

---

# Can Stories Help LLMs Reason?

## Curating Information Space Through Narrative

---

Anonymous Author(s)

Affiliation

Address

email

### Abstract

1 Narrative is widely recognized as a powerful tool for structuring information  
2 and facilitating comprehension of complex ideas in various domains, such as  
3 science communication. This paper investigates whether incorporating narrative  
4 elements can assist Large Language Models (LLMs) in solving complex tasks more  
5 effectively. We propose a novel approach, **Story of Thought (SoT)**, integrating  
6 narrative structures into prompting techniques for problem-solving tasks. This  
7 approach involves constructing narratives around problem statements and creating  
8 a framework to identify and organize relevant information. We hypothesize that this  
9 narrative-based information curation process enhances problem comprehension by  
10 contextualizing critical information and highlighting causal relationships within the  
11 problem space. Our experimental results show that the SoT approach consistently  
12 surpasses Chain of Thought (CoT) and Analogical Reasoning in GPQA tasks,  
13 achieving higher accuracy and better solutions in physics, chemistry, and biology  
14 problem-solving tasks with all tested OpenAI, Meta, and Mistral LLMs.

### 15 1 Introduction

16 Humans have an exceptional ability to understand and reason through narratives. A narrative-driven  
17 approach can enhance the comprehension and retention of complex subjects compared to simple  
18 fact listing Fisher [2021], Abbott [2020], Gottschall [2012]. For example, storytelling effectively  
19 structures information in science communication Dahlstrom [2014], Norris et al. [2005], Martinez-  
20 Conde and Macknik [2017], education Engel et al. [2018], Negrete and Lartigue [2004], and health  
21 communication Dudley et al. [2023], revealing relationships and contextual nuances Zak [2015]. As  
22 shown in Figure 1, the *factual approach* presents information in a concise manner akin to a reference  
23 source, whereas the *narrative approach* conveys information through storytelling to contextualize  
24 facts within a broader setting.

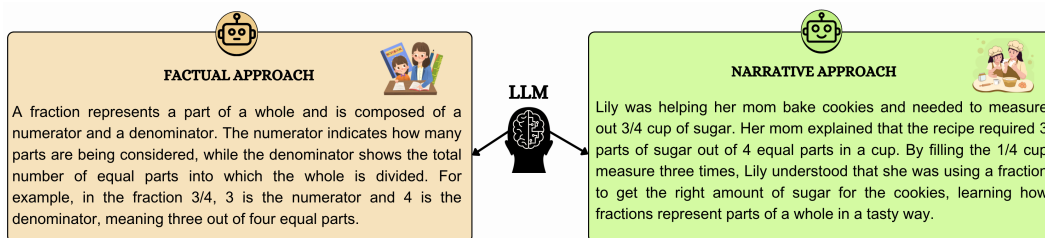


Figure 1: Contrasting approaches to information delivery, illustrated on explaining the concept of fractions: Factual vs. Narrative.

25 To date, large language models (LLMs) struggle with complex problem-solving tasks that require the  
26 ability to integrate, structure, and apply relevant information effectively Qiao et al. [2023], Wang  
27 et al. [2023]. Prompting techniques based on breaking tasks into smaller subtasks, such as Chain-  
28 of-Thought (CoT) Wei et al. [2022] and its more recent adaptations Xia et al. [2024], have led to  
29 considerable improvements in problem-solving benchmarks. The strategies of constructing natural  
30 language rationales Ling et al. [2017], in the CoT context also called reasoning processes, play a  
31 vital role in LLM prompting Ye and Durrett [2024], Min et al. [2022], Wang et al. [2022], Li et al.  
32 [2023].

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

36 **RQ 1:** Can LLMs generate coherent and relevant narratives around problem statements to facilitate  
37 problem comprehension and reasoning?

38 **RQ 2:** Can incorporating narrative elements into LLM prompting techniques improve their perfor-  
39 mance on complex problem-solving tasks?

40 We make the following contributions to the RQs:

41 (i) We introduce a novel method, called **Story of Thought (SoT)**, that aids LLMs to identify and  
42 arrange relevant information for solving complex tasks by incorporating narrative structures into the  
43 prompting process. (ii) We evaluate the effectiveness of SoT on diverse, complex tasks, including  
44 physics, chemistry, and biology problem-solving in GPQA, showing superior performance to existing  
45 task-decomposition-based prompting techniques, such as zero-shot and few-shot chain of thought  
46 and analogical reasoning. (iii) We analyze the impact of individual narrative techniques used on the  
47 generated narrative-based explanation to investigate why they improve LLMs reasoning capabilities.

## 48 2 Related Work

49 **Narrative** *Narrative*, as a noun, refers to a story or a description of a series of events <sup>1</sup>. In other  
50 words, it is a particular way of explaining or understanding events and plays a crucial role in human  
51 communication and cognition Hineline [2018]. The terms “story” and “narrative” are often used  
52 interchangeably. However, there is a subtle difference between them. A “story” is a narrative’s  
53 content or substance, while a “narrative” is the structure or way the story is presented Abbott [2020].  
54 Narrative plays a crucial role in various aspects of human communication and cognition. Bruner  
55 [1991] argues that narrative is a fundamental mode of human thought, allowing individuals to  
56 organize and make sense of their experiences. However, there are also potential disadvantages to  
57 using narrative. One concern is that narratives can oversimplify complex issues or events, leading to  
58 a reductionist understanding Dahlstrom and Ho [2012]. Furthermore, an over-reliance on narrative  
59 structures may limit the exploration of alternative viewpoints or non-linear forms of information  
60 presentation Negrete and Lartigue [2004]. It is essential to balance the benefits of narrative and the  
61 need for a nuanced understanding of the problem space Avraamidou and Osborne [2009].

62 **Narrative and Human Cognition** Research into *narrative transportation* examines how individu-  
63 als become cognitively and emotionally absorbed in stories. This immersive experience enhances  
64 emotional responses and alters attitudes and beliefs by aligning the listener’s brain with the sto-  
65 ryteller’s Oschatz and Marker [2020], Bilandzic et al. [2020]. Neuroimaging techniques such as  
66 functional Magnetic Resonance Imaging (fMRI), Positron Emission Tomography (PET), and Magne-  
67 toencephalography (MEG) show differences in identified active brain regions involved in narrative  
68 comprehension compared to factual processing Sanford and Emmott [2012], Armstrong [2020],  
69 Aboud et al. [2019], Coopmans and Cohn [2022]. Furthermore, presenting information in narratives  
70 can enhance learning and memory, as well as promote engagement and motivation Willingham [2004],  
71 Rowe et al. [2010], Chen et al. [2023], which led to the development of narrative-based educational  
72 strategies Bower and Clark [1969], Mawasi et al. [2020], Norris et al. [2005].

73 **Role of Narrative in Solving Tasks** In problem-solving, narratives can serve as a framework for  
74 organizing and presenting information relevant to the task Jonassen and Hernandez-Serrano [2002],

---

<sup>1</sup><https://dictionary.cambridge.org/>

75 Andrews et al. [2009]. Structuring the problem space as a narrative makes it easier to identify key ele-  
 76 ments, such as characters, goals, obstacles, and potential solutions San Roque et al. [2012]. Narratives  
 77 can also help to break down complex problems into sub-problems, providing a step-by-step ap-  
 78 proach to problem-solving Szurmak and Thuna [2013]. For example, progressive disclosure, analogy,  
 79 and analogical reasoning are powerful narrative techniques that facilitate problem-solving Salvucci  
 80 and Anderson [2001], Gick and Holyoak [1980]. These techniques involve presenting information  
 81 gradually in sub-problems, drawing comparisons and similarities between two seemingly disparate  
 82 concepts, and using these similarities to generate insights or solutions Norris et al. [2005], Holyoak  
 83 and Thagard [1989].

84 **Narrative and LLM Prompting** The intuitive approach to improving reasoning with LLMs is  
 85 prompt engineering (see recent survey in Qiao et al. [2023]. Starting from CoT prompting Wei et al.  
 86 [2022], these techniques leverage LLMs’ strong in-context learning ability, adding intermediate steps  
 87 to generate a reasoning process before answering. While analogical reasoning can be a technique  
 88 used in narrative generation, to our knowledge the narrative technique has never been fully explored  
 89 as a coherent set of interconnected didactic approaches to improve the reasoning abilities of LLMs  
 90 for problem-solving.

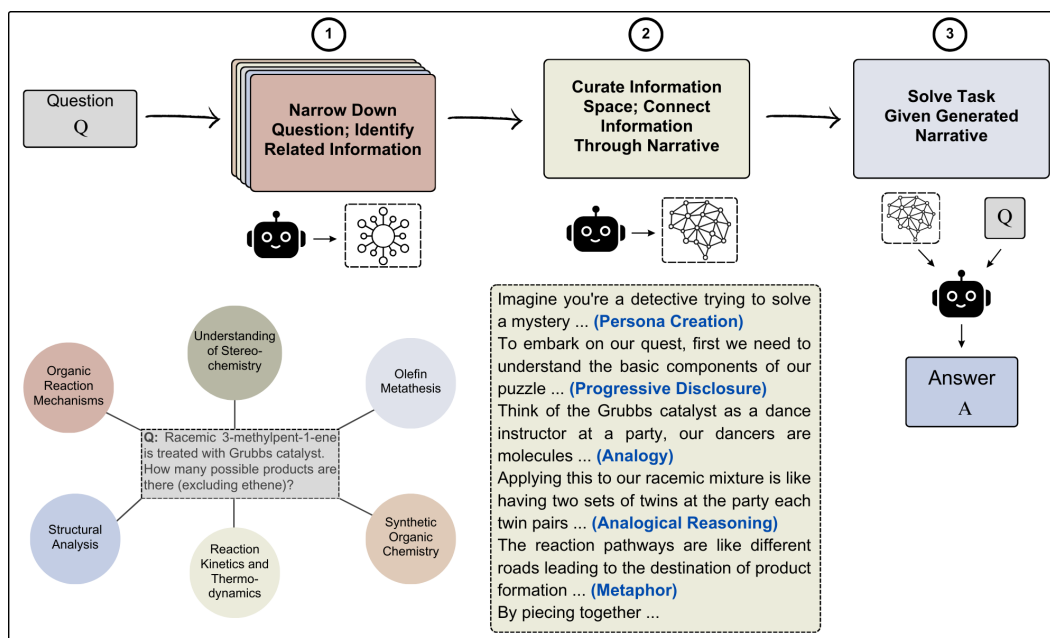


Figure 2: A high-level overview of **Story of Thought** (SoT), consisting of three steps (top): ① Question Clarification, ② Narrative Generation, ③ Solving Task and an actual example of LLM output (bottom) in each step for the GPQA task. The prompt designed for step 2 incorporates the narrative techniques (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 B for the complete prompt for each step. See Appendix A for a complete example.

### 91 3 Methodology

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

96 Figure 2 gives an overview of SoT. It involves three steps using narrative techniques: (i) **Question**  
97 **Clarification** (i.e., acting as an explorer to dissect and clarify complex questions (Section 3.1));  
98 (ii) **Narrative Generation** (i.e., generating detailed narratives from the clarified question components  
99 using different narrative techniques (Section 3.2)); and (iii) **Solving Task** (i.e., leveraging narratives  
100 to prompt the LLMs to solve the tasks (Section 3.3)).

### 101 3.1 Step 1: Question Clarification

102 In the first step, we use the LLM’s ability to explore and clarify the problem. Starting with a  
103 specialized prompt, the LLM breaks down the question into its core components, identifying relevant  
104 subtopics and areas. This detailed analysis is crucial for generating a coherent narrative that thoroughly  
105 addresses the question. The prompt is shown in Appendix B.1.

### 106 3.2 Step 2: Narrative Generation

107 The second step involves generating detailed narratives based on the breakdown and clarification  
108 performed in Step 1 (i.e., Question Clarification described in the previous section). These narratives  
109 provide a structured context for the questions to enhance the LLM’s understanding, reasoning, and  
110 problem-solving abilities. In synthesizing the literature discussed in Section 2, we integrate the  
111 following narrative techniques into our prompt below and task LLM to generate a narrative, based on  
112 the information identified in Step 1 (See Appendix B.2 for designed prompt):

- 113 1. **Progressive Disclosure:** Reveals information gradually, guiding the LLM step-by-step  
114 through the problem-solving process.
- 115 2. **Branching:** Explores different paths or approaches to understanding the problem by provid-  
116 ing multiple perspectives.
- 117 3. **Analogy:** Uses comparisons to familiar concepts or situations to make abstract components  
118 more understandable.
- 119 4. **Analogical Reasoning:** Facilitates understanding by reasoning through similarities between  
120 the problem and known situations.
- 121 5. **Metaphor:** Simplifies complex ideas through metaphorical representation.

### 122 3.3 Step 3: Solving Task

123 In the final step, the LLM uses the narrative generated in Step 2 to solve the original QA task. The  
124 structured and contextual understanding provided by the narrative supports LLM in accessing relevant  
125 aspects of the task. (The prompt is shown in Appendix B.3)

## 126 4 Experimental Setup

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

### 129 4.1 Evaluation Tasks

130 We focus our evaluation on reasoning-intensive tasks spanning multiple domains, including physics,  
131 biology, and chemistry problem-solving. In particular, we utilize the **GPQA** (Diamond set), a  
132 Graduate-level Problem-solving QA dataset Rein et al. [2023], which comprises expert-crafted  
133 multiple-choice questions. These tasks are diverse and extremely challenging, requiring in-depth  
134 reasoning and domain knowledge, making them well-suited for assessing our approach’s ability to  
135 understand complex tasks and contextualize salient information within the problem space.

### 136 4.2 Evaluated Large Language Models

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

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

Table 1: Performance (QA accuracy) of LLMs across prompting methods on GPQA (Diamond set).

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
<b>Zero-shot</b>	34.2	39.5	35.8	36.36	30.6	34.7	<b>28.79</b>	42.42
<b>Zero-shot CoT</b>	40.91	41.92	31.82	35.35	28.1	35.7	24.75	39.39
<b>Analogical Reasoning (3-shot)</b>	40.91	47.47	37.9	26.26	28.1	41.41	16.67	<b>48.48</b>
<b>Ours: Knowledge Identification</b>	40.4	48.99	35.35	37.77	27.77	40.90	20.71	37.88
<b>Ours: Story of Thought (SoT)</b>	<b>43.43</b>	<b>51.01</b>	<b>38.4</b>	<b>38.89</b>	<b>30.8</b>	<b>48.98</b>	22.73	36.36

141 These models were selected to cover a wide spectrum of capabilities and sizes, enabling a compre-  
 142 hensive evaluation of their strengths and limitations. By including models from multiple leading AI  
 143 research organizations, we aim to provide a balanced comparison.

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

### 149 4.3 Prompting Methods Benchmarked

150 We compared our proposed approach against several prompting techniques, including:

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

155 **Zero-shot CoT** Wei et al. [2022]: This technique extends the zero-shot prompting approach by  
 156 encouraging the LLM to explicitly reason through the steps required to arrive at an answer. By  
 157 prompting the model to generate a chain of thought, this method aims to improve the model’s ability  
 158 to solve complex problems by breaking them down into smaller, more manageable steps.

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

165 **Ours: Knowledge Identification:** To measure the effectiveness of our proposed approach, namely  
 166 utilizing narrative in solving tasks, we prompt LLMs to solve the task based solely on the generated  
 167 conceptual knowledge from Step 1 (described in Section 3.1). This allows us to compare the  
 168 model capability in solving tasks using only the identified relevant knowledge versus leveraging this  
 169 knowledge to structure a coherent narrative.

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

## 172 5 Results

### 173 5.1 Benchmark Performance Results

174 The main results of our experiments on the GPQA task are presented in Table 1. We evaluate  
 175 the performance of various prompting methods across eight different LLMs from four major AI  
 176 companies: Meta, Mistral, OpenAI, and Microsoft. Our proposed prompt-driven reasoning approach  
 177 (SoT), consistently outperformed the baseline approaches, including zero-shot prompting, zero-shot  
 178 Chain-of-Thought (CoT) prompting, and analogical reasoning, for six out of the eight LLMs tested.  
 179 This finding highlights the potential of leveraging narrative structures to improve the ability of LLMs  
 180 to understand and reason about the given information in various intensive-reasoning tasks across a  
 181 range of models. In particular, the open-source Llama 3 70B model records the highest accuracy

<sup>2</sup><https://huggingface.co/docs/transformers>

182 using the SoT method, achieving a score of 51.01%. This is the highest accuracy observed among  
 183 all models and methods tested in the study, and, at the time of writing also a state-of-the-art result  
 184 compared to public leaderboards (including e.g., Claude 3 and Gemini 1.5<sup>3</sup>). Furthermore, the  
 185 GPT-4 model shows the most notable improvement in accuracy when the SoT method is employed,  
 186 compared to its zero-shot baseline. Specifically, the accuracy for GPT-4 increased from 34.7% under  
 187 zero-shot conditions to 48.98% with SoT (i.e., an absolute increase of 14.28%, or a relative increase  
 188 of 41% respectively).

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

## 194 5.2 Role of the Narrative Quality/Choice

195 We further investigate the role of the choice of *narrator* model (i.e., the model that generates  
 196 narratives) for problem-solving tasks. In the following experiments, we apply the narratives generated  
 197 by other large and small open-source LLMs to the Phi-3 Mini and Phi-3 Medium models. The results  
 198 of these experiments are presented in Table 2.

199 We observe that the narratives generated by the Llama 3 8B, Llama 3 70B, and Mistral 7B models  
 200 consistently improve the accuracy of both Microsoft models compared to the baseline (i.e., when both  
 201 models use their own generated narratives in Step 2 to solve the tasks, shown in Table 1). The absolute  
 202 improvements range from 1.0% to 2.8%, with the Llama 3 70B model generating the most effective  
 203 narratives. A slight decrease in accuracy is observed with the mixture-of-experts Mixtral 8x7B  
 204 narratives for the Phi-3 Medium model, highlighting the need for careful selection and evaluation of  
 205 narrator models to ensure compatibility and optimal performance.

Table 2: Applying generated narratives by open-source models to Microsoft models to solve the tasks.

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↑)	<b>39.39</b> (+2.79↑)
Mistral 7B	24.24 (+1.51↑)	38.38 (+1.78↑)
Mixtral 8x7B	24.74 (+2.01↑)	35.86 (-0.74↓)

## 206 5.3 Impact of Narrative Elements

207 To measure the impact of each of the five individual narrative techniques, we jointly prompted on the  
 208 performance of open-source Meta models, we ablate the designed prompt in Step 2 (of Section 3.2)  
 209 to apply each of the techniques separately. The results in Table 3 indicate that employing any  
 210 single narrative technique at a time is notably less effective at boosting QA accuracy than utilizing a  
 211 combination of these simultaneously.

Table 3: 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).

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

<sup>3</sup><https://klu.ai/glossary/gpqa-eval>

212 For both models (Llama 3 8B and 7B), the decrease in accuracy is comparably smaller (-3.0%  
 213 to -5.6%) when using only the analogical components of the narrative (*Analogy* and *Analogical*  
 214 *Reasoning*) than when using only the structural instructions (*Progressive Disclosure* or *Branching*)  
 215 which leads to larger (-6.0% to -9.1%) accuracy loss.

216 However, reasoning alone does not perform on par with the full narrative generation listing all the  
 217 techniques. Prompting for *Metaphor* usage only leads to a larger accuracy loss in the 70B model  
 218 (-6.6%) compared to the smaller one (-2.0%). This makes us wonder to which extent the narrative  
 219 techniques are correlated, and to which extent the model can “understand” what it is prompted for,  
 220 which we attempt to analyze in the following subsections.

## 221 5.4 Analyzing Generated Narratives

222 To gain deeper insights into the generated narratives, we designed a prompt (shown below) that  
 223 utilizes our best-performing model (LLama 3 70B) to annotate the number of occurrences of each  
 224 narrative technique for each generated narrative by all models used in our experiments. The intuition  
 225 behind this experiment is that we can better interpret how the model executed the narrative technique  
 226 prompt, by asking it to label if and where the mentioned techniques are used in the text generated.  
 227 Less frequently labeled techniques might be the ones where LLM doesn’t have a clear understanding  
 228 of what it is asked to do. A proportion of the techniques and their correlation can provide us with a  
 229 better picture of LLM’s interpretation of the instruction as well.

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

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
<b>Progressive Disclosure</b>	427	597	191	191	744	570	367	368
<b>Branching</b>	30	56	51	20	72	168	34	61
<b>Analogy</b>	418	425	117	161	498	595	569	499
<b>Analogical Reasoning</b>	205	191	78	108	213	336	276	206
<b>Metaphor</b>	249	316	103	137	811	428	418	291
$\Sigma$	1329	1585	540	617	2338	2097	1664	1425

230 We aim to uncover patterns and variations in the use of narrative techniques across different LLMs.  
 231 Table 4 indicates a comparison of the total number of occurrences for each narrative technique across  
 232 various LLMs.

233 **Variability in Utilization of Narrative Techniques Across Models:** In our designed prompt in  
 234 Step 2 (i.e., Narrative Generation, described in Section 3.2), we task LLMs to generate narrative using  
 235 all of the 5 narrative techniques. However, as Table 4 indicates, not all techniques were employed  
 236 equally. The result reveals that while some techniques like *Analogy* and *Progressive Disclosure* were  
 237 consistently utilized, others such as *Branching* were applied less frequently.

238 We observe a trend across all LLM families where models with larger capacities, such as Llama 3 70B  
 239 and GPT-4, consistently show higher occurrences of narrative techniques compared to their smaller  
 240 counterparts. Furthermore, OpenAI’s models (ChatGPT 3.5 & GPT-4) demonstrate the highest total  
 241 occurrences of narrative techniques, with 2,338 and 2,097, respectively with a notable emphasis on  
 242 *Metaphors* and *Analogies*.

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

## 249 6 Limitations

250 **Contribution limitations.** The occurrences of narrative techniques do not necessarily imply the  
 251 quality or effectiveness of the generated narratives; rather, they provide insights into the models’  
 252 tendencies and preferences in employing these techniques. Therefore, answering the question of

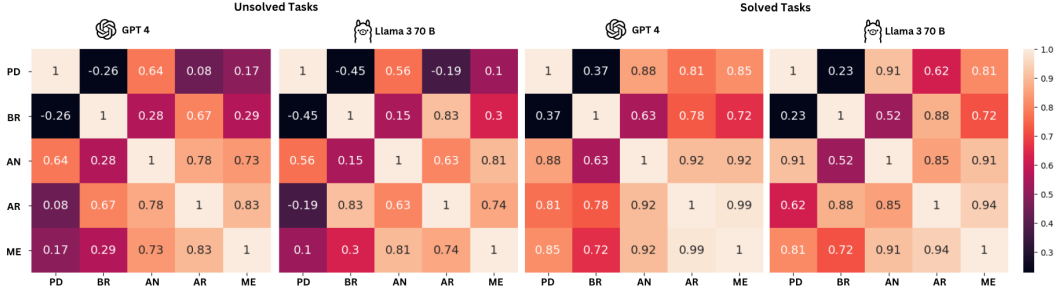


Figure 3: Correlation coefficients among all narrative techniques (**PD** = Progressive Disclosure, **BR** = Branching, **AN** = Analogy, **AR** = Analogical Reasoning, **ME** = Metaphor) used in the SoT approach for GPT-4 and Llama 3 70 B in solved and unsolved tasks.

Table 5: 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.

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

253 why narrative is helping LLMs is more complex and needs to be further investigated by looking into  
 254 different research areas such as cognitive and communication theories.

255 **Method limitations.** This method might not be efficient for tasks such as the MMLU benchmark,  
 256 where the answer to the question depends on the provided options, because part of the information  
 257 necessary to determine the correct answer is contained within the options themselves. To address this,  
 258 we may include the options' information as part of the question, thereby ensuring that all relevant  
 259 information is available for the method to process and derive the correct answer.

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

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

271 **LLM Robustness Limitations (Position of Correct Option).** In the original GPQA dataset  
 272 used for our experiments, the correct answers are always presented as the first option among the  
 273 multiple choices. However, To further evaluate the robustness of the LLMs, we conduct an additional  
 274 experiment where the correct answers are placed in the second option instead. Table 5 presents the  
 275 results of these experiments, comparing the performance of various prompting methods across six  
 276 different open-source LLMs. We observe that most LLMs experience a significant drop in accuracy  
 277 when the correct answer is moved to the second option. However, despite the overall decrease in  
 278 accuracy, our proposed approach, Story of Thought (SoT), consistently outperforms the baseline  
 279 methods for most LLMs. The SoT method achieves the highest accuracy for the Meta Llama 3  
 280 8B, Meta Llama 3 70B, Mistral 8x7B, and Microsoft Phi-3 Medium models, demonstrating its  
 281 effectiveness in enhancing the robustness of LLMs to changes in the problem structure.



## 282 7 Conclusions

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

## 289 References

- 290 H Porter Abbott. *The Cambridge introduction to narrative*. Cambridge University Press, 2020.
- 291 Katherine S Aboud, Stephen K Bailey, Stephanie N Del Tufo, Laura A Barquero, and Laurie E  
292 Cutting. Fairy tales versus facts: Genre matters to the developing brain. *Cerebral Cortex*, 29(11):  
293 4877–4888, 2019.
- 294 Dee H Andrews, Thomas D Hull, and Jennifer A Donahue. *Storytelling as an instructional method:*  
295 *Descriptions and research questions*. 2009.
- 296 Paul B Armstrong. *Stories and the brain: The neuroscience of narrative*. Johns Hopkins University  
297 Press, 2020.
- 298 Lucy Avraamidou and Jonathan Osborne. The role of narrative in communicating science. *Interna-*  
299 *tional Journal of Science Education*, 31(12):1683–1707, 2009.
- 300 Helena Bilandzic, Susanne Kinnebrock, and Magdalena Klingler. The emotional effects of science  
301 narratives: a theoretical framework. *Media and Communication*, 8(1):151–163, 2020.
- 302 Gordon H Bower and Michal C Clark. Narrative stories as mediators for serial learning. *Psychonomic*  
303 *science*, 14(4):181–182, 1969.
- 304 Jerome Bruner. The narrative construction of reality. *Critical inquiry*, 18(1):1–21, 1991.
- 305 Althea Y Chen, Chun-Ching Chen, and Wen-Yin Chen. The design narrative in design learning:  
306 Adjusting the inertia of attention and enhancing design integrity. *The Design Journal*, 26(4):  
307 519–535, 2023.
- 308 Cas W Coopmans and Neil Cohn. An electrophysiological investigation of co-referential processes  
309 in visual narrative comprehension. *Neuropsychologia*, 172:108253, 2022.
- 310 Michael F Dahlstrom. Using narratives and storytelling to communicate science with nonexpert  
311 audiences. *Proceedings of the national academy of sciences*, 111(supplement\_4):13614–13620,  
312 2014.
- 313 Michael F Dahlstrom and Shirley S Ho. Ethical considerations of using narrative to communicate  
314 science. *Science Communication*, 34(5):592–617, 2012.
- 315 Matthew Z Dudley, Gordon K Squires, Tracy M Petroske, Sandra Dawson, and Janesse Brewer. The  
316 use of narrative in science and health communication: a scoping review. *Patient Education and*  
317 *Counseling*, page 107752, 2023.
- 318 Alison Engel, Kathryn Lucido, and Kyla Cook. Rethinking narrative: Leveraging storytelling for  
319 science learning. *Childhood Education*, 94(6):4–12, 2018.
- 320 Walter R Fisher. *Human communication as narration: Toward a philosophy of reason, value, and*  
321 *action*. University of South Carolina Press, 2021.
- 322 Mary L Gick and Keith J Holyoak. Analogical problem solving. *Cognitive psychology*, 12(3):  
323 306–355, 1980.
- 324 Jonathan Gottschall. *The storytelling animal: How stories make us human*. Houghton Mifflin  
325 Harcourt, 2012.

- 326 Philip N Hineline. Narrative: Why it's important, and how it works. *Perspectives on Behavior*  
327 *Science*, 41:471–501, 2018.
- 328 Keith J Holyoak and Paul Thagard. A computational model of analogical problem solving. *Similarity*  
329 *and analogical reasoning*, 242266, 1989.
- 330 David H Jonassen and Julian Hernandez-Serrano. Case-based reasoning and instructional design:  
331 Using stories to support problem solving. *Educational technology research and development*, 50  
332 (2):65–77, 2002.
- 333 Yifei Li, Zeqi Lin, Shizhuo Zhang, Qiang Fu, Bei Chen, Jian-Guang Lou, and Weizhu Chen. Making  
334 language models better reasoners with step-aware verifier. In Anna Rogers, Jordan Boyd-Graber,  
335 and Naoaki Okazaki, editors, *Proceedings of the 61st Annual Meeting of the Association for*  
336 *Computational Linguistics (Volume 1: Long Papers)*, pages 5315–5333, Toronto, Canada, July  
337 2023. Association for Computational Linguistics. doi: 10.18653/v1/2023.acl-long.291. URL  
338 <https://aclanthology.org/2023.acl-long.291>.
- 339 Wang Ling, Dani Yogatama, Chris Dyer, and Phil Blunsom. Program induction by rationale genera-  
340 tion: Learning to solve and explain algebraic word problems. In Regina Barzilay and Min-Yen Kan,  
341 editors, *Proceedings of the 55th Annual Meeting of the Association for Computational Linguistics*  
342 *(Volume 1: Long Papers)*, pages 158–167, Vancouver, Canada, July 2017. Association for Computa-  
343 tional Linguistics. doi: 10.18653/v1/P17-1015. URL <https://aclanthology.org/P17-1015>.
- 344 Susana Martinez-Conde and Stephen L Macknik. Finding the plot in science storytelling in hopes of  
345 enhancing science communication. *Proceedings of the National Academy of Sciences*, 114(31):  
346 8127–8129, 2017.
- 347 Areej Mawasi, Peter Nagy, and Ruth Wylie. Systematic literature review on narrative-based learning  
348 in educational technology learning environments (2007-2017). 2020.
- 349 Sewon Min, Xinxu Lyu, Ari Holtzman, Mikel Artetxe, Mike Lewis, Hannaneh Hajishirzi, and Luke  
350 Zettlemoyer. Rethinking the role of demonstrations: What makes in-context learning work? In  
351 *EMNLP*, 2022.
- 352 Aquiles Negrete and Cecilia Lartigue. Learning from education to communicate science as a good  
353 story. *Endeavour*, 28(3):120–124, 2004.
- 354 Stephen P Norris, Sandra M Guilbert, Martha L Smith, Shahram Hakimelahi, and Linda M Phillips.  
355 A theoretical framework for narrative explanation in science. *Science education*, 89(4):535–563,  
356 2005.
- 357 Corinna Oschatz and Caroline Marker. Long-term persuasive effects in narrative communication  
358 research: A meta-analysis. *Journal of Communication*, 70(4):473–496, 2020.
- 359 Shuofei Qiao, Yixin Ou, Ningyu Zhang, Xiang Chen, Yunzhi Yao, Shumin Deng, Chuanqi Tan, Fei  
360 Huang, and Huajun Chen. Reasoning with language model prompting: A survey. In Anna Rogers,  
361 Jordan Boyd-Graber, and Naoaki Okazaki, editors, *Proceedings of the 61st Annual Meeting of*  
362 *the Association for Computational Linguistics (Volume 1: Long Papers)*, Toronto, Canada, July  
363 2023. Association for Computational Linguistics. doi: 10.18653/v1/2023.acl-long.294. URL  
364 <https://aclanthology.org/2023.acl-long.294>.
- 365 David Rein, Betty Li Hou, Asa Cooper Stickland, Jackson Petty, Richard Yuanzhe Pang, Julien Dirani,  
366 Julian Michael, and Samuel R Bowman. Gpqa: A graduate-level google-proof q&a benchmark.  
367 *arXiv preprint arXiv:2311.12022*, 2023.
- 368 Jonathan P Rowe, Lucy R Shores, Bradford W Mott, and James C Lester. Integrating learning and  
369 engagement in narrative-centered learning environments. In *Intelligent Tutoring Systems: 10th*  
370 *International Conference, ITS 2010, Pittsburgh, PA, USA, June 14-18, 2010, Proceedings, Part II*  
371 *10*, pages 166–177. Springer, 2010.
- 372 Dario D Salvucci and John R Anderson. Integrating analogical mapping and general problem solving:  
373 the path-mapping theory. *Cognitive Science*, 25(1):67–110, 2001.

- 374 Lila San Roque, Lauren Gawne, Darja Hoenigman, Julia Miller, Stef Spronck, Alan Rumsey, Alice  
375 Carroll, and Nicholas Evans. Getting the story straight: Language fieldwork using a narrative  
376 problem-solving task. 2012.
- 377 Anthony J Sanford and Catherine Emmott. *Mind, brain and narrative*. Cambridge University Press,  
378 2012.
- 379 Joanna Szurmak and Mindy Thuna. Tell me a story: The use of narrative as tool for instruction. 2013.
- 380 Xiaoxuan Wang, Ziniu Hu, Pan Lu, Yanqiao Zhu, Jieyu Zhang, Satyen Subramaniam, Arjun R  
381 Loomba, Shichang Zhang, Yizhou Sun, and Wei Wang. Scibench: Evaluating college-level  
382 scientific problem-solving abilities of large language models. *arXiv preprint arXiv:2307.10635*,  
383 2023.
- 384 Xuezhi Wang, Jason Wei, Dale Schuurmans, Quoc Le, Ed Huai hsin Chi, and Denny Zhou. Self-  
385 consistency improves chain of thought reasoning in language models. *ArXiv*, abs/2203.11171,  
386 2022. URL <https://api.semanticscholar.org/CorpusID:247595263>.
- 387 Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Fei Xia, Ed Chi, Quoc V Le, Denny  
388 Zhou, et al. Chain-of-thought prompting elicits reasoning in large language models. *Advances in  
389 neural information processing systems*, 35:24824–24837, 2022.
- 390 Daniel T Willingham. Ask the cognitive scientist the privileged status of story. *American Educator*,  
391 28:43–45, 2004.
- 392 Yu Xia, Rui Wang, Xu Liu, Mingyan Li, Tong Yu, Xiang Chen, Julian McAuley, and Shuai Li. Beyond  
393 chain-of-thought: A survey of chain-of-x paradigms for llms. *arXiv preprint arXiv:2404.15676*,  
394 2024.
- 395 Michihiro Yasunaga, Xinyun Chen, Yujia Li, Panupong Pasupat, Jure Leskovec, Percy Liang,  
396 Ed H Chi, and Denny Zhou. Large language models as analogical reasoners. *arXiv preprint  
397 arXiv:2310.01714*, 2023.
- 398 Xi Ye and Greg Durrett. The unreliability of explanations in few-shot prompting for textual reasoning.  
399 In *Proceedings of the 36th International Conference on Neural Information Processing Systems*,  
400 NIPS '22, Red Hook, NY, USA, 2024. Curran Associates Inc. ISBN 9781713871088.
- 401 Paul J Zak. Why inspiring stories make us react: The neuroscience of narrative. In *Cerebrum: the  
402 Dana forum on brain science*, volume 2015. Dana Foundation, 2015.

## 403 **A Story of Thought (SoT) vs. Chain of Thought (CoT)**

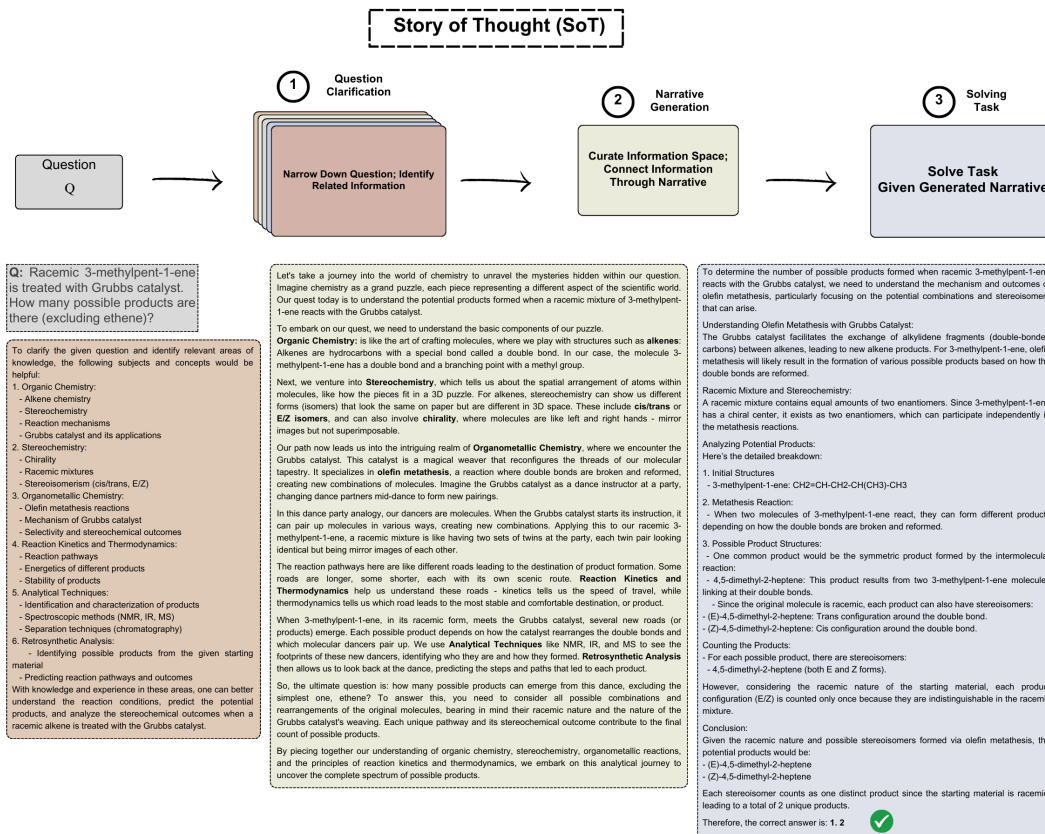


Figure 4: An actual example of SoT.

404 **B Designed Prompts**

405 **B.1 Step 1: Question Clarification**

406 You are an explorer who wants to identify and collect different related and  
407 specialized subject areas to clarify the question. Your goal is to narrow down  
408 the question and provide relevant areas of knowledge and experience you have  
409 that help clarify the question mentioned below. You should not answer the  
410 question.  
411  
412  
413 <question>

415 **B.2 Step 2: Narrative Generation**

416 You are an expert in narrative-based explanations for science communication. Your  
417 goal is to clarify the following question in a narrative way through the  
418 interconnected information provided below to enable a non-expert to comprehend  
419 the question in a more coherent and contextually rich manner. You should not  
420 answer the question.  
421  
422 Make sure to use all of these narrative techniques when clarifying the question  
423 through the interconnected information: Progressive Disclosure, Branching,  
424 Analogy, Analogical Reasoning, and Metaphor.  
425  
426  
427 <question>  
428  
429 <generated information in the previous step>  
430

431 **B.3 Step 3: Solving Task**

432 You are an expert in analyzing narrative-based explanations for solving tasks.  
433 Please answer the following question based on the following narrative-based  
434 clarification:  
435  
436  
437 <question>  
438  
439 Options:  
440 <options>  
441  
442 <generated narrative in the previous step>

444 **B.4 Analyzing Generated Narratives**

445 You are an expert in analyzing narrative-based explanations for science  
446 communication. Your goal is to find out which narrative techniques have been  
447 used in the following narrative-based explanation.  
448  
449 Label the narrative-based explanation using the following narrative-based techniques  
450 :  
451  
452 1. Progressive Disclosure  
453 2. Branching  
454 3. Analogy  
455 4. Analogical Reasoning  
456 5. Metaphor  
457  
458 <generated narrative>  
459