# Focused Chain-of-Thought: Efficient LLM Reasoning via Structured Input Information

**Lukas Struppek**[1*]          **Dominik Hintersdorf**[2,3]          **Hannah Struppek**[4]

**Daniel Neider**[5,6]          **Kristian Kersting**[2,3,7,8]

[1]FAR.AI,
[2]German Research Center for Artificial Intelligence (DFKI),
[3]Technical University of Darmstadt,
[4]University of Kassel,
[5]TU Dortmund University,
[6]TU Center for Trustworthy Data Science and Security, University Alliance Ruhr,
[7]Hessian Center for AI (Hessian.AI),
[8]Centre for Cognitive Science, Technical University of Darmstadt

## Abstract

Recent large language models achieve strong reasoning performance by generating detailed chain-of-thought traces, but this often leads to excessive token use and high inference latency. Existing efficiency approaches typically focus on model-centric interventions, such as reinforcement learning or supervised fine-tuning, to reduce verbosity. In contrast, we propose a training-free, input-centric approach. Inspired by cognitive psychology, we introduce Focused Chain-of-Thought (F-CoT), which separates information extraction from the reasoning process. F-CoT first organizes the essential information from a query into a concise, structured context and then guides the model to reason exclusively over this context. By preventing attention to irrelevant details, F-CoT naturally produces shorter reasoning paths. On arithmetic word problems, F-CoT reduces generated tokens by 2–3× while maintaining accuracy comparable to standard zero-shot CoT. These results highlight structured input as a simple yet effective lever for more efficient LLM reasoning.

## 1 Introduction

Large language models (LLMs) are trained to predict the next token given a sequence of previous ones. Scaling model parameters and training data has substantially improved their performance on mathematical reasoning benchmarks, with recent models continuing to push the state of the art. Many LLMs reveal their internal reasoning by producing an explicit chain-of-thought (CoT) (Wei et al., 2022) – a step-by-step, natural-language rationale that makes reasoning traceable to humans. While CoT outputs increase transparency, they also generate long reasoning traces that are costly in time and computation. Moreover, locating errors within a long CoT is challenging, since the entire trace must typically be checked to identify the mistakes and verify correctness.

The reasoning processes of LLMs are often compared to human logical thinking. Foundational work in cognitive psychology, such as the Active Control of Thought (ACT) framework (Anderson, 1976) models human problem-solving as sequential, resource-efficient processes, beginning with the representation and structuring of information before higher-order reasoning. Modern LLMs exhibit analogous reasoning behavior. However, the stages of information extraction and structuring in LLMs are not clearly distinguished from the subsequent reasoning phase and are often interwoven with it. We hypothesize that this entanglement blurs the boundaries between relevant and irrelevant information, thereby complicating the LLM reasoning process and contributing to the generation of unnecessary tokens.

---

*Work mainly done at DFKI/Technical University of Darmstadt. Contact: *FirstName*@far.ai.
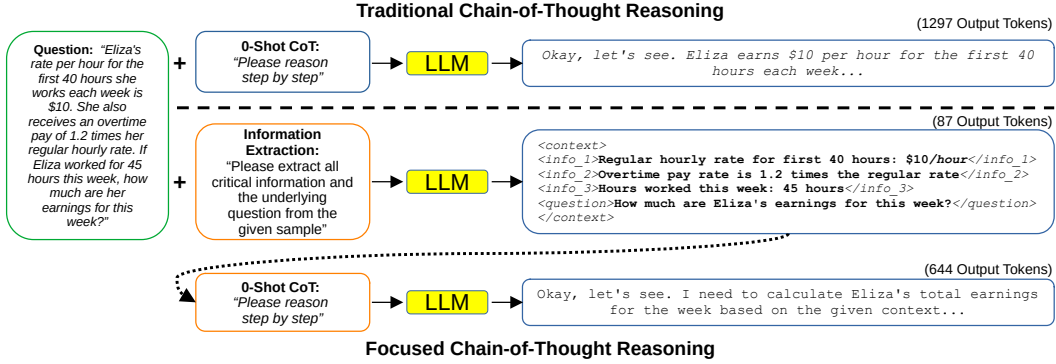
Figure 1: Focused Chain-of-Thought reasoning. The model first extracts key information into an XML-like context block and then performs reasoning based on that block. The context can also be pre-defined by the user or generated automatically by a larger LLM. When queried using only the context, large reasoning models produce significantly shorter reasoning traces compared to standard natural-language inputs. In this particular example, Qwen3 14B produces 43% fewer tokens compared to standard CoT prompting. Shown prompts are abbreviated; see Appx. A.1 and A.5 for full prompts.

In this paper, we introduce *Focused Chain-of-Thought (F-CoT)*, a novel prompting strategy for reasoning tasks. Drawing inspiration from the ACT framework in cognitive psychology, F-CoT explicitly separates information extraction and structuring from the core reasoning process. In the first stage, we prompt an LLM to extract and organize relevant information from a question into a fixed, compact, and structured format. In the second stage, the model receives only this structured representation and performs standard chain-of-thought reasoning. A high-level overview of F-CoT is provided in Fig. 1. By decoupling these two phases, F-CoT substantially reduces the number of generated tokens, accelerating inference by 2–3x compared to standard chain-of-thought prompting, while preserving strong reasoning performance. Crucially, our prompting strategy does not involve instructing or fine-tuning the LLM to shorten its reasoning—providing input information in a structured way is already sufficient to speed up reasoning. This makes F-CoT orthogonal to methods that explicitly encourage token efficiency through prompting strategies (Xu et al., 2025a; Lee et al., 2025), supervised fine-tuning (Yu et al., 2025; Luo et al., 2025), or reward optimization (Aggarwal & Welleck, 2025; Yeo et al., 2025).

## 2 REASONING IN LARGE LANGUAGE MODELS

LLMs are autoregressive models that generate text token by token. While recent models (Dubey et al., 2024; DeepSeek-AI, 2025) have achieved impressive performance across diverse domains, unlocking their full potential requires more than large parameter counts and extensive training data. Their capabilities crucially depend on advanced inference strategies that elicit structured reasoning, i.e., the generation of logically consistent and interpretable intermediate steps rather than shallow text completions. This process is often compared to the slow, deliberative, and analytical (*System 2*) mode of human thinking described by Kahneman (2011).

**Prompt-Based Reasoning.** A common way to trigger such reasoning behavior is to prompt the model to decompose complex problems into a series of intermediate steps, known as chain-of-thought (CoT) prompting (Wei et al., 2022). Notably, LLMs have been shown to reason effectively even in zero-shot settings, generating coherent reasoning steps without any in-context examples. Simply instructing a model to *"think step by step"* can substantially enhance its reasoning performance (Kojima et al., 2022). Building on these insights, subsequent work has proposed multi-path reasoning approaches such as tree-of-thought reasoning (Yao et al., 2023a), self-consistency (Wang et al., 2022), and ReAct (Yao et al., 2023b), which explore multiple reasoning paths before selecting or combining outcomes.

**Dedicated Reasoning LLMs.** An emerging class of specialized reasoning models, including the DeepSeek-R1 (DeepSeek-AI, 2025) and Qwen3 (Qwen Team, 2025) families, employ a structured two-stage generation process. The model first enters a reasoning stage, often triggered by a special

<think> token, followed by a concise summary stage where it synthesizes the final answer. These large reasoning models achieve state-of-the-art performance on logic and math benchmarks but often produce extremely long reasoning traces, sometimes spanning thousands of tokens for relatively simple questions (Chen et al., 2025; Cuadron et al., 2025). Even though these reasoning models perform better than traditional non-reasoning LLMs, their high parameter counts, combined with these lengthy reasoning processes significantly increase inference time and computational cost.

**Test-Time Scaling.** However, this relationship of improved performance due to longer reasoning traces can also be exploited by allocating more inference time, typically by generating additional tokens, commonly referred to as test-time scaling. Various inference-time techniques implement this principle by extending or parallelizing the reasoning process, for example, by using reconsideration tokens (Muennighoff et al., 2025) or self-refinement steps (Madaan et al., 2023; Tian et al., 2025). While such approaches aim to enhance accuracy by scaling up inference compute, the opposite direction, i.e., reducing the number of generated tokens to save time and resources, has received comparatively less attention, even though large reasoning models are known to suffer from overthinking (Chen et al., 2025; Chiang & Lee, 2024).

**Efficient Reasoning.** Existing approaches for reducing token count during reasoning typically intervene at the training or fine-tuning stage. Reinforcement learning methods explicitly reward shorter outputs (Aggarwal & Welleck, 2025; Yeo et al., 2025; Arora & Zanette, 2025; Ye et al., 2025). Data-centric approaches, on the other hand, shorten reasoning traces either by removing redundant and irrelevant tokens (Zhuang et al., 2025; Xia et al., 2025; Yuan et al., 2025; Xiao et al., 2025) or by fine-tuning on shorter chain-of-thought examples (Yu et al., 2025; Luo et al., 2025). Existing training-free approaches aim to reduce the token count by explicitly instructing the model to keep its reasoning short and concise (Xu et al., 2025a; Nayab et al., 2024).

All these methods for increasing reasoning efficiency have in common that they target the model itself, either by explicitly instructing or fine-tuning the model to produce shorter reasoning traces. In contrast, we adopt a fundamentally different perspective that targets the input rather than the model. With our Focused Chain-of-Thought (F-CoT) framework, we demonstrate that presenting information in a structured and compact form naturally leads to shorter reasoning traces and substantial token savings, without any additional instructions or fine-tuning. While this approach is related to chain-of-thought prompting (Wei et al., 2022), as it still encourages sequential reasoning, it fundamentally differs in how the input information is represented. Similarly, it differentiates itself from schema-based prompting (Zhong et al., 2023), which organize outputs in structured formats such as JSON or XML, as F-CoT instead structures inputs to enhance reasoning efficiency. Moreover, while F-CoT has similarities with retrieval-augmented generation (Lewis et al., 2020) as it separates information retrieval from reasoning, it performs self-extraction instead of querying an external database. Notably, while knowledge-graph approaches (Pan et al., 2024; Kau et al., 2024; Ma et al., 2025) also leverage structured factors, F-CoT fundamentally differs by relying solely on self-contained, compact representations for reasoning, rather than externally provided knowledge.

## 3 FROM NATURAL LANGUAGE TO STRUCTURED REPRESENTATIONS

Natural language descriptions in mathematical reasoning tasks often contain the essential facts alongside irrelevant wording. This unstructured form of information can lead LLMs to produce long, verbose reasoning traces. Consider the following arithmetic problem:

> **Prompt:** An apple costs \$2 in the supermarket. A pear costs one and a half times the price of an apple. Eve wants to buy some fruits for herself to eat more healthily. What is the total cost of two apples and two pears?

When given such a question, an LLM typically begins by interpreting the text and directly reasoning about it:

> **LLM Output:** Okay, let me try to figure out this problem. The question is about finding the total cost of two apples and two pears in the supermarket. [...] First, it says an apple costs \$2. [...] Then, a pear costs one and a half times the price of an apple. [...]

This natural-language-based reasoning process generates many unnecessary tokens and scatters key information across verbose text, even when the facts are simple. By contrast, humans are taught in school first to extract and note down the essential facts in a structured format. For example, giving the same problem to a person, one would start by identifying all the necessary information for solving the task:

> **Info 1:** Price per apple: $2
> **Info 2:** Price per pear: $1.5\times$ price per apple

A structured representation allows humans to attend to key information more easily, since it is no longer buried in full sentences and irrelevant information but isolated in a clear and concise format. This compact representation isolates the relevant details, making reasoning easier and reducing cognitive overhead. Inspired by this human strategy, we hypothesize that providing LLMs with similarly structured inputs can shorten reasoning paths, reduce irrelevant attention distraction, and improve efficiency and faithfulness.

## 3.1 FROM TEXT TO STRUCTURED FACTS

To test this hypothesis, we define a fixed, structured context format that explicitly separates factual information from the question:

> $<$**context**$>$
>     $<$**info_1**$>$Price per apple: $2$<$**/info_1**$>$
>     $<$**info_2**$>$Price per pear: 1.5x price per apple$<$**/info_2**$>$
>     $<$**question**$>$Total cost of 2 apples and 2 pears$<$**/question**$>$
> $<$**/context**$>$

This structure contains two components: (1) enumerated information blocks $<$info_k$>$...$<$/info_k$>$, and (2) a $<$question$>$ field specifying the objective. The XML-like format supports automatic parsing, syntax validation, and compatibility with downstream pipelines.

While we use this XML-like format for clarity, our experiments (see Sec. 4.5) show that simpler enumerated lists yield comparable performance, suggesting that the key benefit arises from structure itself, not from the specific syntax.

## 3.2 IMPLEMENTING STRUCTURED REASONING PIPELINES

There are two practical strategies to supply structured information in our F-CoT framework:

1. **Pre-formatting by the user:** The user manually provides the question in the defined structure, allowing the model to focus entirely on reasoning. This minimizes model workload but requires additional user effort.

2. **Two-step prompting:** The model first extracts the context (without reasoning) and then reasons over it. For less capable LLMs, extraction may be unreliable; in such cases, a larger model can generate the context, which a smaller model then uses for reasoning. This hybrid pipeline combines the extraction capabilities of large models with the efficiency of smaller ones, reducing total inference cost. The two-step prompting approach is also depicted in Fig. 1.

Throughout this paper, we follow the second approach, which first queries an LLM to extract and generate the context for a given question. We then provide only this extracted context, without the original question, to the model, asking it to reason about it. Our reasoning prompt extends zero-shot chain-of-thought prompting by explicitly instructing the model to (i) focus solely on the structured context, (ii) reason step by step, and (iii) cite relevant $<$info_k$>$ blocks when using them. These explicit citations make the reasoning process interpretable and facilitate debugging. We present the exact prompts used for F-CoT in Sec. A.1.

# 4 EXPERIMENTS

In the following, we experimentally demonstrate that our proposed F-CoT matches the accuracy of traditional zero-shot CoT reasoning (Kojima et al., 2022) with state-of-the-art reasoning LLMs on standard math benchmarks, while producing far more concise reasoning.

## 4.1 EXPERIMENTAL PROTOCOL

**LLMs and Hyperparameters:** We use *Qwen-3* (0.6B, 4B, 14B, 32B) (Qwen Team, 2025) reasoning LLMs, i.e., models with an explicit thinking mode. For all models, we use sampling parameters recommended by the model developers: temperature 0.6, top-p 0.95, min-p 0.0, and top-k 20. The maximum number of generated tokens is set to 32k, matching the context length supported by Qwen-3 models. When using the LLMs to generate the context, we skip the reasoning process for Qwen-3 models by appending a *no_think* flag to the input prompt.

**Evaluation Benchmarks:** We focus on arithmetic word problems, i.e., problems that combine key information with unrelated details such as descriptions of settings, names, and situations. We primarily rely on three datasets with varying levels of difficulty: *SVAMP*, *GSM-Hard*, and *MATH-500*. The SVAMP dataset (Patel et al., 2021) is comparably easy, containing short and straightforward questions. GSM-Hard is a more challenging version of the GSM8k (Cobbe et al., 2021) test set of grade-school math problems. The hard version (Gao et al., 2023) replaces the original numbers with less common ones, thereby increasing the difficulty of the reasoning process and reducing the risk of test set contamination. In addition, we employ the more demanding MATH-500 benchmark (Hendrycks et al., 2021). While our main focus is on verbose mathematical questions, we also analyze the impact of highly condensed problems using the AIME2024 and AIME2025 datasets (Competitions, 2025). Since these problems are not the central focus of this paper, we report the corresponding results in the appendix and discuss only the key findings in the main text.

**Metrics:** We compute the Pass@5 metric, which evaluates whether at least one of the five outputs generated by a model for a given problem is correct. Additionally, we calculate the average number of tokens generated per question, denoted as *# Tokens*. For settings using our F-CoT reasoning where the LLM generates the context itself, this token count also includes tokens from the extracted context to ensure a fair comparison. For self-generated contexts, we further measure the proportion of valid context blocks generated, i.e., whether the context adheres to a valid XML structure.

**Baselines:** We compare the F-CoT prompting strategy against the standard zero-shot CoT prompting without any contextual information provided (*0-CoT*). As additional baselines, we investigate Plan-and-Solve prompting (*P&S*) (Wang et al., 2023) and Condition-Retrieving Instruction prompting (*CoRe*) (Xu et al., 2025b), two prompting strategies designed to improve 0-CoT prompting by asking the model to first identify crucial conditions and objectives before starting the reasoning process. These methods follow a strategy similar to our F-CoT, but without explicitly separating the two steps or enforcing a fixed structure. Since we do not explicitly instruct the model to reduce its reasoning effort, e.g., by employing reinforcement learning-based fine-tuning to produce fewer tokens, we do not compare against such approaches. Instead, our goal is to investigate how the availability of structured information benefits the reasoning process.

## 4.2 BOOSTING REASONING EFFICIENCY WITH PRE-COMPUTED CONTEXT

We begin by considering the setting in which the user provides the LLM with a pre-computed context. While transforming a question into such a context can be done manually, we use *GPT-5 mini* (version: *2025-08-07*), a state-of-the-art LLM, for this task. These experiments aim to evaluate the impact of structured information on the reasoning process of LLMs. The LLMs' ability to generate context autonomously is assessed in the next subsection.

For context extraction, we instructed GPT-5 to include only information directly relevant to answering the question, while ignoring background details. The model was explicitly asked to refrain from solving the problem or performing calculations, to avoid evaluation biases stemming from pre-computed reasoning steps. Additionally, we provided in-context two examples of correctly extracted contexts to ensure proper formatting. Full prompts are available in Sec. A.2. The prompts for GSM-Hard/SVAMP and MATH differ only in the examples provided, reflecting the differences in
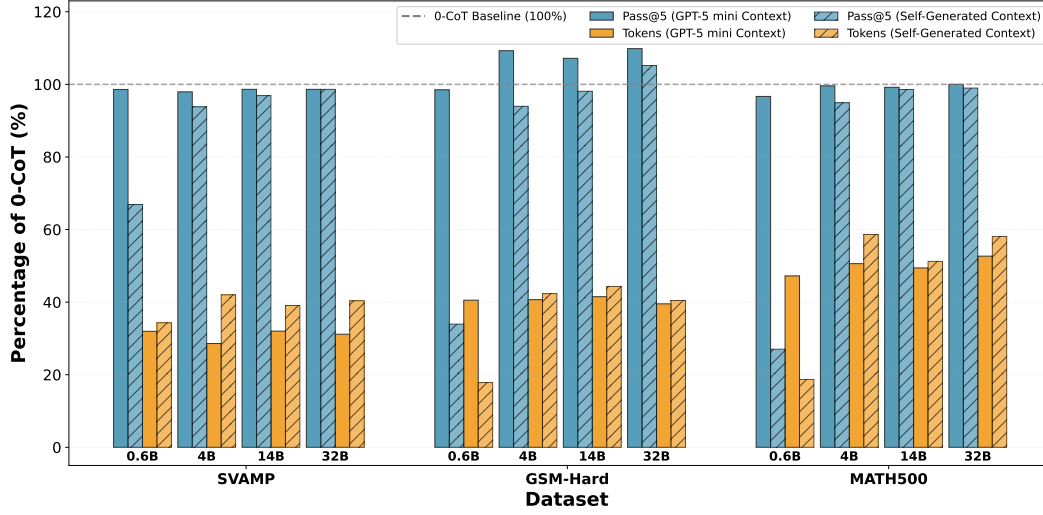
Figure 2: Comparison of 0-CoT and our F-CoT using Qwen3 models of various sizes. For F-CoT, two settings are shown: context pre-computed by GPT-5 mini (solid bars) and generated by the model itself (hatched bars). F-CoT results are expressed relative to 0-CoT. While F-CoT matches 0-CoT performance in most cases, it generates substantially fewer tokens, thereby improving inference efficiency. Detailed numerical results are provided in Appendices C.1 and C.2.

question style. On average, GPT-5 generated contexts for MATH-500 consist of 126 tokens, meaning context generation introduces minimal overhead.

The evaluation results are shown in Fig. 2 (non-hatched bars). We report the F-CoT results as relative values compared to zero-shot CoT (0-CoT). For example, a token usage of 40% indicates that the model required only 40% of the tokens on average compared to 0-CoT. Absolute metrics are additionally provided in Tab. 1 in Sec. C.1.

The most notable finding is that providing reasoning LLMs with context substantially reduces the number of generated tokens (orange bars). For SVAMP questions, token counts drop to roughly one-third of 0-shot CoT across all model sizes. For more challenging GSM-Hard and MATH-500 questions, token counts decrease by about half. For example, Qwen3 32B generates only around 2.3k tokens on MATH-500 questions when provided with the context, compared to roughly 4.4k tokens under 0-shot CoT, while achieving the same Pass@5 scores. This reduction translates to a significant inference speed-up, and we observe similar token savings across all model sizes.

In terms of reasoning accuracy (blue bars), performing reasoning over the provided context yields results comparable to 0-CoT across all three datasets, except for the smallest Qwen3 0.6B model, which exhibits slightly lower performance. Larger models, in contrast, outperform 0-shot CoT on GSM-Hard. When manually analyzing the results, we find that the main differences arise in questions requiring negative results. With 0-CoT, models sometimes misinterpret large numbers due to hallucinated/confabulated typos, producing incorrect answers. When reasoning over a structured context, the models still note unusually large numbers but avoid false interpretations, adhering to the provided information. An example reasoning trace is shown in Sec. B.1. However, since careful prompting could mitigate such errors, 0-CoT performance could likely match structured context results if explicitly instructed to interpret numbers as given.

Prompting the models with *Plan-and-Solve* and *CoRe* does not lead to a noticeable difference, neither in reasoning capabilities nor in the number of generated tokens – both values are comparable to the results with 0-CoT prompting. We, therefore, report the corresponding results only in Sec. C.1.

For AIME questions, which already present condensed information and do not require prior extraction, we still observe a substantial reduction in token usage when using structured context. However, correctness slightly decreases. Given the small dataset size (30 samples) and potential test-set contamination (Golchin & Surdeanu, 2023; Oren et al., 2023), we cannot conclusively claim that

F-CoT harms reasoning in this setting. The observed difference may instead reflect that the model has likely seen the original questions during training, whereas the structured context representation is novel, a phenomenon that also applies to other benchmarks.

## 4.3 CONTEXT GENERATION AND REASONING WITHIN A SINGLE MODEL

So far, we have explored F-CoT by providing the model with external information in the form of pre-computed context generated by another LLM (GPT-5 mini). We now focus on the capabilities of vanilla reasoning LLMs to generate such contexts themselves. In other words, following the concept in Fig. 1, both steps, i.e., context generation and reasoning, are performed by the same model. As before, models see only the generated context during the reasoning step, with no access to the original input. The prompt used for context generation is provided in Sec. A.1, and the relative results are again reported in Fig. 2 (hatched bars). See Sec. C.2 for numerical results. Unlike the context generation with GPT, where the prompt includes some examples of successfully generated contexts, we do not provide such few-shot demonstrations to offer the models more flexibility in their context generation and to obtain a clearer image of their inherent capabilities in generating contexts.

Regarding token counts (including context tokens), we observe values comparable to those in the previous section, again highlighting a substantial reduction compared to 0-shot CoT. In terms of reasoning performance, larger models tend to perform better, often matching the results achieved with GPT-5 pre-computed contexts. However, smaller models, particularly Qwen3 0.6B, exhibit a notable drop in performance. This decline is primarily due to their inability to reliably extract and structure information in the desired format: fewer than 2% of generations produce a valid context. While some of these contexts still contain enough information to solve the question, the model can answer only a limited subset of problems, underperforming compared to both 0-shot CoT and externally computed contexts. Context generation is more stable for Qwen3 14B and larger, where almost all generations produce valid contexts.

These results suggest a practical and cost-efficient strategy for F-CoT: pre-compute the context, which typically requires only a few tokens, using a larger, slower, and more expensive LLM, and then perform the token-intensive reasoning step with a smaller, faster, and more efficient model.

## 4.4 QUANTIFYING AND CHARACTERIZING REASONING DYNAMICS

Next, we examine how F-CoT impacts the reasoning dynamics of LLMs. We focus in this section on the Qwen3 14B model and its results on MATH-500 since the model offers a good balance between reasoning capabilities and size. Moreover, the MATH-500 benchmark offers a diverse range of questions.

We first adapt the overthinking score originally proposed by Cuadron et al. (2025) to quantify overthinking in reasoning LLMs on agentic tasks. In our study, overthinking refers to a model's tendency to prioritize abstract reflection, speculation, or multi-branch planning over performing the concrete mathematical steps that directly advance the solution. The overthinking score, ranging from 0 to 10, measures the degree to which a model emphasizes internal deliberation over concrete problem-solving. Low scores indicate focused, step-by-step reasoning with minimal speculative discussion, while high scores reflect extensive meta-reasoning, frequent method-switching, or premature conclusions without deriving intermediate results.

To compute the score, we modify the original prompt to emphasize mathematical reasoning rather than general agentic behavior. Sample-wise scores are then predicted by GPT-5 Mini, with the full prompt provided in Sec. A.8. The overthinking score for the model using 0-shot CoT is $2.35 \pm 1.5$, indicating mild overthinking. Incorporating the context reduces the score to $1.74 \pm 1.4$, demonstrating improved reasoning efficiency. Moreover, the score decreases in $50.2\%$ of the samples and remains unchanged in $38.4\%$, highlighting consistent reductions in unnecessary reasoning.

We further annotate each output sentence of the Qwen3 model using GPT-5 Mini as *Extraction*, *Reasoning*, or *Filler*. Extraction sentences copy or paraphrase information from the input question without transformation or inference. Reasoning sentences involve logical or arithmetic inference or provide explanations not explicitly stated in the input. Filler sentences include remaining utterances that do not contribute to the solution, such as phrases like *"Wait, let me re-read the problem."* or *"So*

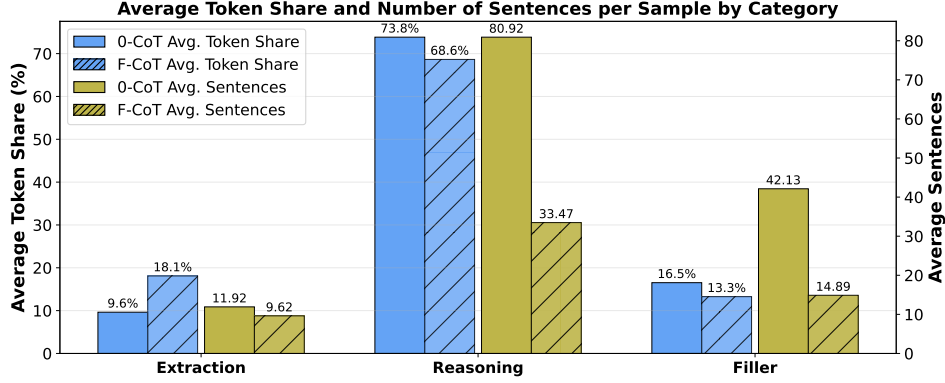**Average Token Share and Number of Sentences per Sample by Category**

Figure 3: Analysis of reasoning traces during the chain-of-thought, where each sentence is classified as *Extraction*, *Reasoning*, or *Filler*. Blue bars indicate the average share of tokens per category, while green bars show the average number of sentences per category. Although the relative token distribution remains largely unchanged, the number of reasoning and filler sentences is substantially reduced when using our F-CoT compared to 0-CoT. Model: Qwen3-14B; Dataset: MATH-500.

*that makes sense."* We then compare the token distribution across the different categories, as well as the average number of sentences classified as each category.

Fig. 3 visualizes the results. When comparing the token distribution across categories (blue bars), no fundamental difference is apparent between standard 0-CoT and our F-CoT. A slight token increase is observed for F-CoT in the extraction category. However, this can be attributed to the model explicitly referencing individual information blocks during reasoning, which are consequently labeled as extraction. In terms of the average number of sentences per category (green bars), both prompting strategies produce roughly the same number of extraction sentences. In contrast, F-CoT more than halves the number of reasoning sentences and substantially reduces the number of filler sentences, which do not directly contribute to solving the problem. We therefore conclude that the overall token reduction achieved by F-CoT primarily results from fewer reasoning and filler sentences, while the relative distribution of tokens across categories remains largely unchanged.

### 4.5 ABLATION STUDY AND SENSITIVITY ANALYSIS

In the final experimental section, we want to clarify the impact of individual design choices for our F-CoT experiments. We again focus on Qwen3 14B and the MATH-500 dataset. We report exact numerical results in Tab. 3 in Sec. C.3

**Q1: Does the prompt format influence the token count?** Previous research (Sclar et al., 2024; Errica et al., 2025) has pointed out the sensitivity of LLMs to their prompt design. To ensure that the structured information is the driving force behind the shorter reasoning traces, we repeat the experiments on MATH-500 using the Qwen3 14B model. We minimize the details in our LLM instruction to use the context to a minimum, followed by the standard 0-shot CoT prompt after the context. The shortened prompt is presented in Listing 9 in Sec. A.5. Removing the explicit instructions on how to use the context increases the average token count to about 3.1k (+747 tokens compared to default F-CoT) while also slightly increasing the pass@5 to 99.2% (+0.6pp).

We repeat this experiment using the same prompt with the standard questions instead of the compact context representation to see whether it also affects these results. Indeed, we observe a token count of 4.2k (−685 tokens compared to default 0-CoT) and also a pass@5 accuracy of 99.2%. Yet, there remains a clear gap of more than 1k tokens between the inference with and without the context, even when using the same prompt type. We also report this prompt in Listing 5 in Sec. A.3.

**Q2: Does the context format influence the results?** We further assess whether the format of the context plays a crucial role. Instead of using the XML format, we explore three alternative design choices that abandon this strict structure. First, we use a simple *enumerated list* of information items followed by the question. Second, we employ an *unnumbered list* to present the information. Finally,

we completely remove any structured representation and simply concatenate all information, followed by the question. The structures of all three context variations are illustrated in Sec. A.6.

We observe no difference in Pass@5 or token count when provided with the context as an enumerated list. While performance remains unchanged for the unnumbered list, the token count increases by 12.6%. A similar effect is observed when concatenating all information without structure. We therefore conclude that the exact format or representation of the extracted information is not crucial for the success of our F-CoT. All forms of structured and compressed information substantially reduce the number of generated tokens. However, introducing more structure, such as numbering individual information pieces, further decreases the token count.

**Q3: Can a model benefit from receiving both the original question and the extracted context?** While we implicitly assume that the provided context contains all information necessary to answer a question, including the original question alongside the extracted context may offer additional benefits. To test this, we repeat the inference process by providing both the structured context and the original question. The model is instructed to focus primarily on the context and refer to the original question only if certain information appears unclear or ambiguous. The adjusted prompt is shown in Listing 13 in Sec. A.5.

When queried with both the context and the original question, the model achieves a slight improvement in Pass@5 (+0.8 percentage points), accompanied by a 19% increase in token count. Nevertheless, the total number of generated tokens remains substantially lower than in the 0-CoT setting, approximately 41% fewer tokens overall. We conclude that this prompting strategy can be beneficial in scenarios where precomputing a precise context is challenging due to unclear phrasing or potentially ambiguous interpretation of the question.

**Q4: Do larger models compute better contexts?** We further examine whether larger models produce more reliable context information than smaller ones. To evaluate this, we generate context for all MATH-500 samples using different Qwen 3 model sizes (ranging from 0.6B to 32B), and then perform the F-CoT reasoning step with the Qwen 3 14B model.

As shown in Tab. 5 in Sec. C.3, larger models ($\geq$14B) indeed produce higher-quality context, leading to improved performance. However, these gains saturate at the 14B scale, and we do not expect further improvements beyond this point. The experiments also show that models with 4B parameters or more consistently generate valid XML-like context, whereas the 1.7B model does so in about 95% of cases. Even though the 0.6B model almost always produces context that is invalid from a strict XML standpoint, the 14B model still achieves a Pass@5 accuracy of 41.2%, indicating that perfect XML format validity is not strictly necessary for correct reasoning.

## 5 Discussion

We finish the paper with a general discussion of the limitations and possible future directions.

### 5.1 Limitations

When automatically creating context windows with LLMs, information loss may occasionally occur during the extraction process. For example, we observed rare cases in which instructions such as "state the answer in percentage values" or other details critical for interpreting the input are omitted, causing the model to output decimal numbers instead (e.g., 0.05 instead of 5%). Such discrepancies can negatively impact performance on math benchmarks.

In other cases, the extracted information may be misinterpreted, which can subsequently affect the downstream reasoning process. Thanks to the structured format of our context, however, such issues can be easily identified and corrected. We also found that transforming already highly condensed information, such as AIME questions, into our context scheme does not improve reasoning correctness and can sometimes even degrade it. Moreover, smaller models, particularly the 0.6B variant, struggle with self-generated context extraction. In these cases, a larger model must be used to pre-compute the context to achieve speedups without significantly harming the smaller model's capabilities.

## 5.2 Future Directions

Until now, we have explored F-CoT in isolation from other prompting techniques. However, combining a structured representation of input information with more advanced prompting strategies could further reduce the number of generated tokens or improve performance. Another promising avenue is to integrate F-CoT with test-time scaling methods such as tree-of-thoughts, potentially accelerating them while enabling richer exploration.

Although our focus is on reasoning-oriented language models, we believe that similar structured information extraction could also benefit multimodal models. For example, when querying a vision-language model with an image, first extracting key visual elements, either in a structured format or as image excerpts, before performing higher-level reasoning may yield efficiency gains similar to those observed in our setting. Another direction is to investigate how structured reasoning can be incorporated directly into model training or fine-tuning, improving the model's inherent ability to identify and extract salient information. Currently, the evaluated models were not trained with such structured inputs, which may limit their native capacity to process them.

We also envision using the context as a dynamic notepad, where the model can store intermediate reasoning steps by updating existing entries or adding new ones. Thanks to the fixed XML structure, context updates can be introduced incrementally. Adjusting the model's prefix to reflect these updates could help concentrate the model's attention on the curated context rather than on thousands of previously generated tokens. However, this approach would likely require additional fine-tuning and would necessitate recomputing cached key–value states whenever the context is modified.

## 6 Conclusion

We present Focused Chain-of-Thought (F-CoT), a simple yet highly effective prompting strategy that explicitly separates information extraction from the reasoning process. By querying reasoning LLMs with structured context rather than raw natural language, F-CoT addresses the issue of excessive token usage, achieving a 2-3x reduction in generated tokens while preserving reasoning performance. Our analysis reveals that this efficiency stems from a reduction in filler and redundant reasoning steps, as structured inputs naturally mitigate the model's tendency to overthink. Ultimately, these results demonstrate that optimizing input representations offers a powerful, training-free alternative to model-centric optimization for efficient deployment.

### References

Pranjal Aggarwal and Sean Welleck. L1: Controlling how long a reasoning model thinks with reinforcement learning. *arXiv preprint arXiv:2503.04697*, 2025.

J.R. Anderson. *Language, Memory, and Thought*. Experimental psychology series. L. Erlbaum Associates, 1976. ISBN 9780898591071.

Daman Arora and Andrea Zanette. Training language models to reason efficiently. *arXiv preprint arXiv:2502.04463*, 2025.

Xingyu Chen, Jiahao Xu, Tian Liang, Zhiwei He, Jianhui Pang, Dian Yu, Linfeng Song, Qiuzhi Liu, Mengfei Zhou, Zhuosheng Zhang, Rui Wang, Zhaopeng Tu, Haitao Mi, and Dong Yu. Do NOT think that much for 2+3=? on the overthinking of long reasoning models. In *Forty-second International Conference on Machine Learning*, 2025.

Cheng-Han Chiang and Hung-yi Lee. Over-reasoning and redundant calculation of large language models. *arXiv preprint arXiv:2401.11467*, 2024.

Karl Cobbe, Vineet Kosaraju, Mohammad Bavarian, Mark Chen, Heewoo Jun, Lukasz Kaiser, Matthias Plappert, Jerry Tworek, Jacob Hilton, Reiichiro Nakano, Christopher Hesse, and John Schulman. Training verifiers to solve math word problems. *arXiv preprint arXiv:2110.14168*, 2021.

American Mathematics Competitions. Aime problems and solutions, 2025. `https://artofproblemsolving.com/wiki/index.php/AIME_Problems_and_Solutions`.

Alejandro Cuadron, Dacheng Li, Wenjie Ma, Xingyao Wang, Yichuan Wang, Siyuan Zhuang, Shu Liu, Luis Gaspar Schroeder, Tian Xia, Huanzhi Mao, et al. The danger of overthinking: Examining the reasoning-action dilemma in agentic tasks. *arXiv preprint arXiv:2502.08235*, 2025.

DeepSeek-AI. Deepseek-r1: Incentivizing reasoning capability in llms via reinforcement learning, 2025.

Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha Letman, Akhil Mathur, Alan Schelten, Amy Yang, Angela Fan, et al. The llama 3 herd of models. *arXiv e-prints*, pp. arXiv–2407, 2024.

Federico Errica, Davide Sanvito, Giuseppe Siracusano, and Roberto Bifulco. What did i do wrong? quantifying llms' sensitivity and consistency to prompt engineering. In *Proceedings of the 2025 Conference of the Nations of the Americas Chapter of the Association for Computational Linguistics: Human Language Technologies*, pp. 1543–1558, 2025.

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

Shahriar Golchin and Mihai Surdeanu. Time travel in llms: Tracing data contamination in large language models. *arXiv preprint arXiv:2308.08493*, 2023.

Dan Hendrycks, Collin Burns, Saurav Kadavath, Akul Arora, Steven Basart, Eric Tang, Dawn Song, and Jacob Steinhardt. Measuring mathematical problem solving with the math dataset. *NeurIPS*, 2021.

Daniel Kahneman. Thinking, fast and slow. *Farrar, Straus and Giroux*, 2011.

Amanda Kau, Xuzeng He, Aishwarya Nambissan, Aland Astudillo, Hui Yin, and Amir Aryani. Combining knowledge graphs and large language models. *arXiv preprint arXiv:2407.06564*, 2024.

Takeshi Kojima, Shixiang Shane Gu, Machel Reid, Yutaka Matsuo, and Yusuke Iwasawa. Large language models are zero-shot reasoners. *Advances in neural information processing systems*, 35: 22199–22213, 2022.

Ayeong Lee, Ethan Che, and Tianyi Peng. How well do llms compress their own chain-of-thought? a token complexity approach. *arXiv preprint arXiv:2503.01141*, 2025.

Patrick Lewis, Ethan Perez, Aleksandra Piktus, Fabio Petroni, Vladimir Karpukhin, Naman Goyal, Heinrich Küttler, Mike Lewis, Wen-tau Yih, Tim Rocktäschel, et al. Retrieval-augmented generation for knowledge-intensive nlp tasks. *Advances in neural information processing systems*, 33: 9459–9474, 2020.

Feng Luo, Yu-Neng Chuang, Guanchu Wang, Hoang Anh Duy Le, Shaochen Zhong, Hongyi Liu, Jiayi Yuan, Yang Sui, Vladimir Braverman, Vipin Chaudhary, et al. Autol2s: Auto long-short reasoning for efficient large language models. *arXiv preprint arXiv:2505.22662*, 2025.

Chuangtao Ma, Yongrui Chen, Tianxing Wu, Arijit Khan, and Haofen Wang. Large language models meet knowledge graphs for question answering: Synthesis and opportunities. *arXiv preprint arXiv:2505.20099*, 2025.

Aman Madaan, Niket Tandon, Prakhar Gupta, Skyler Hallinan, Luyu Gao, Sarah Wiegreffe, Uri Alon, Nouha Dziri, Shrimai Prabhumoye, Yiming Yang, et al. Self-refine: Iterative refinement with self-feedback. *Advances in Neural Information Processing Systems*, 36:46534–46594, 2023.

Niklas Muennighoff, Zitong Yang, Weijia Shi, Xiang Lisa Li, Li Fei-Fei, Hannaneh Hajishirzi, Luke Zettlemoyer, Percy Liang, Emmanuel Candès, and Tatsunori Hashimoto. s1: Simple test-time scaling. *arXiv preprint arXiv:2501.19393*, 2025.

Sania Nayab, Giulio Rossolini, Marco Simoni, Andrea Saracino, Giorgio Buttazzo, Nicolamaria Manes, and Fabrizio Giacomelli. Concise thoughts: Impact of output length on llm reasoning and cost. *arXiv preprint arXiv:2407.19825*, 2024.

Yonatan Oren, Nicole Meister, Niladri S Chatterji, Faisal Ladhak, and Tatsunori Hashimoto. Proving test set contamination in black-box language models. In *The Twelfth International Conference on Learning Representations*, 2023.

Shirui Pan, Linhao Luo, Yufei Wang, Chen Chen, Jiapu Wang, and Xindong Wu. Unifying large language models and knowledge graphs: A roadmap. *IEEE Transactions on Knowledge and Data Engineering*, 36(7):3580–3599, 2024.

Arkil Patel, Satwik Bhattamishra, and Navin Goyal. Are nlp models really able to solve simple math word problems? *arXiv preprint arXiv:2103.07191*, 2021.

Qwen Team. Qwen3 technical report, 2025.

Melanie Sclar, Yejin Choi, Yulia Tsvetkov, and Alane Suhr. Quantifying language models' sensitivity to spurious features in prompt design or: How i learned to start worrying about prompt formatting. In *The Twelfth International Conference on Learning Representations*, 2024.

Xiaoyu Tian, Sitong Zhao, Haotian Wang, Shuaiting Chen, Yunjie Ji, Yiping Peng, Han Zhao, and Xiangang Li. Think twice: Enhancing llm reasoning by scaling multi-round test-time thinking. *arXiv preprint arXiv:2503.19855*, 2025.

Lei Wang, Wanyu Xu, Yihuai Lan, Zhiqiang Hu, Yunshi Lan, Roy Ka-Wei Lee, and Ee-Peng Lim. Plan-and-solve prompting: Improving zero-shot chain-of-thought reasoning by large language models. *arXiv preprint arXiv:2305.04091*, 2023.

Yizhong Wang, Yeganeh Kordi, Swaroop Mishra, Alisa Liu, Noah A Smith, Daniel Khashabi, and Hannaneh Hajishirzi. Self-instruct: Aligning language models with self-generated instructions. *arXiv preprint arXiv:2212.10560*, 2022.

Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Fei Xia, Ed Chi, Quoc V Le, Denny Zhou, et al. Chain-of-thought prompting elicits reasoning in large language models. *Advances in neural information processing systems*, 35:24824–24837, 2022.

Heming Xia, Chak Tou Leong, Wenjie Wang, Yongqi Li, and Wenjie Li. Tokenskip: Controllable chain-of-thought compression in llms. *arXiv preprint arXiv:2502.12067*, 2025.

Yang Xiao, Jiashuo Wang, Ruifeng Yuan, Chunpu Xu, Kaishuai Xu, Wenjie Li, and Pengfei Liu. Limopro: Reasoning refinement for efficient and effective test-time scaling. *arXiv preprint arXiv:2505.19187*, 2025.

Silei Xu, Wenhao Xie, Lingxiao Zhao, and Pengcheng He. Chain of draft: Thinking faster by writing less. *arXiv preprint arXiv:2502.18600*, 2025a.

Xin Xu, Tong Xiao, Zitong Chao, Zhenya Huang, Can Yang, and Yang Wang. Can LLMs solve longer math word problems better? In *The Thirteenth International Conference on Learning Representations*, 2025b.

Shunyu Yao, Dian Yu, Jeffrey Zhao, Izhak Shafran, Tom Griffiths, Yuan Cao, and Karthik Narasimhan. Tree of thoughts: Deliberate problem solving with large language models. *Advances in neural information processing systems*, 36:11809–11822, 2023a.

Shunyu Yao, Jeffrey Zhao, Dian Yu, Nan Du, Izhak Shafran, Karthik Narasimhan, and Yuan Cao. React: Synergizing reasoning and acting in language models. In *International Conference on Learning Representations (ICLR)*, 2023b.

Guanghao Ye, Khiem Duc Pham, Xinzhi Zhang, Sivakanth Gopi, Baolin Peng, Beibin Li, Janardhan Kulkarni, and Huseyin A Inan. On the emergence of thinking in llms i: Searching for the right intuition. *arXiv preprint arXiv:2502.06773*, 2025.

Edward Yeo, Yuxuan Tong, Morry Niu, Graham Neubig, and Xiang Yue. Demystifying long chain-of-thought reasoning in llms. *arXiv preprint arXiv:2502.03373*, 2025.

Bin Yu, Hang Yuan, Haotian Li, Xueyin Xu, Yuliang Wei, Bailing Wang, Weizhen Qi, and Kai Chen. Long-short chain-of-thought mixture supervised fine-tuning eliciting efficient reasoning in large language models. *arXiv preprint arXiv:2505.03469*, 2025.

Hang Yuan, Bin Yu, Haotian Li, Shijun Yang, Christina Dan Wang, Zhou Yu, Xueyin Xu, Weizhen Qi, and Kai Chen. Not all tokens are what you need in thinking. *arXiv preprint arXiv:2505.17827*, 2025.

Wanjun Zhong, Yifan Gao, Ning Ding, Zhiyuan Liu, Ming Zhou, Jiahai Wang, Jian Yin, and Nan Duan. Improving task generalization via unified schema prompt. *AI Open*, 4:120–129, 2023.

Ren Zhuang, Ben Wang, and Shuifa Sun. Accelerating chain-of-thought reasoning: When goal-gradient importance meets dynamic skipping. *arXiv preprint arXiv:2505.08392*, 2025.

# A PROMPT DESIGNS

## A.1 CONTEXT EXTRACTION WITH LLMs WITHOUT EXAMPLES

In our experiment, we used the following prompt shown in Listing 1 to instruct the model to extract context information from input data.

```
1  Please extract all critical information and the underlying question from
       the given sample.
2  The output must follow this format:
3  <context>
4  <info_1>Information 1</info_1>
5  <info_2>Information 2</info_2>
6  <question>The question that needs to be answered based on the information
        provided.</question>
7  </context>
8
9  Guidelines:
10 - Only include information directly relevant to answering the question.
11 - Write information in concise, factual statements.
12 - If quantities, measurements, or mathematical relations are given,
       include them (use LaTeX for numbers and units when appropriate, e.g.,
       $4.5 \\, \\text{inches}$).
13 - Ignore irrelevant background or narrative details.
14 - Restate the question clearly and concisely, keeping only what is
       necessary to answer it.
15 - Do not solve the problem or perform calculations.
16
17 Please extract the information and question from the following sample:
18
19 [ORIGINAL QUESTION]
```

Listing 1: The prompt used to generate context information.

In our experiment, we used the prompt shown in Listing 2 to instruct the model to extract context information from the GSM-hard and SVAMP data. Analogously, the prompt in Listing 3 is used to extract context information from MATH-500 and AIME data.

```
1  Please extract all critical information and the underlying question from
       the given sample.
2  The output must follow this format:
3  <context>
4  <info_1>Information 1</info_1>
5  <info_2>Information 2</info_2>
6  <question>The question that needs to be answered based on the information
       provided.</question>
7  </context>
8
9  Guidelines:
10 - Only include information directly relevant to answering the question.
11 - Write information in concise, factual statements
12 - If quantities, measurements, or mathematical relations are given,
       include them (use LaTeX for numbers and units when appropriate, e.g.,
       $4.5 , \text{inches}$).
13 - Ignore background or story-like details.
14 - Restate the question in a minimal, clear form. Avoid repeating
       irrelevant phrasing.
15 - Do not solve the problem or perform calculations.
16
17 Here is an example:
18 Question: Three of the women at the cocktail party are wearing 4.5-inch
       heels and three are wearing 2.5-inch heels. What is the average
       height of heels at this party?
19
20 Desired Output:
21 <context>
22 <info_1>Three women wearing 4.5-inch heels</info_1>
23 <info_2>Three women wearing 2.5-inch heels</info_2>
24 <question>What is the average heel height?</question>
25 </context>
26
27 Here is another example:
28 Question: Fishio posted her selfie on Instagram. She received 20000 likes
        on the photo after 1 week. Three weeks later, the number of likes
       was 70 times as many as the initial number of likes. If she received
       200000 more new likes recently, how many Instagram likes are there?
29
30 Desired Output:
31 <context>
32 <info_1>Initial likes after 1 week: 20000</info_1>
33 <info_2>Likes after 3 weeks: \(70(20000)\)</info_2>
34 <info_3>Additional $200000$ likes received after the 3-week count</info_3
       >
35 <question>How many Instagram likes are there?</question>
36 </context>
37
38 Please extract the information and question from the following sample:
39
40 [ORIGINAL QUESTION]
```

Listing 2: The prompt used to generate context information from GSM-hard and SVAMP samples. Model: GPT-5 Mini (August 2025 version).

```
1
2  Please extract all critical information and the underlying question from
       the given sample.
3
4  The output must follow this format:
5
6  <context>
7
8  <info_1>Information 1</info_1>
9
10 <info_2>Information 2</info_2>
11
12 <question>The question that needs to be answered based on the information
       provided.</question>
13
14 </context>
15
16 Guidelines:
17
18 - Only include information directly relevant to answering the question.
19
20 - Write information in concise, factual statements
21
22 - If quantities, measurements, or mathematical relations are given,
       include them (use LaTeX for numbers and units when appropriate, e.g.,
       $4.5 , \text{inches}$).
23
24 - Ignore background or story-like details.
25
26 - Restate the question in a minimal, clear form. Avoid repeating
       irrelevant phrasing.
27
28 - Do not solve the problem or perform calculations.
29
30 Here is an example:
31
32 Question: Let \\[f(x) = \\left\\{\n\\begin{array}{cl} ax+3, &\\text{ if }
       x>2, \\\\\nx-5 &\\text{ if } -2 \\le x \\le 2, \\\\\n2x-b &\\text{ if
       } x <-2.\n\\end{array}\n\\right.\\]Find $a+b$ if the piecewise
       function is continuous (which means that its graph can be drawn
       without lifting your pencil from the paper).
33
34 Desired Output:
35
36 <context>
37
38 <info_1>Piecewise function:
39
40 \[
41
42 f(x) =
43
44 \\begin{cases}
45
46 ax + 3, &   ext{if } x > 2 \
47
48 x - 5, &  ext{if } -2 \le x \le 2 \
49
50 2x - b, &   ext{if } x < -2
51
52 \end{cases}
53
54 \]</info_1>
55
```

```
56 <info_2>The function is continuous (can be drawn without lifting the
      pencil)</info_2>
57
58 <question>Find $a+b$.</question>
59
60 </context>
61
62 Here is another example:
63
64 Question: A rectangular band formation is a formation with $m$ band
      members in each of $r$ rows, where $m$ and $r$ are integers. A
      particular band has less than 100 band members. The director arranges
       them in a rectangular formation and finds that he has two members
      left over. If he increases the number of members in each row by 1 and
       reduces the number of rows by 2, there are exactly enough places in
      the new formation for each band member. What is the largest number of
       members the band could have?
65
66 Desired Output:
67
68 <context>
69
70 <info_1>Band has less than 100 members</info_1>
71
72 <info_2>Original formation: $m$ members per row, $r$ rows, with 2 members
       left over</info_2>
73
74 <info_3>New formation: $(m+1)$ members per row, $(r-2)$ rows, exactly
      enough places for all members</info_3>
75
76 <question>What is the largest possible number of band members?</question>
77
78 </context>
79
80 Please extract the information and question from the following sample:
81
82 [ORIGINAL QUESTION]
83 0
```

Listing 3: The prompt used to generate context information from MATH-500 and AIME samples.
Model: GPT-5 Mini (August 2025 version).

### A.3 0-SHOT CHAIN-OF-THOUGHT PROMPTS

```
1 [ORIGINAL QUESTION]
2
3 Please reason step by step, and put your final answer within \\boxed{}.
```
Listing 4: Prompt used for 0-shot chain-of-thought reasoning (0-CoT).

```
1 "You are given ONLY the information in the following question. Use ONLY
     the facts provided by the question to compute the answer.
2
3 [ORIGINAL QUESTION]
4
5 Please reason step by step, and put your final answer within \\boxed{}.
```
Listing 5: The shortened prompt used to elicit reasoning across the standard question.

### A.4 BASELINE PROMPTS

```
1 [ORIGINAL QUESTION]
2
3 Let's first understand the problem, extract relevant variables and their
     corresponding numerals, and make and devise a complete plan. Then,
     let's carry out the plan, calculate intermediate variables (pay
     attention to correct numerical calculation and commonsense), solve
     the problem step by step, and show the answer.
4
5 Please reason step by step, and put your final answer within \\boxed{}.
```
Listing 6: Prompt used for Plan and Solve reasoning (PS).

```
1 [ORIGINAL QUESTION]
2
3 Let's first understand the problem, then list all the known conditions,
     which are formed by numbers or quantitative relationships along with
     their contexts from the problem text, and identify the final goal of
     the problem.
4
5 Please reason step by step, and put your final answer within \\boxed{}.
```
Listing 7: Prompt used for CoRe reasoning.

## A.5 F-CoT Reasoning with Context Blocks

In our F-CoT experiment, we used the following prompts shown in Listing 8 to instruct the model to reason based on the provided context. The other prompts correspond to different variations from the sensitivity analysis.

```
1 You are given ONLY the structured <context> below. Use ONLY the facts
    inside <context> to compute the answer. Do NOT reference or repeat
    the original natural-language question or any outside information.
2
3 Instructions:
4 1) Show step-by-step reasoning between <think> and </think> blocks.
5 2) In your reasoning, explicitly cite the relevant <info_k> entries when
    you use them, including the angle brackets (e.g., 'From <info_1>,
    ...').
6 3) After the reasoning, output the final answer in \\boxed{}.
7
8 <context>
9 [INSERTED CONTEXT]
10 </context>
11
12 Please reason step by step, and put your final answer within \\boxed{}.
```
Listing 8: The standard F-CoT prompt used to elicit reasoning across the context.

```
1 "You are given ONLY the structured <context> below. Use ONLY the facts
    inside <context> to compute the answer.
2
3 <context>
4 [INSERTED CONTEXT]
5 </context>
6
7 Please reason step by step, and put your final answer within \\boxed{}.
```
Listing 9: The shortened prompt from the sensitivity analysis used to elicit reasoning across the context.

```
1 You are given ONLY the structured context below. Use ONLY the provided
    facts to compute the answer. Do NOT reference or repeat the original
    natural-language question or any outside information
2
3 Instructions:
4 1) Show step-by-step reasoning between <think> and </think> blocks.
5 2) In your reasoning, explicitly cite the relevant info_k entries when
    you use them (e.g., 'From info_1, ...').
6 3) After the reasoning, output the final answer in \\boxed{}.
7
8 [INSERTED CONTEXT - VARIATION 1]
9
10 Please reason step by step, and put your final answer within \\boxed{}.
```
Listing 10: Adjusted prompt from the sensitivity analysis for context variation 1.

```
1 You are given ONLY the structured information below. Use ONLY the
     provided facts to compute the answer. Do NOT reference or repeat the
     original natural-language question or any outside information.
2
3 Instructions:
4 1) Show step-by-step reasoning between <think> and </think> blocks.
5 2) In your reasoning, explicitly cite the relevant information when you
     use them.
6 3) After the reasoning, output the final answer in \\boxed{}.
7
8 [INSERTED CONTEXT - VARIATION 2]
9
10 Please reason step by step, and put your final answer within \\boxed{}.
```
Listing 11: Adjusted prompt from the sensitivity analysis for context variation 2.

```
1 You are given ONLY the structured information below. Use ONLY the
     provided facts to compute the answer. Do NOT reference or repeat the
     original natural-language question or any outside information.
2
3 Instructions:
4 1) Show step-by-step reasoning between <think> and </think> blocks.
5 2) After the reasoning, output the final answer in \\boxed{}.
6
7 [INSERTED CONTEXT - VARIATION 3]
8
9 Please reason step by step, and put your final answer within \\boxed{}.
```
Listing 12: Adjusted prompt from the sensitivity analysis for context variation 3.

```
1 You are given the structured <context> below and the original question.
     Use ONLY the facts inside <context> to compute the answer. Refer to
     the original question ONLY if some information in <context> is
     unclear or ambiguous.
2
3 Instructions:
4 1) Show step-by-step reasoning between <think> and </think> blocks.
5 2) In your reasoning, explicitly cite the relevant information when you
     use them.
6 3) After the reasoning, output the final answer in \\boxed{}.
7
8 Original question: [INSERTED QUESTION]
9
10 Context: [INSERTED CONTEXT]
11
12 Please reason step by step, and put your final answer within \\boxed{}.
```
Listing 13: Adjusted prompt from the sensitivity analysis for context and original question.

## A.6 Context Designs

The context design in Listing 14 corresponds to the default context variant. The other context designs are explored in our sensitivity analysis.

```
1 <context>
2     <info_1>Price per apple: $2<info_1>
3     <info_2>Price per pear: 1.5x price per apple<info_2>
4     <question>Total cost of 2 apples and 2 pears</question>
5 </context>
```
Listing 14: Standard F-CoT context with XML-like format.

```
1 Context:
2 info_1: Price per apple: $2
3 info_2: Price per pear: 1.5x price per apple
4 question: Total cost of 2 apples and 2 pears</question>
```
Listing 15: Context alternative design 1 - enumerated list.

```
1 **Information**
2 - Price per apple: $2
3 - Price per pear: 1.5x price per apple
4
5 **Question**
6 Total cost of 2 apples and 2 pears</question>
```
Listing 16: Context alternative design 2 - unnumbered list.

```
1 Price per apple: $2. Price per pear: 1.5x price per apple. Question:
    Total cost of 2 apples and 2 pears</question>
```
Listing 17: Context alternative design 3 - concatenated information.

## A.7 SENTENCE TYPE ANNOTATION

```
1
2  You are an annotator distinguishing between *extraction* and *reasoning*
       in model-generated solutions.
3
4  - **Extraction** = A sentence that copies or paraphrases information
       directly from the input question without transformation or inference.
5
6  - **Reasoning** = A sentence that performs logical inference, arithmetic,
        explanation, or deduction not explicitly given in the input.
7
8  - **Filler** = A sentence that does not contain extraction or reasoning,
       such as generic phrases ("Let's solve step by step", "Now I will
       calculate").
9
10 Instructions:
11
12 1. I will give you a question and a step-by-step solution.
13
14 2. Split the solution into individual sentences.
15
16 3. For each sentence, output one of the following labels: `Extraction`, `
       Reasoning`, or `Filler`.
17
18 4. Return the result as a JSON list of objects, where each object has two
        keys: `"sentence"` and `"label"`.
19
20 ---
21
22 Question: [ORIGINAL QUESTION]
23
24 Solution: [MODEL GENERATION]
25
26 Output format:
27
28 [
29
30 {"sentence": "...", "label": "..."},
31
32 {"sentence": "...", "label": "..."}
33
34 ]
```

Listing 18: Prompt used annotating the sentence type.

## A.8 Overthinking Score

```
1
2  You are an AI judge focused on detecting when a model prefers its
       internal reasoning chain over performing concrete mathematical
       reasoning steps that advance toward a solution.
3
4  <INTERACTION>
5
6  User Question: [ORIGINAL QUESTION]
7
8  Model Solution: [MODEL GENERATION]
9
10 </INTERACTION>
11
12 Analyze the <INTERACTION> and determine if the model is preferring
       internal reasoning and reflection over actually performing concrete
       steps to solve the problem.
13
14 How could this be detected?
15
16 <CORE PRINCIPLE>
17 - Analysis Paralysis:
18   The model repeatedly analyzes strategies or discusses what could be
        done (e.g., substitution, factoring, induction) without executing any
         mathematical steps.
19 - Rogue Reasoning:
20   The model jumps between multiple methods (e.g., algebraic -> geometric
        -> numerical) without completing any of them or checking intermediate
         results.
21 - Premature Disengagement:
22   The model concludes the problem without verifying or deriving the
        answer, often due to overconfidence ("The answer must be X") or self-
        dismissal ("This can't be solved").
23 </CORE PRINCIPLE>
24
25 <SCORING SYSTEM (0-10)>
26 0-3: Clear, step-by-step reasoning
27 - Progresses directly toward the solution through algebraic, numerical,
        or logical steps.
28 - A brief plan or method selection is fine if quickly followed by
        computation.
29 - Checks and verifies results.
30 - Repeats calculations or re-derives expressions carefully.
31 - Stays focused on the mathematical task.
32
33 4-7: Some overthinking but still problem-solving
34 - Includes long reasoning or multiple reflections before computing, but
        eventually performs concrete steps.
35 - Might discuss multiple solution paths once but settles on one.
36 - Verification is present but delayed.
37 - Shows minor hesitation or digression but remains mostly focused.
38
39 8-10: Strong overthinking / minimal problem interaction
40 - Engages in extended speculation about possible methods or properties
        without any real computation.
41 - Switches between multiple unrelated strategies.
42 - Concludes without a derivation.
43 - Produces long internal dialogue instead of concrete reasoning steps.
44 - Gets stuck in unproductive self-reflection.
45 </SCORING SYSTEM>
46
47 <ANALYSIS STEPS>
48 1. Analysis Paralysis
49    - Does the model spend a long time discussing potential approaches?
```

```
50      - Are there few or no concrete equations, simplifications, or
         calculations?

51
52  2. Rogue Reasoning
53      - Does it jump between many approaches without applying or finishing
         any?
54      - Are intermediate results ignored or overwritten by new plans?

55
56  3. Premature Disengagement
57      - Does it state a final answer without derivation or checking?
58      - Does it give up or claim it cannot solve the problem?
59  </ANALYSIS STEPS>

60
61  <EXAMPLES>
62  Example 1 - Focused Reasoning (Good):
63  Problem: Solve x^2 - 5x + 6 = 0.
64  Model:
65  "Let's compute the discriminant: D = 25 - 24 = 1.
66  Then x = (5 \pm 1)/2, so x = 2 or x = 3."
67  Score: 0 - The model works step-by-step, no unnecessary reflection.

68
69  Example 2 - Brief Planning (Good):
70  Model:
71  "This looks quadratic, so I can either factor or use the quadratic
         formula. Let's try factoring: (x-2)(x-3) = 0, giving roots 2 and 3."
72  Score: 2 - Some internal planning, but quickly applied and solved.

73
74  Example 3 - Mild Overthinking:
75  Model:
76  "We might approach this by factoring, or perhaps by completing the square
         . Factoring works if the product and sum fit correctly... Hmm, 2×3 =
         6 and 2+3 = 5, that works. So, yes, x=2,3."
77  Score: 5 - Some unnecessary deliberation, but correct and eventually
         concrete.

78
79  Example 4 - Analysis Paralysis:
80  Model:
81  "This problem could be tackled in multiple ways - factoring, completing
         the square, or graphing. Maybe I should first recall the discriminant
          formula, but perhaps the sum-product relation is easier.
         Alternatively, we might think geometrically about the roots of a
         parabola... On the other hand..."
82  Score: 9 - Prolonged reflection with no actual computation.

83
84  Example 5 - Rogue Reasoning:
85  Model:
86  "Let's test some integer values: x=0 gives 6, x=1 gives 2, x=2 gives 0 so
          that works. But maybe we can confirm using the discriminant method
         ... or perhaps symmetry arguments from the vertex form... though also
          we might estimate roots by graph intersection."
87  Score: 8 - Multiple scattered strategies without systematic follow-
         through.

88
89  Example 6 - Premature Disengagement:
90  Model:
91  "The answer should probably be around 2 or 3 because it looks symmetric."
92  Score: 10 - Concludes without derivation or verification.
93  </EXAMPLES>

94
95  <IMPORTANT>
96  Format your response as:
97  <answer>
98  {

99
100 "overthinking_score": "[0-10]",
```

```
101  "reasoning": "Explain your reasoning for the score, in one line suitable
         for JSON."
102  }
103  </answer>
104
105  Always surround your answer with <answer> and </answer> tags.
106  Take your time to understand the reasoning trace and analyze it carefully
         .
107  Think step-by-step about whether the model overthinks (reflects
         excessively) instead of reasoning concretely toward a mathematical
         solution.
108  </IMPORTANT>
```

Listing 19: The prompt used to compute the overthinking score. Model: GPT-5 Mini (August 2025 version).

## B QUALITATIVE RESULTS

### B.1 GSM-HARD RESULTS

```
1 <think>
2 Okay, let's see. Janet has ducks that lay 16 eggs each day. She eats
      three for breakfast every morning. Then she bakes muffins using
      4933828... wait, that number seems really big. Maybe that's a typo?
      Let me check the problem again.
3
4 "Janet's ducks lay 16 eggs per day. She eats three for breakfast every
      morning and bakes muffins for her friends every day with 4933828. She
       sells the remainder at the farmers' market daily for $2 per fresh
      duck egg. How much in dollars does she make every day at the farmers'
       market?"
5
6 Hmm, baking muffins with 4933828... that number is way too large. Maybe
      it's supposed to be 493,382.8? Or maybe 493.3828? Or perhaps it's a
      different number? Wait, maybe it's a typo and the number is supposed
      to be 493, but written as 4933828? Let me think. If the number is
      493,382.8, that would be like 493 thousand eggs, which is way more
      than the 16 eggs she has. That doesn't make sense. If she only has 16
       eggs a day, she can't use 493 thousand eggs to bake muffins. So
      maybe there's a mistake in the number? Maybe it's 4.933828? Or maybe
      it's 49 eggs? Wait, the original problem says "bakes muffins for her
      friends every day with 4933828." Maybe that's supposed to be 49 eggs?
       Let me check again.
7
8 Wait, the original problem says "bakes muffins for her friends every day
      with 4933828." Maybe that's a typo and should be 49 eggs? Or maybe
      4.933828? But the number is written as 4933828. Maybe it's a misinput
      . Let me check if that number is related to something else.
      Alternatively, maybe the number is correct, but it's in a different
      unit? Wait, that seems unlikely.
9
10 [...]
11
12 But given that the original number is 4933828, which is 4,933,828, and if
       it's a misinput of 4.933828, then the answer is approximately $16
      .13.
13
14 However, since the problem might have a different intended number, like
      493,382.8 divided by 100,000, which is 4.933828, I think that's the
      best I can do.
15
16 So, I'll go with $16.13 as the answer, but I'm not fully confident due to
       the ambiguity in the number provided.
17
18 **Final Answer**
19 \boxed{16.13}
```

Listing 20: Generated output for Qwen3 14B queried with a GSM-Hard Example and 0-CoT. The total number of generated tokens is 6,401. Parts of the output are shortened, indicated by ... .

```
1  Okay, let's see. The problem is asking how much money she makes per day
       at the farmers' market. The context gives me some information about
       eggs.
2
3  First, from info_1, ducks lay 16 eggs per day. So total eggs available
       each day are 16.
4
5  Then, info_2 says she eats 3 eggs for breakfast every morning. So that's
       3 eggs used up by her.
6
7  Next, info_3 mentions she uses 4,933,828 eggs daily to bake muffins. Wait
       , that number seems really large. Let me check. Oh, maybe that's a
       typo? Because if ducks only lay 16 eggs a day, how can she use over 4
        million eggs for muffins? That doesn't make sense. But the problem
       says to use only the given facts. Maybe there's a mistake in the
       numbers? Or maybe the numbers are correct? Wait, maybe the info_3 is
       a typo. Let me think again. If the numbers are as given, then perhaps
        the problem is designed with those numbers. But that would mean that
        the total eggs used would be 3 (breakfast) + 4,933,828 (muffins) =
       4,933,831 eggs. But the total eggs laid per day are only 16. So that
       would mean she's using way more eggs than are available. That's
       impossible. But since the problem says to use only the given facts,
       maybe I need to proceed regardless. Wait, maybe the info_3 is a
       mistake, like maybe it's supposed to be 4.933828 eggs? But the number
        is written as 4,933,828, which is 4 million. Hmm. Alternatively,
       maybe the problem is designed to have a negative number of eggs sold?
        But that would be odd. Alternatively, maybe the numbers are correct
       and the answer is that she can't sell any eggs because she uses more
       than she has. But the problem says to compute the answer based on the
        given info. Let me check again.
8
9  [...]
10
11 This is confusing. But the user says to use only the facts inside the
       context. So I need to proceed with the numbers given.
12
13 [...]
14
15 This is a tricky situation. Given the instructions, I have to use the
       given numbers, even if they are inconsistent. So the answer would be
       the calculation based on the numbers, even if it's negative.
       Therefore, the final answer would be \boxed{-9867630}. But I'm not
       sure if that's what the problem expects. However, given the strict
       instructions, this is the only possible answer.
16 </think>
17
18 [...]
19 ### Final Answer
20 $$
21 \boxed{-9867630}
22 $$
```

Listing 21: Generated output for Qwen3 14B queried with a GSM-Hard Example and F-CoT. The total number of generated tokens is 2,820. Parts of the output are shortened, indicated by ... .

# C  QUANTITATIVE RESULTS

## C.1  EXPERIMENTS WITH GPT-5 MINI PRE-COMPUTED CONTEXT INFORMATION

Tab. 1 summarizes the performance of multiple reasoning LLMs across five benchmarks, comparing standard zero-shot CoT (0-CoT) with our F-CoT approach, as well as Plan-and-Solve (PS) and CoRe prompting strategies. For F-CoT, the context is pre-computed using GPT-5 Mini, ensuring structured information is provided to the model before reasoning.

Table 1: Performance of reasoning LLMs on five benchmarks using context pre-computed by GPT-5 Mini.

| Model | Prompting | SVAMP | | GSM-Hard | | MATH-500 | | AIME2024 | | AIME2025 | |
|---|---|---|---|---|---|---|---|---|---|---|---|
| | | Pass@5 | # Tokens | Pass@5 | # Tokens | Pass@5 | # Tokens | Pass@5 | # Tokens | Pass@5 | # Tokens |
| **Qwen 3 0.6B** | 0-CoT | 95.67% | 1,358 | 65.88% | 5,235 | 90.20% | 5,643 | 30.00% | 16,907 | 23.33% | 14,904 |
| | PS | 96.00% | 1,397 | 65.43% | 4,989 | 88.40% | 5,522 | 30.00% | 16,653 | 26.67% | 14,289 |
| | CoRe | 96.33% | 1,380 | 66.34% | 4,854 | 91.00% | 5,140 | 16.67% | 17,833 | 23.33% | 14,035 |
| | F-CoT | 94.33% | 434 | 64.90% | 2,122 | 87.20% | 2,665 | 26.67% | 9,615 | 23.33% | 8,460 |
| **Qwen 3 4B** | 0-CoT | 97.00% | 1,871 | 70.36% | 5,674 | 98.80% | 5,099 | 86.67% | 12,640 | 83.33% | 16,293 |
| | PS | 96.67% | 1,706 | 71.49% | 5,318 | 98.60% | 5,100 | 83.33% | 12,924 | 80.00% | 16,115 |
| | CoRe | 97.00% | 1,844 | 71.95% | 5,509 | 98.40% | 5,005 | 83.33% | 13,541 | 76.67% | 16,695 |
| | F-CoT | 95.00% | 535 | 76.88% | 2,308 | 98.40% | 2,580 | 83.33% | 9,422 | 73.33% | 11,287 |
| **Qwen 3 14B** | 0-CoT | 97.00% | 1,487 | 71.87% | 4,429 | 99.40% | 4,931 | 86.67% | 12,697 | 80.00% | 16,185 |
| | PS | 97.00% | 1,457 | 72.86% | 4,285 | 99.40% | 4,697 | 90.00% | 13,276 | 83.33% | 15,701 |
| | CoRe | 96.67% | 1,481 | 73.31% | 4,253 | 98.80% | 4,692 | 90.00% | 14,255 | 80.00% | 16,086 |
| | F-CoT | 95.67% | 476 | 77.03% | 1,838 | 98.60% | 2,437 | 83.33% | 10,593 | 76.67% | 12,788 |
| **Qwen 3 32B** | 0-CoT | 97.00% | 1,377 | 70.13% | 4,516 | 99.00% | 4,384 | 90.00% | 12,343 | 86.67% | 15,137 |
| | PS | 96.00% | 1,395 | 70.96% | 4,238 | 99.00% | 4,397 | 90.00% | 12,358 | 86.67% | 14,644 |
| | CoRe | 96.33% | 1,464 | 71.27% | 4,387 | 99.20% | 4,688 | 90.00% | 12,505 | 83.33% | 15,456 |
| | F-CoT | 95.67% | 429 | 77.03% | 1,785 | 99.00% | 2,309 | 83.33% | 8,964 | 76.67% | 11,143 |

## C.2  SELF-COMPUTED CONTEXT INFORMATION

Tab. 2 summarizes the performance of multiple reasoning LLMs across five benchmarks, comparing F-CoT using self-generated context with F-CoT using context pre-computed by GPT-5 Mini.

Table 2: Performance of reasoning LLMs on five benchmarks using context computed by the models themselves.

| Model | Prompting | SVAMP | | | GSM-Hard | | | MATH-500 | | | AIME2024 | | | AIME2025 | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | | Pass@5 | # Tokens | Valid Context | Pass@5 | # Tokens | Valid Context | Pass@5 | # Tokens | Valid Context | Pass@5 | # Tokens | Valid Context | Pass@5 | # Tokens | Valid Context |
| **Qwen 3 0.6B** | 0-CoT | 95.67% | 1,358 | N/A | 65.88% | 5,235 | N/A | 90.20% | 5,643 | N/A | 30.00% | 16,907 | N/A | 23.33% | 14,904 | N/A |
| | F-CoT (Provided) | 94.33% | 434 | 100% | 64.90% | 2,122 | 100% | 87.20% | 2,665 | 100% | 26.67% | 9,615 | 100% | 23.33% | 8,460 | 100% |
| | F-CoT (Self-Generated) | 64.00% | 466 | 0.00% | 22.37% | 934 | 0.01% | 24.40% | 1,055 | 1.20% | 3.33% | 2,851 | 3.33% | 3.33% | 2,087 | 0.00% |
| **Qwen 3 4B** | 0-CoT | 97.00% | 1,871 | N/A | 70.36% | 5,674 | N/A | 98.80% | 5,099 | N/A | 86.67% | 12,640 | N/A | 83.33% | 16,293 | N/A |
| | F-CoT (Provided) | 95.00% | 535 | 100% | 76.88% | 2,308 | 100% | 98.40% | 2,580 | 100% | 83.33% | 9,422 | 100% | 73.33% | 11,287 | 100% |
| | F-CoT (Self-Generated) | 91.00% | 786 | 98.67% | 66.11% | 2,403 | 98.94% | 93.80 | 2,990 | 99.80% | 60.00% | 9,253 | 100% | 66.67% | 12,439 | 100.0% |
| **Qwen 3 14B** | 0-CoT | 97.00% | 1,487 | N/A | 71.87% | 4,429 | N/A | 99.40% | 4,931 | N/A | 86.67% | 12,697 | N/A | 80.00% | 16,185 | N/A |
| | F-CoT (Provided) | 95.67% | 476 | 100% | 77.03% | 1,838 | 100% | 98.60% | 2,437 | 100% | 83.33% | 10,593 | 100% | 76.67% | 12,788 | 100% |
| | F-CoT (Self-Generated) | 94.00% | 581 | 99.67% | 70.51% | 1,963 | 99.55% | 98.00% | 2,526 | 100% | 80.00% | 9,389 | 100% | 76.67% | 11,695 | 100% |
| **Qwen 3 32B** | 0-CoT | 97.00% | 1,377 | N/A | 70.13% | 4,516 | N/A | 99.00% | 4,384 | N/A | 90.00% | 12,343 | N/A | 86.67% | 15,137 | N/A |
| | F-CoT (Provided) | 95.67% | 429 | 100% | 77.03% | 1,785 | 100% | 99.00% | 2,309 | 100% | 83.33% | 8,964 | 100% | 76.67% | 11,143 | 100% |
| | F-CoT (Self-Generated) | 95.67% | 556 | 100% | 73.77% | 1,827 | 99.92% | 98.00% | 2,547 | 100% | 86.67% | 9,395 | 100% | 80.00% | 11,040 | 100% |

## C.3  ABLATION STUDY AND SENSITIVITY ANALYSIS

We conduct a series of ablation and sensitivity experiments to assess the robustness of F-CoT under various prompt and context conditions. Tab. 3 shows the effect of different reasoning prompt designs and context representations on the Qwen3 14B model using GPT-5 Mini pre-computed context. Importantly, the information and sentences between the various formats are identical, and only their representation (XML, list, concatenated) is changed.

We observe that while the default F-CoT prompt achieves a substantial reduction in token count compared to 0-CoT, adjusting the prompt or including the original question alongside the pre-computed context can slightly improve accuracy with minimal impact on efficiency. Alternative context formats (enumerated list, unnumbered list, or concatenated text) produce similar accuracy, confirming that F-CoT is largely robust to the exact structure of the extracted information.

Tab. 4 investigates the effect of in-context examples (ICL) for self-generated context extraction on Qwen3 14B. Including few-shot examples has a negligible impact on Pass@5 and token efficiency, though it slightly decreases the proportion of valid context blocks.

Finally, Tab. 5 evaluates how model size affects context quality and reasoning performance. Smaller models (0.6B) struggle to produce valid contexts, resulting in a dramatic drop in accuracy. Larger models consistently generate valid contexts and achieve high Pass@5 scores, highlighting the benefit of using more capable LLMs for context extraction in F-CoT. Overall, these analyses demonstrate that F-CoT is robust across prompt variations, context formats, and model scales, while retaining its efficiency advantages.

Table 3: Sensitivity Analysis performed on Qwen3 14B with GPT5-mini precomputed context

| Setting | MATH-500 | |
| --- | --- | --- |
| | Pass@5 | # Tokens |
| 0-CoT – Default Prompt (Listing 4) | 99.40% | 4,931 |
| 0-CoT – Adjusted Prompt (Listing 5) | 99.20% | 4,246 |
| F-CoT – Default Prompt (Listing 8) | 98.60% | 2,437 |
| F-CoT – Adjusted Prompt (Listing 9) | 99.20% | 3,184 |
| F-CoT – Context Variation 1 (Listing 15) | 98.80% | 2,453 |
| F-CoT – Context Variation 2 (Listing 16) | 98.80% | 2,744 |
| F-CoT – Context Variation 3 (Listing 17) | 98.40% | 2,707 |
| F-CoT – Context + Original Question | 99.20% | 2,898 |

Table 4: Sensitivity Analysis Context Extraction Prompt on Qwen-14B (Self-generated context)

| Setting | MATH-500 | | |
| --- | --- | --- | --- |
| | Pass@5 | # Tokens | Valid Context |
| F-CoT – Prompt w/o ICL (Listing 1) | 97.00% | 2,526 | 100% |
| F-CoT – Prompt w/ ICL (Listing 3) | 96.80% | 2,563 | 98.20% |

Table 5: Sensitivity Analysis for Context Extraction with Models of Different Sizes on Qwen-14B.

| Setting | MATH-500 | | |
| --- | --- | --- | --- |
| | Pass@5 | # Tokens | Valid Context |
| F-CoT – Qwen3-0.6B Context | 41.2% | 1,472 | 0.2% |
| F-CoT – Qwen3-1.7B Context | 91.4% | 2,378 | 95.2% |
| F-CoT – Qwen3-4B Context | 92.8% | 2,496 | 99.8% |
| F-CoT – Qwen3-8B Context | 94.8% | 2,547 | 99.6% |
| F-CoT – Qwen3-14B Context | 97.2% | 2,625 | 99.8% |
| F-CoT – Qwen3-32B Context | 97.0% | 2,542 | 99.8% |