

STOC-ToT: Stochastic Tree-of-Thought with Constrained Decoding for Complex Reasoning in Multi-Hop Question Answering

Anonymous EMNLP submission

Abstract

Multi-hop question answering (MHQA) requires a model to retrieve and integrate information from multiple passages to answer a complex question. Recent systems leverage the power of large language models and integrate evidence retrieval with reasoning prompts (e.g., chain-of-thought reasoning) for the MHQA task. However, the complexities in the question types (bridge v.s. comparison questions) and the reasoning types (sequential v.s. parallel reasonings) require more novel and fine-grained prompting methods to enhance the performance of MHQA under the zero-shot setting. In this paper, we propose STOC-ToT, a stochastic tree-of-thought reasoning prompting method with constrained decoding for MHQA and conduct a detailed comparison with other reasoning prompts on different question types and reasoning types. Specifically, we construct a tree-like reasoning structure by prompting the model to break down the original question into smaller sub-questions to form different reasoning paths. In addition, we prompt the model to provide a probability estimation for each reasoning path at each reasoning step. At answer time, we conduct constrained decoding on the model to generate more grounded answers and reduce hallucination. Experiments comparing STOC-ToT with on two MHQA datasets and five large language models showed that STOC-ToT outperforms other reasoning prompts by a significant margin.

1 Introduction

Question answering (QA) is a fundamental task in natural language processing (NLP) that involves designing systems capable of understanding human language questions and providing accurate and relevant answers. With the recent advancement of large language models (LLMs) that demonstrated superior reasoning ability (Brown et al., 2020), researchers have been focusing more on complex QA tasks, such as multi-hop question answering

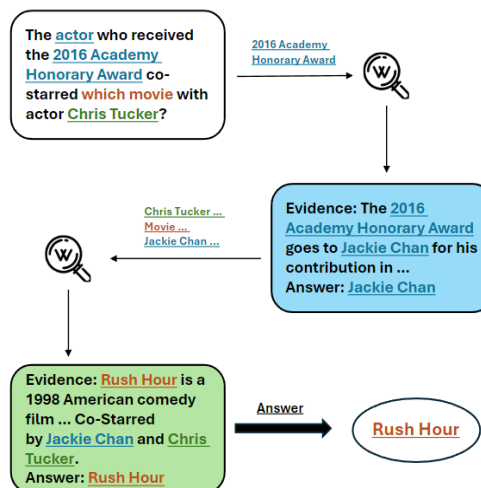


Figure 1: An example of the MHQA question. This question has two hops that require the model to reason about before answering the final question.

(MHQA). MHQA is more challenging as it requires models to understand complicated questions, perform multiple reasoning steps, and gather evidence across documents. Figure 1 shows an example of a two-hop MHQA question. To answer that question in Figure 1, the QA model needs to first figure out who is the actor that received the 2016 Academy Honorary Award. Then based on the answer to the previous question, the QA model needs to further answer a second question about which movie the actor co-starred with Chris Tucker.

State-of-the-art methods for MHQA are fully-supervised methods that often follow a retrieve-and-read framework, including a passage retrieving module that gathers relative evidence from documents and a reading comprehension module to reason about the evidence (Zhu et al., 2021; Li et al., 2022). Other methods include beam-search (Zhang et al., 2023) and label-smoothing (Yin et al., 2023). However, these methods often require extensive pre-training or fine-tuning and do not generalize well to other datasets.

Large language models (LLMs), on the other hand, show remarkable reasoning ability and rich knowledge of general-domain questions. Many LLMs can answer simple and straightforward questions that do not require complex reasoning without any supervision involved but often fail to deal with complex questions requiring multiple reasoning steps. To tackle the problem, researchers have developed many prompting techniques to improve LLM’s reasoning ability, such as chain-of-thought (CoT) (Wei et al., 2022), self-consistency CoT (Sc-CoT) (Wang et al., 2023), and tree-of-thought (ToT) prompting (Yao et al., 2023a).

CoT has been shown effective across tasks requiring extensive, step-by-step reasoning, such as math calculation and reading comprehension. However, there could be various possible reasoning paths for many complex multi-hop questions, and CoT models cannot "turn back" when they have made a mistake along their reasoning paths. Sc-CoT further improves on CoT by proposing different chains of thought, thus expanding the reasoning space. However, there is no local reasoning expansion within each chain, and the "majority voting" strategy often fails in open-domain tasks where the output space is unlimited. ToT, designed to maintain different reasoning paths along its reasoning process, is more suitable for dealing with complex question types. However, the intermediate reasoning steps in NLP generation tasks are much less constrained and require more than a simple rule-based evaluation. The complexities in the question types (bridge v.s. comparison questions in Table 1), as well as the reasoning types (sequential v.s. parallel reasonings in Table 2), require more novel and fine-grained prompting methods to enhance the reasoning ability of LLMs.

To tackle the challenges and design a more reliable reasoning method for open-domain NLP tasks, we propose STOC-ToT, a stochastic ToT-based framework that instructs the model to generate different reasoning paths from the same question and assign probability scores to reasoning paths to effectively avoid reasoning dead-ends. To the best of our knowledge, our work is the first to adapt the tree-of-thought reasoning prompting to natural language tasks that require complex reasoning, such as MHQA. We provide an example overview of our framework in Figure 2. Specifically, we construct a tree-like reasoning structure by prompting the model to break down the original question into

smaller sub-questions to form different reasoning paths. We evaluate the validity of each reasoning path on three levels of aspects and arrive at a model-given probability score. At answer time, we innovatively propose to use constrained decoding in the answering process to reduce hallucination by forcing the model to generate grounded answers from evidence and letting models give concise and exact answers. Ultimately, we arrive at the best answer by choosing the path with the highest aggregated probability score. Experiments on two benchmarking MHQA datasets demonstrate that STOC-ToT significantly improves the reasoning ability of LLMs in complex reasoning scenarios, especially with GPT-4, improving Exact Match accuracy by 7%, and F1 score by 7.8 points on the HotpotQA dataset over the original tree-of-thought prompting. Our contributions are as follows:

- We propose STOC-ToT, which constructs a stochastic reasoning tree in complex reasoning scenarios. We introduce stochastic estimations on different reasoning paths, which helps the model have a more reliable reasoning process than previous reasoning prompting methods.
- We innovatively propose to use constrained decoding in the answering process. This step reduces model hallucination by forcing the model to generate grounded answers from evidence and letting models give concise and exact answers.
- We evaluate the effectiveness of STOC-ToT by conducting experiments on two MHQA datasets. We observe substantial improvements over other reasoning prompting methods, with STOC-ToT surpassing all other selected reasoning prompting baselines on 5 tested models.

2 Related Work

Multi-Hop Question Answering Multi-hop Question Answering (MHQA) is a challenging task requiring models to reason over different evidence across documents to answer a complex multi-hop question. Many high-quality MHQA datasets have been developed, including HotpotQA (Yang et al., 2018), WikiHop (Welbl et al., 2018), MuSiQue (Trivedi et al., 2022), and others. Among these, HotpotQA is the task’s most representative and widely used dataset. Previous state-of-the-art MHQA models often follow a two-stage pipeline: a retriever that extracts evidence from the documents, and a reader that reasons about the evidence

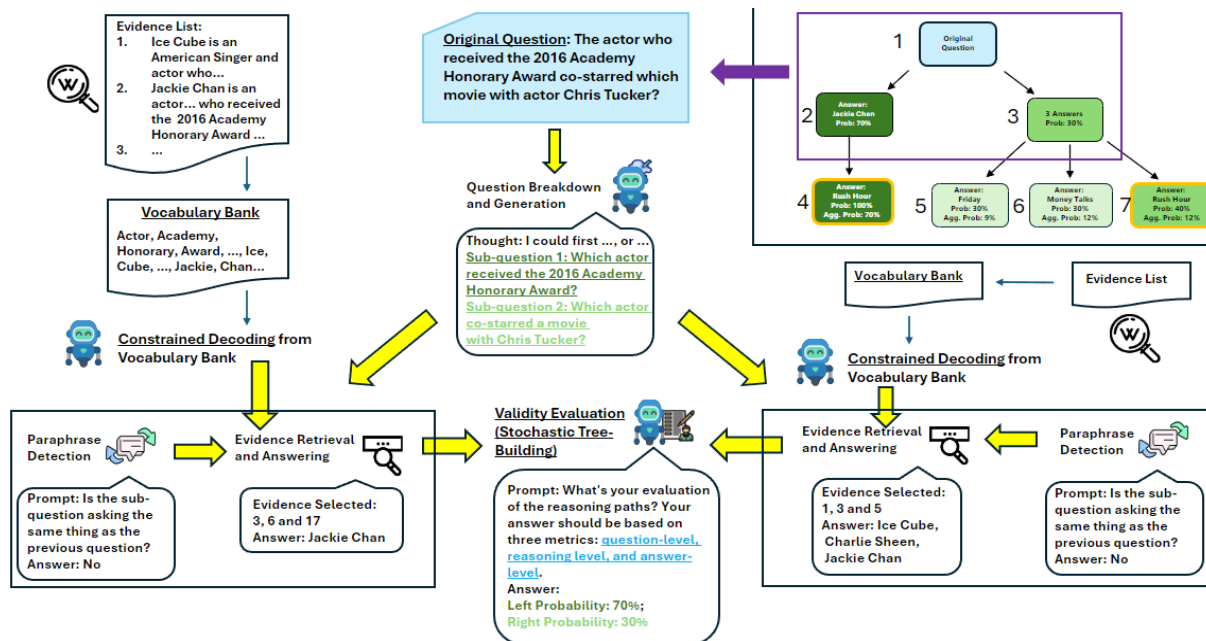


Figure 2: Overview of our framework, with the example in Figure 1. The top-right Corner shows the overall structure of the constructed tree, with each node’s label on the left. **Darker green** in the nodes means a **higher evaluated probability of the reasoning path**. The original Question is colored in blue. We chose the first round of our tree-building process as an example in the purple block.

165 to arrive at an answer (Zhu et al., 2021; Li et al.,
 166 2022). Other methods include beam-search (Zhang
 167 et al., 2023) and label-smoothing (Yin et al., 2023).
 168 Some LLM-based frameworks (Yao et al., 2023b;
 169 Gou et al., 2024) were also evaluated on the task of
 170 MHQA, but their performance fell short compared
 171 with supervised methods.

172 **Reasoning Prompting of LLMs** Various prompt
 173 engineering methods have been developed (Wei
 174 et al., 2022; Wang et al., 2023; Yao et al., 2023a;
 175 Besta et al., 2024; Sel et al., 2024; Chen et al.,
 176 2023), aiming to improve large language models’
 177 reasoning ability across various tasks and domains.
 178 Chain-of-thought (CoT) prompting (Wei et al.,
 179 2022) prompts the large language models (LLMs)
 180 to divide their reasoning process into smaller
 181 steps when solving a question, forming a chain
 182 of thoughts. Chain-of-thought self-consistency
 183 prompting (Wang et al., 2023) improves on the CoT
 184 method by proposing different reasoning chains
 185 and ensembles on the final result. Tree-of-thought
 186 (ToT) prompting method (Yao et al., 2023a) ac-
 187 tively maintains a tree of thoughts, where each
 188 thought is a coherent language sequence that serves
 189 as an intermediate step toward problem-solving.
 190 Graph-of-thought (Besta et al., 2024) further im-
 191 proves ToT by constructing a Directed Graph in-

192 stead of a tree. LLMs can loop over a thought to
 193 refine it and aggregate thoughts or chains.

194 **Constrained Decoding** Constrained decoding is
 195 the technique that asks the models to generate out-
 196 puts following a given set of rules. The most
 197 common way of conducting constrained generation
 198 uses beam search (Och and Ney, 2004) in decoding
 199 time. Before the LLM era, works on constrained
 200 decoding focused on task-specific sequence-to-
 201 sequence models that span across many fields, such
 202 as machine translation (Hokamp and Liu, 2017;
 203 Post and Vilar, 2018), named entity recognition
 204 (Lester et al., 2020), and dialogue generation (Bal-
 205 akrishnan et al., 2019). Recently, Microsoft intro-
 206 duced Guidance¹, which allows users of various
 207 large language models to control their outputs given
 208 a human-defined vocabulary or rules.

3 Method

3.1 Task Formation

210 Given a multi-hop question Q and background cor-
 211 pus of evidence P , the goal of our framework is
 212 to output the answer A to question Q , drawing its
 213 reasoning with the support of multiple evidence
 214 passages p_1, p_2, \dots retrieved from corpus P .
 215

¹<https://github.com/guidance-ai/guidance>

3.2 STOC-TOT Framework

For each of the questions Q , multiple reasoning lines and, thus, multiple ways of breaking down the question could exist. However, not every reasoning line would lead us to the right answer, and they take us to dead ends. To avoid such reasoning dead-ends, we build a stochastic reasoning tree to represent the possible reasoning lines and the probability of each reasoning line taking us to the right answer. We achieve this by proposing a self-interactive framework that automatically builds the reasoning tree given a multi-hop question. Figure 2 shows our framework with an example question.

In our reasoning process, we first prompt the model to propose different possible sub-questions to solve at each reasoning step. Each sub-question corresponds to one possible reasoning path and is presented as a node in the tree. We then ask the model to answer the generated sub-questions. To prevent hallucination and make the model more focused on the given question and evidence, we build a vocabulary bank using words from the evidence list and the original question and instruct the model to do constrained decoding from the vocabulary bank when generating its answers. After answering every sub-question generated from the same question in the previous reasoning level, we prompt the model to evaluate each reasoning path and estimate how likely the reasoning path would lead us to the right answer. This probability estimation would be assigned to the corresponding node in the tree. After the reasoning process finishes, each reasoning path would have an aggregated probability calculated from nodes along the path.

Formally, given a question Q , we instruct the model to generate sub-questions q_1, q_2, \dots, q_n , and build a tree structure with the original question Q as the root node and each question q_i as subsequent nodes. The tree would expand as each sub-question q_i has its sub-question q_j , and the reasoning paths are thus represented as branches in the tree structure. From the original question Q and the evidence list $E = e_1, e_2, \dots, e_n$, we build a vocabulary bank $V = [w_1, w_2, \dots, w_n], w_i \in Q, w_j \in E$. We then prompt the model to generate their answer a_1, a_2, \dots, a_n using only $w_i \in V$. We describe the details of our framework below.

Example-Based Sub-Question Generation Our framework starts with the sub-question generation module, which generates sub-questions q_1, q_2, \dots, q_n using the question Q_g from the pre-

vious reasoning level. The sub-questions are generated based on both the model’s reasoning ability and the model’s semantic understanding of the question Q_g . An example is given in Figure 2, where the sub-questions from nodes 2 and 3 were generated using the question from node 1. However, we cannot guarantee that each sub-question asked is a good sub-question, and sometimes, the generated sub-question merely repeats the previous question. We introduce the paraphrase detection module and pass on the generated sub-questions to reduce redundancy and improve question quality.

Paraphrase Detection Answering repetitive questions often leads to low-quality answers and time-consuming steps. Following the sub-question generation module, we introduce the paraphrase detection module to reduce redundancy and improve question quality. In this module, we prompt the model and ask it to distinguish informative questions from questions that merely repeat what is already stated at the previous reasoning level. If a sub-question is a paraphrase, we instruct the model to stop generating sub-questions from the current question. In other words, we prune the low-quality sub-branch of the tree that could otherwise be generated. By pruning these branches, we effectively improve the efficiency of our framework.

Evidence Retrieval and Answering We then move on to answering the question after our paraphrase detection module. Our evidence retrieval and answering module focuses on retrieving evidence and generating answers to the given sub-question. We also pass in the full evidence list provided and prompt the model to give out an answer to the given sub-question. The evidence retrieval and answering module selects relative evidence from an evidence pool for each sub-question and uses words only from the vocabulary bank to generate its final answer. We will discuss details of constrained decoding in Section 3.3. The generated sub-answer and the answered sub-question are then passed on to the sub-question generation module at the next level to continue the reasoning process.

Validity Estimation Not each sub-question asked is a good sub-question, and not each reasoning path is reasonable. After every sub-question q_i generated from the same question Q_g has been answered, we prompt the model to provide a probability estimation p_i for each (q_i, a_i) pair. This probability is the model’s evaluation of going down the

317	correct reasoning path. Specifically, this probability is obtained by prompting the model to consider the following three aspects:	word $w_j \notin V$, we manually set the score to negative infinity to prevent the model from generating them. Thus, every answer generated will only use words from the evidence list.	364
318			365
319			366
320	• Question Level: Is the question semantically clear and answerable?		367
321			
322	• Reasoning Level: Is the reasoning line coherent when considering previous levels?	Prompt-based Constrained Decoding For closed-source LLMs (e.g., GPT models), since we do not have access to their decoding function, we had to instruct the GPT models using prompts to do constrained decoding. We provide our prompt template used in Appendix A.	368
323			369
324	• Answer Level: Does the evidence fully support the answer to the question?		370
325			371
326	As shown in Figure 2, we conduct validity estimation for sub-questions and sub-answers in nodes 2 and 3 since the sub-questions were generated from the same question in node 1.		372
327			373
328			
329	At the leaf node of our tree, we would have a final question q_f , along with a final answer A to the original question Q , and also an aggregated probability $p_{final} = \prod_i p_i$, with each p_i being the probability of the nodes along the reasoning path. We assign p_{final} to the leaf node, representing the aggregated probability of answer A being the correct answer to Q .	4 Experimental Setup	374
330			
331		Dataset We compare STOC-ToT with baseline methods on the HotpotQA dataset (Yang et al., 2018) and the MuSiQue dataset (Trivedi et al., 2022), both of which are widely used MHQA datasets across state-of-the-art MHQA baselines. The experiments are conducted under the distractor setting, where we provide the model with an evidence pool containing both golden and irrelevant evidence. The model needs to find the golden evidence to answer the question correctly. We randomly selected 200 examples from each dataset as our evaluation set.	375
332			376
333			377
334			378
335			379
336			380
337			381
338	3.3 Constrained Decoding		382
339	One challenge for generative LLMs in the task of question answering is hallucination. LLMs often fail to pay attention to the golden evidence and hallucinate their own reference even when large amounts of evidence exist. To alleviate the problem of LLM hallucination during evidence selection and answer generation, we innovatively propose to use constrained decoding in the answering process to reduce hallucination by forcing the model to generate grounded answers from evidence and let models give concise and exact answers. As shown in Figure 2, we conduct constrained decoding by asking the model to generate words from the vocabulary bank, consisting of words taken only from the original question and the evidence list provided. More formally, we construct a vocabulary bank $V = w_1, w_2, \dots, w_i$ from all words in the provided evidence sentences. We conduct a simple filtering by removing common English stop words. We then instruct the model’s evidence retrieval and answering module to construct its answers using words only from the given vocabulary V .	Baselines We included three baselines:	383
340			384
341		• Vanilla Prompting with no examples provided. We only provide the model with questions and evidence and instruct it to output the answer.	385
342			386
343		• Chain-of-Thought (CoT) prompting (Wei et al., 2022) with a standard input-output (IO) prompt. We design the prompt with one in-context example, which presents the whole reasoning chain, including all intermediate steps.	387
344			388
345		• Tree-of-Thought prompting (Yao et al., 2023a) with slight modifications to adapt to the MHQA task. We largely followed the original framework and used majority voting on the reasoning lines to decide the final answer.	389
346			390
347			391
348			392
349			393
350			394
351			395
352			396
353			397
354			398
355			399
356			400
357			401
358			402
359			403
360			404
361	Code-based Constrained Decoding For open-source LLMs (e.g., Llama), we build our logit processor at the decoding time. Specifically, for every	We recognize that there are LLM-based retrieval augmented generation frameworks (Yao et al., 2023b; Gou et al., 2024) that were also evaluated on HotpotQA. However, we excluded them from our baselines as they used outside knowledge bases, which are under a different testing scenario.	405
362			406
363			407
		4.1 Implementation	407
		We experiment with the baselines and our model utilizing five LLMs: GPT-3.5-turbo (Brown et al.,	408
			409

Table 1: Performance comparison of STOC-ToT and baseline methods on the HotpotQA dataset.

Prompting Method	GPT3.5		GPT4		LLaMa2(13B)		LLaMa2(70B)		LLaMa3(8B)	
	EM	F1	EM	F1	EM	F1	EM	F1	EM	F1
Zero-Shot Vanilla	34.0	45.0	51.0	65.0	25.5	36.5	30.5	41.0	27.5	40.7
Chain-of-Thought	35.5	47.3	52.0	66.8	30.5	42.5	33.5	45.0	32.5	44.6
Tree-of-Thought	36.5	49.5	55.0	68.5	29.5	41.3	35.5	47.3	30.5	37.5
STOC-ToT	45.5	56.2	62.0	76.3	31.0	43.0	43.0	56.3	33.0	44.5
w/o constrained decoding	40.5	53.5	59.5	73.0	31.0	43.0	40.5	53.5	32.0	44.3

Table 2: Performance comparison of STOC-ToT and baseline methods on the MusiQue dataset.

Prompting Method	GPT3.5		GPT4		LLaMa2(13B)		LLaMa3(8B)	
	EM	F1	EM	F1	EM	F1	EM	F1
Zero-Shot Vanilla	17.0	28.8	31.5	41.2	9.5	16.0	12.0	19.2
Chain-of-Thought	18.0	29.7	32.5	44.2	11.0	17.5	12.5	21.6
Tree-of-Thought	20.5	32.0	35.0	47.3	11.0	17.2	12.0	20.6
STOC-ToT	26.5	38.0	42.0	55.3	11.5	18.0	14.5	22.0
w/o constrained decoding	24.0	35.5	38.5	51.0	11.5	18.0	14.0	22.0

2020) and GPT-4(OpenAI et al., 2024) from OpenAI, LLaMa 2-13B (Touvron et al., 2023), LLaMa 2-70B, and LLaMa 3-8B from MetaAI. Due to the lengthy running time, LLaMa 2-70B was not tested on the MusiQue dataset. For all models, We set the temperature to 0.5, top_k to 1.0, and maximum number of iterations to 5.

4.2 Evaluation Metric

Following the metrics in (Yang et al., 2018), we use Exact Match and F1 score as two evaluation metric. For an answer a given by our framework, the Exact Match score equals 1 if the answer span matches the golden answer exactly and 0 otherwise. The F1 metric measures the average overlap between the prediction and ground truth answers.

5 Results

5.1 Overall Results

We compare STOC-ToT with LLM baselines on the HotpotQA dataset and the MusiQue dataset and present our results in Tables 1 and 2. The backbone LLMs in our experiments include GPT3.5, GPT4, Llama2-13B, Llama2-70B, and Llama3-8B. Due to time constraints, we only tested with Llama2-70B on the HotpotQA dataset. On the HotpotQA dataset, STOC-ToT attains an on-average increase in performance of over 6 % compared with vanilla prompting on GPT models, and the improvement goes up to 11% when we further implement STOC-ToT with constrained decoding. On the more challenging MusiQue dataset, we still see an increase

in performance of STOC-ToT compared with the other baselines, most notably on GPT4, where we observe an 11.5% EM improvement (from 31.50 to 42.0).

Comparison with Tree-of-Thought STOC-ToT surpasses the original Tree-of-Thought prompting by 7% with the GPT4 model on both tested datasets. For LLMs with inferior reasoning ability, such as LLaMa2-8B, we still observe a performance improvement, even on the harder MusiQue dataset. These results suggest that STOC-ToT is more effective at forming and selecting reliable reasoning paths under complex reasoning scenarios.

Constrained Decoding Even though the LLM’s reasoning ability can be improved by reasoning prompting, such techniques have little help in preventing hallucination. However, STOC-ToT implements constrained decoding, which makes the model much more grounded to evidence when answering the question, effectively addressing hallucination issues and improving the overall performance of our framework.

5.2 Ablation Study

Sensitivity to Demonstration Question Type

We study the effect on STOC-ToT performance when different types of demonstration questions are provided in the prompt template. The HotPotQA dataset specified two types of questions. The "Bridge" question contains a "bridge entity" that connects the question and the final answer. In

Table 3: Performance of STOC-ToT with different prompt types on the HotpotQA dataset in terms of EM score. "Com" represents comparison questions, and "Bri" represents bridge questions.

Model Variant	GPT3.5		GPT4		LLaMa2(13B)		LLaMa2(70B)		LLaMa3(8B)	
	Com	Bri	Com	Bri	Com	Bri	Com	Bri	Com	Bri
Prompt: Comparison	58.8	41.0	76.5	57.2	38.2	31.9	58.8	41.0	44.1	33.7
Prompt: Bridge	55.9	43.4	73.5	59.0	35.3	32.5	55.9	42.2	41.2	34.9

Table 4: Question Type Examples. On the left side, the bridging entity is highlighted in red, and the final question is highlighted in orange. On the right side, entities that are being compared are highlighted in blue.

Bridge Question	Comparison Question
What distinction is held by the former NBA player who was a member of the Charlotte Hornets during their 1992-93 season and was head coach for the WNBA team Charlotte Sting?	Were Scott Derrickson and Ed Wood of the same nationality?

Table 5: Reasoning Type Examples. On the left side, the entity in red needs to be found before solving the question in orange. On the right side, questions with parallel reasoning contain parts (highlighted in blue) that can be solved in arbitrary order.

Sequential Reasoning	Parallel Reasoning
The football manager who recruited David Beckham managed Manchester United during what time-frame?	What distinction is held by the former NBA player who was a member of the Charlotte Hornets during their 1992-93 season and was head coach for the WNBA team Charlotte Sting?

contrast, the "Comparison" question requires the model to compare two entities of the same type. Of the 200 questions in our evaluation set, 34 are comparison questions, and 166 are bridge questions. Examples of bridge and comparison questions are in Table 4.

We examined STOC-ToT performance under the two different question types, each with a different prompt template: one containing only a comparison question as an example and the other containing only a bridge question as an example. We provide the content of our templates in Appendix A. Results are shown in Table 3. We observe that the difference in prompt templates influences the performance of our framework under different question types by a small margin. The comparison questions are generally easier to solve, and STOC-ToT performs better on comparison questions than on bridge questions. STOC-ToT will handle comparison questions better if the prompt template contains comparison questions and vice versa.

Question and Reasoning Types We examine STOC-ToT, Tree-of-Thought prompting, and Chain-of-Thought prompting by comparing their performance under different question-type settings. Detailed results are shown in Figure 3(a). STOC-ToT performs better at both Bridge Questions and Sequential Questions, suggesting that STOC-ToT can avoid reasoning dead-ends and is better at form-

ing intermediate reasoning lines.

We also conduct an in-depth analysis of the reasoning types in the existing MHQA datasets by randomly selecting 100 questions from our testing set. The questions are roughly divided into two categories: 1) tree-like parallel reasoning and 2) chain-like sequential reasoning. Questions with parallel reasoning contain two or more reasoning paths that can be solved arbitrarily. Questions with sequential reasoning follow a strict reasoning chain, and all the sub-questions must be solved to form the correct reasoning process. All comparison questions are parallel reasoning, but some bridge questions contain parallel reasoning. Examples of sequential and parallel reasoning questions are in Table 5. Out of the selected 100 questions, 59 questions were Sequential and 41 questions were Parallel. Results are shown in Figure 3(b). STOC-ToT performs better on both reasoning types, especially on questions containing parallel reasoning. This suggests that STOC-ToT’s stochastic way of forming the tree is very effective when solving questions containing multiple reasoning paths.

Performance and Hops As the number of hops increases in a question, the reasoning line gets more complex and varied. Figure 4 shows the performances of different prompting techniques on questions in the MusiQue dataset with differ-

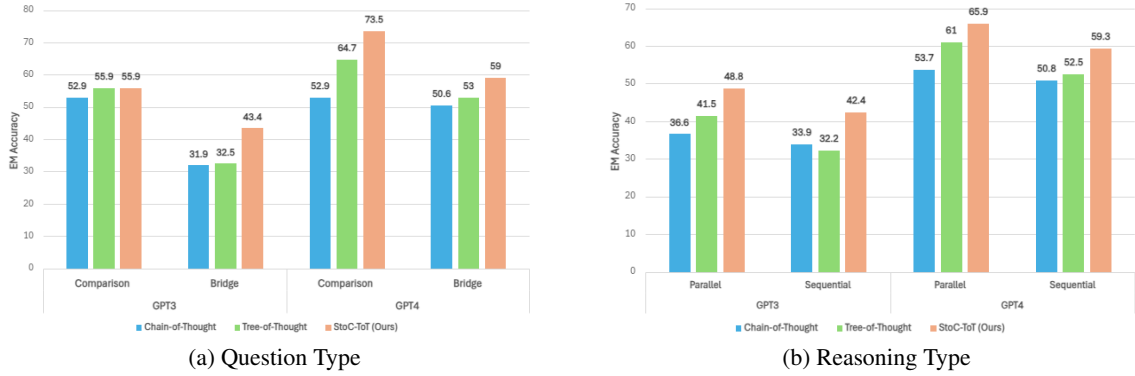


Figure 3: Performance comparison of Chain-of-Thought, Tree-of-Thought, and STOC-ToT on questions of different question types (Left) and reasoning types (Right). Experiments were done on the HotpotQA dataset.

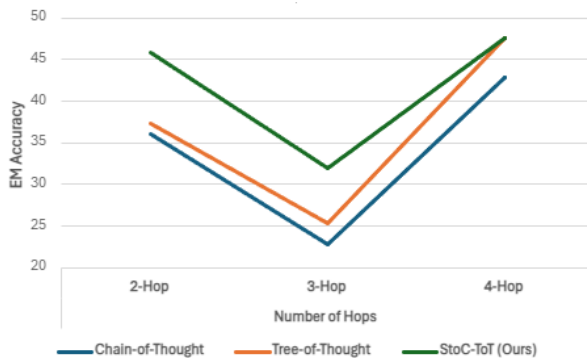


Figure 4: Performance comparison of CoT, ToT, and STOC-ToT on different number of hops in the question. Experiments done in the MusiQue dataset.

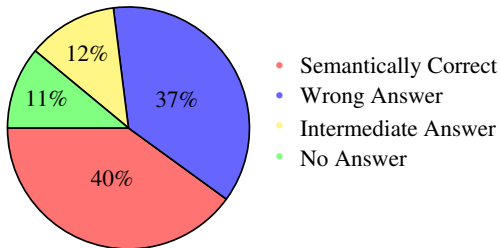


Figure 5: Ratio of different categories in error cases, on the HotpotQA dataset.

ent numbers of hops. STOC-ToT performs best in all categories, demonstrating our framework’s superior ability to deal with complex reasoning scenarios. This ablation study was conducted only on GPT4, as other models performed poorly on 3-hop and 4-hop scenarios, regardless of the reasoning prompting technique used.

Error Analysis We conduct a detailed analysis of the errors made by our framework on GPT3 and GPT4, and present our results in Figure 5. We categorize the errors into four types: (1) **No Answer**:

our framework did not come up with an answer for the question due to not finishing the reasoning process; (2) **Intermediate Answer**: our framework came up with an answer for one of the intermediate hops instead of for the final question; (3) **Wrong Answer**: our framework came up with an answer that is neither the final answer nor one of the intermediate answers; (4) **Semantically Correct**: our framework came up with the right answer, but did not have an exact match with the final answer. Appendix B shows examples of each error category. Large amounts of error cases were correct answers with extra wording or hallucination errors, signaling potential improvements over our constrained decoding scheme. Reasoning process errors, including no answer and intermediate answer, make up only 25% of the total error cases. This result shows that our framework is capable of building a robust reasoning process for complex questions.

6 Conclusion

This paper proposes STOC-ToT, a stochastic tree-of-thought reasoning framework with constrained generation for multi-hop question answering. STOC-ToT is specialized in dealing with complex reasoning scenarios in natural language tasks. Experiments on two benchmark datasets show that our framework outperforms previous reasoning prompting techniques with multiple Large Language Models. Detailed analysis shows that our framework is capable of building a robust reasoning process given different types of questions. Further research can aim to enhance the reliability of our framework by proposing better validity evaluation schemes and more effective methods for improving groundedness and preventing hallucination.

573
574
575
576
577
578
579
580
581
582
583
584
585
586
587
588
589
590
591
592
593
594
595
596
597
598
599
600
601
602
603
604
605
606
607
608
609
610
611
612
613
614
615
616
617
618
619
620
621
622
623

Limitations

Our framework relies on initiating multiple model instances and requires multiple prompts per round. The repetitive callings impose heavy time costs for our framework, even after implementing our paraphrase module. Another limitation comes from how we generated sub-questions. Currently, we directly prompt the model to generate sub-questions. A more complex standard can be used to increase the quality of the sub-questions generated. Also, more extensive experiments should be provided, including experimenting on other different datasets and case studies.

Ethics Statement

This research adhered to the ethical standards and best practices outlined in the ACL Code of Ethics. Language Models can sometimes produce illogical or inaccurate reasoning paths, so their outputs should be cautiously used. The outputs are only examined to understand how a model arrives at its answers and investigate why it makes certain errors. All experiments used publicly available datasets from previously published works and did not involve ethical or privacy issues.

References

Anusha Balakrishnan, Jinfeng Rao, Kartikeya Upasani, Michael White, and Rajen Subba. 2019. [Constrained decoding for neural NLG from compositional representations in task-oriented dialogue](#). In *Proceedings of the 57th Conference of the Association for Computational Linguistics, ACL 2019, Florence, Italy, July 28- August 2, 2019, Volume 1: Long Papers*, pages 831–844. Association for Computational Linguistics.

Maciej Besta, Nils Blach, Ales Kubicek, Robert Gerstenberger, Michal Podstawski, Lukas Gianinazzi, Joanna Gajda, Tomasz Lehmann, Hubert Niewiadomski, Piotr Nyczyk, and Torsten Hoefler. 2024. [Graph of thoughts: Solving elaborate problems with large language models](#). In *Thirty-Eighth AAAI Conference on Artificial Intelligence, AAAI 2024, Thirty-Sixth Conference on Innovative Applications of Artificial Intelligence, IAAI 2024, Fourteenth Symposium on Educational Advances in Artificial Intelligence, EAAI 2014, February 20-27, 2024, Vancouver, Canada*, pages 17682–17690. AAAI Press.

Tom B. Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared 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 M. Ziegler, Jeffrey Wu,

Clemens Winter, Christopher 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](#). *CoRR*, abs/2005.14165. 624
625
626
627
628
629

Wenhu Chen, Xueguang Ma, Xinyi Wang, and William W. Cohen. 2023. [Program of thoughts prompting: Disentangling computation from reasoning for numerical reasoning tasks](#). *Transactions on Machine Learning Research*. 630
631
632
633
634

Zhibin Gou, Zhihong Shao, Yeyun Gong, yelong shen, Yujia Yang, Nan Duan, and Weizhu Chen. 2024. [CRITIC: Large language models can self-correct with tool-interactive critiquing](#). In *The Twelfth International Conference on Learning Representations*. 635
636
637
638
639

Chris Hokamp and Qun Liu. 2017. [Lexically constrained decoding for sequence generation using grid beam search](#). In *Proceedings of the 55th Annual Meeting of the Association for Computational Linguistics, ACL 2017, Vancouver, Canada, July 30 - August 4, Volume 1: Long Papers*, pages 1535–1546. Association for Computational Linguistics. 640
641
642
643
644
645
646

Brian Lester, Daniel Pressel, Amy Hemmeter, Sagnik Ray Choudhury, and Srinivas Bangalore. 2020. [Constrained decoding for computationally efficient named entity recognition taggers](#). In *Findings of the Association for Computational Linguistics: EMNLP 2020, Online Event, 16-20 November 2020*, volume EMNLP 2020 of *Findings of ACL*, pages 1841–1848. Association for Computational Linguistics. 647
648
649
650
651
652
653
654

Xin-Yi Li, Weixian Lei, and Yubin Yang. 2022. [From easy to hard: Two-stage selector and reader for multi-hop question answering](#). *ICASSP 2023 - 2023 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pages 1–5. 655
656
657
658
659

Franz Josef Och and Hermann Ney. 2004. [The alignment template approach to statistical machine translation](#). *Comput. Linguistics*, 30(4):417–449. 660
661
662

OpenAI, Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman, Diogo Almeida, Janko Altenschmidt, Sam Altman, Shyamal Anadkat, Red Avila, Igor Babuschkin, Suchir Balaji, Valerie Balcom, Paul Baltescu, Haiming Bao, Mohammad Bavarian, Jeff Belgum, Irwan Bello, Jake Berdine, Gabriel Bernadett-Shapiro, Christopher Berner, Lenny Bogdonoff, Oleg Boiko, Madelaine Boyd, Anna-Luisa Brakman, Greg Brockman, Tim Brooks, Miles Brundage, Kevin Button, Trevor Cai, Rosie Campbell, Andrew Cann, Brittany Carey, Chelsea Carlson, Rory Carmichael, Brooke Chan, Che Chang, Fotis Chantzis, Derek Chen, Sully Chen, Ruby Chen, Jason Chen, Mark Chen, Ben Chess, Chester Cho, Casey Chu, Hyung Won Chung, Dave Cummings, Jeremiah Currier, Yunxing Dai, Cory Decareaux, Thomas Degry, Noah Deutsch, Damien Deville, Arka Dhar, David Dohan, Steve Dowling, Sheila Dunning, Adrien Ecoffet, Atty Eleti,

682	Tyna Eloundou, David Farhi, Liam Fedus, Niko Felix,	745
683	Simón Posada Fishman, Juston Forte, Isabella Ful-	746
684	ford, Leo Gao, Elie Georges, Christian Gibson, Vik	747
685	Goel, Tarun Gogineni, Gabriel Goh, Rapha Gontijo-	748
686	Lopes, Jonathan Gordon, Morgan Grafstein, Scott	749
687	Gray, Ryan Greene, Joshua Gross, Shixiang Shane	
688	Gu, Yufei Guo, Chris Hallacy, Jesse Han, Jeff Harris,	
689	Yuchen He, Mike Heaton, Johannes Heidecke, Chris	
690	Hesse, Alan Hickey, Wade Hickey, Peter Hoeschele,	
691	Brandon Houghton, Kenny Hsu, Shengli Hu, Xin	
692	Hu, Joost Huizinga, Shantanu Jain, Shawn Jain,	
693	Joanne Jang, Angela Jiang, Roger Jiang, Haozhun	
694	Jin, Denny Jin, Shino Jomoto, Billie Jonn, Hee-	
695	woo Jun, Tomer Kaftan, Łukasz Kaiser, Ali Ka-	
696	mali, Ingmar Kanitscheider, Nitish Shirish Keskar,	
697	Tabarak Khan, Logan Kilpatrick, Jong Wook Kim,	
698	Christina Kim, Yongjik Kim, Jan Hendrik Kirchn-	
699	er, Jamie Kiros, Matt Knight, Daniel Kokotajlo,	
700	Łukasz Kondraciuk, Andrew Kondrich, Aris Kon-	
701	stantinidis, Kyle Kopic, Gretchen Krueger, Vishal	
702	Kuo, Michael Lampe, Ikai Lan, Teddy Lee, Jan	
703	Leike, Jade Leung, Daniel Levy, Chak Ming Li,	
704	Rachel Lim, Molly Lin, Stephanie Lin, Mateusz	
705	Litwin, Theresa Lopez, Ryan Lowe, Patricia Lue,	
706	Anna Makanju, Kim Malfacini, Sam Manning, Todor	
707	Markov, Yaniv Markovski, Bianca Martin, Katie	
708	Mayer, Andrew Mayne, Bob McGrew, Scott Mayer	
709	McKinney, Christine McLeavey, Paul McMillan,	
710	Jake McNeil, David Medina, Aalok Mehta, Jacob	
711	Menick, Luke Metz, Andrey Mishchenko, Pamela	
712	Mishkin, Vinnie Monaco, Evan Morikawa, Daniel	
713	Mossing, Tong Mu, Mira Murati, Oleg Murk, David	
714	Mély, Ashvin Nair, Reiichiro Nakano, Rajeev Nayak,	
715	Arvind Neelakantan, Richard Ngo, Hyeonwoo Noh,	
716	Long Ouyang, Cullen O’Keefe, Jakub Pachocki, Alex	
717	Paino, Joe Palermo, Ashley Pantuliano, Giambat-	
718	tista Parascandolo, Joel Parish, Emy Parparita, Alex	
719	Passos, Mikhail Pavlov, Andrew Peng, Adam Perel-	
720	man, Filipe de Avila Belbute Peres, Michael Petrov,	
721	Henrique Ponde de Oliveira Pinto, Michael, Poko-	
722	rny, Michelle Pokrass, Vitchyr H. Pong, Tolly Pow-	
723	ell, Alethea Power, Boris Power, Elizabeth Proehl,	
724	Raul Puri, Alec Radford, Jack Rae, Aditya Ramesh,	
725	Cameron Raymond, Francis Real, Kendra Rimbach,	
726	Carl Ross, Bob Rotsted, Henri Roussez, Nick Ry-	
727	der, Mario Saltarelli, Ted Sanders, Shibani Santurkar,	
728	Girish Sastry, Heather Schmidt, David Schnurr, John	
729	Schulman, Daniel Selsam, Kyla Sheppard, Toki	
730	Sherbakov, Jessica Shieh, Sarah Shoker, Pranav	
731	Shyam, Szymon Sidor, Eric Sigler, Maddie Simens,	
732	Jordan Sitkin, Katarina Slama, Ian Sohl, Benjamin	
733	Sokolowsky, Yang Song, Natalie Staudacher, Fel-	
734	ipe Petroski Such, Natalie Summers, Ilya Sutskever,	
735	Jie Tang, Nikolas Tezak, Madeleine B. Thompson,	
736	Phil Tillet, Amin Tootoonchian, Elizabeth Tseng,	
737	Preston Tuggle, Nick Turley, Jerry Tworek, Juan Fe-	
738	lipe Cerón Uribe, Andrea Vallone, Arun Vijayvergiya,	
739	Chelsea Voss, Carroll Wainwright, Justin Jay Wang,	
740	Alvin Wang, Ben Wang, Jonathan Ward, Jason Wei,	
741	CJ Weinmann, Akila Welihinda, Peter Welinder, Ji-	
742	ayi Weng, Lilian Weng, Matt Wiethoff, Dave Willner,	
743	Clemens Winter, Samuel Wolrich, Hannah Wong,	
744	Lauren Workman, Sherwin Wu, Jeff Wu, Michael	
	Wu, Kai Xiao, Tao Xu, Sarah Yoo, Kevin Yu, Qim-	745
	ing Yuan, Wojciech Zaremba, Rowan Zellers, Chong	746
	Zhang, Marvin Zhang, Shengjia Zhao, Tianhao	747
	Zheng, Juntang Zhuang, William Zhuk, and Barret	748
	Zoph. 2024. Gpt-4 technical report .	749
	Matt Post and David Vilar. 2018. Fast lexically con-	750
	strained decoding with dynamic beam allocation for	751
	neural machine translation . In <i>Proceedings of the</i>	752
	<i>2018 Conference of the North American Chapter</i>	753
	<i>of the Association for Computational Linguistics:</i>	754
	<i>Human Language Technologies, NAACL-HLT 2018,</i>	755
	<i>New Orleans, Louisiana, USA, June 1-6, 2018, Vol-</i>	756
	<i>ume 1 (Long Papers)</i> , pages 1314–1324. Association	757
	for Computational Linguistics.	758
	Bilgehan Sel, Ahmad Tawaha, Vanshaj Khattar, Ruoxi	759
	Jia, and Ming Jin. 2024. Algorithm of thoughts:	760
	Enhancing exploration of ideas in large language	761
	models . In <i>Forty-first International Conference on</i>	762
	<i>Machine Learning</i> .	763
	Hugo Touvron, Louis Martin, Kevin Stone, Peter Al-	764
	bert, Amjad Almahairi, Yasmine Babaei, Nikolay	765
	Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti	766
	Bhosale, Dan Bikel, Lukas Blecher, Cristian Canton	767
	Ferrer, Moya Chen, Guillem Cucurull, David Esiobu,	768
	Jude Fernandes, Jeremy Fu, Wenyin Fu, Brian Fuller,	769
	Cynthia Gao, Vedanuj Goswami, Naman Goyal, An-	770
	thony Hartshorn, Saghar Hosseini, Rui Hou, Hakan	771
	Inan, Marcin Kardas, Viktor Kerkez, Madian Khabsa,	772
	Isabel Kloumann, Artem Korenev, Punit Singh Koura,	773
	Marie-Anne Lachaux, Thibaut Lavril, Jenya Lee, Di-	774
	ana Liskovich, Yinghai Lu, Yuning Mao, Xavier Mar-	775
	tinet, Todor Mihaylov, Pushkar Mishra, Igor Moly-	776
	bog, Yixin Nie, Andrew Poulton, Jeremy Reizen-	777
	stein, Rashi Rungta, Kalyan Saladi, Alan Schelten,	778
	Ruan Silva, Eric Michael Smith, Ranjan Subrama-	779
	nian, Xiaoqing Ellen Tan, Binh Tang, Ross Tay-	780
	lor, Adina Williams, Jian Xiang Kuan, Puxin Xu,	781
	Zheng Yan, Iliyan Zarov, Yuchen Zhang, Angela Fan,	782
	Melanie Kambadur, Sharan Narang, Aurelien Rod-	783
	riguez, Robert Stojnic, Sergey Edunov, and Thomas	784
	Scialom. 2023. Llama 2: Open foundation and fine-	785
	tuned chat models .	786
	Harsh Trivedi, Niranjan Balasubramanian, Tushar Khot,	787
	and Ashish Sabharwal. 2022. 9835 musique: Multi-	788
	hop questions via single-hop question composition .	789
	<i>Trans. Assoc. Comput. Linguistics</i> , 10:539–554.	790
	Xuezhi Wang, Jason Wei, Dale Schuurmans, Quoc V.	791
	Le, Ed H. Chi, Sharan Narang, Aakanksha Chowd-	792
	hery, and Denny Zhou. 2023. Self-consistency	793
	improves chain of thought reasoning in language	794
	models . In <i>The Eleventh International Conference</i>	795
	<i>on Learning Representations, ICLR 2023, Kigali,</i>	796
	<i>Rwanda, May 1-5, 2023</i> . OpenReview.net.	797
	Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten	798
	Bosma, Brian Ichter, Fei Xia, Ed H. Chi, Quoc V. Le,	799
	and Denny Zhou. 2022. Chain-of-thought prompting	800
	elicits reasoning in large language models . In <i>Ad-</i>	801
	<i>vances in Neural Information Processing Systems 35:</i>	802

803 *Annual Conference on Neural Information Process-*
804 *ing Systems 2022, NeurIPS 2022, New Orleans, LA,*
805 *USA, November 28 - December 9, 2022.*

806 Johannes Welbl, Pontus Stenetorp, and Sebastian Riedel.
807 2018. [Constructing datasets for multi-hop reading](#)
808 [comprehension across documents](#). *Trans. Assoc.*
809 *Comput. Linguistics*, 6:287–302.

810 Zhilin Yang, Peng Qi, Saizheng Zhang, Yoshua Ben-
811 gio, William W. Cohen, Ruslan Salakhutdinov, and
812 Christopher D. Manning. 2018. HotpotQA: A dataset
813 for diverse, explainable multi-hop question answer-
814 ing. In *Conference on Empirical Methods in Natural*
815 *Language Processing (EMNLP)*.

816 Shunyu Yao, Dian Yu, Jeffrey Zhao, Izhak Shafran,
817 Tom Griffiths, Yuan Cao, and Karthik Narasimhan.
818 2023a. [Tree of thoughts: Deliberate problem solving](#)
819 [with large language models](#). In *Advances in Neural*
820 *Information Processing Systems 36: Annual Confer-*
821 *ence on Neural Information Processing Systems 2023,*
822 *NeurIPS 2023, New Orleans, LA, USA, December 10*
823 *- 16, 2023.*

824 Shunyu Yao, Jeffrey Zhao, Dian Yu, Nan Du, Izhak
825 Shafran, Karthik R. Narasimhan, and Yuan Cao.
826 2023b. [React: Synergizing reasoning and acting](#)
827 [in language models](#). In *The Eleventh International*
828 *Conference on Learning Representations, ICLR 2023,*
829 *Kigali, Rwanda, May 1-5, 2023*. OpenReview.net.

830 Zhangyue Yin, Yuxin Wang, Xiannian Hu, Yiguang Wu,
831 Hang Yan, Xinyu Zhang, Zhao Cao, Xuanjing Huang,
832 and Xipeng Qiu. 2023. [Rethinking label smoothing](#)
833 [on multi-hop question answering](#). In *Chinese*
834 *Computational Linguistics - 22nd China National*
835 *Conference, CCL 2023, Harbin, China, August 3-5,*
836 *2023, Proceedings*, volume 14232 of *Lecture Notes*
837 *in Computer Science*, pages 72–87. Springer.

838 Jiahao Zhang, Haiyang Zhang, Dongmei Zhang, Yong
839 Liu, and Shen Huang. 2023. [Beam retrieval: General](#)
840 [end-to-end retrieval for multi-hop question answer-](#)
841 [ing](#). *CoRR*, abs/2308.08973.

842 Fengbin Zhu, Wenqiang Lei, Chao Wang, Jianming
843 Zheng, Soujanya Poria, and Tat-Seng Chua. 2021.
844 [Retrieving and reading: A comprehensive sur-](#)
845 [vey on open-domain question answering](#). *CoRR*,
846 abs/2101.00774.

847
848
849
850
851
852
853
854
855
856
857
858
859
860
861
862
863
864
865
866
867
868
869
870
871
872
873
874
875
876
877
878
879
880
881
882
883
884
885
886
887
888
889
890
891
892
893
894
895

A Prompt Templates

Sub Question Generation Template The prompt template containing one comparison question and one bridge question is given below:

prompt: Break a question into high-quality sub-questions that are easier to answer. Here are two examples as guidelines:

"Question: Are Tokyo and Busan in the same country? Thought 1: I could either find which country Tokyo is located in, or which country Busan is located in. Sub Question 1-1: Which country is Tokyo located in? Sub Question 1-2: Which country is Busan located in?"

"Question: Tokyo is located in the country that has what colors present on its national flag? Thought 1: I need to first find out which country Tokyo is located in. Sub Question 1-1: Which country is Tokyo located in?"

Only give out your thought process and current-level sub-questions. Do not give out answers to your questions. Question: *Given Question*. Thought 1:

Prompt-based Constrained Generation Template The prompt template at answering time is given below:

prompt: Given a question and a list of evidence that may of help, give your answer directly, using words only from the vocabulary bank, without any explanations.

Question: *Given Question*. Evidence as reference: *Given Evidence*. Vocabulary Bank: *Given Vocabulary*. Answer:

B Examples of the Error Cases

•Type-2: Intermediate Answer

Question:
Where does the hotel and casino located in which Bill Cosby's third album was recorded?

Answer given by STOC-TOT on GPT4:
Las Vegas.

Golden Answer:
Las Vegas Strip in Paradise.

•Type-3: Wrong Answer

Question:
Aside from the Apple Remote, what other device can control the program Apple Remote was originally designed to interact with?

Answer given by STOC-TOT on GPT4:
siri remote and devices with netsupport manager software
Golden Answer:
keyboard function keys

•Type-4: Semantically Correct

Question:
Roger O. Egeberg was Assistant Secretary for Health and Scientific Affairs during the administration of a president that served during what years?

Answer given by STOC-TOT on GPT4:
1969 to 1974
Golden Answer:
1969 until 1974