When an LLM initially produces an answer, there may be some  
69 inaccuracies due to the model’s limitations in understanding the full context or providing detailed  
70 reasoning. However, the model has the ability to correct or refine its output by "thinking twice." One  
71 way this can be achieved is by prompting the model to reverse the roles of questions and answers. By  
72 taking the original answer and transforming it back into a question, followed by asking the LLM to  
73 generate an updated or revised answer, researchers can often obtain a more accurate and thoughtful  
74 response.

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

## 82 2.5 Memorizing Transformers and Self-Reflective Retrieval-Augmented Generation 83 (SELF-RAG)

84 Furthermore, exploring memorizing transformers and their approaches could provide additional  
85 strategies for extending Transformer architectures using kNN [11]. These methods focus on enhancing  
86 the model’s ability to retain and recall information across longer contexts, potentially offering practical  
87 solutions for memory limitations.

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

## 93 3 Framework design

94 We propose a novel Flow of Thoughts framework that can dynamically process and reason about  
95 long passages using self-RAG based generation combined with knowledge graph construction.  
96 This framework enables large language models to effectively manage and utilize information from  
97 extensive texts while maintaining coherent logic chains. The framework consists of four main  
98 components:

### 99 3.1 Preprocessing

100 First, the LLM constructs a knowledge graph by splitting input long passages into syntactically and  
101 semantically independent units that fit within the model’s context window.

102 For academic papers, this involves breaking down the text into logical sections such as abstract,  
103 introduction, related work, methodology, and results. Each section is further subdivided if needed to  
104 ensure it fits within context limits while preserving meaningful relationships. For reading comprehen-  
105 sion tasks, the framework segments passages into coherent paragraphs that maintain narrative flow  
106 and topical consistency while staying within context window constraints.

107 The model can also use predefined knowledge graphs as well.

### 108 3.2 Generating methods

109 The framework then employs a multi-strategy approach by having the LLM generate diverse solution  
110 methods tailored to the specific problem type. For example, when tackling sorting problems, the  
111 LLM generates multiple sorting algorithms like bubble sort, quick sort, and merge sort.

112 For each method, the LLM also performs an error analysis by identifying potential pitfalls and edge  
113 cases that could arise during execution for future refinement.

### 114 3.3 Self-verification and refining

115 On each branch of the solution, the framework feeds the problem statement, generated methods, and  
116 data corpus to the LLM. The model evaluates passage relevance and uses the generated methods to  
117 solve the problem. The LLM will also check if the generated thought follows the generated methods  
118 and the problem statement using self-verification [10]. To achieve this, we provide the LLM with the  
119 method description and generated output, we ask the LLM to self-verify the output and the method  
120 and capture the mistakes as we go.

121 Then, we ask the LLM to iteratively check the list of potential mistakes that was generated in the  
122 previous steps.

123 If new errors are found, we may also add the new errors to the checking process.

124 For instance, when answering questions about technical approaches used in a project, the LLM  
125 selectively focuses on relevant sections like "Related Work" while appropriately discounting less  
126 relevant sections like conclusions or experimental results. This targeted attention mechanism helps  
127 maintain solution quality while managing computational resources effectively.

128 Ideally, in this step, the LLM should take control of the thought generation and refinement process  
 129 and should be able to determine the steps required to solve the problem. However, in this research,  
 130 we don't have sufficient time to implement such a protocol, which can be done in the future. 6.5

### 131 3.4 Aggregation

132 The final phase aggregates the final thoughts from each independent branch to reach the final answer.  
 133 Unlike Graph of Thoughts approaches that merge partial solutions, our framework uses a majority  
 134 voting system to determine the final answer.

---

#### 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  $Q$ .  
**Require:** Predefined root thought node that contains the problem statement  $problem\_root$   
 $methods \leftarrow LM(\text{How to solve the problem } Q, \text{ number of methods})$   $\triangleright$  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**  
    $thoughts \leftarrow [problem\_root]$   $\triangleright$  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**  
       $thought \leftarrow thoughts.pop()$   
      **for** each  $operation \in thought.child\_operations$  **do**  
         $operation.execute()$   $\triangleright$  Execute the operation with various types, such as retrieval, generation, etc. We can capture the mistakes and refine the solution on execution steps.  
        **for** each  $child\_thoughts \in operation.child\_thoughts$  **do**  
          **if**  $child\_thoughts.is\_executable()$  **then**  
             $thoughts.add(child\_thoughts)$   
          **end if**  
        **end for**  
      **end for**  
   **end while**  
    $solutions.add(thoughts)$   
**end for**  
**return**  $majority(solutions)$

---

135 The framework is designed to be a general framework that can be applied to any problem that can be  
 136 solved by LLMs.

### 137 3.5 Novelty of the solution

138 The Flow of Thoughts framework introduces a new approach to processing and extracting relevant  
 139 information from a large data corpus that might contain the necessary information to solve the target  
 140 problem using Self-RAG and a Graph of Thoughts structure, significantly enhancing the efficiency,  
 141 relevance, and quality of interactions with long-form content.



Figure 1: Flow of Thoughts framework using graph visualization, the pink nodes are the nodes from the knowledge graph or definition of the problem, the yellow nodes are generated by the LLM, the purple nodes are refinement of the solution generated by the LLM. In the final layer, the LLM will aggregate the solution and generate the final answer (the cyan node)

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

### 146 3.6 Efficiency

147 Compared with a normal Graph of Thoughts, the framework proposes independent approaches to  
 148 solve the problem and solves it automatically with self-checking for common mistakes and refines the  
 149 solution along each step. This uses generation of LLM more effectively since each path is independent  
 150 and the LLM can generate more diverse solutions compared to the normal Graph of Thoughts or Tree  
 151 of Thoughts.

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

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

### 158 3.7 Effectiveness

159 The framework is designed to be a general framework that can be applied to most of the problems  
 160 that can be solved by LLMs. All the prompts and methods are generated by the LLM itself; we only  
 161 need to provide few-shot examples to guide the LLM to generate the prompts and methods in our  
 162 setting. In later experiments, we observed that our scheme is competitive with Graph of Thoughts  
 163 and even outperforms it in some cases.

## 164 4 Experiment results

165 To test the performance of the Flow of Thoughts framework, we have implemented basic functions  
 166 for the Flow of Thoughts in revised Graph of Thoughts framework with implementation of basic  
 167 input and output (IO), Chain-of-Thought (CoT), Tree-of-Thought (ToT), Graph-of-Thought (GoT),  
 168 and Flow-of-Thought (FoT). All the prompts are given with 3-shot examples. For the flow of thoughts  
 169 framework, the shots are given three simple tasks with potential approaches.

170 We use black-box models GPT-4o-mini and Claude-3.5-Sonnet with API calls for the experiment. For  
 171 both models, we set temperature to 0.7 with all other parameters to default.

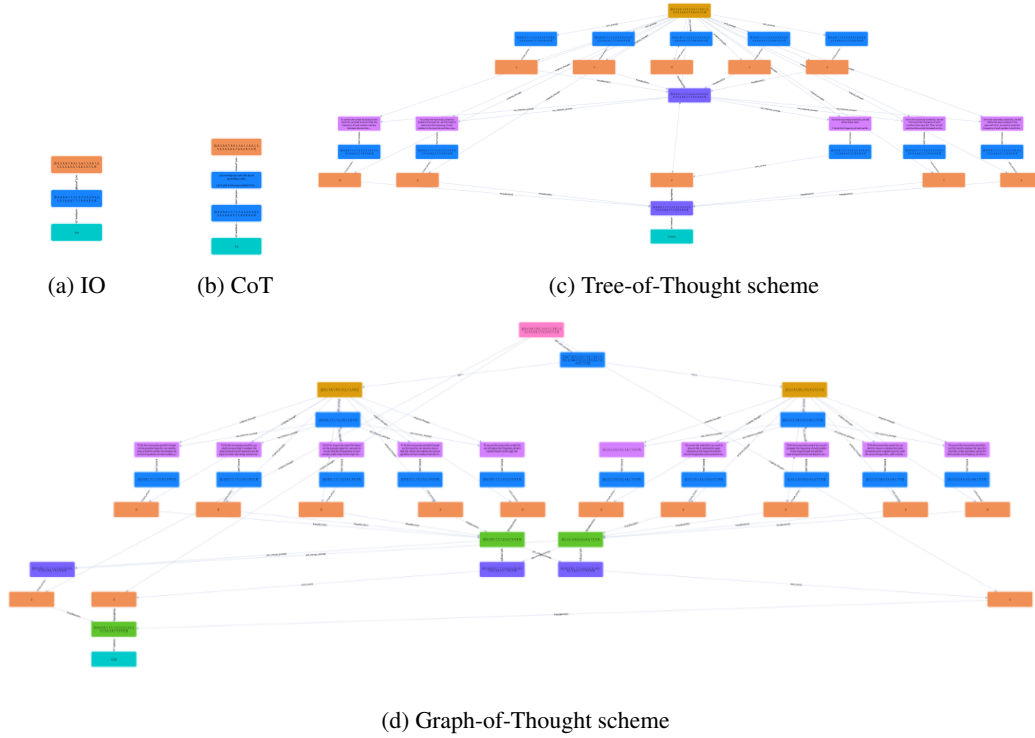


Figure 2: Other framework structures defined in the experiment, only used for reference, detailed structure can be found in each run of the experiment

#### 4.1 Sorting

The sorting task is a basic task that can be solved by Graph of Thoughts. We give the LLM a list of 32 random numbers and let each scheme sort the numbers. The number of trials is 100. The number of errors is calculated by the number of trials in which the LLM cannot sort the numbers correctly. Dataset can be found in appendix 7.1.1.

All the prompts are given with 3-shot examples. For the flow of thoughts framework, the shots are given three simple tasks with potential approaches.

Model	GPT-4o-mini	Claude-3.5-Sonet
<b>IO</b>	2.14	2.03
<b>CoT</b>	1.88	1.86
<b>ToT</b>	1.20	1.32
<b>GoT</b>	1.10	1.32
<b>FoT</b>	1.13	<b>1.09</b>

Table 1: Average number of incorrect responses in sorting 32 integers

#### 4.2 Set Intersection

The set intersection task is a basic task that also can be solved by Graph of Thoughts. The task is defined as follows: Given two sets of 32 numbers, find the intersection of the two sets.

We use the same dataset as the Graph of Thoughts and compare the performance with Graph of Thoughts. Dataset can be found in appendix 7.1.2.

The experiment is conducted with 100 trials. The number of errors is calculated by the number of misclassified elements in the intersection.

Model	GPT-4o-mini	Claude-3.5-Sonet
<b>IO</b>	3.02	2.88
<b>CoT</b>	2.46	1.89
<b>ToT</b>	1.93	1.97
<b>GoT</b>	1.24	1.56
<b>FoT</b>	<b>1.13</b>	1.20

Table 2: Average number of incorrect responses in set intersection

### 186 4.3 Reading Comprehension

187 For the reading comprehension task, we use the RACE dataset [5] to test the performance of the Flow  
188 of Thoughts framework. The dataset contains 28,000+ passages and 100,000+ questions, which is a  
189 large dataset for reading comprehension tasks. We use the subset of the top 100 longest passages to  
190 test the performance of the Flow of Thoughts framework. (min passage length 3850 words, details of  
191 the passages can be found in the appendix)

192 The experiment is conducted with 100 passages. Each passage weights 1 point and we calculate the  
193 total score of the model by summing up the scores of each passage. The best score obtained is 85.5,  
194 whereas the SOTA score is 91.4 by ALBERT (Ensemble)[4].

Model	GPT-4o-mini	Claude-3.5-Sonet
<b>IO</b>	68.5	65.25
<b>CoT</b>	72.33	69.5
<b>ToT</b>	76.25	75.33
<b>GoT</b>	78.5	83.25
<b>FoT</b>	84.33	<b>85.5</b>

Table 3: Score of the reading comprehension task on RACE-min dataset

## 195 5 Analysis

### 196 5.1 Method generation

197 By default, we let the LLM generate 3 methods for each task all by themselves. However, the model  
198 may not be able to generate 3 practical methods that LLM can execute.

199 For sorting task, in some rare cases, the LLM will generate methods using python code. However,  
200 in our setting, the LLM does not have the capability to execute python code and the methods  
201 will be processed as a normal tree of thoughts branch and don’t contribute to the diversity of the  
202 solutions. This case is also observed in the set intersection task, where the LLM will also use  
203 `set.intersection()` to solve the problem.

204 For reading comprehension task, instead of generating 3 methods, the LLM will sometimes generate  
205 1 or 2 methods only or even generate methods that are not related to the task. For instance, the LLM  
206 will produce Chunking and Visualizing as methods for the reading comprehension task, where you  
207 break down large pieces of information into smaller, more manageable sections (chunks). However,  
208 the LLM does not have the capability to generate images or diagrams for the passages that can be fed  
209 to itself.

### 210 5.2 Self-checking process

211 The self-checking process is a crucial part of the Flow of Thoughts framework. It allows the LLM to  
212 check the correctness of the methods and refine the solution. However, the self-checking process is  
213 not perfect and the LLM may not be able to check the correctness of the methods.

214 For sorting task, we observe that the LLM cannot count the number of elements in the list correctly  
215 when the number of elements is large. The LLM will fail to count the number of missing elements

216 in the list and snowball the error. For example, the LLM will generate 31 as the number of missing  
217 elements in the list when the correct number of elements is 32. The LLM sometimes fails to check  
218 the length of the output list and uses illusions to fill the gaps in the list.

## 219 **6 Limitations and potential future work**

### 220 **6.1 Approach generation**

221 The generated methods may still not be approachable for LLM to execute (Shor’s algorithm to  
222 factorize large primes)

223 Among all the generated methods, most of them are approachable for LLM to execute and will not  
224 affect the majority voting procedure of the final answer. However, this statement may not hold for  
225 more complex problems. The model may need access to external knowledge or API to execute some  
226 of the methods. One further improvement is to let the LLM access external knowledge or API to  
227 execute some of the methods as described in ToolLLM [8].

### 228 **6.2 High expense for thoughts generation**

229 During the experiment, we found that the generation of the graph is expensive. The generation of the  
230 graph is  $O(nm)$  where  $n$  is the number of approaches and  $m$  is the number of checking steps for each  
231 method. (Maybe we can share and check some common mistakes to save costs for improvements)  
232 We tried to reuse some of the common mistakes to save costs for improvements. However, the  
233 performance is still non-linear especially when the number of mistakes is large.

### 234 **6.3 Flexibility for Divide and Conquer approach**

235 The final Flow of Thoughts framework should allow the LLM to generate additional methods to  
236 divide the problem into sub-problems. However, the current framework is not flexible for this feature  
237 due to time constraints for this project. During the sorting process in Flow of Thoughts scheme, the  
238 LLM cannot divide the problem into sub-problems as normal Graph of Thoughts do, which might  
239 potentially hinder the performance of the framework for longer array size.

240 In future work, it’s possible to design a protocol to let the LLM decide whether to divide the problem  
241 into sub-problems or generate additional filtering functions to filter out the invalid or incorrect  
242 solutions.

243 Our current Flow of Thoughts framework also doesn’t fully support the self-retrieval process. The  
244 LLM will not be able to retrieve the information from the previous steps if that information is not  
245 stored in previous steps. Additional protocol is needed to enable the self-retrieval process as LLM  
246 needs.

### 247 **6.4 Ablation study**

248 Another future work is to do an ablation study to see the contribution of each component to the final  
249 performance. During the experiment, we found that each component has different contributions to  
250 the final performance in different tasks. For example, the method generation is more important for  
251 the sorting task than the set intersection task since many proposed methods for set intersection are  
252 not practical for LLM to execute. More experiments could be done to see the contribution of each  
253 component to the final performance.

254 However, due to time restrictions of this project, we cannot do more tasks to see the contribution of  
255 each component to the final performance, which can be done in the future with more task designs.

### 256 **6.5 Potential improvements that can be made**

257 Due to time constraints, we only implemented the basic version of the Flow of Thoughts framework.  
258 There are many potential improvements that might be done to improve the performance of the  
259 framework.



260 Instead of finite steps generations, we can let the LLM continue to generate thoughts until we find a  
 261 solution that LLM agrees with with high confidence by adding certainty attribute [6] to each thought  
 262 and using priority queue to check the most promising thoughts first.

## 263 References

- 264 [1] Akari Asai, Zeqiu Wu, Yizhong Wang, Avirup Sil, and Hannaneh Hajishirzi. Self-rag: Learning to retrieve,  
 265 generate, and critique through self-reflection, 2023.
- 266 [2] Maciej Besta, Nils Blach, Ales Kubicek, Robert Gerstenberger, Michal Podstawski, Lukas Gianinazzi,  
 267 Joanna Gajda, Tomasz Lehmann, Hubert Niewiadomski, Piotr Nyczyk, and Torsten Hoeffer. Graph of  
 268 thoughts: Solving elaborate problems with large language models. *Proceedings of the AAAI Conference on*  
 269 *Artificial Intelligence*, 38(16):17682–17690, March 2024.
- 270 [3] Grégoire Delétang, Anian Ruoss, Paul-Ambroise Duquenne, Elliot Catt, Tim Genewein, Christopher  
 271 Mattern, Jordi Grau-Moya, Li Kevin Wenliang, Matthew Aitchison, Laurent Orseau, Marcus Hutter, and  
 272 Joel Veness. Language modeling is compression, 2024.
- 273 [4] Yufan Jiang, Shuangzhi Wu, Jing Gong, Yahui Cheng, Peng Meng, Weiliang Lin, Zhibo Chen, and Mu li.  
 274 Improving machine reading comprehension with single-choice decision and transfer learning, 2020.
- 275 [5] Guokun Lai, Qizhe Xie, Hanxiao Liu, Yiming Yang, and Eduard Hovy. Race: Large-scale reading  
 276 comprehension dataset from examinations. *arXiv preprint arXiv:1704.04683*, 2017.
- 277 [6] Stephanie Lin, Jacob Hilton, and Owain Evans. Teaching models to express their uncertainty in words,  
 278 2022.
- 279 [7] Nelson F. Liu, Kevin Lin, John Hewitt, Ashwin Paranjape, Michele Bevilacqua, Fabio Petroni, and Percy  
 280 Liang. Lost in the middle: How language models use long contexts, 2023.
- 281 [8] Yujia Qin, Shihao Liang, Yining Ye, Kunlun Zhu, Lan Yan, Yaxi Lu, Yankai Lin, Xin Cong, Xiangru  
 282 Tang, Bill Qian, Sihan Zhao, Lauren Hong, Runchu Tian, Ruobing Xie, Jie Zhou, Mark Gerstein, Dahai Li,  
 283 Zhiyuan Liu, and Maosong Sun. Toolllm: Facilitating large language models to master 16000+ real-world  
 284 apis, 2023.
- 285 [9] Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Ed Huai hsin Chi, F. Xia, Quoc Le, and Denny  
 286 Zhou. Chain of thought prompting elicits reasoning in large language models. *ArXiv*, abs/2201.11903,  
 287 2022.
- 288 [10] Yixuan Weng, Minjun Zhu, Fei Xia, Bin Li, Shizhu He, Shengping Liu, Bin Sun, Kang Liu, and Jun Zhao.  
 289 Large language models are better reasoners with self-verification, 2023.
- 290 [11] Yuhuai Wu, Markus N. Rabe, DeLesley Hutchins, and Christian Szegedy. Memorizing transformers, 2022.
- 291 [12] Muru Zhang, Ofir Press, William Merrill, Alisa Liu, and Noah A. Smith. How language model hallucina-  
 292 tions can snowball, 2023.

## 293 7 Appendix

### 294 7.1 Prompts used in solving the problem

#### 295 7.1.1 Sorting

296 Dataset for sorting can be found in [https://github.com/Trance-0/Flow-of-Thoughts/](https://github.com/Trance-0/Flow-of-Thoughts/blob/main/flow-of-thoughts/test/sorting_test/sorting_032.csv)  
 297 [blob/main/flow-of-thoughts/test/sorting\\_test/sorting\\_032.csv](https://github.com/Trance-0/Flow-of-Thoughts/blob/main/flow-of-thoughts/test/sorting_test/sorting_032.csv)

298 Prompt for sorting can be found in [https://github.com/Trance-0/Flow-of-Thoughts/blob/](https://github.com/Trance-0/Flow-of-Thoughts/blob/main/flow-of-thoughts/test/sorting_test/sort_prompt.json)  
 299 [main/flow-of-thoughts/test/sorting\\_test/sort\\_prompt.json](https://github.com/Trance-0/Flow-of-Thoughts/blob/main/flow-of-thoughts/test/sorting_test/sort_prompt.json)

### 300 **7.1.2 Set Intersection**

301 Dataset for set intersection can be found in [https://github.com/Trance-0/](https://github.com/Trance-0/Flow-of-Thoughts/blob/main/flow-of-thoughts/test/set_intersections/set_intersection_032.csv)  
302 [Flow-of-Thoughts/blob/main/flow-of-thoughts/test/set\\_intersections/set\\_](https://github.com/Trance-0/Flow-of-Thoughts/blob/main/flow-of-thoughts/test/set_intersections/set_intersection_032.csv)  
303 [intersection\\_032.csv](https://github.com/Trance-0/Flow-of-Thoughts/blob/main/flow-of-thoughts/test/set_intersections/set_intersection_032.csv)

304 Prompt for set intersection can be found in [https://github.com/Trance-0/](https://github.com/Trance-0/Flow-of-Thoughts/blob/main/flow-of-thoughts/test/set_intersections/set_intersection_prompt.json)  
305 [Flow-of-Thoughts/blob/main/flow-of-thoughts/test/set\\_intersections/set\\_](https://github.com/Trance-0/Flow-of-Thoughts/blob/main/flow-of-thoughts/test/set_intersections/set_intersection_prompt.json)  
306 [intersection\\_prompt.json](https://github.com/Trance-0/Flow-of-Thoughts/blob/main/flow-of-thoughts/test/set_intersections/set_intersection_prompt.json)

### 307 **7.1.3 Reading Comprehension**

308 Passages used for reading comprehension test can be found in github repository  
309 [https://github.com/Trance-0/Flow-of-Thoughts/tree/main/flow-of-thoughts/](https://github.com/Trance-0/Flow-of-Thoughts/tree/main/flow-of-thoughts/test/reading_comprehension/RACE_min)  
310 [test/reading\\_comprehension/RACE\\_min](https://github.com/Trance-0/Flow-of-Thoughts/tree/main/flow-of-thoughts/test/reading_comprehension/RACE_min)

311 Prompt for approach generation can be found in [https://github.com/Trance-0/](https://github.com/Trance-0/Flow-of-Thoughts/blob/main/flow-of-thoughts/test/reading_comprehension/reading_comprehension_prompt.json)  
312 [Flow-of-Thoughts/blob/main/flow-of-thoughts/test/reading\\_comprehension/](https://github.com/Trance-0/Flow-of-Thoughts/blob/main/flow-of-thoughts/test/reading_comprehension/reading_comprehension_prompt.json)  
313 [reading\\_comprehension\\_prompt.json](https://github.com/Trance-0/Flow-of-Thoughts/blob/main/flow-of-thoughts/test/reading_comprehension/reading_comprehension_prompt.json)