# Attention Reveals More Than Tokens: Training-Free Long-Context Reasoning with Attention-guided Retrieval

**Yuwei Zhang, Jayanth Srinivasa[2], Gaowen Liu[2], Jingbo Shang[1*]**
University of California, San Diego[1]    Cisco[2]
{yuz163, jshang}@ucsd.edu
{jasriniv, gaoliu}@cisco.com

## Abstract

Large Language Models (LLMs) often exhibit substantially shorter effective context lengths than their claimed capacities, especially when handling complex reasoning tasks that require integrating information from multiple parts of a long context and performing multi-step reasoning. Although Chain-of-Thought (CoT) prompting has shown promise in reducing task complexity, our empirical analysis reveals that it does not fully resolve this limitation. Through controlled experiments, we identify poor recall of implicit facts as the primary cause of failure, which significantly hampers reasoning performance. Interestingly, we observe that the internal attention weights from the generated CoT tokens can effectively ground implicit facts, even when these facts are not explicitly recalled. Building on this insight, we propose a novel training-free algorithm, AT-TRIEVAL, which leverages attention weights to retrieve relevant facts from the long context and incorporates them into the reasoning process. Additionally, we find that selecting context tokens from CoT tokens further improves performance. Our results demonstrate that ATTRIEVAL enhances long-context reasoning capability notably on both synthetic and real-world QA datasets with various models.

## 1 Introduction

Recent advancements in long-context language models have unlocked the ability to process much larger input sequences (Zaheer et al., 2020; Gu and Dao, 2023; Peng et al., 2023b; Chen et al., 2023b,c; Jin et al., 2024; Wang et al., 2024), achieving near perfect recall on retrieval tasks such as *needle-in-a-haystack* (gkamradt, 2023). However, real-world applications—including multi-hop question answering (Yang et al., 2018; Trivedi et al., 2022), document-level reasoning (Mou et al., 2021; Dasigi et al., 2021), and multi-turn conversational
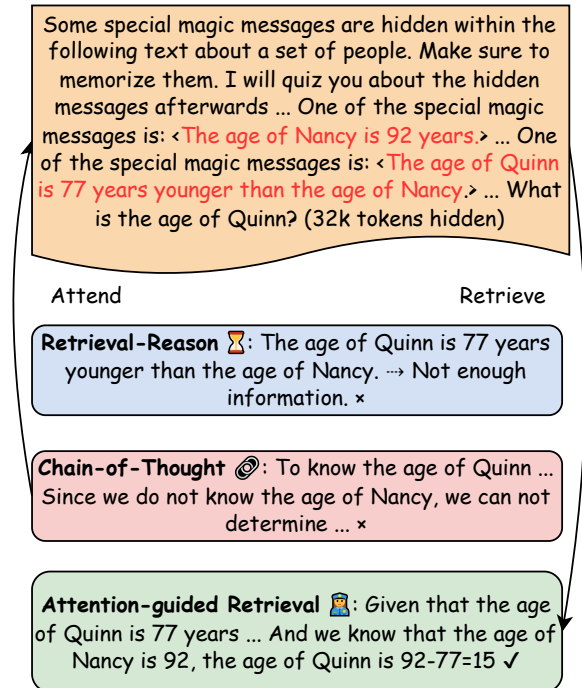
---

\* Corresponding authors.



Figure 1: Both Retrieval-Reason (agentic framework) and Chain-of-Though (CoT) might suffer from poor recall of implicit facts. Our proposed ATTRIEVAL leverage the internal attention weights to resolve this issue.

agents (Wu et al., 2024) demand more than verbatim fact extraction, sometimes requiring aggregating information from scattered evidence into coherent conclusions. While existing models excel at locating explicit statements, their performance degrades significantly as context length increases for tasks requiring reasoning, even when all necessary facts are present in the input (Hsieh et al., 2024; Kuratov et al., 2024; Ling et al., 2025; Bai et al., 2024; Zhang et al., 2024b). This discrepancy reveals a critical gap: strong single-hop retrieval capabilities do not inherently enable robust reasoning.

Chain-of-Thought (CoT) reasoning (Wei et al., 2022) offers a promising framework for complex tasks by decomposing reasoning into retrieval and inference steps that receives few attention in pre-

vious benchmarking. The step-by-step CoT reasoning turns multi-hop questions into single-hop retrieval tasks that are easier to be solved by long-context models. Yet, we observe that even with CoT, performance degrades sharply as context length increases. We hypothesize that this stems from failures to retrieve implicit facts—information critical for reasoning but lacking explicit surface cues as illustrated by Figure 1. To test this, we introduce **Deduction**, a diagnostic benchmark requiring models to (1) retrieve numerical facts from long contexts and (2) perform arithmetic reasoning. By analyzing responses for both fact recall and final accuracy, we find that the performance is mostly bottlenecked by the missed implicit (or second-hop) facts, not faulty arithmetic.

Notably, agentic frameworks have been explored to improve CoT in the literature by explicitly prompt the LLMs to retrieve-then-reason. For instance, Zhang et al. (2024c) proposed a multi-agent framework that distributes the long-context across multiple agents and then aggregates information through model collaboration. Zhang et al. (2024a) proposed an automatic attention steering framework that utilizes prompt-based method to elicit the model to generate useful facts and "steer" the attention weights. Chen et al. (2023a) proposed memory maze that summarizes the long-context into a hierarchical structure and then perform tree search during inference time. However, neither of them solve the implicit fact retrieval problem inherently (Figure 2a) and might introduce laborious prompt engineering efforts or rely on strong close-source LLMs. Furthermore, CoT reasoning inherently outperforms agentic workflows by leveraging LLM's native generation of coherent, self-contained reasoning paths while maintaining computational efficiency and scalability.

In this paper, we first make the observation that the internal attention weights often highlight the overlooked implicit facts, suggesting a disconnect between latent retrieval signals (attention) and explicit generation. Inspired by these findings, we propose Attention-guided Retrieval (ATTRIEVAL), a training-free framework that enhances long-context reasoning without compromising short-context performance or requiring laborious prompt engineering. ATTRIEVAL operates in three key stages: (1) The input context is partitioned into discrete facts, which are then ranked by their attention weights from intermediate CoT tokens. (2) To counter the dominance of "atten-tion sink" tokens, we filter out facts that appear in the top-$k$ attended positions for an excessive proportion of CoT tokens. (3) We introduce a cross-evaluation framework to identify retriever tokens from the generated CoT sequence by measuring the KL-divergence between model predictions with and without the context. The final retrieved facts are reintegrated into the context, enabling the model to reason over both explicit and previously overlooked implicit information.

Our works makes the following contributions:

- We introduce Deduction, a controlled benchmark for long-context reasoning, and identify retrieval failures—particularly for latent facts—as the primary bottleneck in existing methods.
- We demonstrate that attention weights encode latent factual relevance even when generated tokens fail to reference them explicitly, challenging the assumption that token outputs fully reflect model "knowledge".
- ATTRIEVAL provides the first training-free solution that leverages attention patterns to bridge the gap between retrieval and reasoning, achieving state-of-the-art performance across both synthetic and realistic QA benchmarks (e.g., +47% accuracy on Deduction and +11% accuracy on MuSiQue on 32K context length).

## 2 Preliminary

We formally define *long-context reasoning* task in this section.

**Definition 1** (Long-Context Reasoning)**.** *Let $Q$ be a question (e.g., a natural-language query), and let $\hat{I} = \{i_1, i_2, \ldots, i_r\}$ be a set of* informative *(or* relevant*) facts needed to correctly answer $Q$. Let $\hat{N} = \{n_1, n_2, \ldots, n_s\}$ be a set of* noisy *(or irrelevant) facts. Define the* long *context $C$ as the union of these two sets: $C = \hat{I} \cup \hat{N}$. Suppose there is an (ideal) reasoning function $R : \mathcal{Q} \times \mathcal{I} \to \mathcal{A}$, where $\mathcal{Q}$ is the space of all possible questions, $\mathcal{I}$ is the space of all possible informative-fact sets, and $\mathcal{A}$ is the space of all possible answers.*

*The* long-context reasoning problem *is to construct a function $\hat{R} : \mathcal{Q} \times \mathcal{C} \to \mathcal{A}$ that approximates $R$ when presented with the full long context $C$, i.e., $\hat{R}(Q, C) \approx R(Q, \hat{I})$.*

From a probabilistic point of view, long-context reasoning requires the model to be able to "filter out" (or marginalize) the noise $\hat{N}$ in posterior distribution:

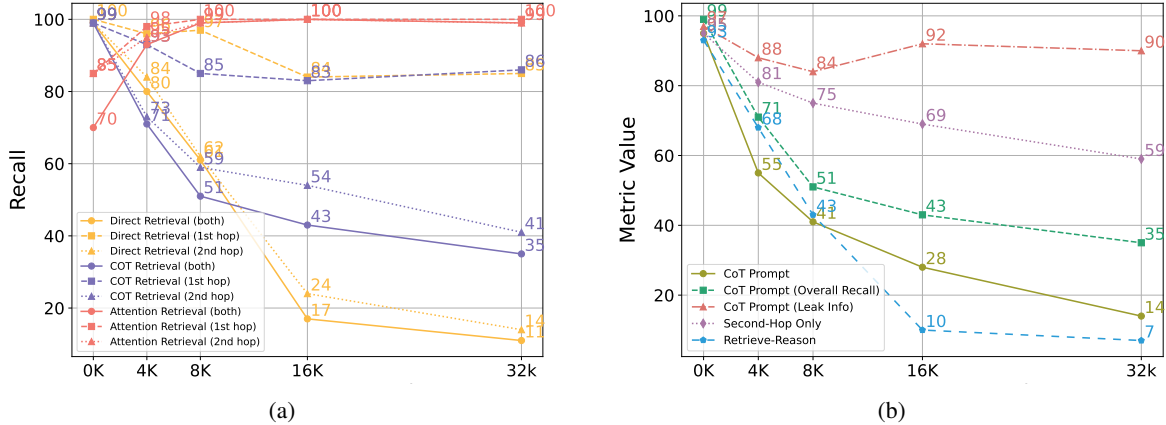$$P(A = a|Q, C) \approx P(A = a|Q, \hat{I}) \quad (1)$$

2

Figure 2: Analysis on CoT tokens, including: (a) recall with various retrieval methods; (b) accuracy with various prompts and questions. See section 3 for more details.

## 3 Analysis on CoT Tokens

A natural strategy for improving long-context reasoning is Chain-of-Thought (CoT) prompting (Wei et al., 2022), which enables models to strategically search through extended contexts (Yu et al., 2023; Li et al., 2024a,c). The generated reasoning chain can decompose long-context reasoning into two subtasks: *retrieval* and *reasoning*. This process can be formulated as follows:

$$P(A, Y|Q, C) = \underbrace{P(Y|Q, C)}_{retrieval} \underbrace{P(A|Y, Q, C)}_{reasoning}$$

(2)

where $Y$ represents retrieved facts during *retrieval* phase. While models can dynamically alternate between retrieval and reasoning to iteratively refine outputs, our experiments in Figure 2b demonstrate that CoT alone fails to mitigate the performance degradation in long-context scenarios. We identify two distinct failure modes: (1) search errors, where retrieved facts $Y$ are incomplete or misaligned with $Q$; or (2) reasoning errors, where the model misapplies logical rules despite accurate retrieval. To dissect these issues, we first quantify the relative impact of each error type through empirical analysis. We then reveal a critical insight: transformers' internal attention mechanisms exhibit stronger grounding to contextually relevant facts compared to explicit CoT-generated retrieval tokens. This finding suggests inherent limitations in relying solely on CoT's discrete search phase for long-context understanding.

### 3.1 Which is the Devil? Search or Reasoning?

To systematically diagnose the interplay between search and reasoning errors, we require a benchmark where both retrieval validity (whether all necessary facts are recalled) and reasoning validity (whether logic is correctly applied) can be unambiguously evaluated. Existing long-context datasets often conflate these two aspects, as their open-ended questions and implicit grounding in context make it difficult to isolate failure modes.

To address this, we introduce **Deduction**, a diagnostic benchmark featuring synthetic reasoning tasks with explicit ground-truth retrieval requirements. Each task embeds a set of atomic facts (*e.g.*, "Nancy's age is 92") within a long, distractor-filled context, followed by a deterministic question (*e.g.*, "What is Quinn's age?") solvable only by recalling all relevant facts (*e.g.*, "Quinn is 77 years younger than Nancy") and applying basic arithmetic. Crucially, our design ensures both controlled retrieval evaluation for both explicit and implicit facts and deterministic reasoning. See Appendix A for details about dataset creation.

**Observations** As illustrated in Figure 2, four key patterns emerge: (1) First-hop recall (retrieving explicit facts like "Nancy's age is 92") remains robust (80–90% across 4K–32K contexts), while second-hop recall (implicit dependencies like "Quinn's age depends on Nancy") drops sharply as sequence length increases, narrowing the gap between overall recall and second-hop recall (Figure 2a). (2) Despite being widely employed in agentic workflows, directly prompting the model retrieve useful information amplify the incomplete retrieve issue compared with a more natural CoT prompt (Figure 2a). (3) Final answer accuracy lags behind recall by 15–20% (Figure 2b), indicating that even when models retrieve partial facts, they might fail to synthesize them into correct answers. (3) Retrieval is the primary bottleneck. When explicitly prompted for second-hop facts ("What is Nancy's

age?"), retrieval success improves by 35% (Figure 2b, Second-Hop Only), confirming that models can reason accurately if retrieval is guaranteed. (4) Appending ground-truth facts post-context (Leak Info) restores 85% of 0K baseline performance (Figure 2b), yet a residual 15% accuracy gap persists, likely due to attention dispersion over long sequences. These results underscore that while reasoning errors occur, the dominant failure mode is retrieval: models struggle to retrieve implicit, interdependent facts from long contexts. The compounding effect of partial retrieval and flawed logic explains the steep performance decline in multi-hop tasks.

## 3.2 Can Attention Weights Retrieve Latent Facts?

While CoT prompting struggles to surface implicit facts through its explicit token generations, we find that the model's internal attention patterns reveal richer evidence of factual grounding. We hypothesize that this discrepancy arises because generated tokens represent high-level discretizations of latent states, potentially obscuring the model's sensitivity to specific input features. By contrast, attention weights provide continuous-valued signals that better preserve these fine-grained associations. This observation motivates our central investigation: *Does the model internally attend to factual evidence that remains implicit in its generations?* Through quantitative analysis of attention patterns (Figure 3), we demonstrate that the model allocates substantial attention to second-hop factual relationships, even when these fail to surface in CoT generations. To formalize this analysis, let $t \in \{1, 2, \ldots, T\}$ denote the positions of the generated tokens and $i \in \{1, 2, \ldots, N\}$ denote the positions of the input tokens. For a given layer $l$, we normalize the attention weights $A_{t,i}^{(l)}$ so that they satisfy $\sum_{i=1}^{N} A_{t,i}^{(l)} = 1$. For a statement spanning input tokens indexed by $I_{\text{stmt}} \subseteq \{1, 2, \ldots, N\}$, we compute the aggregated attention score for each layer and generated token as

$$H_{\text{stmt}}(l, t) = \sum_{i \in I_{\text{stmt}}} A_{t,i}^{(l)}. \qquad (3)$$

Our case study in Figure 3 reveals two key patterns: (1) Early generated tokens exhibit heightened attention to both first-hop and second-hop factual statements (red boxes), despite the CoT ultimately failing to verbalize the latter, and (2) While first-hop attention resurfaces in later tokens, second-hop

---

**Algorithm 1** Attention-Guided Retrieval (ATTRIEVAL)

**Require:** Input context $\mathcal{X}$, generated CoT tokens $\{t_1, \ldots, t_T\}$, layers $\mathcal{L}$, top-$k$ threshold, frequency threshold $\tau$, min tokens $m$, max facts $n$

**Ensure:** Retrieved facts $\mathcal{F}_{\text{retrieved}}$
1: **Stage 1: Multi-Layer Attention Aggregation**
2: **for** each generated token $t \in \{1, \ldots, T\}$ **do**
3:     **for** each input token $i \in \mathcal{X}$ **do**
4:         Compute $\bar{A}_{t,i}$ via Equation 4
5:     **end for**
6: **end for**
7: **Stage 2: Common Facts Filtering**
8: Segment $\mathcal{X}$ into facts $\{c\}$ via punctuation
9: **for** each generated token $t$ **do**
10:     Identify top-$k$ tokens: $\mathcal{T}_t$ via Equation 5
11: **end for**
12: **for** each fact $c$ **do**
13:     Compute frequency: $f(c)$ via Equation 6
14: **end for**
15: Filter sinks: $\mathcal{F}_{\text{filtered}} \leftarrow \{c : f(c) < \tau\}$
16: **Stage 3: Fact Scoring & Selection**
17: **for** each fact $c \in \mathcal{F}_{\text{filtered}}$ **do**
18:     Aggregate fact score: $s(c)$ via Equation 7
19: **end for**
20: Sort facts by $s(c)$, filter length $\geq m$ tokens
21: Return $\mathcal{F}_{\text{retrieved}} \leftarrow$ top-$n$ facts

---

attention remains suppressed. We further show in Figure 5 that the rankings of attention weights spent on the statement tokens are usually high. Notably, second-hop statements achieve comparable ranking positions to first-hop statements during initial generated tokens. Nonetheless, it remains challenging to extract these statements from the overall input, as high attention weights are also assigned to other irrelevant tokens, such as those at the beginning of the prompt and the most recent tokens (Xiao et al., 2023; Han et al., 2023).

## 4 Methodology

Inspired by the previous observations that the attention weights perform better at grounding in the long-context setting, we now introduce a novel algorithm that improves long-context reasoning without any additional training or extensive prompt engineering. Intuitively, the proposed algorithm performs attention-based retrieval based on the generated CoT tokens, and then incorporate them for

Figure 3: Proportion of attention from generated tokens to the input prompt across layers.

reasoning.

Formally, given a pre-defined set of layers $\mathcal{L}$, we first aggregate the attention over heads and layers,

$$\bar{A}_{t,i} = \frac{1}{|\mathcal{L}|} \sum_{l \in \mathcal{L}} \left( \frac{1}{H} \sum_{h=1}^{H} A_{t,i}^{(l,h)} \right). \quad (4)$$

The input sequence is segmented into discrete facts $\{c\}$ based on punctuations. Each input token $i$ is mapped to its corresponding fact $c(i)$. For each generated token $t$, we identify the top-$k$ input tokens with the highest aggregated attention scores, denoting their indices by the set $\mathcal{T}_t$.

$$\mathcal{T}_t = \arg \operatorname{top-}_i k(\bar{A}_{t,i}) \quad (5)$$

We then define the frequency of a fact $c$ as

$$f(c) = \frac{1}{T} \sum_{t=1}^{T} \mathbb{I}\left\{ c \in \{c(i) : i \in \mathcal{T}_t\} \right\}, \quad (6)$$

where $\mathbb{I}\{\cdot\}$ is the indicator function. Facts with $f(c) \geq \tau$ (where $\tau$ is a threshold) are filtered as potential attention sinks (Xiao et al., 2023)—frequently attended tokens that provide little informational value. For remaining facts, we compute a relevance score by first averaging the aggregated attention over all generated tokens for each input token, and then averaging these scores over all tokens belonging to fact $c$. Concretely, if $I_c = \{i : c(i) = c\}$, then the fact score is defined as

$$s(c) = \frac{1}{|I_c|} \sum_{i \in I_c} \left( \frac{1}{T} \sum_{t=1}^{T} \bar{A}_{t,i} \right). \quad (7)$$

These scores $s(c)$ provide a measure of relevance between facts and generated tokens. We then take the top-$n$ facts while filtering out those with less than $m$ tokens. These facts are then incorporated

into the context for generating the final answers. We can also select a subset of tokens to calculate final score as illustrated in the next paragraph. The prompt we used for this procedure requires minimal design. See Appendix D for prompts and Algorithm 1 for algorithm procedure.

**Cross-Evaluation for Token Selection.** As shown in Figure 3, we observe that there exist two kinds of tokens: *retriever tokens* that cite the context and spread more attention on the ground truth facts; *reasoner tokens* that focuses on reasoning with previously cited context. We hypothesize that the *retriever tokens* might be better at retrieving relevant information from the context. Therefore, in this section, we propose a simple method to automatically detect *retriever tokens* via cross-evaluation (shown in Algorithm 2). Given context $C$ and question $Q$, the model evaluates the generate CoT tokens with both a long prompt $\mathcal{P}_L(C, Q)$ and a short prompt $\mathcal{P}_S(Q)$. The token-wise KL divergence $D_{KL}(P_L^{(t)} \| P_S^{(t)})$ identifies tokens where contextual information most significantly alters their predictions. We then simply take the top-$s$ tokens as the selected *retriever tokens* for Equation 7.

## 5 Main Results

### 5.1 Experimental Setting

In this paper, we mainly study three open-source models `meta-llama/Llama-3.2-3B-Instruct`, `meta-llama/Llama-3.1-8B-Instruct` and `Qwen/Qwen2.5-3B-Instruct`. In our preliminary experiments, we found that attention weights from higher layers can better ground to the context, we thus always choose the last $1/4$ layers as $\mathcal{L}$. We fix $k = 50$, $\tau = 0.99$ for filtering common facts, and minimum fact length $m = 3$. We always choose $n = 10$ facts for all the experiments. For token selection strategies, we consistently use $s = 10$.

5

| Model | Method | 0K | 4K | 8K | 16K | 32K | Overall |
|---|---|---|---|---|---|---|---|
| | | | | Deduction | | | |
| Llama-3.2-3B-Instruct | CoT | 95 | 55 | 41 | 28 | 14 | 47 |
| | ATTRIEVAL | 97 | 76 | 75 | 63 | 57 | 74 |
| | ATTRIEVAL-kl | 96 | 79 | 80 | 77 | 61 | 79 |
| Llama-3.1-8B-Instruct | CoT | 99 | 96 | 93 | 74 | 70 | 86 |
| | ATTRIEVAL | 99 | 97 | 99 | 92 | 81 | 94 |
| | ATTRIEVAL-kl | 100 | 98 | 96 | 91 | 83 | 94 |
| Qwen2.5-3B-Instruct | CoT | 96 | 46 | 50 | 34 | 28 | 51 |
| | ATTRIEVAL | 99 | 59 | 60 | 56 | 40 | 63 |
| | ATTRIEVAL-kl | 99 | 55 | 55 | 48 | 56 | 63 |
| | | | | MuSiQue | | | |
| Llama-3.2-3B-Instruct | CoT | 78 | - | 25 | 24 | 18 | 36 |
| | ATTRIEVAL | 71 | - | 38 | 33 | 23 | 41 |
| | ATTRIEVAL-kl | 74 | - | 43 | 33 | 29 | 45 |
| Llama-3.1-8B-Instruct | CoT | 78 | - | 61 | 45 | 29 | 53 |
| | ATTRIEVAL | 87 | - | 64 | 55 | 45 | 63 |
| | ATTRIEVAL-kl | 79 | - | 68 | 58 | 41 | 62 |
| Qwen2.5-3B-Instruct | CoT | 69 | - | 42 | 33 | 23 | 42 |
| | ATTRIEVAL | 72 | - | 48 | 39 | 18 | 44 |
| | ATTRIEVAL-kl | 64 | - | 45 | 36 | 22 | 42 |
| | | | | HotpotQA | | | |
| Llama-3.2-3B-Instruct | CoT | 71 | 68 | 62 | 57 | 58 | 63 |
| | ATTRIEVAL | 73 | 58 | 66 | 61 | 59 | 63 |
| | ATTRIEVAL-kl | 68 | 67 | 66 | 61 | 55 | 63 |
| Llama-3.1-8B-Instruct | CoT | 71 | 71 | 71 | 70 | 68 | 70 |
| | ATTRIEVAL | 74 | 72 | 69 | 72 | 69 | 71 |
| | ATTRIEVAL-kl | 69 | 71 | 70 | 70 | 67 | 69 |
| Qwen2.5-3B-Instruct | CoT | 64 | 61 | 60 | 54 | 54 | 59 |
| | ATTRIEVAL | 66 | 58 | 52 | 58 | 48 | 56 |
| | ATTRIEVAL-kl | 63 | 60 | 58 | 55 | 50 | 57 |

Table 1: Main results with color annotations. Green numbers exceed CoT; red numbers are lower than CoT. For MuSiQue, the context already exceeds 4k tokens.

| Tokens From | Select Strategy | 0K | 4K | 8K | 16K | 32K | Overall |
|---|---|---|---|---|---|---|---|
| CoT | None | 95 | 55 | 41 | 28 | 14 | 47 |
| CoT | All | 97 | 76 | 75 | 63 | 57 | 74 |
| Random | All | **99** | 61 | 60 | 55 | 48 | 65 |
| CoT | First-s | 98 | 66 | 56 | 49 | 39 | 62 |
| CoT | Random-s | 95 | 76 | 69 | 62 | 60 | 72 |
| CoT | KL top-s | 96 | **79** | **80** | **77** | **61** | **79** |

Table 2: Analysis on the generated tokens used in ATTRIEVAL to calculate attention matrices and the retriever token selection strategy. Studied dataset and model are Deduction and `meta-llama/Llama-3.2-3B-Instruct`.

All the experiments are done on a single A100 GPU.

## 5.2 Evaluation Dataset

We evaluate on both synthetic and realistic QA datasets. For synthetic QA, we evaluate with Deduction dataset with 2 main entities and 6 distraction entities. For realistic QA, we choose HotpotQA (Yang et al., 2018) and MuSiQue (Trivedi et al., 2022) since they mainly aim for multi-hop QA task. We follow RULER (Hsieh et al., 2024) for dataset creation. Furthermore, we evaluate on a more challenging benchmark dataset BABI-LONG (Kuratov et al., 2024) that requires the algorithm to be sensitive to the order of facts presented in the long context. For RULER datasets, we evaluate up to 32K context length, and 16K for BABI-LONG due to limited computational resource. We

employ the same evaluation metric as proposed in each benchmark dataset. Due to limited computation resource, we evaluate 100 examples for each of the length and dataset.

## 5.3 Comparison on Benchmark Datasets

Our experiments evaluate the performance of AT-TRIEVAL and its variant ATTRIEVAL-kl across three open-source models and four datasets. The results summarized in Table 1 and Figure 4, demonstrate consistent improvements over the baseline CoT prompting, particularly in tasks requiring complex reasoning and long-context understanding. We conclude several key findings in the following, and show case studies in the Appendix:

**Superiority of ATTRIEVAL over CoT** On Deduction, ATTRIEVAL and ATTRIEVAL-kl significantly outperform CoT across all models. For in-
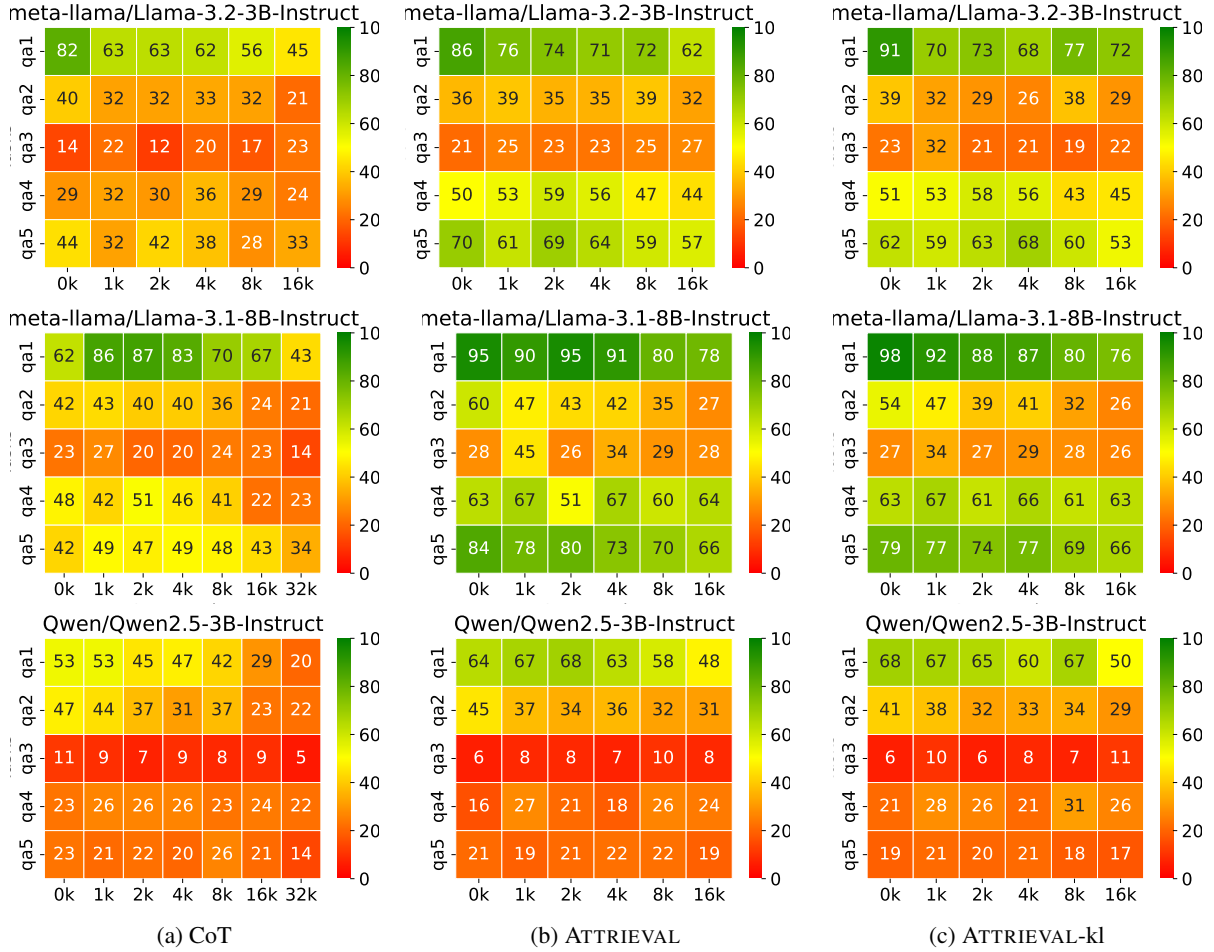
Figure 4: BABILONG results. Greener colors represent higher scores.

---

**Algorithm 2** Cross-Evaluation Token Selection

**Require:** Context $C$, question $Q$, token count $s$, Model $M$

**Ensure:** Selected retriever tokens $\mathcal{T}_{\text{retrieve}}$

1: Generate token distributions:
2: $\quad P_L^{(1:T)} \leftarrow M(\mathcal{P}_L(C,Q))$ ▷ Long prompt
3: $\quad P_S^{(1:T)} \leftarrow M(\mathcal{P}_S(Q))$ ▷ Short prompt
4: Compute token-wise divergence:
5: **for** each token $t \in \{1, \dots, T\}$ **do**
6: $\quad D_{\text{KL}}^{(t)} \leftarrow D_{\text{KL}}\left(P_L^{(t)} \parallel P_S^{(t)}\right)$
7: **end for**
8: $\mathcal{T}_{\text{retriever}} \leftarrow \arg\underset{t}{\text{top-}s}(D_{\text{KL}}^{(t)})$

---

stance, Llama-3.2-3B-Instruct with ATTRIEVAL-kl achieves an overall score of 79 (vs. CoT's 47), while Qwen2.5-3B-Instruct improves from 51 (CoT) to 63 (ATTRIEVAL). Gains are especially pronounced at longer context lengths (16K–32K), where ATTRIEVAL-kl mitigates performance degradation (e.g., Llama-3.2-3B-Instruct at 32K: 61 *vs.* CoT's 14). On MuSiQue, ATTRIEVAL-based methods exhibit stronger robustness. Llama-3.1-8B-Instruct with ATTRIEVAL achieves an overall score of 63, surpassing its CoT counterpart by 21 points. Even smaller models like Qwen2.5-3B-Instruct show improvements (ATTRIEVAL: 44 vs. CoT: 42). On HotpotQA, ATTRIEVAL-based methods perform modest. Llama-3.1-8B-Instruct with AT-TRIEVAL achieves a 71 overall score (*vs.* CoT's 59), but smaller models like Qwen2.5-3B-Instruct show narrower margins (ATTRIEVAL: 56 vs. CoT: 59). We also notice that larger models (e.g., Llama-3.1-8B-Instruct) consistently outperform smaller counterparts when paired with ATTRIEVAL, highlighting synergies between method efficacy and model capacity. For example, on MuSiQue, Llama-3.1-8B-Instruct with ATTRIEVAL scores 63, far exceeding Qwen2.5-3B-Instruct's 44. We also show case studies in Appendix C. And refer to Figure 2a for the recall analysis of retrieved facts.

**Effectiveness of ATTRIEVAL-kl Variant** The AT-TRIEVAL-kl variant consistently matches or exceeds the base ATTRIEVAL method. For example, on Deduction with Llama-3.2-3B-Instruct, AT-

TRIEVAL-kl achieves a 79 overall score (*vs.* AT-TRIEVAL's 74), driven by superior performance at 16K (77 vs. 63). This suggests that integrating our proposed token selection strategy enhances retrieval performance.

### 5.4 How does Token Selection Affect Performance?

We study the effect of varying the generated tokens used for calculating attention scores and the retriever token selection strategy. Specifically, we first treat a paragraph of 150 random words as if they are generated CoT tokens and proceed with ATTRIEVAL algorithm. Surprisingly, as shown in Table 2, we found that even though the tokens are completely irrelevant with the context, they can still improve the performance over the vanilla CoT. We further study the effect of token selection strategy. We propose several variants against our proposed KL-divergence based selection. "First-s" means we only select the first $s$ tokens in the sequence. "Random-s" means we select random $s$ tokens. From Table 2, we found that our proposed strategy performs the best among others. However, we do notice that on Qwen models, our strategy does not perform better than using all tokens, and we hypothesize that this is because Qwen models tend to generate longer CoT and selecting more tokens could help.

## 6 Related Works

**Long-context Reasoning** Architectural innovations, such as modified positional encodings (Chen et al., 2023b; Peng et al., 2023b; Jin et al., 2024; Chen et al., 2023c), sparse attention mechanisms (Zaheer et al., 2020; Lou et al., 2024), RNN-like models (Gu and Dao, 2023; Peng et al., 2023a) have enabled efficient processing of extended sequences while mitigating computational costs, as surveyed in Wang et al. (2024). However, challenges persist in multi-hop reasoning, where models exhibit sensitivity to noisy contexts (Bai et al., 2024; Hsieh et al., 2024; Kuratov et al., 2024; Ling et al., 2025; Zhang et al., 2024b; Wu et al., 2024). To tackle this problem, recent research often employ fine-tuning based approaches that either focus on collecting complex long-context training data (Li et al., 2024b; An et al., 2024; Chen et al., 2024) or training the model to retrieve and cite the context before generating the answers (Li et al., 2024a,c; Yu et al., 2023). While effective, these approaches face two key limitations: collecting high-quality long-context data is prohibitively expensive, and excessive specialization risks degrading performance on short-context tasks. On the other hand, training-free agentic workflows are proposed improve long-context capability (Zhang et al., 2024c; Chen et al., 2023a; Zhang et al., 2024a). This work argues that these approaches does not inherently solve implicit fact retrieval problem.

**Attention-guided Retrieval** Unlike traditional retrieval-augmented generation (RAG) pipelines that rigidly separate retrieval and generation stages, recent approaches leverage attention mechanisms to dynamically guide retrieval process. Notably, Jiang et al. (2022) unifies retrieval and geenration in a single Transformer. Jiang et al. (2023); Li et al. (2022) uses attention distribution guide or trigger retrieval. Wu et al. (2022); Borgeaud et al. (2022) introduce memory banks into Transformers via cross-attention. This work makes the observation that attention from CoT tokens can improve reasoning capability over long-context.

## 7 Discussion and Conclusion

This work starts by making several key observations on the Chain-of-Thought (CoT) of long-context reasoning tasks: (1) CoT struggles with multi-hop reasoning mainly due to incomplete retrieval of implicit facts; (2) attention patterns from intermediate CoT tokens consistently highlight relevant facts, even when those facts remain unmentioned in generated text. We then present AT-TRIEVAL, a novel training-free framework that enhances long-context reasoning by grounding retrieval in the latent signals of transformer attention mechanisms. By identifying and reintegrating these retrieved facts, ATTRIEVAL mitigates the performance degradation of LLMs on tasks requiring multi-hop reasoning over extended contexts. Our results on Deduction, BABILong, and real-world benchmarks like MuSiQue demonstrate its broad applicability and robustness. This work advances the understanding of how attention mechanisms can be harnessed to align context retrieval and reasoning, offering a lightweight yet effective solution. Our work also spurs two potential future directions: (1) iteratively combine CoT generation and attention-guided retrieval based on the model uncertainty; (2) utilize attention weights from generated CoT as a supervision signal to better finetune long-context model.

## Limitations

Despite its effectiveness on various tasks and models, we point the following limitations of AT-TRIEVAL: (1) it requires two steps of response generation—one for acquiring attention matrix and the other for answer generation—which approximately doubles the inference costs. Future work could explore when to early stop the first-round generation and start retrieval. (2) ATTRIEVAL can be effectively applied on applications with shorter CoT. However, when it is applied on long-form generation tasks, ATTRIEVAL should be applied iteratively during the generation. (3) ATTRIEVAL still does not completely solve long-context performance degradation. There is still a minor issue that the reasoning steps can be distracted by excessive attentions spread on previous sequence. Future work could explore context reduction guided by attention weights.

## Ethics Consideration

This paper only studies datasets in English language.

## References

Shengnan An, Zexiong Ma, Zeqi Lin, Nanning Zheng, and Jian-Guang Lou. 2024. Make your llm fully utilize the context. *Preprint*, arXiv:2404.16811.

Yushi Bai, Shangqing Tu, Jiajie Zhang, Hao Peng, Xiaozhi Wang, Xin Lv, Shulin Cao, Jiazheng Xu, Lei Hou, Yuxiao Dong, et al. 2024. Longbench v2: Towards deeper understanding and reasoning on realistic long-context multitasks. *arXiv preprint arXiv:2412.15204*.

Sebastian Borgeaud, Arthur Mensch, Jordan Hoffmann, Trevor Cai, Eliza Rutherford, Katie Millican, George Bm Van Den Driessche, Jean-Baptiste Lespiau, Bogdan Damoc, Aidan Clark, et al. 2022. Improving language models by retrieving from trillions of tokens. In *International conference on machine learning*, pages 2206–2240. PMLR.

Howard Chen, Ramakanth Pasunuru, Jason Weston, and Asli Celikyilmaz. 2023a. Walking down the memory maze: Beyond context limit through interactive reading. *arXiv preprint arXiv:2310.05029*.

Longze Chen, Ziqiang Liu, Wanwei He, Yunshui Li, Run Luo, and Min Yang. 2024. Long context is not long at all: A prospector of long-dependency data for large language models. *arXiv preprint arXiv:2405.17915*.

Shouyuan Chen, Sherman Wong, Liangjian Chen, and Yuandong Tian. 2023b. Extending context window

of large language models via positional interpolation. *arXiv preprint arXiv:2306.15595*.

Yukang Chen, Shengju Qian, Haotian Tang, Xin Lai, Zhijian Liu, Song Han, and Jiaya Jia. 2023c. Longlora: Efficient fine-tuning of long-context large language models. *arXiv preprint arXiv:2309.12307*.

Pradeep Dasigi, Kyle Lo, Iz Beltagy, Arman Cohan, Noah A Smith, and Matt Gardner. 2021. A dataset of information-seeking questions and answers anchored in research papers. *arXiv preprint arXiv:2105.03011*.

gkamradt. 2023. Llmtest needle in a haystack - pressure testing llms. https://github.com/gkamradt/LLMTest_NeedleInAHaystack. GitHub repository for evaluating long-context retrieval capabilities of LLMs.

Albert Gu and Tri Dao. 2023. Mamba: Linear-time sequence modeling with selective state spaces. *arXiv preprint arXiv:2312.00752*.

Chi Han, Qifan Wang, Wenhan Xiong, Yu Chen, Heng Ji, and Sinong Wang. 2023. Lm-infinite: Simple on-the-fly length generalization for large language models. *arXiv preprint arXiv:2308.16137*.

Cheng-Ping Hsieh, Simeng Sun, Samuel Kriman, Shantanu Acharya, Dima Rekesh, Fei Jia, Yang Zhang, and Boris Ginsburg. 2024. Ruler: What's the real context size of your long-context language models? *arXiv preprint arXiv:2404.06654*.

Zhengbao Jiang, Luyu Gao, Jun Araki, Haibo Ding, Zhiruo Wang, Jamie Callan, and Graham Neubig. 2022. Retrieval as attention: End-to-end learning of retrieval and reading within a single transformer. *arXiv preprint arXiv:2212.02027*.

Zhengbao Jiang, Frank F Xu, Luyu Gao, Zhiqing Sun, Qian Liu, Jane Dwivedi-Yu, Yiming Yang, Jamie Callan, and Graham Neubig. 2023. Active retrieval augmented generation. *arXiv preprint arXiv:2305.06983*.

Hongye Jin, Xiaotian Han, Jingfeng Yang, Zhimeng Jiang, Zirui Liu, Chia-Yuan Chang, Huiyuan Chen, and Xia Hu. 2024. Llm maybe longlm: Self-extend llm context window without tuning. *arXiv preprint arXiv:2401.01325*.

Yuri Kuratov, Aydar Bulatov, Petr Anokhin, Ivan Rodkin, Dmitry Sorokin, Artyom Sorokin, and Mikhail Burtsev. 2024. Babilong: Testing the limits of llms with long context reasoning-in-a-haystack. *arXiv preprint arXiv:2406.10149*.

Huayang Li, Pat Verga, Priyanka Sen, Bowen Yang, Vijay Viswanathan, Patrick Lewis, Taro Watanabe, and Yixuan Su. 2024a. Alr$^2$: A retrieve-then-reason framework for long-context question answering. *arXiv preprint arXiv:2410.03227*.

Siheng Li, Cheng Yang, Zesen Cheng, Lemao Liu, Mo Yu, Yujiu Yang, and Wai Lam. 2024b. Large language models can self-improve in long-context reasoning. *arXiv preprint arXiv:2411.08147*.

Xiang Lisa Li, Ari Holtzman, Daniel Fried, Percy Liang, Jason Eisner, Tatsunori Hashimoto, Luke Zettle-moyer, and Mike Lewis. 2022. Contrastive decoding: Open-ended text generation as optimization. *arXiv preprint arXiv:2210.15097*.

Yanyang Li, Shuo Liang, Michael R Lyu, and Li-wei Wang. 2024c. Making long-context language models better multi-hop reasoners. *arXiv preprint arXiv:2408.03246*.

Zhan Ling, Kang Liu, Kai Yan, Yifan Yang, Weijian Lin, Ting-Han Fan, Lingfeng Shen, Zhengyin Du, and Jiecao Chen. 2025. Longreason: A synthetic long-context reasoning benchmark via context expansion. *arXiv preprint arXiv:2501.15089*.

Chao Lou, Zixia Jia, Zilong Zheng, and Kewei Tu. 2024. Sparser is faster and less is more: Efficient sparse attention for long-range transformers. *arXiv preprint arXiv:2406.16747*.

Xiangyang Mou, Chenghao Yang, Mo Yu, Bingsheng Yao, Xiaoxiao Guo, Saloni Potdar, and Hui Su. 2021. Narrative question answering with cutting-edge open-domain qa techniques: A comprehensive study. *Transactions of the Association for Computational Linguistics*, 9:1032–1046.

Bo Peng, Eric Alcaide, Quentin Anthony, Alon Albalak, Samuel Arcadinho, Stella Biderman, Huanqi Cao, Xin Cheng, Michael Chung, Matteo Grella, et al. 2023a. Rwkv: Reinventing rnns for the transformer era. *arXiv preprint arXiv:2305.13048*.

Bowen Peng, Jeffrey Quesnelle, Honglu Fan, and Enrico Shippole. 2023b. Yarn: Efficient context window extension of large language models. *arXiv preprint arXiv:2309.00071*.

Harsh Trivedi, Niranjan Balasubramanian, Tushar Khot, and Ashish Sabharwal. 2022. Musique: Multi-hop questions via single-hop question composition. *Transactions of the Association for Computational Linguistics*, 10:539–554.

Xindi Wang, Mahsa Salmani, Parsa Omidi, Xiangyu Ren, Mehdi Rezagholizadeh, and Armaghan Eshaghi. 2024. Beyond the limits: A survey of techniques to extend the context length in large language models. *arXiv preprint arXiv:2402.02244*.

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. *Advances in neural information processing systems*, 35:24824–24837.

Di Wu, Hongwei Wang, Wenhao Yu, Yuwei Zhang, Kai-Wei Chang, and Dong Yu. 2024. Longmemeval: Benchmarking chat assistants on long-term interactive memory. *arXiv preprint arXiv:2410.10813*.

Yuhuai Wu, Markus N Rabe, DeLesley Hutchins, and Christian Szegedy. 2022. Memorizing transformers. *arXiv preprint arXiv:2203.08913*.

Guangxuan Xiao, Yuandong Tian, Beidi Chen, Song Han, and Mike Lewis. 2023. Efficient streaming language models with attention sinks. *arXiv preprint arXiv:2309.17453*.

Zhilin Yang, Peng Qi, Saizheng Zhang, Yoshua Bengio, William W. Cohen, Ruslan Salakhutdinov, and Christopher D. Manning. 2018. HotpotQA: A dataset for diverse, explainable multi-hop question answering. In *Conference on Empirical Methods in Natural Language Processing (EMNLP)*.

Wenhao Yu, Hongming Zhang, Xiaoman Pan, Kaixin Ma, Hongwei Wang, and Dong Yu. 2023. Chain-of-note: Enhancing robustness in retrieval-augmented language models. *arXiv preprint arXiv:2311.09210*.

Manzil Zaheer, Guru Guruganesh, Kumar Avinava Dubey, Joshua Ainslie, Chris Alberti, Santiago Ontanon, Philip Pham, Anirudh Ravula, Qifan Wang, Li Yang, et al. 2020. Big bird: Transformers for longer sequences. *Advances in neural information processing systems*, 33:17283–17297.

Qingru Zhang, Xiaodong Yu, Chandan Singh, Xiaodong Liu, Liyuan Liu, Jianfeng Gao, Tuo Zhao, Dan Roth, and Hao Cheng. 2024a. Model tells itself where to attend: Faithfulness meets automatic attention steering. *Preprint*, arXiv:2409.10790.

Xinrong Zhang, Yingfa Chen, Shengding Hu, Zihang Xu, Junhao Chen, Moo Hao, Xu Han, Zhen Thai, Shuo Wang, Zhiyuan Liu, and Maosong Sun. 2024b. ∞Bench: Extending long context evaluation beyond 100K tokens. In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 15262–15277, Bangkok, Thailand. Association for Computational Linguistics.

Yusen Zhang, Ruoxi Sun, Yanfei Chen, Tomas Pfister, Rui Zhang, and Sercan Ö Arik. 2024c. Chain of agents: Large language models collaborating on long-context tasks. *arXiv preprint arXiv:2406.02818*.

## A  Deduction: A Diagnostic Benchmark for Long-context Reasoning

We first create 6 problem types including: fruit price, person age, car speed, city population, book length and planet temperature. For each of them, we generate 15 entities as candidates. Then a statement is generated by first randomly sample a problem type and then a subset of entities is randomly sampled. Unique values are then assigned to these entities using a controlled random number generation process that ensures non-duplicative values. The statements further encodes relationships

Figure 5: Ranking of tokens most attended in the statements. The example shows a failure case.
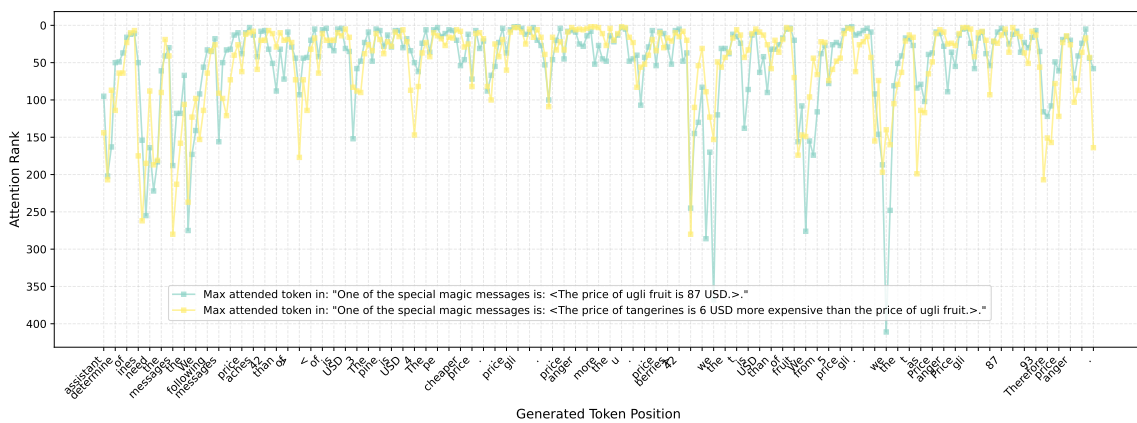


Figure 6: Ranking of tokens most attended in the statements. The example shows a success case.

Figure 7: Proportion of attention from generated tokens to the input prompt across layers.

between entities by formulating independent and pairwise conditions based on templated statements, which describe direct values or comparative differences (e.g., "more expensive" or "older"). To increase the complexity and challenge of inference, additional distractor conditions involving extra entities are optionally introduced. Finally, a question is generated to prompt for numeric responses. Finally statements and distractors are randomly inserted to a haystack (gkamradt, 2023) to flexibly extend the context length.

## B  Statement Ranking from the Attention

In subsection 3.2, we mention Figure 5 about the ranking of statement tokens. We now illustrate how we get the rankings. Specifically, we calculate the ranking by first identifying the most attended token within the statement $i^*(t) = \arg\max_{i \in I_{\text{stmt}}} A_{t,i}^{(l)}$, and then determine its rank among all the input tokens

$$r(t) = 1 + \sum_{i=1}^{N} \mathbf{1}\left(A_{t,i}^{(l)} > A_{t,i^*(t)}^{(l)}\right). \quad (8)$$

Apart from a failure case, we also show the heatmap plots and ranking plot for a successful retrieval case in Figure 7 and Figure 6.

## C  Case Study of Retrieved Facts

We show the cases where CoT suffers from implicit retrieval while ATTRIEVAL can successfully retrieve the ground truth context.

## D  Prompt Used to Integrate Facts

We use the same prompt for Deduction and two QA datasets. The template is as follows:

*"{anything before question}*

*Some clauses extracted from the context that might be related:*
*{clauses}*

*{anything after question starts}"*

For BABILONG, we employ a slightly different template since these tasks are sensitive fact order.

*"{anything before question}*

*Some clauses are extracted from the context that might be related:*
*{clauses}*

*Notice that the clause indices represents the order of them appearing in the context. Larger clause indices indicate that they appear later in the context. The answer to the question is sensitive to the order in the context. The clauses only serve as a hint, please check the original context for exact information.*

*{anything after question starts}"*

For BABILONG, we employ a slightly different template since these tasks are sensitive fact order.

We reorder the retrieved facts according to the order they appear in the context.

12

| Ground Truth Facts or Answer | Original CoT | Retrieved Facts |
|---|---|---|
| The length of Fables of Yldora is 40 pages shorter than the length of Saga of Lynora. The length of Saga of Lynora is 92 pages. | To answer this question, we need to refer to one of the special magic messages provided in the text. The message is:\n\n"One of the special magic messages is: <The length of Fables of Yldora is 40 pages.>."\n\nFrom this message, we can infer that the length of Fables of Yldora is 40 pages. | clause 136: One of the special magic messages is: <The length of Fables of Yldora is 40 pages shorter than the length of Saga of Lynora ... clause 460: One of the special magic messages is: <The length of Saga of Lynora is 92 pages. |
| Marie de' Medici | Based on the document, Philippe, Duke of Orléans (Document 53) was the younger son of Louis XIII of France and his wife, Anne of Austria. His mother, Anne of Austria, was the daughter of King Philip II of Spain and Elisabeth of Valois, making Elisabeth of Valois the grandmother of Philippe, Duke of Orléans. | clause 185:Document 54:Leonora Dori Leonora Dori Galiga (19 May 1568 – 8 July 1617) was a French courtier of Italian origin, an influential favourite of the French regent Marie de' Medici, mother of King Louis XIII of France ... clause 189:Document 56:Louise Élisabeth d'Orléans Louise Élisabeth d'Orléans was one of the seven daughters of Philippe II, Duke of Orléans, and his wife, Françoise Marie de Bourbon. |
| from 1986 to 2013 | The football manager who recruited David Beckham managed Manchester United between 1945 and 1969, and again for the second half of the 1970–71 season. The answer is: Sir Matt Busby managed Manchester United between 1945 and 1969, and again for the second half of the 1970–71 season. | clause 4:Document 2:\nAlex Ferguson\nSir Alexander Chapman Ferguson, CBE (born 31 December 1941) is a Scottish former football manager and player who managed Manchester United from 1986 to 2013. |

Table 3: Cases where CoT is unable to retrieve the ground truth facts but ATTRIEVAL can successfully retrieve.