# Can Stories Help LLMs Reason?
# Curating Information Space Through Narrative

**Vahid Sadiri Javadi** [Ψ]    **Johanne R. Trippas** [⌖]    **Yash Kumar Lal** [ψ]    **Lucie Flek** [Ψ]

[Ψ] University of Bonn, [⌖] RMIT University, [ψ] Stony Brook University

{vahid.sadirij, lflek}@uni-bonn.de,
j.trippas@rmit.edu.au, ylal@cs.stonybrook.edu

## Abstract

Narratives are widely recognized as a powerful tool for structuring information and facilitating comprehension of complex ideas in various domains such as science communication. This paper investigates whether incorporating narrative elements can assist Large Language Models (LLMs) in solving complex problems more effectively. We propose a novel approach, **Story of Thought (SoT)**, integrating narrative structures into prompting techniques for problem solving. This approach involves constructing narratives around problem statements and creating a framework to identify and organize relevant information. Our experiments show that using various LLMs with SoT consistently surpasses using them with other techniques on physics, chemistry, math, and biology questions in both the GPQA and JEEBench datasets. The narrative-based information curation process in SoT enhances problem comprehension by contextualizing critical in-domain information and highlighting causal relationships within the problem space.

## 1 Introduction

Humans have an exceptional ability to understand and reason through narratives. A narrative-driven approach can enhance the comprehension and retention of complex subjects compared to simple fact listing (Fisher, 2021; Abbott, 2020; Gottschall, 2012). For example, storytelling effectively structures information in science communication (Dahlstrom, 2014; Norris et al., 2005; Martinez-Conde and Macknik, 2017), education (Engel et al., 2018; Negrete and Lartigue, 2004), and health communication (Dudley et al., 2023), revealing relationships and contextual nuances (Zak, 2015). While *narrative approach* contextualizes facts within a daily life scenario (story) with a planned structure, a *factual approach* conveys information in a concise in-domain manner.

To date, large language models (LLMs) struggle with complex problem-solving tasks that require the ability to integrate, structure, and apply relevant information effectively (Qiao et al., 2023; Wang et al., 2023). Prompting techniques based on breaking tasks into smaller subtasks, such as Chain-of-Thought (CoT) (Wei et al., 2022) and its more recent adaptations (Xia et al., 2024), have led to considerable improvements in problem-solving benchmarks. The strategies of constructing natural language rationales (Ling et al., 2017), in the CoT context also called reasoning processes, play a vital role in LLM prompting (Ye and Durrett, 2024; Min et al., 2022; Wang et al., 2022; Li et al., 2023).

Inspired by the effectiveness of narrative in *(i)* identifying and explaining abstract concepts and *(ii)* organizing the information flow coherently, we explore integrating narrative elements into prompt-driven reasoning. The main research questions addressed in this work are:

**RQ1:** Can LLMs generate coherent and relevant narratives around problem statements to facilitate comprehension and reasoning?
**RQ2:** Can incorporating narrative elements into prompting techniques improve model performance on complex problem-solving tasks?

We make the following contributions: *(i)* We introduce a novel method, **Story of Thought (SoT)**, that aids LLMs to identify and arrange relevant information for solving complex tasks by incorporating narrative structures into the prompting process, *(ii)* We evaluate the effectiveness of SoT on GPQA and JEEBench datasets of complex problems, showing superior performance to existing prompting techniques with SotA models, and *(iii)* We analyze the impact of narrative techniques to generate narrative-based explanations and investigate why they improve LLMs' reasoning abilities.

## 2 Related Work

Bruner (1991) posit that narratives are a fundamental mode of human thought, allowing individuals to convey complex concepts in a more understandable manner. Presenting information through narratives can enhance learning and memory, promote engagement and motivation (Willingham, 2004; Chen et al., 2023). The development of narrative-based educational strategies (Bower and Clark, 1969; Mawasi et al., 2020; Norris et al., 2005) paved the way for using them as a framework for organizing information for problem solving. The use of narratives can break down complex problems into sub-problems, providing a step-by-step approach to answering a question (Szurmak and Thuna, 2013). Sadiri Javadi et al. (2024) use different narratives techniques to satisfy diverse requirements for conversational information-seeking systems.

There are a plethora of datasets focusing on answering questions about given contexts. Reading comprehension datasets (Khashabi et al., 2018; Welbl et al., 2018; Williams et al., 2018; Mihaylov et al., 2018) explicitly evaluate a system's ability to answer questions that need information from multiple sentences in a passage. NarrativeQA (Kočiský et al., 2018) provides a dataset of 1,567 narratives and associated QA pairs as written by human annotators. ROCStories (Mostafazadeh et al., 2016) is a collection of 5 sentence short stories over which numerous datasets such as TellMeWhy (Lal et al., 2021) have been built to facilitate answering questions about narratives. However, none of these datasets use narratives as a tool of understanding, or relate to problem solving.

Problem solving datasets focus on mathematics, physics or other scientific domains. GSM8K (Cobbe et al., 2021) is a dataset of 8.5K high quality linguistically diverse grade school math word problems created by human problem writers. SciQ (Welbl et al., 2017) is built using a novel method for obtaining high-quality, domain-targeted multiple choice questions from crowd workers, and contains 13.7K multiple choice science exam questions. ScienceQA (Lu et al., 2022) adds multimodal context to collected elementary and high school science questions. While there has been rapid progress on these tasks, prior work has not integrated educational strategies such as narratives to tackle them, a setting which is likely to be used in the real world. MedMCQA (Pal et al., 2022) contains MCQ questions designed to address real-world medical entrance exam questions. Such datasets have been used extensively as yardsticks to measure the progress of NLP techniques.

The strength of modern LLMs, coupled with the paradigm of prompting, has driven up performance on problem solving tasks. In-context learning through few-shot examples has been used to teach LLMs about new tasks using a small number of examples. Chain of thought prompting (Wei et al., 2022) nudges LLMs to generate intermediate steps to mimic an explicit reasoning process before answering a question. Similarly, Tree of Thoughts (ToT) (Yao et al., 2023) and Graph of Thoughts (GoT) (Besta et al., 2024) induce intermediate reasoning structures, trees and graphs respectively, to decide on an answer. However, despite the fact that narratives have been used as a way to simplify problems, they have never been explored to improve the problem solving abilities of LLMs.

## 3 Methodology: Story of Thought

We introduce **Story of Thought** (SoT), a novel prompt-driven reasoning approach that generates narrative-based clarification to guide LLMs' reasoning process. Inspired by the narrative format, the SoT approach leverages the cognitive benefits of storytelling, such as contextual understanding and relational reasoning, that can help LLMs identify and maintain the information structure.

Figure 1 gives an overview of SoT. It involves three steps using narrative techniques: *(i)* **Question clarification** (i.e., acting as an explorer to dissect and clarify complex questions (Section 3.1)); *(ii)* **Narrative Generation** (i.e., generating detailed narratives from the clarified question components using different narrative techniques (Section 3.2)); and *(iii)* **Problem Solving** (i.e., leveraging narratives to prompt the LLMs to solve the tasks (Section 3.3)). We describe the exact prompts used in each step in Appendix C.

### 3.1 Step 1: Question Clarification

In the first step, we use the LLM's ability to explore and clarify the problem. Starting with a specialized prompt, the LLM breaks down the question into its core components, identifying relevant subtopics and areas. This detailed analysis is crucial for generating a coherent narrative that thoroughly addresses the question.
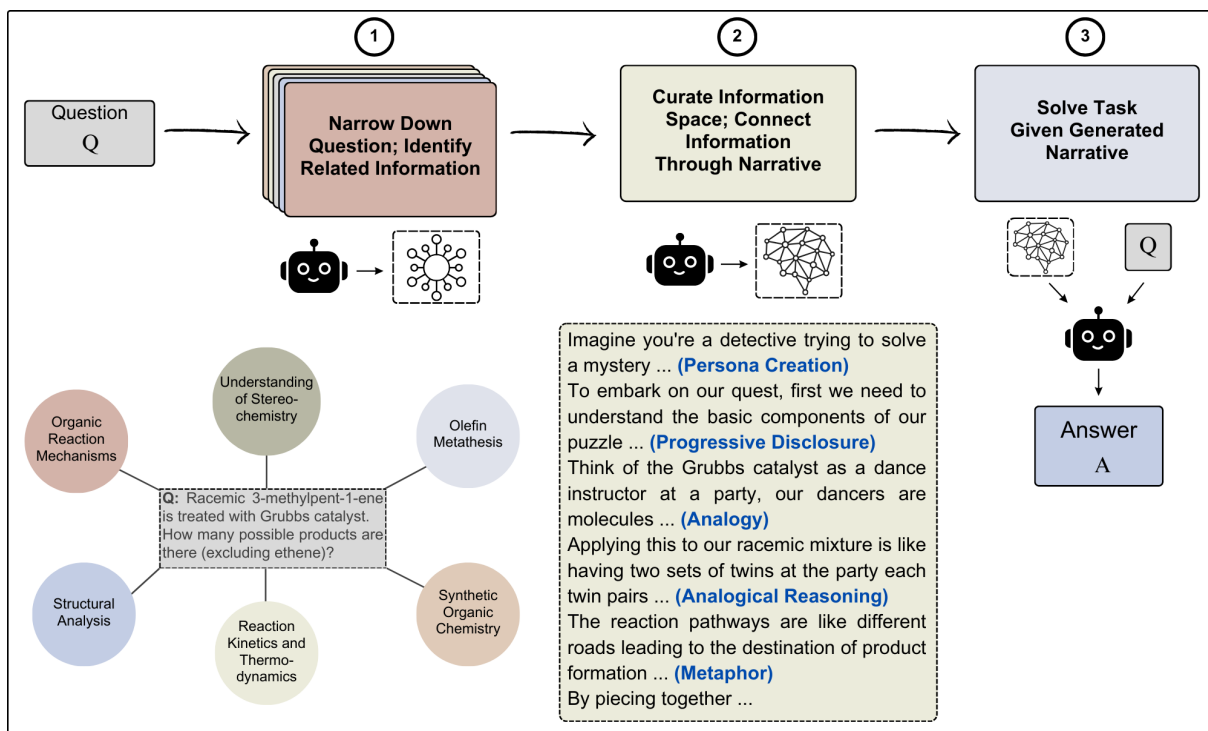
Figure 1: A high-level overview of **Story of Thought** (SoT), consisting of three steps (top): ① Question Clarification, ② Narrative Generation, ③ Solving Task and an actual example of LLM output (bottom) in each step for the GPQA task. The prompt designed for step 2 incorporates the narrative techniques (highlighted in blue) such as *analogical reasoning*, which identifies similarities between the target concept (information being conveyed) and a more familiar concept (*analogy*) and *progressive disclosure* which reveals information gradually throughout the narrative, rather than presenting it all at once. See Appendix C for prompts for each step and Appendix E for an example.

## 3.2 Step 2: Narrative Generation

The second step involves generating detailed narratives based on the breakdown and clarification performed in Step 1 (question clarification). These narratives provide a structured context for the questions to enhance the LLM's understanding, reasoning, and problem-solving abilities. Sadiri Javadi et al. (2024) discuss different narrative techniques required in conversational information seeking systems. We integrate the below subset of these techniques into our prompt and task LLMs to generate a narrative, based on the information from Step 1:

1. **Progressive Disclosure**: Reveals information gradually, guiding the LLM step-by-step through the problem-solving process.

2. **Branching**: Explores different paths or approaches to understanding the problem by providing multiple perspectives.

3. **Analogy**: Uses comparisons to familiar concepts or situations to make abstract components more understandable.

4. **Analogical Reasoning**: Facilitates understanding by reasoning through similarities between the problem and known situations.

5. **Metaphor**: Simplifies complex ideas through metaphorical representation.

## 3.3 Step 3: Problem Solving

In the final step, the LLM uses the narrative generated in Step 2 to solve the original QA task. The structured and contextual understanding provided by the narrative supports LLM in accessing relevant aspects of the task.

## 4 Experimental Setup

To comprehensively evaluate the effectiveness of our proposed approach, we conduct experiments across a diverse set of tasks and models, employing various prompting techniques for comparison.

### 4.1 Evaluation Tasks

We focus our evaluation on reasoning-intensive tasks spanning multiple domains, including physics, biology, and chemistry problem-solving. In particular, we utilize the **GPQA** (Diamond set) (Rein et al., 2024) and **JEEBench** (Arora et al., 2023). GPQA is a Graduate-level Problem-solving QA dataset which comprises expert-crafted

| Prompting Method | Meta | | Mistral | | OpenAI | | Microsoft | |
|---|---|---|---|---|---|---|---|---|
| | Llama 3 8B | Llama 3 70B | Mistral 7B | Mixtral 8x7B | ChatGPT 3.5 | GPT 4 | Phi-3 Mini | Phi-3 Medium |
| **Zero-shot** | 34.2 | 39.5 | 35.8 | 36.36 | 30.6 | 34.7 | **28.79** | 42.42 |
| **Zero-shot CoT** | 40.91 | 41.92 | 31.82 | 35.35 | 28.1 | 35.7 | 24.75 | 39.39 |
| **Tree of Thoughts** | 34.34 | 43.43 | 29.79 | 32.82 | 24.24 | 42.42 | 18.68 | 31.81 |
| **Graph of Thoughts** | 33.83 | 43.43 | 28.78 | 30.30 | 23.23 | 40.90 | 19.69 | 28.78 |
| **Analogical Reasoning (3-shot)** | 40.91 | 47.47 | 37.9 | 26.26 | 28.1 | 41.41 | 16.67 | **48.48** |
| **Ours: Knowledge Identification** | 40.4 | 48.99 | 35.35 | 37.77 | 27.77 | 40.90 | 20.71 | 37.88 |
| **Ours: Story of Thought (SoT)** | **43.43** | **51.01** | **38.4** | **38.89** | **30.8** | **48.98** | 22.73 | 36.36 |

Table 1: On GPQA (Diamond set), Story of Thought (SoT) consistently outperforms other techniques. We present the performance (QA accuracy) of different methods with various LLMs on GPQA Diamond set.

multiple-choice questions. It contains 448 multiple-choice questions written by domain experts in biology, physics, and chemistry of high quality and difficulty. We use the Diamond subset, which contains 198 questions on which both expert annotators agree. JEEBench contains 515 challenging pre-engineering mathematics, physics and chemistry problems from the highly competitive IIT JEE-Advanced exam.

These problems span the different sciences and are extremely challenging, requiring in-depth reasoning and domain knowledge, making them well-suited for assessing our approach's ability to understand complex tasks and contextualize salient information within the problem space.

### 4.2 Benchmarking Models

To evaluate the performance of our approach across a wide range of Large Language Models, we experiment with the following LLM families:

1. **Meta**: Llama 3 8B & Llama 3 70B
2. **Mistral**: Mistral 7B & Mixtral 8x7B
3. **OpenAI**: GPT-3.5-turbo & GPT 4
4. **Microsoft**: Phi 3 Medium & Phi 3 Mini

These models were selected to cover a wide spectrum of capabilities, sizes and families, enabling a comprehensive evaluation of their strengths and limitations. More details on the implementation can be found in Appendix B.

### 4.3 Methods Studied

We compared our proposed approach against several prompting techniques:

**Zero-shot Prompting**: This method, similar to our approach (SoT), does not rely on labeled examples. Instead, LLMs are prompted to solve tasks based solely on their pre-trained knowledge without any context provided. This approach serves as a baseline, demonstrating the LLMs'

ability to solve problems without explicit guidance.

**Zero-shot CoT** (Wei et al., 2022): This technique prompts the LLM to explicitly reason through the steps required to arrive at an answer. By prompting the model to generate a chain of thought, this method aims to improve the model's ability to solve complex problems by breaking them down into smaller, more manageable steps.

**Tree of Thoughts** (Yao et al., 2023): This method systematically explores multiple reasoning paths instead of a single linear progression. In ToT, a tree-structured solution to a problem is generated by breaking it down into sub-problems. This approach enables the model to consider a broader set of potential solutions by evaluating each branch for correctness before proceeding further.

**Graph of Thoughts** (Besta et al., 2024): This technique extends the Tree of Thoughts (ToT) approach by allowing for a more flexible and non-hierarchical representation of problem-solving steps. In this method, the reasoning steps are treated as nodes, and the connections between them are edges that represent logical relationships or dependencies. In our experiments, we adopt the same setup described in their original work.

**Analogical Reasoning** (Yasunaga et al., 2023): This approach leverages analogies to help the model draw parallels between known concepts and the task at hand. By providing analogical examples, the model is guided to understand and apply similar reasoning patterns to new problems. In our experiment, we allow the LLMs to self-generate three exemplars for each question (akin to the prompt described in their paper). This enables them to identify relevant examples and adapt their reasoning accordingly.

**Ours: Knowledge Identification**: To mea-

sure the effectiveness of our proposed approach, namely utilizing narrative in solving tasks, we prompt LLMs to solve the task based solely on the generated conceptual knowledge from Step 1 (described in Section 3.1). This allows us to compare the model capability in solving tasks using only the identified relevant knowledge versus leveraging this knowledge to structure a coherent narrative.

**Ours: Story of Thought (SoT)**: This approach represents the core of our proposed method, where we leverage the generated narratives from Step 2 (described in Section 3.2) to solve the given tasks.

# 5 Results

Our proposed SoT approach that incorporates narrative structures improves over almost all previous prompting approaches across two different problem-solving datasets. This highlights the potential of using narratives to improve the ability of LLMs to understand and reason about the given information in various intensive reasoning tasks.

## 5.1 Performance on GPQA

We present the results of our experiments on GPQA (Diamond) are presented in Table 1. For this task, SoT is the best method to use with six of eight models. The open-source Llama 3 70B model records the highest accuracy using the SoT method, achieving a score of 51.01%. This is the highest accuracy observed among all models and methods tested in the study. Furthermore, the GPT-4 model shows the most notable improvement in accuracy when the SoT method is employed, compared to its zero-shot baseline. Specifically, the accuracy for GPT-4 increased from 34.7% under zero-shot conditions to 48.98% with SoT (i.e., an absolute increase of 14.28%, or a relative increase of 41% respectively).[1]

Interestingly, all reasoning strategies lead to an accuracy drop for the comparably smaller Phi-3 Mini model, and all CoT strategies except Analogical Reasoning also lead to the accuracy drop of the Phi-3 Medium model compared to its zero-shot baseline. We hypothesize that this is due to the low quality of the generated explanations (whether CoT steps or SoT narrative), as further studied in §6.1.

---

[1]We also find that Llama 3 70B with SoT outperforms zero-shot o1-preview which uses CoT style reasoning internally. `https://openai.com/index/learning-to-reason-with-llms/`.

Figure 2 presents the performance of different models when using SoT across the different problem domains in GPQA. We note that, on average, models improve the most on Biology problems when using SoT. We hypothesize that this is because it is easier to simplify information for graduate level biology problems that can be used by models to come up with a solution.

## 5.2 Performance on JEEBench

Table 2 presents detailed experimental results on JEEBench. Our proposed Story of Thought (SoT) method consistently improves the performance of seven out of the eight LLMs. Using SoT, Llama 3 70B performance surpasses even the GPT models. It obtains the highest scores in all subjects and question types (Except Single-Correct), with an overall aggregate score of 0.453. This is a significant improvement on the previous SOTA, which was a strong GPT4 model used with both CoT and Self-Consistency.

Across models, the results highlight the effectiveness of Story of Thought (SoT) in enhancing model performance on complex, multi-disciplinary benchmarks like JEEBench, setting new SOTA results in several categories. The improvements are particularly notable in the subject categories and question types where the other methods struggle.

In Figure 2, we present subject-wise performance of different models on JEEBench. On average, model performance is highest on Chemistry problems when using SoT. This is in contrast to findings on GPQA and could occur due to the difference in degree of difficulty of problems in the two datasets (graduate level vs high school level). Regardless, improvements on Biology problems are not far behind those for Chemistry.

# 6 Analysis of SoT Aspects

## 6.1 Role of the Narrative Quality/Choice

The choice of *narrator* model (i.e., the model that generates narratives) can impact the problem-solving resuls. In the following experiments, we apply the narratives generated by other large and small open-source LLMs to the Phi-3 Mini and Phi-3 Medium models. The results of these experiments are presented in Table 3.

We observe that the **narratives** generated by most models **consistently improve the accuracy** of both Microsoft models compared to the baseline (i.e., when both models use their own generated

| | Chemistry | Mathematics | Physics | Integer | Single-Correct | Multi-Correct | Numeric | Total |
|---|---|---|---|---|---|---|---|---|
| GPT-4+CoT+SC@8* | 0.463 | 0.308 | 0.449 | 0.293 | **0.618** | 0.410 | 0.234 | 0.389 |
| Llama 3 8B | 0.143 | 0.082 | 0.089 | 0.061 | 0.127 | 0.148 | 0.044 | 0.102 |
| Llama 3 8B+CoT | 0.127 | 0.101 | 0.116 | **0.11** | 0.145 | 0.149 | 0.036 | 0.112 |
| Ours: Llama 3 8B+SoT | **0.154** | **0.195** | **0.172** | 0.072 | **0.259** | **0.324** | 0.028 | **0.173** |
| Llama 3 70B | 0.324 | 0.189 | 0.274 | 0.171 | 0.345 | 0.316 | 0.131 | 0.25 |
| Llama 3 70B+CoT | 0.264 | 0.228 | 0.268 | 0.159 | 0.291 | 0.317 | 0.175 | 0.249 |
| Ours: Llama 3 70B+SoT | **0.554** | **0.329** | **0.471** | **0.446** | **0.42** | **0.485** | **0.462** | **0.453** |
| Mistral 7B | 0.119 | 0.079 | 0.091 | 0.049 | 0.109 | 0.159 | 0.022 | 0.094 |
| Mistral 7B+CoT | 0.106 | 0.123 | 0.059 | 0.073 | 0.118 | 0.165 | 0.022 | 0.102 |
| Ours: Mistral 7B+SoT | **0.2** | **0.177** | **0.201** | **0.11** | **0.245** | **0.224** | **0.146** | **0.19** |
| Mixtral 8x7B | 0.22 | 0.151 | 0.167 | 0.122 | 0.218 | 0.261 | 0.058 | 0.176 |
| Mixtral 8x7B+CoT | 0.237 | 0.142 | 0.152 | 0.061 | 0.209 | 0.27 | 0.08 | 0.173 |
| Ours: Mixtral 8x7B+SoT | **0.253** | **0.251** | **0.274** | **0.268** | **0.309** | 0.277 | **0.182** | **0.257** |
| ChatGPT 3.5 | **0.228** | **0.146** | 0.173 | 0.073 | **0.318** | **0.249** | 0.029 | **0.177** |
| ChatGPT 3.5+CoT | 0.17 | 0.111 | 0.167 | 0.11 | 0.173 | 0.206 | 0.051 | 0.142 |
| Ours: ChatGPT 3.5+SoT | 0.189 | 0.128 | **0.189** | 0.073 | 0.291 | 0.204 | 0.051 | 0.161 |
| GPT 4 | 0.423 | 0.212 | 0.352 | 0.207 | 0.455 | 0.383 | 0.153 | 0.309 |
| GPT 4+CoT | 0.468 | 0.280 | 0.335 | 0.256 | **0.473** | 0.448 | 0.175 | 0.350 |
| Ours: GPT 4+SoT | **0.535** | **0.294** | **0.413** | **0.378** | 0.4 | **0.453** | **0.321** | **0.395** |
| Phi-3 Mini | **0.256** | 0.12 | **0.199** | 0.146 | 0.255 | 0.224 | 0.08 | 0.18 |
| Phi-3 Mini+CoT | **0.256** | 0.137 | 0.171 | 0.134 | 0.209 | 0.216 | **0.139** | 0.181 |
| Ours: Phi-3 Mini+SoT | 0.224 | **0.209** | 0.181 | **0.183** | **0.282** | **0.234** | 0.124 | **0.207** |
| Phi-3 Medium | 0.298 | 0.193 | 0.165 | 0.146 | 0.255 | 0.286 | 0.139 | 0.218 |
| Phi-3 Medium+CoT | 0.253 | 0.195 | 0.199 | 0.171 | 0.236 | 0.274 | 0.139 | 0.214 |
| Ours: Phi-3 Medium+SoT | **0.279** | **0.203** | **0.224** | **0.232** | **0.273** | 0.263 | **0.153** | **0.231** |

Table 2: On JEEBench, Story of Thought (SoT) outperforms previous SOTA as well as other methods. We present the aggregate score by subject as well as question type and present the overall aggregate score. The best overall scores are highlighted in **blue** while the best score by method for a model is in **bold**. * reported in (Arora et al., 2023).



Figure 2: Performance of **Story of Thought** (SoT) on GPQA and JEEBench across various LLMs and domains.

| Narrative Generator | Solver Models | |
|---|---|---|
| | Phi-3 Mini | Phi-3 Medium |
| Llama 3 8B | 23.74 (+1.01↑) | 37.88 (+1.28↑) |
| Llama 3 70B | 25.25 (+2.52↑) | **39.39** (+2.79↑) |
| Mistral 7B | 24.24 (+1.51↑) | 38.38 (+1.78↑) |
| Mixtral 8x7B | 24.74 (+2.01↑) | 35.86 (-0.74↓) |

Table 3: Applying generated narratives by open-source models to Microsoft models to solve the tasks.

narratives in Step 2 to solve the tasks, shown in Table 1). The absolute improvements range from 1.0% to 2.8%, with the Llama 3 70B model generating the most effective narratives. A slight decrease in accuracy is observed with the mixture-of-experts Mixtral 8x7B narratives for the Phi-3 Medium model, highlighting the need for careful selection and evaluation of narrator models to ensure compatibility and optimal performance. Larger models

| Narrative Technique | Meta | | Mistral | | OpenAI | | Microsoft | |
|---|---|---|---|---|---|---|---|---|
| | Llama 3 8B | Llama 3 70B | Mistral 7B | Mixtral 8x7B | ChatGPT 3.5 | GPT 4 | Phi-3 Mini | Phi-3 Medium |
| **Progressive Disclosure** | 427 | 597 | 191 | 191 | 744 | 570 | 367 | 368 |
| **Branching** | 30 | 56 | 51 | 20 | 72 | 168 | 34 | 61 |
| **Analogy** | 418 | 425 | 117 | 161 | 498 | 595 | 569 | 499 |
| **Analogical Reasoning** | 205 | 191 | 78 | 108 | 213 | 336 | 276 | 206 |
| **Metaphor** | 249 | 316 | 103 | 137 | 811 | 428 | 418 | 291 |
| $\sum$ | 1329 | 1585 | 540 | 617 | 2338 | 2097 | 1664 | 1425 |

Table 4: Comparing Generated Narratives - Total Number of Occurrences for each Narrative Techniques (Evaluator: Llama 3 70B)

generate narratives that break down problems to make them more easily solvable. Unsurprisingly, there is larger room for improving the problem solving abilities of smaller models.

## 6.2 Impact of Narrative Elements

To measure the impact of each of the narrative techniques we jointly prompted on the performance of open-source Meta models, we ablate the designed prompt in Step 2 (of Section 3.2) to apply each of the techniques separately. The results in Table 5 indicate that **employing any single narrative technique at a time is notably less effective at boosting QA accuracy than utilizing a combination of these simultaneously**.

| Narrative Technique | Meta | |
|---|---|---|
| | Llama 3 8B | Llama 3 70B |
| **Progressive Disclosure** | 34.85 (-8.58↓) | 44.95 (-6.06↓) |
| **Branching** | 34.34 (**-9.09**↓) | 44.95 (-6.06↓) |
| **Analogy** | 39.39 (-4.04↓) | 46.46 (-4.55↓) |
| **Analogical Reasoning** | 40.4 (-3.03↓) | 45.45 (-5.56↓) |
| **Metaphor** | 41.41 (-2.02↓) | 44.44 (-6.57↓) |
| **All** | 43.43 | 51.01 |

Table 5: Comparing accuracy when using a single narrative technique. The values in parentheses represent the decrease in accuracy percentage points compared to a combination of multiple narrative techniques simultaneously (shown in Table 1).

For both models (Llama 3 8B and 70B), the decrease in accuracy is comparably smaller (-3.0% to -5.6%) when using only the analogical components of the narrative (*Analogy* and *Analogical Reasoning*) than when using only the structural instructions (*Progressive Disclosure* or *Branching*) which leads to larger (-6.0% to -9.1%) accuracy loss. However, reasoning alone does not perform on par with the full narrative generation listing all the techniques. Prompting for *Metaphor* usage only leads to a larger accuracy loss in the 70B model (-6.6%) compared to the smaller one (-2.0%). This

makes it difficult to determine both how the narrative techniques relate to each other and whether the model truly comprehends the prompts it receives. We disentangle and study the two going forward.

## 6.3 Analyzing Generated Narratives

To gain deeper insights into the generated narratives, we designed a prompt (shown below) that utilizes our best-performing model (Llama 3 70B) to annotate the number of occurrences of each narrative technique for each generated narrative by all models used in our experiments. We can better interpret how the model executed the narrative technique prompt, by asking it to label if and where the mentioned techniques are used in the text generated. Less frequently labeled techniques might be the ones where LLM doesn't have a clear understanding of what it is asked to do. A proportion of the techniques and their correlation can provide us with a better picture of LLM's interpretation of the instruction as well. We detail the instruction given to the LLM in Appendix C.

We aim to uncover patterns and variations in the use of narrative techniques across different LLMs. Table 4 compares the total number of occurrences for each narrative technique across various LLMs.

**Variability in Utilization of Narrative Techniques Across Models:** In our designed prompt in Step 2 (i.e., Narrative Generation), LLMs generate narrative using all 5 narrative techniques. However, as Table 4 indicates, not all techniques were employed equally. This reveals that while some techniques like *Analogy* and *Progressive Disclosure* were consistently utilized, others such as *Branching* were applied less frequently.

We observe a trend across all LLM families where models with larger capacities, such as Llama 3 70B and GPT-4, consistently show higher occurrences of narrative techniques compared to their smaller counterparts. Furthermore, OpenAI's mod-

| Similarity Metric | BertScore | | ROUGE-L | | BLEU | |
|---|---|---|---|---|---|---|
| | SoT Reasoning | CoT Reasoning | SoT Reasoning | CoT Reasoning | SoT Reasoning | CoT Reasoning |
| **Llama 3 8B** | **0.28** | 0.06 | **0.19** | 0.11 | **6.57** | 0.19 |
| **Llama 3 70B** | **0.3** | 0.04 | **0.2** | 0.1 | **8.18** | 0.06 |
| **Mistral 7B** | 0.27 | **0.33** | 0.18 | **0.2** | **8.12** | 4.65 |
| **Mixtral 8x7B** | 0.3 | **0.34** | 0.19 | **0.21** | **8.92** | 8.14 |
| **ChatGPT 3.5** | **0.3** | 0.24 | **0.19** | 0.16 | 6.1 | **6.07** |
| **GPT 4** | 0.31 | **0.34** | 0.19 | **0.2** | **8.84** | 6.73 |
| **Phi-3 Mini** | 0.27 | **0.31** | 0.17 | **0.19** | **6.54** | 6.36 |
| **Phi-3 Medium** | 0.3 | **0.35** | 0.2 | **0.21** | 7.13 | **8.4** |

Table 6: Comparison of generated Story of Thought (SoT) and Chain of Thought (CoT) reasoning with Human Explanations on the GPQA (Diamond set) using BERTScore, ROUGE-L, and BLEU metrics across various large language models. Bold values indicate the reasoning approach that is more similar to human explanations for each model and metric pair.



Figure 3: Correlation coefficients among all narrative techniques (**PD** = Progressive Disclosure, **BR** = Branching, **AN** = Analogy, **AR** = Analogical Reasoning, **ME** = Metaphor) used in the SoT approach for GPT 4 and Llama 3 70B in solved and unsolved tasks.

els (ChatGPT 3.5 & GPT-4) demonstrate the highest total occurrences of narrative techniques, with 2,338 and 2,097, respectively with a notable emphasis on *Metaphors* and *Analogies*.

**Correlation Among Narrative Techniques:** To further investigate the dynamics of narrative techniques, we compute correlations between the frequencies of narrative techniques across solved and unsolved tasks, as shown in Figure 3. This analysis aims to uncover if the models consistently use certain narrative techniques together or vary significantly. Our initial results indicate diverse correlation patterns, suggesting that the effectiveness of narrative techniques in solving tasks across various LLMs needs to be further analyzed.

### 6.4 Analyzing SoT Reasoning

Table 6 compares the similarity of Story of Thought (SoT) and Chain of Thought (CoT) reasoning outputs to human explanations for different language models on the GPQA (Diamond set) dataset, using BertScore, ROUGE-L, and BLEU.

The differences between ROUGE-L values are insignificant and do not display any clear trends. However, according to BLEU scores, using SoT results in explanations closer to humans and the differences are more pronounced.

As per BertScore (an embedding-based similarity metric), Llama 3 models (8B and 70B) explanations are more similar to human ones when using SoT reasoning across all three metrics. However, Mistral models (7B and 8x7B), GPT-4, and Phi-3 Mini generate explanations more similar to human explanations when using CoT reasoning across all metrics. The semantic similarity of narratives generated by Llama 3 70B to human explanations combined with their effect of improving smaller models indicates that these narratives present information about the problems in a simplified manner.

## 7 Conclusion

Inspired by findings from human cognitive processes explored in didactics research, in this work, we propose to use narrative techniques in LLM prompting.We present strong evidence on public benchmark datasets that narrative techniques have the potential to notably enhance the reasoning abilities of LLMs in complex problem-solving tasks. By incorporating narrative structures, which mimic human cognitive processes of organizing and interpreting information, LLMs can achieve higher levels of performance and provide more contextually enriched responses.

## Limitations

**Contribution limitations.** The occurrences of narrative techniques do not necessarily imply the quality or effectiveness of the generated narratives; rather, they provide insights into the models' tendencies and preferences in employing these techniques. Therefore, answering the question of why narrative is helping LLMs is more complex and needs to be further investigated by looking into different research areas such as cognitive and communication theories.

**Dataset limitations.** So far, we used only GPQA and JEEBench problems as the most challenging set of problem-solving benchmarks we were aware of. Other comparable benchmarks, such as MGSM, are much closer to human or superhuman accuracy already without reasoning prompts and will be explored in future work.

**Analysis limitations.** We used Llama 70 B to respectively analyze the narratives. The intuition behind this experiment is that we can better interpret how the model executed the narrative technique prompt, by asking it to label if and where the mentioned techniques are used in the text generated. An alternative would be a thorough human assessment and further analysis of the impact on downstream performance, both of which we pursue in ongoing follow-up experiments. (We also previously prompted the LLMs in Step 2 to explain each of these five narrative techniques to make sure the concepts are understood before generating the narrative.)

## References

H Porter Abbott. 2020. *The Cambridge introduction to narrative*. Cambridge University Press.

Daman Arora, Himanshu Singh, and Mausam. 2023. Have LLMs advanced enough? a challenging problem solving benchmark for large language models. In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pages 7527–7543, Singapore. Association for Computational Linguistics.

Maciej Besta, Nils Blach, Ales Kubicek, Robert Gerstenberger, Lukas Gianinazzi, Joanna Gajda, Tomasz Lehmann, Michał Podstawski, Hubert Niewiadomski, Piotr Nyczyk, and Torsten Hoefler. 2024. Graph of Thoughts: Solving Elaborate Problems with Large Language Models. *Proceedings of the AAAI Conference on Artificial Intelligence*, 38(16):17682–17690.

Gordon H Bower and Michal C Clark. 1969. Narrative stories as mediators for serial learning. *Psychonomic science*, 14(4):181–182.

Jerome Bruner. 1991. The narrative construction of reality. *Critical inquiry*, 18(1):1–21.

Althea Y Chen, Chun-Ching Chen, and Wen-Yin Chen. 2023. The design narrative in design learning: Adjusting the inertia of attention and enhancing design integrity. *The Design Journal*, 26(4):519–535.

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. 2021. Training verifiers to solve math word problems. *arXiv preprint arXiv:2110.14168*.

Michael F Dahlstrom. 2014. Using narratives and storytelling to communicate science with nonexpert audiences. *Proceedings of the national academy of sciences*, 111(supplement_4):13614–13620.

Matthew Z Dudley, Gordon K Squires, Tracy M Petroske, Sandra Dawson, and Janesse Brewer. 2023. The use of narrative in science and health communication: a scoping review. *Patient Education and Counseling*, page 107752.

Alison Engel, Kathryn Lucido, and Kyla Cook. 2018. Rethinking narrative: Leveraging storytelling for science learning. *Childhood Education*, 94(6):4–12.

Walter R Fisher. 2021. *Human communication as narration: Toward a philosophy of reason, value, and action*. University of South Carolina Press.

Jonathan Gottschall. 2012. *The storytelling animal: How stories make us human*. Houghton Mifflin Harcourt.

Daniel Khashabi, Snigdha Chaturvedi, Michael Roth, Shyam Upadhyay, and Dan Roth. 2018. Looking beyond the surface: A challenge set for reading comprehension over multiple sentences. In *Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long Papers)*, pages 252–262, New Orleans, Louisiana. Association for Computational Linguistics.

Tomáš Kočiský, Jonathan Schwarz, Phil Blunsom, Chris Dyer, Karl Moritz Hermann, Gábor Melis, and Edward Grefenstette. 2018. The NarrativeQA reading comprehension challenge. *Transactions of the Association for Computational Linguistics*, 6:317–328.

Yash Kumar Lal, Nathanael Chambers, Raymond Mooney, and Niranjan Balasubramanian. 2021. TellMeWhy: A dataset for answering why-questions in narratives. In *Findings of the Association for Computational Linguistics: ACL-IJCNLP 2021*, pages 596–610, Online. Association for Computational Linguistics.

Yifei Li, Zeqi Lin, Shizhuo Zhang, Qiang Fu, Bei Chen, Jian-Guang Lou, and Weizhu Chen. 2023. Making language models better reasoners with step-aware verifier. In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 5315–5333, Toronto, Canada. Association for Computational Linguistics.

Wang Ling, Dani Yogatama, Chris Dyer, and Phil Blunsom. 2017. Program induction by rationale generation: Learning to solve and explain algebraic word problems. In *Proceedings of the 55th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 158–167, Vancouver, Canada. Association for Computational Linguistics.

Pan Lu, Swaroop Mishra, Tony Xia, Liang Qiu, Kai-Wei Chang, Song-Chun Zhu, Oyvind Tafjord, Peter Clark, and Ashwin Kalyan. 2022. Learn to explain: Multimodal reasoning via thought chains for science question answering. In *The 36th Conference on Neural Information Processing Systems (NeurIPS)*.

Susana Martinez-Conde and Stephen L Macknik. 2017. Finding the plot in science storytelling in hopes of enhancing science communication. *Proceedings of the National Academy of Sciences*, 114(31):8127–8129.

Areej Mawasi, Peter Nagy, and Ruth Wylie. 2020. Systematic literature review on narrative-based learning in educational technology learning environments (2007-2017).

Todor Mihaylov, Peter Clark, Tushar Khot, and Ashish Sabharwal. 2018. Can a suit of armor conduct electricity? a new dataset for open book question answering. In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*, pages 2381–2391, Brussels, Belgium. Association for Computational Linguistics.

Sewon Min, Xinxi Lyu, Ari Holtzman, Mikel Artetxe, Mike Lewis, Hannaneh Hajishirzi, and Luke Zettlemoyer. 2022. Rethinking the role of demonstrations: What makes in-context learning work? In *EMNLP*.

Nasrin Mostafazadeh, Nathanael Chambers, Xiaodong He, Devi Parikh, Dhruv Batra, Lucy Vanderwende, Pushmeet Kohli, and James Allen. 2016. A corpus and cloze evaluation for deeper understanding of commonsense stories. In *Proceedings of the 2016 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 839–849, San Diego, California. Association for Computational Linguistics.

Aquiles Negrete and Cecilia Lartigue. 2004. Learning from education to communicate science as a good story. *Endeavour*, 28(3):120–124.

Stephen P Norris, Sandra M Guilbert, Martha L Smith, Shahram Hakimelahi, and Linda M Phillips. 2005. A theoretical framework for narrative explanation in science. *Science education*, 89(4):535–563.

Ankit Pal, Logesh Kumar Umapathi, and Malaikannan Sankarasubbu. 2022. Medmcqa: A large-scale multi-subject multi-choice dataset for medical domain question answering. In *Proceedings of the Conference on Health, Inference, and Learning*, volume 174 of *Proceedings of Machine Learning Research*, pages 248–260. PMLR.

Shuofei Qiao, Yixin Ou, Ningyu Zhang, Xiang Chen, Yunzhi Yao, Shumin Deng, Chuanqi Tan, Fei Huang, and Huajun Chen. 2023. Reasoning with language model prompting: A survey. In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 5368–5393, Toronto, Canada. Association for Computational Linguistics.

David Rein, Betty Li Hou, Asa Cooper Stickland, Jackson Petty, Richard Yuanzhe Pang, Julien Dirani, Julian Michael, and Samuel R. Bowman. 2024. GPQA: A graduate-level google-proof q&a benchmark. In *First Conference on Language Modeling*.

Vahid Sadiri Javadi, Johanne R Trippas, and Lucie Flek. 2024. Unveiling information through narrative in conversational information seeking. In *Proceedings of the 6th ACM Conference on Conversational User Interfaces*, CUI '24, New York, NY, USA. Association for Computing Machinery.

Joanna Szurmak and Mindy Thuna. 2013. Tell me a story: The use of narrative as tool for instruction.

Xiaoxuan Wang, Ziniu Hu, Pan Lu, Yanqiao Zhu, Jieyu Zhang, Satyen Subramaniam, Arjun R Loomba, Shichang Zhang, Yizhou Sun, and Wei Wang. 2023. Scibench: Evaluating college-level scientific problem-solving abilities of large language models. *arXiv preprint arXiv:2307.10635*.

Xuezhi Wang, Jason Wei, Dale Schuurmans, Quoc Le, Ed Huai hsin Chi, and Denny Zhou. 2022. Self-consistency improves chain of thought reasoning in language models. *ArXiv*, abs/2203.11171.

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.

Johannes Welbl, Nelson F. Liu, and Matt Gardner. 2017. Crowdsourcing multiple choice science questions. In *Proceedings of the 3rd Workshop on Noisy User-generated Text*, pages 94–106, Copenhagen, Denmark. Association for Computational Linguistics.

Johannes Welbl, Pontus Stenetorp, and Sebastian Riedel. 2018. Constructing datasets for multi-hop reading comprehension across documents. *Transactions of the Association for Computational Linguistics*, 6:287–302.

Adina Williams, Nikita Nangia, and Samuel Bowman. 2018. A broad-coverage challenge corpus for sentence understanding through inference. In *Proceedings of the 2018 Conference of the North American*

*Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long Papers)*, pages 1112–1122, New Orleans, Louisiana. Association for Computational Linguistics.

Daniel T Willingham. 2004. Ask the cognitive scientist the privileged status of story. *American Educator*, 28:43–45.

Yu Xia, Rui Wang, Xu Liu, Mingyan Li, Tong Yu, Xiang Chen, Julian McAuley, and Shuai Li. 2024. Beyond chain-of-thought: A survey of chain-of-x paradigms for llms. *arXiv preprint arXiv:2404.15676*.

Shunyu Yao, Dian Yu, Jeffrey Zhao, Izhak Shafran, Thomas L. Griffiths, Yuan Cao, and Karthik Narasimhan. 2023. Tree of Thoughts: Deliberate problem solving with large language models. *Preprint*, arXiv:2305.10601.

Michihiro Yasunaga, Xinyun Chen, Yujia Li, Panupong Pasupat, Jure Leskovec, Percy Liang, Ed H Chi, and Denny Zhou. 2023. Large language models as analogical reasoners. *arXiv preprint arXiv:2310.01714*.

Xi Ye and Greg Durrett. 2024. The unreliability of explanations in few-shot prompting for textual reasoning. In *Proceedings of the 36th International Conference on Neural Information Processing Systems*, NIPS '22, Red Hook, NY, USA. Curran Associates Inc.

Paul J Zak. 2015. Why inspiring stories make us react: The neuroscience of narrative. In *Cerebrum: the Dana forum on brain science*, volume 2015. Dana Foundation.

## A Robustness of LLM Predictions

In the original GPQA dataset used for our experiments, the correct answers are always presented as the first option among the multiple choices. However, To further evaluate the robustness of the LLMs, we conduct an additional experiment where the correct answers are placed in the second option instead. Table 7 presents the results of these experiments, comparing the performance of various prompting methods across six different open-source LLMs. We observe that most LLMs experience a significant drop in accuracy when the correct answer is moved to the second option. However, despite the overall decrease in accuracy, our proposed approach, Story of Thought (SoT), consistently outperforms the baseline methods for most LLMs. The SoT method achieves the highest accuracy for the Meta Llama 3 8B, Meta Llama 3 70B, Mistral 8x7B, and Microsoft Phi-3 Medium models, demonstrating its effectiveness in enhancing the robustness of LLMs to changes in the problem structure.

## B Model Implementation Details

All experiments, except for those involving OpenAI models, were conducted on local machines equipped with GPUs. The models were run locally on a GPU setup without quantization using the Hugging Face Transformer library[2]. For OpenAI's GPT-3.5-turbo and GPT 4 models, we use the OpenAI API to generate outputs. Across all models, we use a temperature of 1.0 and a maximum number of tokens of 8,000 and report the accuracy.

## C Prompts Used in Story of Thought

We describe the prompts used for each stage in pipeline.

### C.1 Question Clarification

```
You are an explorer who wants to identify and
    collect different related and specialized
    subject areas to clarify the question. Your
    goal is to narrow down the question and
    provide relevant areas of knowledge and
    experience you have that help clarify the
    question mentioned below. You should not
    answer the question.

<question>
```

### C.2 Narrative Generation

```
You are an expert in narrative-based
    explanations for science communication.
    Your goal is to clarify the following
    question in a narrative way through the
    interconnected information provided below
    to enable a non-expert to comprehend the
    question in a more coherent and
    contextually rich manner. You should not
    answer the question.

Make sure to use all of these narrative
    techniques when clarifying the question
    through the interconnected information:
    Progressive Disclosure, Branching, Analogy,
    Analogical Reasoning, and Metaphor.

<question>

<generated information in the previous step>
```

### C.3 Problem Solving

```
You are an expert in analyzing narrative-based
    explanations for solving tasks. Please
    answer the following question based on the
    following narrative-based clarification:

<question>

Options:
<options>

<generated narrative in the previous step>
```

### C.4 Analyzing Generated Narratives

```
You are an expert in analyzing narrative-based
    explanations for science communication.
    Your goal is to find out which narrative
    techniques have been used in the following
    narrative-based explanation.

Label the narrative-based explanation using the
    following narrative-based techniques:
1. Progressive Disclosure
2. Branching
3. Analogy
4. Analogical Reasoning
5. Metaphor

<generated narrative>
```

## D Performance on JEEBench

## E Story of Thought (SoT) vs. Chain of Thought (CoT)

---

| Prompting Method | Meta | | Mistral | | Microsoft | |
|---|---|---|---|---|---|---|
| | Llama 3 8B | Llama 3 70B | Mistral 7B | Mixtral 8x7B | Phi-3 Mini | Phi-3 Medium |
| **Zero-shot** | 30.81 (-3.39↓) | 31.31 (-8.19↓) | 19.7 (-16.1↓) | 18.18 (-18.18↓) | 29.8 (+1.01↑) | 21.72 (-20.7↓) |
| **Zero-shot CoT** | 27.27 (-13.64↓) | 33.33 (-8.59↓) | **22.73** (-9.09↓) | 17.17 (-18.18↓) | 32.32 (+7.57↑) | 21.21 (-18.18↓) |
| **Analogical Reasoning** | 27.78 (-13.13↓) | 40.91 (-6.56↓) | 10.61 (-27.29↓) | 19.19 (-7.07↓) | **35.86** (+19.19↑) | 16.67 (-31.81↓) |
| **Ours: Knowledge Identification** | 32.32 (-8.08↓) | 42.4 (-6.59↓) | 16.67 (-18.68↓) | 14.65 (-23.12↓) | 28.28 (+7.57↓) | 23.26 (-14.62↓) |
| **Ours: Story of Thought (SoT)** | **34.85** (-8.58↓) | **45.4** (-5.61↓) | 20.2 (-18.2↓) | **20.2** (-18.69↓) | 27.7 (+4.97↑) | **25.75** (-10.85↓) |

Table 7: Performance of various LLMs across different prompting methods on GPQA (Diamond set). Correct answers are presented in the second option. Values in parentheses indicate the change in accuracy compared to the original setting in Table 1 where the correct answer was in the first option.

| | Chemistry | Mathematics | Physics | Integer | Single-Correct | Multi-Correct | Numeric | Total |
|---|---|---|---|---|---|---|---|---|
| GPT-4+CoT+SC@8* | 0.463 | 0.308 | 0.449 | 0.293 | <u>**0.618**</u> | 0.410 | 0.234 | 0.389 |
| Llama 3 8B | 0.143 | 0.082 | 0.089 | 0.061 | 0.127 | 0.148 | 0.044 | 0.102 |
| Llama 3 8B+CoT | 0.127 | 0.101 | 0.116 | **0.11** | 0.145 | 0.149 | 0.036 | 0.112 |
| Llama 3 8B+Analogical Reasoning (3-shot) | 0.139 | 0.111 | 0.128 | **0.11** | 0.145 | 0.165 | **0.058** | 0.124 |
| Ours: Llama 3 8B+Knowledge Identification | 0.199 | 0.099 | 0.134 | 0.073 | 0.227 | 0.171 | 0.058 | 0.137 |
| Ours: Llama 3 8B+SoT | **0.154** | **0.195** | **0.172** | 0.072 | **0.259** | **0.324** | 0.028 | **0.173** |
| Llama 3 70B | 0.324 | 0.189 | 0.274 | 0.171 | 0.345 | 0.316 | 0.131 | 0.25 |
| Llama 3 70B+CoT | 0.264 | 0.228 | 0.268 | 0.159 | 0.291 | 0.317 | 0.175 | 0.249 |
| Llama 3 70B+Analogical Reasoning (3-shot) | 0.314 | 0.24 | 0.295 | 0.195 | 0.318 | 0.349 | 0.19 | 0.276 |
| Ours: Llama 3 70B+Knowledge Identification | 0.317 | 0.226 | 0.254 | 0.195 | 0.345 | 0.323 | 0.146 | 0.26 |
| Ours: Llama 3 70B+SoT | **0.554** | **0.329** | **0.471** | **0.446** | 0.42 | **0.485** | **0.462** | **0.453** |
| Mistral 7B | 0.119 | 0.079 | 0.091 | 0.049 | 0.109 | 0.159 | 0.022 | 0.094 |
| Mistral 7B+CoT | 0.106 | 0.123 | 0.059 | 0.073 | 0.118 | 0.165 | 0.022 | 0.102 |
| Mistral 7B+Analogical Reasoning (3-shot) | 0.157 | 0.084 | 0.116 | 0.073 | 0.155 | 0.169 | 0.029 | 0.114 |
| Ours: Mistral 7B+Knowledge Identification | 0.109 | 0.055 | 0.063 | 0.037 | 0.091 | 0.117 | 0.022 | 0.073 |
| Ours: Mistral 7B+SoT | **0.2** | **0.177** | **0.201** | **0.11** | **0.245** | **0.224** | **0.146** | **0.19** |
| Mixtral 8x7B | 0.22 | 0.151 | 0.167 | 0.122 | 0.218 | 0.261 | 0.058 | 0.176 |
| Mixtral 8x7B+CoT | 0.237 | 0.142 | 0.152 | 0.061 | 0.209 | 0.27 | 0.08 | 0.173 |
| Mixtral 8x7B+Analogical Reasoning (3-shot) | 0.202 | 0.155 | 0.197 | 0.122 | 0.191 | **0.281** | 0.066 | 0.179 |
| Ours: Mixtral 8x7B+Knowledge Identification | 0.184 | 0.129 | 0.144 | 0.122 | 0.155 | 0.237 | 0.044 | 0.15 |
| Ours: Mixtral 8x7B+SoT | **0.253** | **0.251** | **0.274** | **0.268** | **0.309** | 0.277 | **0.182** | **0.257** |
| ChatGPT 3.5 | **0.228** | **0.146** | 0.173 | 0.073 | **0.318** | **0.249** | 0.029 | **0.177** |
| ChatGPT 3.5+CoT | 0.17 | 0.111 | 0.167 | 0.11 | 0.173 | 0.206 | 0.051 | 0.142 |
| ChatGPT 3.5+Analogical Reasoning (3-shot) | 0.208 | 0.125 | 0.148 | 0.098 | 0.2 | 0.216 | **0.073** | 0.156 |
| Ours: ChatGPT 3.5+Knowledge Identification | 0.155 | 0.141 | 0.167 | **0.122** | 0.209 | 0.188 | **0.073** | 0.151 |
| Ours: ChatGPT 3.5+SoT | 0.189 | 0.128 | **0.189** | 0.073 | 0.291 | 0.204 | 0.051 | 0.161 |
| GPT 4 | 0.423 | 0.212 | 0.352 | 0.207 | 0.455 | 0.383 | 0.153 | 0.309 |
| GPT 4+CoT | 0.468 | 0.280 | 0.335 | 0.256 | **0.473** | 0.448 | 0.175 | 0.350 |
| GPT 4+Analogical Reasoning (3-shot) | 0.479 | 0.286 | 0.396 | 0.305 | 0.4 | 0.43 | 0.307 | 0.371 |
| Ours: GPT 4+Knowledge Identification | 0.481 | 0.287 | 0.386 | 0.293 | 0.373 | 0.452 | 0.314 | 0.373 |
| Ours: GPT 4+SoT | **0.535** | **0.294** | **0.413** | **0.378** | 0.4 | **0.453** | **0.321** | **0.395** |
| Phi-3 Mini | **0.256** | 0.12 | **0.199** | 0.146 | 0.255 | 0.224 | 0.08 | 0.18 |
| Phi-3 Mini+CoT | **0.256** | 0.137 | 0.171 | 0.134 | 0.209 | 0.216 | **0.139** | 0.181 |
| Phi-3 Mini+Analogical Reasoning (3-shot) | 0.205 | 0.159 | 0.195 | 0.146 | 0.264 | 0.218 | 0.088 | 0.182 |
| Ours: Phi-3 Mini+Knowledge Identification | 0.168 | 0.091 | 0.106 | 0.073 | 0.136 | 0.181 | 0.044 | 0.118 |
| Ours: Phi-3 Mini+SoT | 0.224 | **0.209** | 0.181 | **0.183** | **0.282** | **0.234** | 0.124 | **0.207** |
| Phi-3 Medium | 0.298 | 0.193 | 0.165 | 0.146 | 0.255 | 0.286 | 0.139 | 0.218 |
| Phi-3 Medium+CoT | 0.253 | 0.195 | 0.199 | 0.171 | 0.236 | 0.274 | 0.139 | 0.214 |
| Phi-3 Medium+Analogical Reasoning (3-shot) | 0.258 | 0.181 | 0.173 | 0.159 | 0.218 | 0.276 | 0.117 | 0.202 |
| Ours: Phi-3 Medium+Knowledge Identification | 0.288 | 0.163 | 0.205 | 0.207 | 0.236 | 0.235 | 0.161 | 0.211 |
| Ours: Phi-3 Medium+SoT | **0.279** | **0.203** | **0.224** | **0.232** | **0.273** | **0.263** | **0.153** | **0.231** |

Table 8: On JEEBench, Story of Thought (SoT) outperforms previous SOTA as well as other methods. We present the aggregate score by subject as well as question type and present the overall aggregate score. * denotes SOTA results taken from the original paper (Arora et al., 2023).

# Story of Thought (SoT)

**1 — Question Clarification**

**2 — Narrative Generation**

**3 — Solving Task**

Question Q

Narrow Down Question; Identify Related Information

Curate Information Space; Connect Information Through Narrative

Solve Task Given Generated Narrative

---

**Q:** Racemic 3-methylpent-1-ene is treated with Grubbs catalyst. How many possible products are there (excluding ethene)?

To clarify the given question and identify relevant areas of knowledge, the following subjects and concepts would be helpful:
1. Organic Chemistry:
   - Alkene chemistry
   - Stereochemistry
   - Reaction mechanisms
   - Grubbs catalyst and its applications
2. Stereochemistry:
   - Chirality
   - Racemic mixtures
   - Stereoisomerism (cis/trans, E/Z)
3. Organometallic Chemistry:
   - Olefin metathesis reactions
   - Mechanism of Grubbs catalyst
   - Selectivity and stereochemical outcomes
4. Reaction Kinetics and Thermodynamics:
   - Reaction pathways
   - Energetics of different products
   - Stability of products
5. Analytical Techniques:
   - Identification and characterization of products
   - Spectroscopic methods (NMR, IR, MS)
   - Separation techniques (chromatography)
6. Retrosynthetic Analysis:
   - Identifying possible products from the given starting material
   - Predicting reaction pathways and outcomes
With knowledge and experience in these areas, one can better understand the reaction conditions, predict the potential products, and analyze the stereochemical outcomes when a racemic alkene is treated with the Grubbs catalyst.

---

Let's take a journey into the world of chemistry to unravel the mysteries hidden within our question. Imagine chemistry as a grand puzzle, each piece representing a different aspect of the scientific world. Our quest today is to understand the potential products formed when a racemic mixture of 3-methylpent-1-ene reacts with the Grubbs catalyst.

To embark on our quest, we need to understand the basic components of our puzzle.
**Organic Chemistry:** is like the art of crafting molecules, where we play with structures such as **alkenes**: Alkenes are hydrocarbons with a special bond called a double bond. In our case, the molecule 3-methylpent-1-ene has a double bond and a branching point with a methyl group.

Next, we venture into **Stereochemistry**, which tells us about the spatial arrangement of atoms within molecules, like how the pieces fit in a 3D puzzle. For alkenes, stereochemistry can show us different forms (isomers) that look the same on paper but are different in 3D space. These include **cis/trans** or **E/Z isomers**, and can also involve **chirality**, where molecules are like left and right hands - mirror images but not superimposable.

Our path now leads us into the intriguing realm of **Organometallic Chemistry**, where we encounter the Grubbs catalyst. This catalyst is a magical weaver that reconfigures the threads of our molecular tapestry. It specializes in **olefin metathesis**, a reaction where double bonds are broken and reformed, creating new combinations of molecules. Imagine the Grubbs catalyst as a dance instructor at a party, changing dance partners mid-dance to form new pairings.

In this dance party analogy, our dancers are molecules. When the Grubbs catalyst starts its instruction, it can pair up molecules in various ways, creating new combinations. Applying this to our racemic 3-methylpent-1-ene, a racemic mixture is like having two sets of twins at the party, each twin pair looking identical but being mirror images of each other.

The reaction pathways here are like different roads leading to the destination of product formation. Some roads are longer, some shorter, each with its own scenic route. **Reaction Kinetics and Thermodynamics** help us understand these roads - kinetics tells us the speed of travel, while thermodynamics tells us which road leads to the most stable and comfortable destination, or product.

When 3-methylpent-1-ene, in its racemic form, meets the Grubbs catalyst, several new roads (or products) emerge. Each possible product depends on how the catalyst rearranges the double bonds and which molecular dancers pair up. We use **Analytical Techniques** like NMR, IR, and MS to see the footprints of these new dancers, identifying who they are and how they formed. **Retrosynthetic Analysis** then allows us to look back at the dance, predicting the steps and paths that led to each product.

So, the ultimate question is: how many possible products can emerge from this dance, excluding the simplest one, ethene? To answer this, you need to consider all possible combinations and rearrangements of the original molecules, bearing in mind their racemic nature and the nature of the Grubbs catalyst's weaving. Each unique pathway and its stereochemical outcome contribute to the final count of possible products.

By piecing together our understanding of organic chemistry, stereochemistry, organometallic reactions, and the principles of reaction kinetics and thermodynamics, we embark on this analytical journey to uncover the complete spectrum of possible products.

---

To determine the number of possible products formed when racemic 3-methylpent-1-ene reacts with the Grubbs catalyst, we need to understand the mechanism and outcomes of olefin metathesis, particularly focusing on the potential combinations and stereoisomers that can arise.

Understanding Olefin Metathesis with Grubbs Catalyst:
The Grubbs catalyst facilitates the exchange of alkylidene fragments (double-bonded carbons) between alkenes, leading to new alkene products. For 3-methylpent-1-ene, olefin metathesis will likely result in the formation of various possible products based on how the double bonds are reformed.

Racemic Mixture and Stereochemistry:
A racemic mixture contains equal amounts of two enantiomers. Since 3-methylpent-1-ene has a chiral center, it exists as two enantiomers, which can participate independently in the metathesis reactions.

Analyzing Potential Products:
Here's the detailed breakdown:

1. Initial Structures
   - 3-methylpent-1-ene: CH2=CH-CH2-CH(CH3)-CH3

2. Metathesis Reaction:
   - When two molecules of 3-methylpent-1-ene react, they can form different products depending on how the double bonds are broken and reformed.

3. Possible Product Structures:
   - One common product would be the symmetric product formed by the intermolecular reaction:
   - 4,5-dimethyl-2-heptene: This product results from two 3-methylpent-1-ene molecules linking at their double bonds.
     - Since the original molecule is racemic, each product can also have stereoisomers:
   - (E)-4,5-dimethyl-2-heptene: Trans configuration around the double bond.
   - (Z)-4,5-dimethyl-2-heptene: Cis configuration around the double bond.

Counting the Products:
- For each possible product, there are stereoisomers:
   - 4,5-dimethyl-2-heptene (both E and Z forms).

However, considering the racemic nature of the starting material, each product configuration (E/Z) is counted only once because they are indistinguishable in the racemic mixture.

Conclusion:
Given the racemic nature and possible stereoisomers formed via olefin metathesis, the potential products would be:
- (E)-4,5-dimethyl-2-heptene
- (Z)-4,5-dimethyl-2-heptene

Each stereoisomer counts as one distinct product since the starting material is racemic, leading to a total of 2 unique products.

Therefore, the correct answer is: **1. 2** ✅

---

# Chain of Thought (CoT)

To determine the number of possible products formed when racemic 3-methylpent-1-ene is treated with the Grubbs catalyst, let's consider the following steps:

Step 1: Understanding the Reaction
Grubbs catalyst is used in olefin metathesis reactions. In this process, alkenes undergo a redistribution of alkylidene groups. Essentially, two alkenes exchange their substituents to form new alkenes.

Step 2: Structure of Racemic 3-Methylpent-1-ene
The structure of 3-methylpent-1-ene is:
$$ \text{CH}_2=CH-\text{CH}_2-\text{CH}(\text{CH}_3)-\text{CH}_3 $$
Racemic indicates that we have both enantiomers in equal amounts.

Step 3: Possible Metathesis Outcomes
When two molecules of 3-methylpent-1-ene undergo metathesis, the reaction can yield several different products depending on how the alkylidene fragments rearrange.

Step 4: Identifying Different Metathesis Combinations
Consider the double bond positions and the substituents:

1. The reaction involves splitting each alkene into two fragments, then recombining these fragments.
2. Each molecule of 3-methylpent-1-ene can split into two fragments:
   - CH_2=CH-
   - -CH_2-CH(CH_3)-CH_3

Step 5: Combining Fragments
These fragments can recombine in several ways:

1. Symmetric recombination (self-metathesis):
   - Combining two identical fragments.
   - This leads to the same product as the starting material.

2. Asymmetric recombination (cross-metathesis):
   - Combining different fragments from different molecules.
Let's enumerate the products:
- Combining CH_2=CH- with itself:
   - Leads back to the starting material, which is counted as one unique product (1).
- Combining CH_2=CH- with -CH_2-CH(CH_3)-CH_3:
   - Leads to new product: CH_2=CH-CH_2-CH=CH-CH_3
- Combining -CH_2-CH(CH_3)-CH_3 with itself:
   - Leads to the product: CH_3-CH=CH-CH_2-CH(CH_3)-CH_3

Step 6: Counting Products
1. From symmetric recombination:
   - Starting material: 1 product.
2. From asymmetric recombination:
   - Different new alkenes can be formed when fragments from different molecules combine.

Summary of Products:
- The unique alkenes can be:
   - The starting material (racemic 3-methylpent-1-ene)
   - The recombination product CH_2=CH-CH_2-CH=CH-CH_3
   - The recombination product from two -CH_2-CH(CH_3)-CH_3 fragments.

After accounting for all possibilities, we should recognize that symmetric recombinations do not change the count.

Therefore, the total number of unique products, excluding ethene (produced as a byproduct), is: 4
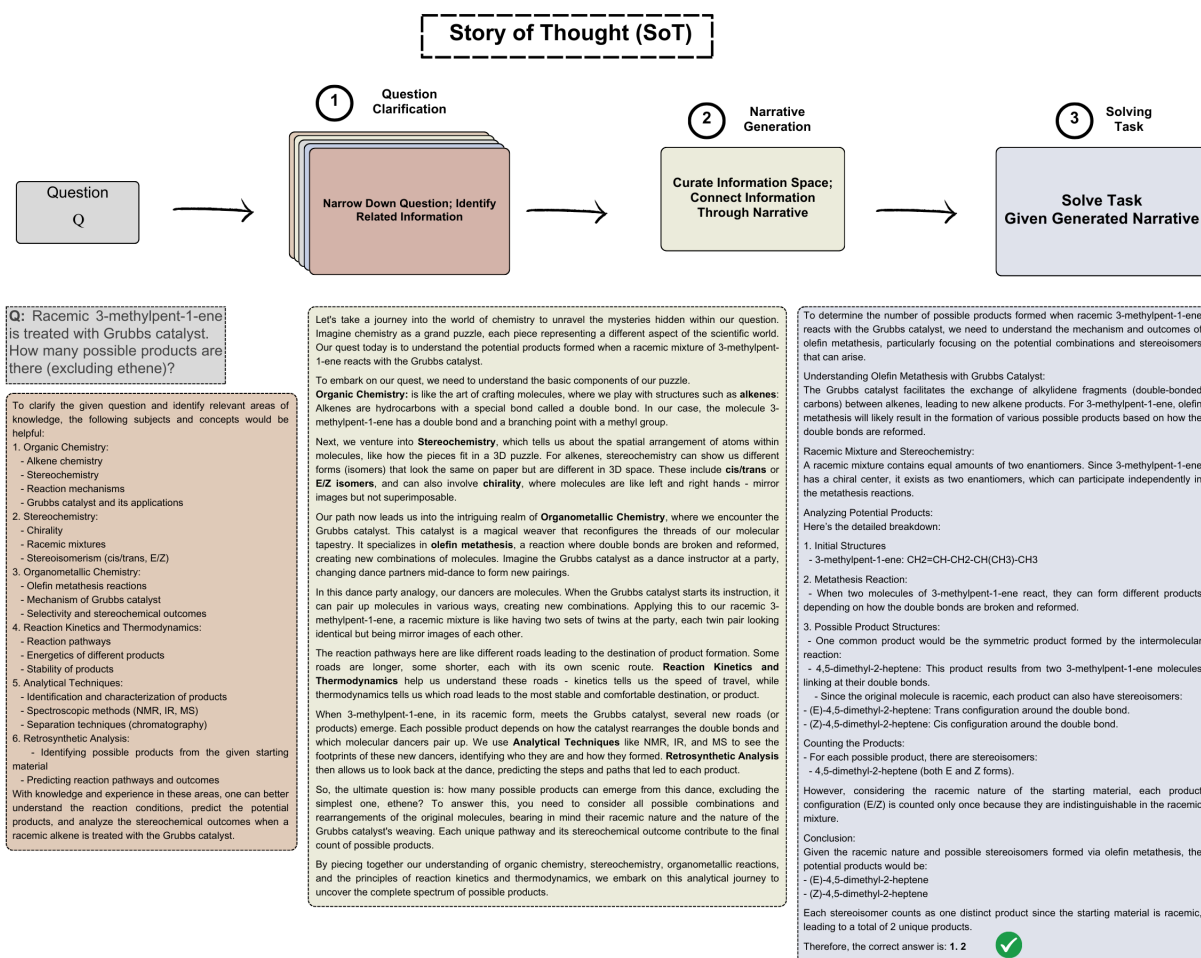
So, the correct answer is **4**. ❌

Figure 4: An actual example of SoT.