

Chain of Evidences and Evidence to Generate: Prompting for Context Grounded and Retrieval Augmented Reasoning

Anonymous ACL submission

Abstract

While chain-of-thoughts (CoT) prompting has revolutionized how LLMs perform reasoning tasks, its current methods and variations (e.g. Self-consistency, ReACT, Reflexion, Tree-of-Thoughts (ToT), Cumulative Reasoning (CR) etc..) suffer from limitations like limited context grounding, hallucination/inconsistent output generation, and iterative sluggishness. To overcome these challenges, we introduce a novel mono/dual-step zero-shot prompting framework built upon two unique strategies **Chain of Evidences (COE)** and **Evidence to Generate (E2G)**. Instead of unverified reasoning claims, our innovative approaches leverage the power of "evidence for decision making" by first focusing exclusively on the thought sequences explicitly mentioned in the context which then serve as extracted evidence, guiding the LLM's output generation process with greater precision and efficiency. This simple yet potent approach unlocks the full potential of chain-of-thoughts prompting, facilitating faster, more reliable, and contextually aware reasoning in LLMs. Our framework consistently achieves remarkable results across various knowledge-intensive reasoning and generation tasks, surpassing baseline approaches with state-of-the-art LLMs. For instance, (i) on the LogiQA benchmark using GPT-4, COE achieves a new state-of-the-art accuracy of 53.8%, surpassing CoT by 18%, ToT by 11%, and CR by 9%; (ii) CoE with PaLM-2 outperforms the variable-shot performance of Gemini Ultra by 0.9 F1 points, achieving an F1 score of 83.3 on DROP.

1 Introduction

Retrieval-augmented or context-based generation serves as a mean for leveraging relevant information, empowering large language models (LLMs) to reduce the factual errors in their generation (Asai et al., 2023a,b). However, despite the expansion in model and data size, LLMs struggle in contextual

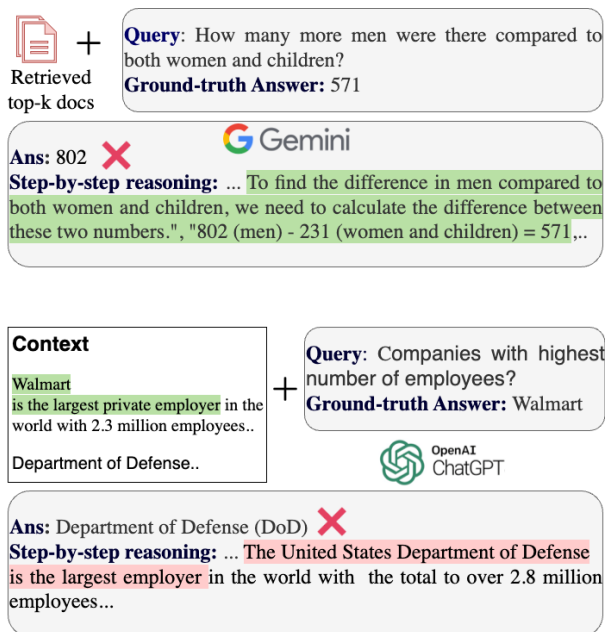


Figure 1: CoT & variants falter in context-aware reasoning. Top: Overwhelming long-text complexity leads models' failure even when it generates partially/fully correct reasoning (in green). Bottom: Ungrounded internal reasoning fails to grasp context, confusing "DoD" (ungrounded private org in red) vs Walmart (in green).

reasoning. This challenge is further amplified when dealing with retrieved information that are often long and imperfect text with distractive contents.

To bolster LLM's reasoning capabilities, the Chain-of-Thought (CoT) prompting paradigm has emerged as a potent tool (Wei et al., 2022). Subsequent methods, including Self-consistency (SC; Wang et al., 2022), ReACT (Yao et al., 2022), Reflexion (Shinn et al., 2023), Tree of Thoughts (ToT; Yao et al., 2023), and Cumulative Reasoning (CR; Zhang et al., 2023b), generalize CoT with various multi-objective, ensemble-based, or tool-augmented, and trial & error approaches but do not address the complexities of context-grounded or retrieval augmented generations (RAG). We highlight two of their pivotal bottlenecks: (i) CoT focuses solely on expanding steps without verifying

060 hypotheses; (ii) excessively long retrieved text can
061 lead to incorrect conclusions even with valid CoT
062 reasonings (example in Figure 1).

063 To tackle, multi-step reasoning prompting
064 (Wang et al., 2023a; Zhao et al., 2023; Trivedi et al.,
065 2023; Fu et al., 2022; Creswell et al., 2022; Li et al.,
066 2023) has emerged as a promising alternative to
067 CoT that breaks down complex problems into se-
068 ries of reasoning substeps. However, in large, they
069 are very complicated, do not address context-based
070 or long text retrieval augmented generation, require
071 rigorous verification in every-step, and importantly
072 employ different intermediate prompts (e.g., ra-
073 tionale selection & inference/premise derivation)
074 that require w/ k-shot annotated in-context exem-
075 plars, which are often very difficult to form (Ya-
076 sunaga et al., 2024). Therefore, unlocking CoT’s
077 true potential for RAG and context-aware reason-
078 ing remains unanswered. To address, in this pa-
079 per, we propose a simple verification-free zero-
080 shot prompting framework for context-aware and
081 retrieval augmented reasoning.

082 Ours framework consists of two unique and real-
083 time prompting strategies particularly tailored for
084 context-aware reasoning. First, single-step **Chain-**
085 **of-Evidences (COE)**: to address the problem of
086 ungrounded reasoning hypotheses, our designed
087 prompt asks for specific thought sequences that are
088 explicitly mentioned in the context. We call these
089 series of intermediate reasoning steps w/ directly
090 extracted rationales from the given context as *ev-*
091 *idence* (as in human decision making). Our key
092 distinction from existing CoT approaches is that
093 instead of mere "thinking step-by-step" (Kojima
094 et al., 2022) our prompt instruction asks for "step-
095 by-step reasoning w/ evidence & explanation".

096 Second, dual-step **Evidence to Generate**
097 **(E2G)**: to facilitate LLMs’ answering the query
098 properly even w/ retrieval augmented long-text con-
099 texts, we split the task into steps. In the first step
100 (E), we adopt prompts similar to COE and generate
101 both the *Answer & Evidence*. Then in next step
102 (G), we pass only the *Evidence* as context for a
103 second round of COE to LLM. G Step *Answer* is
104 predicted as the final answer. In contrast to com-
105 plex long original context in E step, the *Evidence*
106 is a concise short text that directly answer the input
107 query, G step is very fast, and simpler for the model
108 to generate answer.

109 In experiments with different LLMs, we show
110 that our prompts consistently outperform existing
111 approaches in a diverse set of eight context-driven

112 tasks, including natural QA, complex multi-hop,
113 long-form QA, fact checking, dialog generation,
114 and reading comprehension tasks. Since, even with
115 such techniques, it is non-trivial to comprehend
116 why and how this works and how to setup the
117 prompt to function correctly, cost-effectively, and
118 robustly. To this end, we perform case studies, ana-
119 lyze different alternatives and reveal the strengths
120 and weaknesses of our approach. We will release
121 our prompts and outputs on these benchmarks as a
122 new instruction tuning dataset for future research.

123 2 Related Works and Preliminaries

124 2.1 Prompting LLMs

125 Various prompting paradigms have been studied
126 in literature toward enhancing reasoning in LLMs.
127 In Section 1, we provide a (non-exhaustive) list
128 of CoT approaches. Among others, search-based
129 (Pryzant et al., 2023; Lu et al., 2021), Program-
130 aided LLM generation (Liu et al., 2023a; Gao et al.,
131 2023; Jung et al., 2022; Zhu et al., 2022), self gen-
132 eration of prompts (He et al., 2023; Yasunaga et al.,
133 2023; Sun et al., 2022; Kim et al., 2022; Li et al.,
134 2022), self evaluation based approaches (Madaan
135 et al., 2023; Xie et al., 2023; Kim et al., 2023; Paul
136 et al., 2023) have been studied. Other works have
137 also been extended w/ more complex multi-step
138 reasoning procedure (e.g., using a different fine-
139 tuned model (Zelikman et al., 2022; Nye et al.,
140 2021; Lester et al., 2021)) or for domain specific
141 applications (Parvez et al., 2023, 2021; Ouyang
142 et al., 2022; Sanh et al., 2021; Wei et al., 2021).

143 2.2 Chain-of-Thoughts (CoT) Prompting

144 Chain-of-thoughts (CoT; (Wei et al., 2022)) is a
145 prompting framework that guides LLMs to pro-
146 duce intermediate reasoning steps towards the fi-
147 nal answer, enhancing its reasoning. Original ver-
148 sion of CoT employs a few-shot version by pro-
149 viding multiple exemplars of the reasoning process
150 (question–reasoning–answer), leveraging LLMs’
151 in-context learning abilities. However, due to the re-
152 quirement of labeled exemplars, it quickly evolved
153 with a 0-shot instance (Kojima et al., 2022). 0-
154 shot CoT prompts LLMs with a general instruction
155 like “think step by step” to produce intermediate
156 reasoning steps (See Figure 2).

157 3 Our Prompting Framework

158 In this section, we develop our prompting frame-
159 work for context-grounding and retrieval aug-

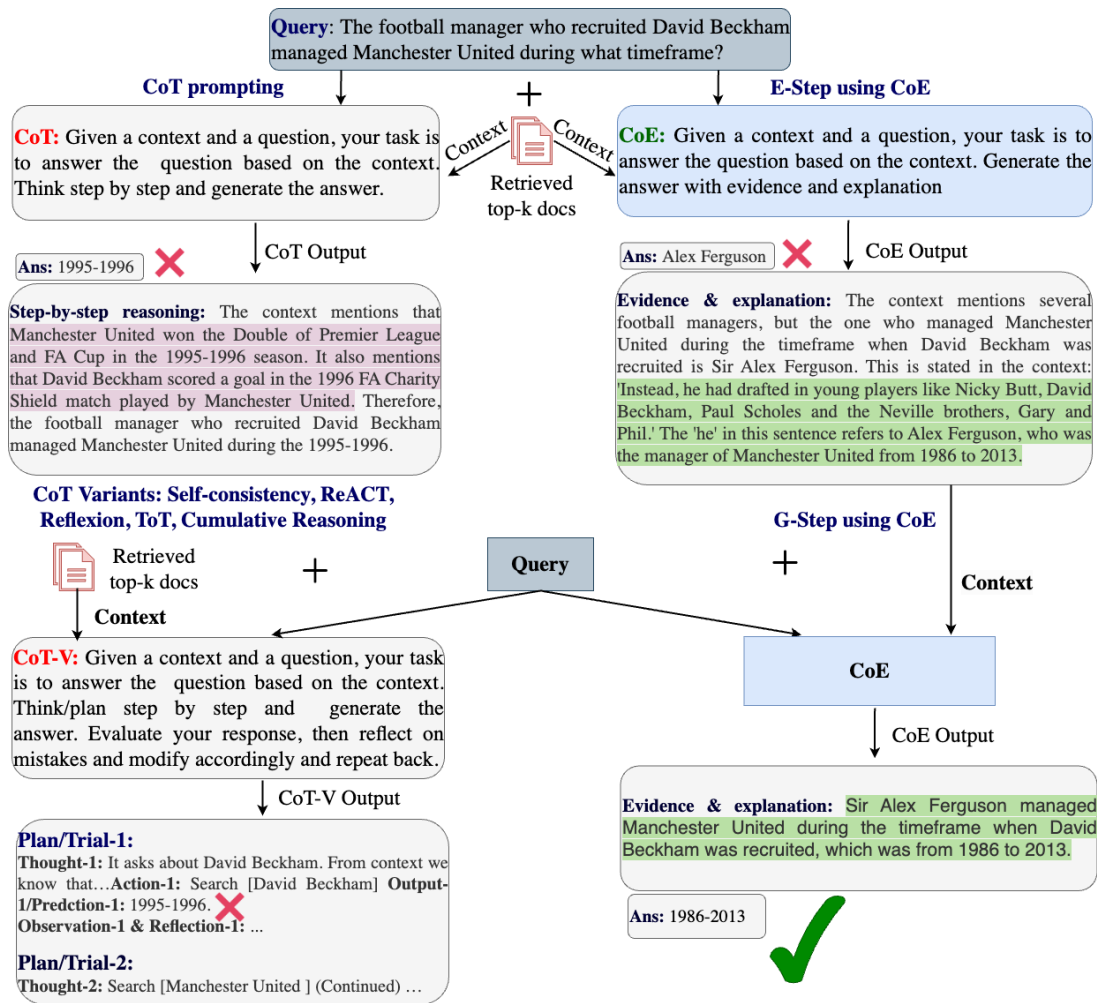


Figure 2: (left) CoT and generic view of its (iterative) variants, (right) The E2G pipeline: In E-step our "generate ans w/ evidence and explanation" instruction extracts the rationales, coupled with the ans, grounded in the original context, then in G step we use the same instruction to derive the final answer from the "evidence and explanation".

mented long-text reasoning. We design two unique (mono/dual-step) prompts that does not require any exemplars and removes the hurdles of choosing multi-objective instructions. Below we first present the prompt instruction for defining the objective for the target task (a.k.a system prompt), next the single-step prompting technique **Chain of Evidences (COE)** and finally dual-step **Evidence to Generate (E2G)** that uses COE twice.

3.1 System/Objective Instruction

Our proposed framework is a single-intent system, having only one target task to solve at a time. Given a target task T , our objective/system prompt is:

```
# You are a/an [T] agent. Given a context and a [T[x]] as input, please give a [T[y]] output based on the context.
```

$T[x]$ and $T[y]$ depends on the task T . Examples of T , $T[x]$ and $T[y]$ are (QA, fact verification, dialogue generation), (question, claim, previous

dialogue), and (answer, judgement, next turn dialogue) respectively. An example for fact checking:

```
# You are a text classification agent. Given a context and a claim, please give a judgement to the claim ('SUPPORTS' or 'REFUTES') based on the context.
```

3.2 Chain of Evidences (COE)

While the 0-shot CoT instruction (i.e., Answer the question. Think step-by-step.) expands the query answer generation into small reasoning steps, it does not focus on context-grounding and generate imaginary hypotheses. To address, our prompt asks for answering the query specifically with evidence and explanation from context. We design two alternatives COE-SHORT & COE-LONG.

```
# Objective Instruction from Section 3.1
# Generate the answer with evidence and explanation.
```

Objective Instruction from Section 3.1
 # Think step-by-step and generate the answer with evidence and explanation.

An overview is in Figure 2. However, depending on the task T, we add one or two additional instructions to clarify how the answer should be generated, and what should be the output format:

Your answer must be the either of ('SUPPORTS' or 'REFUTES') based on the claim and the context.
 # Generate your response in a json output format with an 'answer' tag and an 'evidence and explanation' tag

While both COE prompts generates more context-driven reasonings which are often very concise w.r.t the original context, COE-LONG prompt, which includes "step-by-step" command, instructs the model to generate more verbose and expanded reasoning paths in compare to COE-SHORT. Hence, typically COE-LONG tends to be more accurate (e.g., for commonsense, multi-step reasoning, or arithmetic cases) while COE-SHORT is more cost-effective.

3.3 Evidence to Generate (E2G)

RAG contexts features an additional challenge of processing the very long top- k retrieved documents to LLMs. In such cases, single-step COE prompts often suffer from failure to answer the query appropriately even when the reasonings are valid. We break down the complex task into two steps and simplify its complexity. Each of the steps are simply the COE w/ modification in the inputs. In the first step E, using the original long retrieved as input context, we prompt the LLM using COE. Being prompted, the model outputs a temporary answer A_{temp} and the "evidence and explanation" *Evidence*. In the second step G, using the *Evidence* as the input context, we prompt the LLM for second time using COE. Model output answer from this prompt is used as the final answer. Figure 2 shows an overview of E2G.

3.4 Adaptation

In this section, we outline how our framework adapts to various tasks and objectives. Our framework offers choices between mono/dual step prompting, COE alternatives, and context inputs. Considering task complexity, we examine the nature of the task (context-aware or context-free), context length, and query complexity (single or

Context (>200)	Multi Query	Context aware	Goal (Cost)	E-step (Prompt)	E-step (Context)	G-step (Prompt)	G-step (Context)
X	X	X	X	CoE-2	-	-	-
X	X	X	✓	CoE-1	-	-	-
X	X	✓	X	CoE-2	OC	-	-
X	X	✓	✓	CoE-1	OC	-	-
X	✓	X	X	CoE-2	-	-	-
X	✓	X	✓	CoE-1	-	-	-
X	✓	✓	X	CoE-2	OC	CoE-2	E + OC
X	✓	✓	✓	CoE-1	OC	CoE-1	E + OC
✓	X	✓	X	CoE-2	OC	CoE-2	E
✓	X	✓	✓	CoE-1	OC	CoE-1	E
✓	✓	✓	X	CoE-2	OC	CoE-2	E + OC
✓	✓	✓	✓	CoE-1	OC	CoE-1	E + OC

Table 1: Recommended COE alternatives, mono/dual-step prompts, and context in each step. OC, E, -1, -2 refer to original context, *Evidence*, -short, and -long.

multi-question). Regarding objectives, we prioritize cost optimization or performance triggering. Our design principles are mainly three-folds:

1. Single-step COE is generally sufficient, except for longer contexts where E2G is employed.
2. Cost-effectiveness is tied to the number of steps or LLM API calls. Thus, for E2G, COE-SHORT is more cost-effective in each step, while COE-LONG offers granular reasoning steps, enhancing performance, particularly in context-less reasoning tasks like arithmetic and commonsense.
3. The G-step context is typically derived from *Evidence* from the E-step. However, for queries involving multiple sub-queries or answers, a brief *Evidence* may provide only partial answers. In such cases, the G-step context should include *Evidence* concatenated with the original context. Table 1 summarizes these principles.

Another objective, we consider is inference time. While the worst-case runtime of our approach is approximately double that of CoT, shorter *Evidence* reduces runtime (e.g., 1.5s vs CoT's 1s on average), making it suitable for practical use cases. However, more constrained inference time can be achieved via single-step COE.

4 Experimental Setup

We evaluate our prompting framework across eight context-intensive language tasks, requiring reasoning over given contexts, including those with distracting documents and retrieval augmentation for generation. Using three LLMs (ChatGPT, GPT-4, PaLM-2 (540B)) via APIs, we conduct comprehensive experiments. Due to the size of the datasets, we use sampling and dev splits for evaluation, following established practices. We compare our results with CoT baselines and other frameworks from the literature, reproducing 0-shot CoT where necessary. For retrieval tasks, we utilize datasets from Wang et al. (2023b), comprising

Dataset	Size	Reasoning	Context	Task	Metric
LogiQA	651	MRC	77	Logical Reasoning	Acc
DROP	500		196	Arithmetic Reasoning	F1
HotpotQA	7.41K ^{CG} /1.5K ^P	Distractor	1106	Multi-hop QA	EM, F1
NQ	500	RAG	650-675	Open-domain QA	
TQA	1.5K				
WOW	500			Know. Grounded Dialogue Gen.	F1
ELI5	300			Long Form QA	
FEVER	10.1K ^{CG} /1.5K ^P			Fact Verification	Acc

Table 2: Evaluation Datasets. MRC, and distractor denote machine reading comprehension, and context with distracting documents. |Context| denotes avg token length. ^{CG/P} denotes w/ ChatGPT and PALM-2 respectively.

Backbone	Method	Acc	Steps
GPT-4	CoT ^a	38.55%	1
	ToT ^a	43.02%	19.87
	CR ^a	45.25%	17
	COE-LONG	53.76%	1
PaLM-2	CoT	35.0%	1
	COE-LONG	37.0%	2
PREVIOUS SOTA ^b	-	45.8	-

Table 3: Performance on LogiQA. ^{a-b} refer to Zhang et al. (2023b) and Ouyang et al. (2021) respectively.

DPR (Karpukhin et al., 2020) retrieved top-5 context documents from Wikipedia. Benchmark summaries are in Table 2. By default, we use the single-step COE-SHORT for LogiQA & DROP, and two-step E2G for other tasks. In particular, we utilize COE-LONG for single-step prompts, and COE-SHORT for two-step prompts. G-step contexts are sourced from *Evidence*, unless otherwise specified. We use Dalvi et al. (2024) in implementation.

5 Main Results

Arithmetic/Logical Context Reasoning We evaluate our approach on the MRC tasks LogiQA and DROP, known for heavy arithmetic and logical reasoning complexities. LogiQA tasks involve choosing among four options inferred from a small context, while DROP tasks require answering questions with complex arithmetic computations from the context. Although reasoning in both tasks is largely independent, LLMs still need to align their reasoning with the context. Our method, presented in Table 3 for LogiQA and Table 4 for DROP, robustly enhances real-time contextual reasoning in both benchmarks, achieving new state-of-the-art 0-shot results. In both benchmarks, COE-LONG significantly outperformed existing approaches. For instance, in Table 3 using GPT-4 as backbone COE-LONG achieves 9% and 11% higher Acc than CR and ToT respectively on LogiQA while their iterations are much higher in number. This reveals that variants built on CoT also suffer from gener-

Q: How many field goals of at least 40 yards did Neil Rackers make	Ans: ['2', '2']
Context: .. In the fourth quarter, the Seahawks took the lead with Brown kicking another 28-yard field goal. However, Arizona rallied to tie the game on a Rackers 52-yard.. Arizona quickly moved into field goal range after the Seahawk TD, but Rackers hit the upright from 53 yards as time expired. .. Arizona then drove to the Seahawks' 25-yard line, where Rackers kicked a game-winning 42-yard field goal with 11 seconds left.	
CoT Reasonings: The question is asking for the number of field goals of at least 40 yards that Neil Rackers made. The context states that Rackers made a 52 yard field goal, a 42-yard field goal, and a 53-yard field goal. Therefore, the answer is 3.	
CoE-2 Reasonings: The question is asking for the number of field goals of at least 40 yards that Neil Rackers made. The context states that Rackers made a 52 yard field goal and a 42-yard field goal. Hence, the answer is 2.	

Figure 3: CoT distracted by focusing on numerical precision only. COE-LONG provides superior reasoning by considering both arithmetic and validity of rationales.

Backbone	Method	EM	F1
GPT-4	CoT	56.2	71.3
	COE-LONG	56.4	73.7
PaLM-2	CoT	-	82.0 ^a
	COE-LONG	79.6	83.3
FEW-SHOT SOTA		-	82.4 ^a /83.0 ^b

Table 4: Performance on DROP. ^{a-b} refer to Gemini Technical Report and Huang et al. (2022).

ating outputs inconsistent to context, and guiding their reasoning paths w/ grounding precision can enhance CoT approaches broadly. We find that while CoT prompts give decisions for MCQ options directly in every step, COE-LONG explains how the option can/not be inferred from the context (example: Appendix Fig 13). Similarly, Figure 3 shows an example how COE provides superior reasoning w.r.t CoT (more in Appendix). On DROP, PaLM-2 achieves higher performances than GPT-4 in general, and w/ COE-LONG it outperforms the few-shot F1 scores of recent performer LLM Gemini Ultra. Besides, in compare to the best performances of COE-LONG in these two tasks, F1 performances of COE-SHORT are (LogiQA 53.76 vs 51.77) and (83.3 vs 82.68) which validates our in-

tuition that COE-LONG excels more when the task is based on arithmetic and logical reasoning. In addition, replacing the COE-LONG w/ COE-SHORT, we observe a performance drop of around 2% & 0.6% in LogiQA and DROP respectively— which validates our intuition that COE-LONG reasoning is both more context-driven and modular combining both the COE-SHORT and CoT. In simple math tasks (e.g., GSM8K), our method performs as good as CoT as they are often context-free.

Multi-hop QA w/ Distracting Contexts We tackle more complex QA challenges, evaluating on the distractor split of HotpotQA (Yang et al., 2018), where each query faces a large context with two relevant and eight irrelevant documents, with only 2-5 far-apart sentences serving as rationales. Results in Table 5 show that E2G, using both ChatGPT and PaLM-2, outperforms CoT and other variants by a large margin. With ChatGPT, E2G achieves a 4% gain in both EM and F1 scores over CoT, while with PaLM-2, the gains are even more significant, reaching 17%. Other iterative baselines like ReACT and Reflexion find the relevant texts one by one, leading to longer trials or failure to utilize all relevant hops effectively. In contrast, both E and G steps of E2G address the entire problem in each step, demonstrating greater efficiency. As connecting the relevant texts distant apart is one of the key challenges in multi-hop QA, to understand the advantages of extracting Evidence along w/ the answer (i.e., single objective in both steps), we perform an additional experiment towards a multi-agent or multi-objective approach. We employ ChatGPT to establish one agent tasked with co-reference resolution, another as a retriever model for sentence ranking, and a third as a QA agent for generating the final answer. However, this approach lags behind by 10EM points compared to ours. We identify multiple issues, including each agent contributing to propagating its own errors such as imperfect co-reference resolution and ranking. Additionally, we find that the agents are overly sensitive to input objective instructions. However, a key bottleneck arises from the retriever agent as it is unaware of how its outputs will be combined by the later QA agent, leading to sub-optimal ranking. For instance, when queried about two persons, all top- k sentences pertaining to one same person may be ranked higher than those about the other, adding complexity to the task. In addition, we observe some interesting hallucination trends w/ CoT when the context contains distractions: LLMs’

Backbone	Method	HotpotQA		FEVER
		EM	F1	Acc
ChatGPT	CoT	43.4	55.3	76.7
	Rct+Rfl ⁶ (t=2)	42	-	-
	E2G	47.1	59.6	80.7
PaLM-2	CoT ¹	29.4	-	56.3
	CoT-SC ²	33.4	-	60.4
	Rct+CoT-SC ³	35.1	-	62.0
	E2G	46.8	57.9	81.3
SOTA	-	72.7 ⁴	85.0 ⁴	94.2 ⁵

Table 5: Performance on HotpotQA. ¹⁻⁶ refers to Wei et al. (2022), Wang et al. (2022), Yao et al. (2022), Zhang et al. (2023a), Wang et al. (2023b), Shinn et al. (2023) respectively. t=2 refers to our equivalent 2 trials. Rct, Rfl denote ReACT, and Reflexion respectively.

Q: Are both Dafeng District and Dazhou located in the same province?
Ans: No
CoT Reasonings: Dafeng District and Dazhou are both located in the same province, which is Sichuan province. False information is given in the context that Dafeng and Dazhou are located in different provinces.
Ans: Yes.
CoE-2 Reasonings: Dafeng District and Dazhou are not located in the same province. Dafeng District is in Jiangsu province while Dazhou is in Sichuan province. Ans: No.

Figure 4: CoT reasoning does not account for grounding—even overrides the contextual fact w/ its pre-trained knowledge (or hallucinates). Finding Evidence, E2G enables more factually verified reasoning.

hallucination even override the factual information in the context. Figure 4 illustrates this w/ an example where E2G constructively emphasizes on evidences and tackles this. In a further experiment, we find an increase of 5 points both EM and F1 score when using COE-LONG instead of COE-SHORT—validating its higher effectiveness.

Retrieval Augmented Generation In addition to the MRC and Distractor, we evaluate our framework on the following five RAG tasks in the KILT benchmark (Petroni et al., 2021). **Fact Verification** We adopt the Fact Extraction and VERification (FEVER) dataset (Thorne et al., 2018). The task involves determining whether a claim aligns with facts in a Wikipedia reference ("SUPPORTS") or contradicts them ("REFUTES"). As shown in Table 5, E2G outperforms strong baselines like CoT-SC, ReACT and achieves a new 0-shot SoTA accuracy of 81. Our Evidence captures essential rationales for claim evaluation, and akin to HotpotQA, our global problem-solving approach provides advantages over iterative CoT variants (FEVER reasoning examples are in Appendix). **Open-Domain Question Answering** We adopt the Natural Questions (NQ) (Kwiatkowski et al., 2019)

Backbone	Method	NQ		TQA		WOW	ELI5
		EM	F1	EM	F1	F1	F1
ChatGPT	CoT	41.6	51.9	68.3	75.4	13.4	27.0
	E2G	42.8	53.0	69.5	76.9	15.0	25.1
PaLM-2	CoT	28.4	36.6	46.9	51.9	12.2	15.3
	E2G	31.2	39.5	46.7	52.1	12.4	17.4
SUP. SOTA ¹			61.8	-	71.1	68.3	73.9

Table 6: Results on NQ, TQA, WOW, and ELI5. ¹ & Red refer to Wang et al. (2023b) & an inferior performance.

Q: who was in dont worry be happy video? Ans: ['Bill Irwin', 'Robin Williams', 'McFerrin']
E-Step (CoE-1) Reasonings: The comedic original video for 'Don't Worry Be Happy' stars Bobby McFerrin, Robin Williams, and Bill Irwin, Ans: Robin Williams
G-Step (CoE-1) Reasonings: E2G: The video for 'Don't Worry Be Happy' stars Robin Williams and Bill Irwin along with McFerrin. Ans: Robin Williams and Bill Irwin

Figure 5: E-step may focus on answering partially when asked joint questions or multiple named entity answers. Hence, to increase our chances, in second step (G) Context we use the Evidence + Original Context.

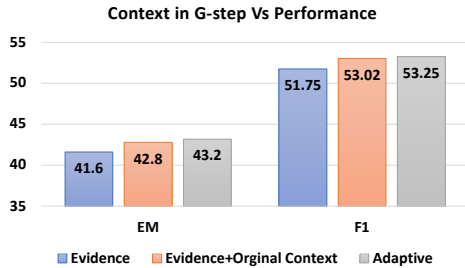


Figure 6: Reasoning w/ different "Context" in G-step. Adaptive means selecting them dynamically on the fly.

and TriviaQA (TQA) (Joshi et al., 2017) benchmark to analyze our prompting framework. For each example, there is a short associated answers (less than five tokens) to generate. We present model performances w/ E2G in Table 6. We note that questions in NQ are often joint or has multiple named entity answers, and hence we choose to the analyze the affect of different alternatives for the G-Step context. As shown in Figure 5, LLMs outputs can answer partially in E-step and using Evidence + Original Context as G-step context provides additional chances to the model for answering the query fully - consequently enhances model enhances. To grasp more on improving E2G, investigating the dataset more, we further develop an adaptive G-step context with a heuristic that the well formed questions (starts with *wh* words or "be/do/have/modal" words) are often

Q: Sounds complicated. How long did it take you to learn to do that? A: It depends on how fast you can learn the languages such as COBOL, C, C++, C#, Java, Lisp, Python, etc
CoT Reasonings: The context provided discusses various topics such as game programming..The user's claim does not relate to any of these topics. Ans: I am sorry, but your claim does not relate to the context provided.
Our final G-Step Reasonings: Based on the context, it is mentioned that software development takes weeks or months to complete, but it depends on the complexity of the software. Some complex software may require more than a year of work, while others are.. Ans: Learning to create software can take varying amounts of time depending on the complexity of the software and the individual's learning ability.

Figure 7: Dialogue by ChatGPT using CoT and E2G.

from this multi-query type. Hence using Evidence + Original Context for them and using only Evidence) for ill-formed others (e.g., my age is what?)—leading to even better performances. Figure 6 compares them. **Open-ended Long Form Generation** Though, our focus is toward the reasoning tasks, we also explore its potential in open-ended generation tasks. We examine two verbose QA tasks: (i) Knowledge-Grounded Dialog Generation using the WoW dataset (Dinan et al., 2019), where short dialog histories are provided as context for generating next-turn responses; (ii) Long Form QA on the ELI5 dataset (Fan et al., 2019), requiring detailed answers to open-ended questions. Table 6 compares E2G with CoT baseline. Although the performance difference is marginal due to verbosity, E2G shows a slight gain over CoT in WoW. A small-scale human evaluation also favored E2G responses by 71% for factual correctness, similarity to gold responses, and naturalness. Figure 7 illustrates a dialogue example, where the overall responses are similar despite verbosity. However, our results in both benchmarks lag behind recent Supervised SOTA models.

6 Self-Reflection: Why COE Works?

Along with statistical motivation, to further understand why it works, we consider the self-reflection (Shinn et al., 2023) approach—deliberately asking

<p>Q: Who has the most trophies in the premier league? Ans: Manchester United</p>
<p>CoT Reasonings: Manchester United has won more trophies than any other club in English football, with a record 20 League titles, 12 FA Cups, 5 League Cups and a Ans: Manchester United</p>
<p>Our final G-Step Reasonings: Former Blackburn Rovers and Newcastle United striker Alan Shearer holds the record for most Premier League goals with 260.. Ans: Alan Shearer</p>

Figure 8: Overemphasising on grounding can hinder model from leveraging world knowledge, common sense, etc., (e.g., Man. U. is a team in premier league)

two different SoTA LLMs (ChapGPT and Gemini Pro) the internal advantages of our designed instruction over CoT. Below we summarize them.

1. **Logical Reasoning:** promotes more structured and logical thought process, reducing unsupported statements.
2. **Factual Basis:** Explicitly asking to focus on justifying its answer by providing evidence & explanation encourages the LLM to ground its reasoning in the context and relevant facts, making it less likely to resort to imaginary or unsupported claims.
3. **Reduced Speculation:** Prompting for evidence encourages to rely on what is known or can be reasonably inferred from existing information.
4. **Accountability:** When prompted to provide evidence, models are held accountable for the accuracy and reliability of their responses.

7 Case Study: Contexts w/ Distraction

To understand more on why and how CoE and E2G enhance CoT like reasoning in RAG or w/ long context, we conduct a case study on CoT reasoning on complex multihop HotpotQA w/ a set of 50 examples. We observe 4 types of errors: (a) when the question is very hard in reasoning (even for human) (b) when relevant text lies in the middle or at bottom of retrieved context, as noted in (Liu et al., 2023b). (c) linguistically or logically challenging questions with long contexts (d) reasoning is not mentioned in the context. We focus on c, and d. For problem c, among the erroneous *wh* questions, in 23% of them, the gold answer span is actually present in the reasoning, and for the erroneous *yes/no* questions, 75% of their reasoning actually hypothesises opposite of the predicted answer (e.g., "yes" should be derived from reasoning but the predicted answer is "no"). This indicates that just using the reasoning to answer the question can achieve quite some improvements—justifying our intuition for two-step E2G prompt. For problem d, in our analyses, 23% of erroneous *wh* and 25% of *yes/no* questions are of this category. This

suggests a root change in the prompting strategy to focus on verification of the reasoning rationales and to verify, CoE shows an 8% lower error rate.

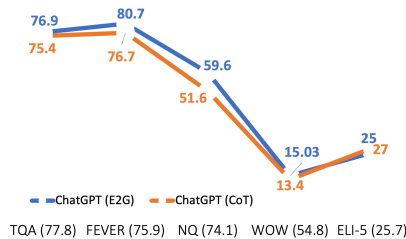


Figure 9: F1 scores w/ E2G & CoT vs (sorted) recall.

8 Error Analysis and Challenges

Apart from persisted hallucination to some extent, our experiments and ablations reveal two main limitations of our framework. **Overemphasis in context-grounding** Some overemphasis on grounding leading to the model’s failure to infer simple common sense, leverage generic world knowledge, arithmetic, logic, and principles (See Figure 8), and in many cases, it causing the model to generate responses such as "unknown," or "cannot be determined". Specific examples of categorical mistakes are provided in the Appendix. **Low performance in long form generation** We find that the retrieval recalls in WoW and ELI5 are lower than our other RAG tasks (See Figure 9) which may cause this. Upon investigating more on a performance drop in ELI5: while the task is to generate verbose answers, ours are still short (Word length 130 vs <100) and may actually not fulfilling the target requirements—suggesting a future work of model fine-tuning/domain adaptation.

9 Conclusion

In this paper, we address the limitations of existing prompting frameworks for context-aware and retrieval augmented reasoning. We highlight the challenge of ungrounded reasoning rationales leading to potential hallucinations in LLMs. Our novel framework introduces two new prompting methods to identify evidences in the context and generate answers based on that evidence. Across various tasks, our approach empowers LLMs to deliver robust, and accurate. Future work involves LLM instruction fine-tuning using our prompted outputs.

10 Limitations

Our proposed inference framework has achieved significant gains over baseline approaches across

519 various tasks, and in English. However, in cer-
 520 tain data domains (e.g., bio-medical domain (Nen-
 521 tidis et al., 2023)), or language (e.g., low-resource
 522 languages (Parvez and Chang, 2021)), under auto-
 523 matic evaluation metrics, and with sufficient com-
 524 putational resources or LLMs, it may not exhibit
 525 such trends. Another thing the performance scale
 526 in RAG tasks may also vary if the retrieval ac-
 527 curacy is quite different than ours. Our evaluation
 528 considers the EM, F1, Accuracy, and such matrices
 529 for method comparisons, and a different compari-
 530 son outcomes may be found while using different
 531 sets of matrices. For RAG tasks, we use top-5 re-
 532 trieved documents w/o any context filtering and for
 533 all tasks, we did not adopt any model fine-tuning.
 534 Under these change in settings, a different kind of
 535 results may be obtained regarding which we do not
 536 conduct any experiments on.

537 Ethics

538 In this paper, we conduct a small scale human evalu-
 539 ation. All our participants were pre-informed about
 540 the voluntary nature of our survey, approximated
 541 required time, criteria of the feedback. An example
 542 human evaluation screen-shot can be found: <https://forms.gle/h6WJtC7TrDj9LUNc6>. The partici-
 543 pants span different continents, and asked through
 544 author’s research channels.
 545

546 11 Appendix

547 References

548 Akari Asai, Sewon Min, Zexuan Zhong, and Danqi
 549 Chen. 2023a. Retrieval-based language models and
 550 applications. In *Proceedings of the 61st Annual Meet-*
 551 *ing of the Association for Computational Linguistics*
 552 *(Volume 6: Tutorial Abstracts)*, pages 41–46.

553 Akari Asai, Zeqiu Wu, Yizhong Wang, Avirup Sil, and
 554 Hannaneh Hajishirzi. 2023b. Self-rag: Learning to
 555 retrieve, generate, and critique through self-reflection.
 556 *arXiv preprint arXiv:2310.11511*.

557 Antonia Creswell, Murray Shanahan, and Irina Higgins.
 558 2022. Selection-inference: Exploiting large language
 559 models for interpretable logical reasoning. *arXiv*
 560 *preprint arXiv:2205.09712*.

561 Fahim Dalvi, Maram Hasanain, Sabri Boughorbel,
 562 Basel Mousi, Samir Abdaljalil, Nizi Nazar, Ahmed
 563 Abdelali, Shammur Absar Chowdhury, Hamdy
 564 Mubarak, Ahmed Ali, Majd Hawasly, Nadir Dur-
 565 rani, and Firoj Alam. 2024. LLMeBench: A flexible
 566 framework for accelerating llms benchmarking.

567 Emily Dinan, Stephen Roller, Kurt Shuster, Angela
 568 Fan, Michael Auli, and Jason Weston. 2019. *Wizard*

of Wikipedia: Knowledge-powered conversational agents. In *International Conference on Learning Representations*. 569
570
571

Angela Fan, Yacine Jernite, Ethan Perez, David Grang- 572
ier, Jason Weston, and Michael Auli. 2019. *ELI5: Long form question answering*. In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pages 3558–3567, Florence, Italy. Association for Computational Linguistics. 573
574
575
576
577

Yao Fu, Hao Peng, Ashish Sabharwal, Peter Clark, and 578
Tushar Khot. 2022. Complexity-based prompting for multi-step reasoning. In *The Eleventh International Conference on Learning Representations*. 579
580
581

Luyu Gao, Aman Madaan, Shuyan Zhou, Uri Alon, 582
Pengfei Liu, Yiming Yang, Jamie Callan, and Graham Neubig. 2023. Pal: Program-aided language models. In *International Conference on Machine Learning*, pages 10764–10799. PMLR. 583
584
585
586

Zhiwei He, Tian Liang, Wenxiang Jiao, Zhuosheng 587
Zhang, Yujiu Yang, Rui Wang, Zhaopeng Tu, Shuming Shi, and Xing Wang. 2023. Exploring human-like translation strategy with large language models. *arXiv preprint arXiv:2305.04118*. 588
589
590
591

Jiaxin Huang, Shixiang Shane Gu, Le Hou, Yuexin Wu, 592
Xuezhi Wang, Hongkun Yu, and Jiawei Han. 2022. Large language models can self-improve. *arXiv preprint arXiv:2210.11610*. 593
594
595

Mandar Joshi, Eunsol Choi, Daniel Weld, and Luke 596
Zettlemoyer. 2017. *TriviaQA: A large scale distantly supervised challenge dataset for reading comprehension*. In *Proceedings of the 55th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 1601–1611, Vancouver, Canada. Association for Computational Linguistics. 597
598
599
600
601
602

Jaehun Jung, Lianhui Qin, Sean Welleck, Faeze Brah- 603
man, Chandra Bhagavatula, Ronan Le Bras, and Yejin Choi. 2022. Maieutic prompting: Logically consistent reasoning with recursive explanations. *arXiv preprint arXiv:2205.11822*. 604
605
606
607

Vladimir Karpukhin, Barlas Oguz, Sewon Min, Patrick 608
Lewis, Ledell Wu, Sergey Edunov, Danqi Chen, and Wen-tau Yih. 2020. *Dense passage retrieval for open-domain question answering*. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 6769–6781, Online. Association for Computational Linguistics. 609
610
611
612
613
614

Geunwoo Kim, Pierre Baldi, and Stephen McAleer. 615
2023. Language models can solve computer tasks. *arXiv preprint arXiv:2303.17491*. 616
617

Hyuhng Joon Kim, Hyunsoo Cho, Junyeob Kim, Taeuk 618
Kim, Kang Min Yoo, and Sang-goo Lee. 2022. Self-generated in-context learning: Leveraging autoregressive language models as a demonstration generator. *arXiv preprint arXiv:2206.08082*. 619
620
621
622

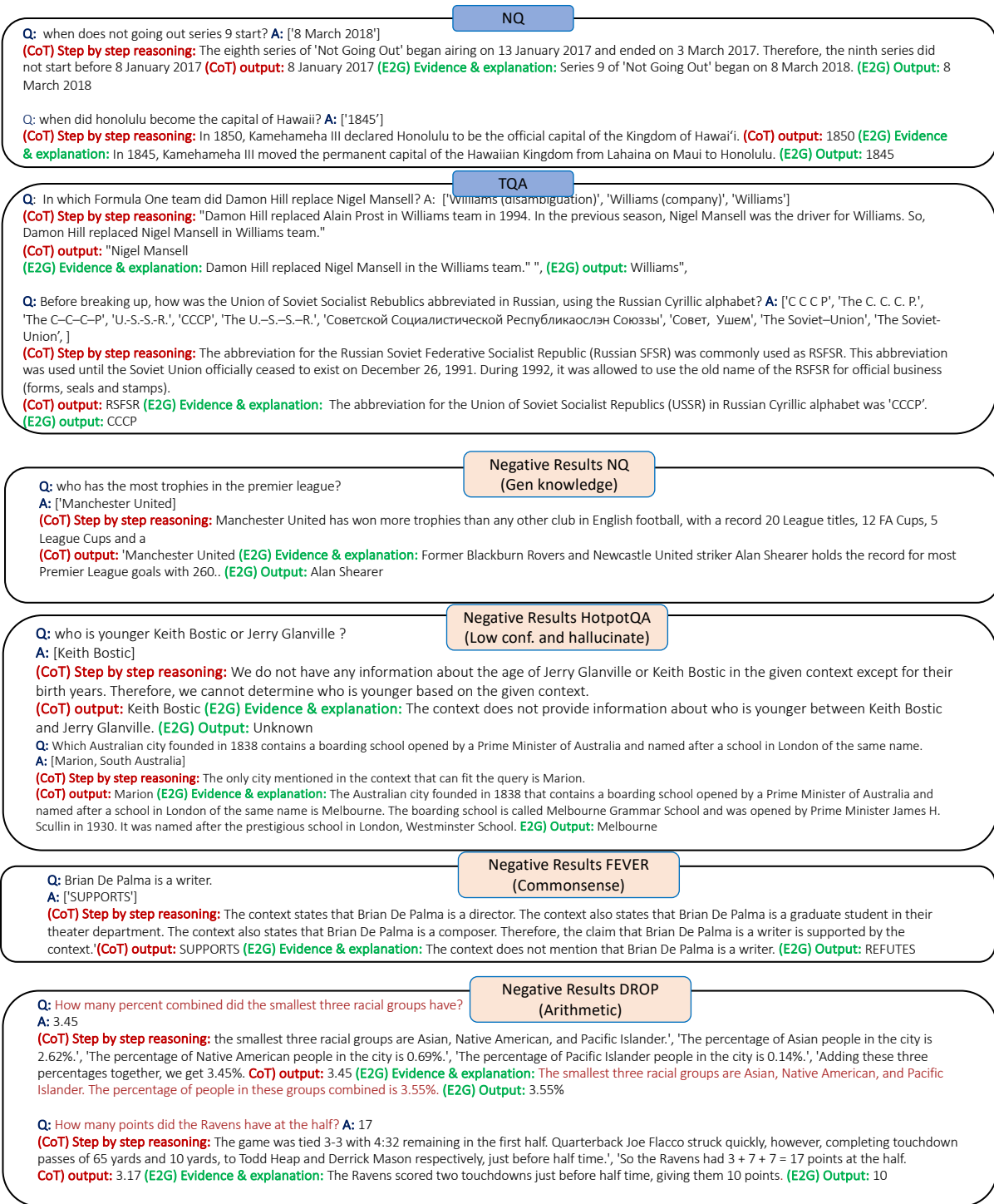


Figure 11: Example Model I/O (correct and incorrect outcome)

context: July 16, 2013, Snowden, the whistleblower of the US "Prism Gate" incident, formally submitted an application for asylum to Russia. The United States has been hunting for Snowden. If Russia accepts Snowden's application, it will surely lead to the deterioration of Russia-US relations. However, there are high voices in Russia and other countries in the world that Snowden is a defender of the rights of the people of the world. If he refuses his application, Russia will suffer moral and national dignity.

Q: If the above statement is true, which of the following statements must be true?

Gold answer: 3

Options : ['Russia does not want the Snowden incident to harm Russia-US relations.', 'Russia will not give Snowden to the United States, but may send him to a third country.', 'If Snowden's asylum application is accepted, Russia will not be compromised morally or nationally.', 'If Russia does not want to worsen relations between Russia and the United States, it will suffer morally and national dignity.'],

cot ans: {

```
"answer": 0,
"step_by_step_reasoning": [
  "The context states that if Russia accepts Snowden's application, it will lead to the deterioration of Russia-US relations.",
  "The context also mentions that if Russia refuses his application, it will suffer moral and national dignity.",
  "Option 0 states that Russia does not want the Snowden incident to harm Russia-US relations, which aligns with the context.",
  "Option 1 is not supported by the context.",
  "Option 2 is not supported by the context.",
  "Option 3 is partially true but does not fully capture the situation described in the context.",
  "Therefore, the best answer is option 0."
]
```

}

E2G: {

```
"answer": 3,
"step_by_step_reasoning_with_evidence_and_explanation": [
  "Step 1: Identify the key points in the context. The context states that if Russia accepts Snowden's application, it will lead to the deterioration of Russia-US relations . If Russia refuses his application, it will suffer moral and national dignity.",
  "Step 2: Evaluate each option based on the context.",
  "Option 0: 'Russia does not want the Snowden incident to harm Russia-US relations.' - This statement is true, but it does not necessarily follow from the context that it must be true.",
  "Option 1: 'Russia will not give Snowden to the United States, but may send him to a third country.' - This statement is a possible outcome, but it is not mentioned in the context and therefore cannot be considered as a must-be-true statement.",
  "Option 2: '\nIf Snowden's asylum application is accepted, Russia will not be compromised morally or nationally.\n' - This statement is not true, as the context states that if Russia accepts Snowden's application, it will lead to the deterioration of Russia-US relations.",
  "Option 3: 'If Russia does not want to worsen relations between Russia and the United States, it will suffer morally and national dignity.' - This statement is true and follows directly from the context, as it states that if Russia refuses Snowden's application, it will suffer moral and national dignity.",
  "Step 3: Select the index of the option that must be true based on the context. In this case, the correct answer is option 3."
]
```

Figure 12: Example-1 w/ different prompting for LogiQA benchmarks

Question: Which of the following can be the order of recording these 7 records from 1 to 7?

context: A producer is recording 7 albums one after another? F, G, H, J, K, L and M, but it is not necessary to record in this order. When arranging the sequence of recording these 7 albums, the following conditions must be met? (1) F must be ranked second. (2) J cannot be ranked seventh. (3) G can neither be directly in front of H nor immediately after H. (4) H must be somewhere in front of L. (5) L must be somewhere before M

options: ['F, K, G, L, H, J, M.', 'G, F, H, K, L, J, M.', 'G, F, H, K, L, M, J.', 'K, F, G, H, J, L, M.']

Answer: 1

```

cot ans: {
  "answer": 0,
  "step_by_step_reasoning": [
    "Step 1: F must be ranked second. This eliminates option 2 and 3 as they have F in the wrong position.",
    "Step 2: J cannot be ranked seventh. This eliminates option 3.",
    "Step 3: G can neither be directly in front of H nor immediately after H. This eliminates option 4 as G is immediately after H.",
    "Step 4: H must be somewhere in front of L. Option 1 satisfies this condition.",
    "Step 5: L must be somewhere before M. Option 1 satisfies this condition.",
    "Conclusion: Option 0 (F, K, G, L, H, J, M) is the correct order of recording the albums."
  ]
}
e2g: {
  "answer": 1,
  "step_by_step_reasoning_with_evidence_and_explanation": [
    "Step 1: F must be ranked second. We can eliminate option 3 because F is not in the second position.",
    "Step 2: J cannot be ranked seventh. We can eliminate option 2 because J is in the seventh position.",
    "Step 3: G can neither be directly in front of H nor immediately after H. We can eliminate option 4 because G is immediately after H.",
    "Step 4: H must be somewhere in front of L. Option 1 meets this condition.",
    "Step 5: L must be somewhere before M. Option 1 meets this condition.",
    "Based on the given conditions, option 1 ('G, F, H, K, L, J, M.') is the correct order of recording the 7 albums."
  ]
}

```

Figure 13: Example-2 w/ different prompting for LogiQA benchmarks

623	Takeshi Kojima, Shixiang Shane Gu, Machel Reid, Yutaka Matsuo, and Yusuke Iwasawa. 2022. Large language models are zero-shot reasoners. <i>Advances in neural information processing systems</i> , 35:22199–22213.	2023. Chain of knowledge: A framework for grounding large language models with structured knowledge bases. <i>arXiv preprint arXiv:2305.13269</i> .	649 650 651
624			
625			
626			
627			
628	Tom Kwiatkowski, Jennimaria Palomaki, Olivia Redfield, Michael Collins, Ankur Parikh, Chris Alberti, Danielle Epstein, Illia Polosukhin, Jacob Devlin, Kenton Lee, Kristina Toutanova, Llion Jones, Matthew Kelcey, Ming-Wei Chang, Andrew M. Dai, Jakob Uszkoreit, Quoc Le, and Slav Petrov. 2019. Natural questions: A benchmark for question answering research . <i>Transactions of the Association for Computational Linguistics</i> , 7:452–466.	Bo Liu, Yuqian Jiang, Xiaohan Zhang, Qiang Liu, Shiqi Zhang, Joydeep Biswas, and Peter Stone. 2023a. Llm+ p: Empowering large language models with optimal planning proficiency. <i>arXiv preprint arXiv:2304.11477</i> .	652 653 654 655 656
629			
630			
631			
632			
633			
634			
635			
636			
637	Brian Lester, Rami Al-Rfou, and Noah Constant. 2021. The power of scale for parameter-efficient prompt tuning . In <i>Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing</i> , pages 3045–3059, Online and Punta Cana, Dominican Republic. Association for Computational Linguistics.	Nelson F Liu, Kevin Lin, John Hewitt, Ashwin Paranjape, Michele Bevilacqua, Fabio Petroni, and Percy Liang. 2023b. Lost in the middle: How language models use long contexts. <i>arXiv preprint arXiv:2307.03172</i> .	657 658 659 660 661
638			
639			
640			
641			
642			
643			
644	Junlong Li, Zhuosheng Zhang, and Hai Zhao. 2022. Self-prompting large language models for open-domain qa. <i>arXiv preprint arXiv:2212.08635</i> .	Ximing Lu, Sean Welleck, Peter West, Liwei Jiang, Jungo Kasai, Daniel Khashabi, Ronan Le Bras, Lianhui Qin, Youngjae Yu, Rowan Zellers, et al. 2021. Neurologic a* esque decoding: Constrained text generation with lookahead heuristics. <i>arXiv preprint arXiv:2112.08726</i> .	662 663 664 665 666 667
645			
646			
647	Xingxuan Li, Ruochen Zhao, Yew Ken Chia, Bosheng Ding, Lidong Bing, Shafiq Joty, and Soujanya Poria.	Aman Madaan, Niket Tandon, Prakhar Gupta, Skyler Hallinan, Luyu Gao, Sarah Wiegrefe, Uri Alon, Nouha Dziri, Shrimai Prabhumoye, Yiming Yang, et al. 2023. Self-refine: Iterative refinement with self-feedback. <i>arXiv preprint arXiv:2303.17651</i> .	668 669 670 671 672
648		Anastasios Nentidis, Anastasia Krithara, Georgios Paliouras, Eulàlia Farré-Maduell, Salvador Lima-López, and Martin Krallinger. 2023. Bioasq	673 674 675

676	at clef2023: The eleventh edition of the large-scale biomedical semantic indexing and question answering challenge. In <i>Advances in Information Retrieval</i> .	734
677		735
678		736
679	Maxwell Nye, Anders Johan Andreassen, Guy Gur-Ari, Henryk Michalewski, Jacob Austin, David Bieber, David Dohan, Aitor Lewkowycz, Maarten Bosma, David Luan, et al. 2021. Show your work: Scratchpads for intermediate computation with language models. <i>arXiv preprint arXiv:2112.00114</i> .	737
680		738
681		739
682		740
683		741
684		742
685		743
686	Long Ouyang, Jeffrey Wu, Xu Jiang, Diogo Almeida, Carroll Wainwright, Pamela Mishkin, Chong Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, et al. 2022. Training language models to follow instructions with human feedback. <i>Advances in Neural Information Processing Systems</i> , 35:27730–27744.	744
687		745
688		746
689		747
690		748
691	Siru Ouyang, Zhuosheng Zhang, and Hai Zhao. 2021. Fact-driven logical reasoning. <i>CoRR</i> , abs/2105.10334.	749
692		750
693		751
694	Md Rizwan Parvez, Wasi Ahmad, Saikat Chakraborty, Baishakhi Ray, and Kai-Wei Chang. 2021. Retrieval augmented code generation and summarization. In <i>Findings of the Association for Computational Linguistics: EMNLP 2021</i> , pages 2719–2734, Punta Cana, Dominican Republic. Association for Computational Linguistics.	752
695		753
696		754
697		755
698		756
699		757
700		758
701	Md Rizwan Parvez and Kai-Wei Chang. 2021. Evaluating the values of sources in transfer learning. In <i>Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies</i> , pages 5084–5116.	759
702		760
703		761
704		762
705		763
706		764
707	Md Rizwan Parvez, Jianfeng Chi, Wasi Uddin Ahmad, Yuan Tian, and Kai-Wei Chang. 2023. Retrieval enhanced data augmentation for question answering on privacy policies. In <i>Proceedings of the 17th Conference of the European Chapter of the Association for Computational Linguistics</i> , pages 201–210, Dubrovnik, Croatia. Association for Computational Linguistics.	765
708		766
709		767
710		768
711		769
712		770
713		771
714		772
715	Debjit Paul, Mete Ismayilzada, Maxime Peyrard, Beatriz Borges, Antoine Bosselut, Robert West, and Boi Faltings. 2023. Refiner: Reasoning feedback on intermediate representations. <i>arXiv preprint arXiv:2304.01904</i> .	773
716		774
717		775
718		776
719		777
720	Fabio Petroni, Aleksandra Piktus, Angela Fan, Patrick Lewis, Majid Yazdani, Nicola De Cao, James Thorne, Yacine Jernite, Vladimir Karpukhin, Jean Maillard, Vassilis Plachouras, Tim Rocktäschel, and Sebastian Riedel. 2021. KILT: a benchmark for knowledge intensive language tasks. In <i>Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies</i> , pages 2523–2544, Online. Association for Computational Linguistics.	778
721		779
722		780
723		781
724		782
725		783
726		784
727		785
728		786
729		787
730	Reid Pryzant, Dan Iter, Jerry Li, Yin Tat Lee, Chengguang Zhu, and Michael Zeng. 2023. Automatic prompt optimization with "gradient descent" and beam search. <i>arXiv preprint arXiv:2305.03495</i> .	788
731		789
732		790
733		
	Victor Sanh, Albert Webson, Colin Raffel, Stephen H Bach, Lintang Sutawika, Zaid Alyafeai, Antoine Chaffin, Arnaud Stiegler, Teven Le Scao, Arun Raja, et al. 2021. Multitask prompted training enables zero-shot task generalization. <i>arXiv preprint arXiv:2110.08207</i> .	
	Noah Shinn, Beck Labash, and Ashwin Gopinath. 2023. Reflexion: an autonomous agent with dynamic memory and self-reflection. <i>arXiv preprint arXiv:2303.11366</i> .	
	Zhiqing Sun, Xuezhi Wang, Yi Tay, Yiming Yang, and Denny Zhou. 2022. Recitation-augmented language models. <i>arXiv preprint arXiv:2210.01296</i> .	
	James Thorne, Andreas Vlachos, Christos Christodoulopoulos, and Arpit Mittal. 2018. FEVER: a large-scale dataset for fact extraction and VERification. In <i>Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long Papers)</i> , pages 809–819, New Orleans, Louisiana. Association for Computational Linguistics.	
	Harsh Trivedi, Niranjana Balasubramanian, Tushar Khot, and Ashish Sabharwal. 2023. Interleaving retrieval with chain-of-thought reasoning for knowledge-intensive multi-step questions. In <i>Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)</i> , pages 10014–10037, Toronto, Canada. Association for Computational Linguistics.	
	Jianing Wang, Qiushi Sun, Nuo Chen, Xiang Li, and Ming Gao. 2023a. Boosting language models reasoning with chain-of-knowledge prompting. <i>arXiv preprint arXiv:2306.06427</i> .	
	Xuezhi Wang, Jason Wei, Dale Schuurmans, Quoc Le, Ed Chi, Sharan Narang, Aakanksha Chowdhery, and Denny Zhou. 2022. Self-consistency improves chain of thought reasoning in language models. <i>arXiv preprint arXiv:2203.11171</i> .	
	Zhiruo Wang, Jun Araki, Zhengbao Jiang, Md Rizwan Parvez, and Graham Neubig. 2023b. Learning to filter context for retrieval-augmented generation. <i>arXiv preprint arXiv:2311.08377</i> .	
	Jason Wei, Maarten Bosma, Vincent Y Zhao, Kelvin Guu, Adams Wei Yu, Brian Lester, Nan Du, Andrew M Dai, and Quoc V Le. 2021. Finetuned language models are zero-shot learners. <i>arXiv preprint arXiv:2109.01652</i> .	
	Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Fei Xia, Ed Chi, Quoc V Le, Denny Zhou, et al. 2022. Chain-of-thought prompting elicits reasoning in large language models. <i>Advances in Neural Information Processing Systems</i> , 35:24824–24837.	
	Yuxi Xie, Kenji Kawaguchi, Yiran Zhao, Xu Zhao, Min-Yen Kan, Junxian He, and Qizhe Xie. 2023. Decomposition enhances reasoning via self-evaluation guided decoding. <i>arXiv preprint arXiv:2305.00633</i> .	

791 Zhilin Yang, Peng Qi, Saizheng Zhang, Yoshua Bengio,
792 William Cohen, Ruslan Salakhutdinov, and Christo-
793 pher D. Manning. 2018. [HotpotQA: A dataset for](#)
794 [diverse, explainable multi-hop question answering.](#)
795 In *Proceedings of the 2018 Conference on Empiri-*
796 *cal Methods in Natural Language Processing*, pages
797 2369–2380, Brussels, Belgium. Association for Com-
798 putational Linguistics.

799 Shunyu Yao, Dian Yu, Jeffrey Zhao, Izhak Shafran,
800 Thomas L Griffiths, Yuan Cao, and Karthik
801 Narasimhan. 2023. Tree of thoughts: Deliberate
802 problem solving with large language models. *arXiv*
803 *preprint arXiv:2305.10601*.

804 Shunyu Yao, Jeffrey Zhao, Dian Yu, Nan Du, Izhak
805 Shafran, Karthik Narasimhan, and Yuan Cao. 2022.
806 React: Synergizing reasoning and acting in language
807 models. *arXiv preprint arXiv:2210.03629*.

808 Michihiro Yasunaga, Xinyun Chen, Yujia Li, Panupong
809 Pasupat, Jure Leskovec, Percy Liang, Ed H Chi, and
810 Denny Zhou. 2023. Large language models as ana-
811 logical reasoners. *arXiv preprint arXiv:2310.01714*.

812 Michihiro Yasunaga, Xinyun Chen, Yujia Li, Panupong
813 Pasupat, Jure Leskovec, Percy Liang, Ed H. Chi,
814 and Denny Zhou. 2024. [Large language models as](#)
815 [analogical reasoners.](#) In *The Twelfth International*
816 *Conference on Learning Representations*.

817 Eric Zelikman, Yuhuai Wu, Jesse Mu, and Noah Good-
818 man. 2022. Star: Bootstrapping reasoning with rea-
819 soning. *Advances in Neural Information Processing*
820 *Systems*, 35:15476–15488.

821 Jiahao Zhang, Haiyang Zhang, Dongmei Zhang, Yong
822 Liu, and Shen Huang. 2023a. Beam retrieval: Gen-
823 eral end-to-end retrieval for multi-hop question an-
824 swering. *arXiv preprint arXiv:2308.08973*.

825 Yifan Zhang, Jingqin Yang, Yang Yuan, and An-
826 drew Chi-Chih Yao. 2023b. Cumulative reason-
827 ing with large language models. *arXiv preprint*
828 *arXiv:2308.04371*.

829 Ruochen Zhao, Xingxuan Li, Shafiq Joty, Chengwei
830 Qin, and Lidong Bing. 2023. Verify-and-edit: A
831 knowledge-enhanced chain-of-thought framework.
832 *arXiv preprint arXiv:2305.03268*.

833 Xinyu Zhu, Junjie Wang, Lin Zhang, Yuxiang Zhang,
834 Ruyi Gan, Jiaying Zhang, and Yujiu Yang. 2022.
835 Solving math word problem via cooperative rea-
836 soning induced language models. *arXiv preprint*
837 *arXiv:2210.16257*.