

Sketch-of-Thought: Efficient LLM Reasoning with Adaptive Cognitive-Inspired Sketching

Simon A. Aytes¹ Jinheon Baek¹ Sung Ju Hwang^{1,2}
KAIST¹ DeepAuto.ai²
{saytes, jinheon.baek, sungju.hwang}@kaist.ac.kr

Abstract

Recent advances in large language models (LLMs) have enabled strong reasoning capabilities through Chain-of-Thought (CoT) prompting, which elicits step-by-step problem solving, but often at the cost of excessive verbosity in intermediate outputs, leading to increased computational overhead. We propose *Sketch-of-Thought* (SoT), a prompting framework that integrates cognitively inspired reasoning paradigms with linguistic constraints to reduce token usage while preserving reasoning accuracy. SoT is designed as a flexible, modular approach and is instantiated with three paradigms—*Conceptual Chaining*, *Chunked Symbolism*, and *Expert Lexicons*—each tailored to distinct reasoning tasks and selected dynamically at test-time by a lightweight routing model. Across 18 reasoning datasets spanning multiple domains, languages, and modalities, SoT achieves token reductions of up to 84% with minimal accuracy loss. In tasks such as mathematical and multi-hop reasoning, it even improves accuracy while shortening outputs.

1 Introduction

Large language models (LLMs) have become central to a wide range of complex reasoning tasks across diverse domains, such as mathematics, science, and commonsense inference (Bubeck et al., 2023; Zhao et al., 2024). Even without dedicated training for reasoning, these models often exhibit emergent capabilities when prompted to decompose problems into intermediate steps (Wei et al., 2023). Chain-of-Thought (CoT) prompting (Wei et al., 2023) exemplifies this approach by encouraging step-by-step natural language reasoning, which has been shown to significantly improve performance on tasks such as logical inference and numerical problem solving (Sprague et al., 2024).

Despite its benefits, CoT often produces verbose outputs that dramatically increase token usage

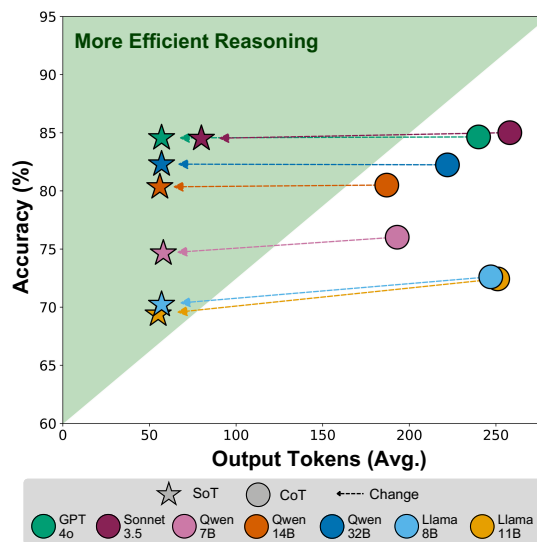


Figure 1: A comparison of accuracy and token usage in Chain-of-Thought (CoT) (Wei et al., 2023) and the proposed Sketch-of-Thought (SoT). Average scores for model performance across 18 datasets. Shaded region represents more efficient reasoning.

and computational overhead, making it less suitable for latency- or budget-constrained deployment scenarios (Nayab et al., 2025; Arora and Zanette, 2025). More sophisticated strategies, such as Self-Consistency (Wang et al., 2023b), Tree-of-Thoughts (Yao et al., 2023), and Graph-of-Thoughts (Besta et al., 2024), further expand the reasoning process via structured exploration, but tend to exacerbate inefficiencies in token usage.

To tackle these limitations, we introduce *Sketch-of-Thought* (SoT), a prompting framework that rethinks how language models externalize reasoning. Inspired by cognitive science, particularly the use of symbolic *sketches* as efficient mental intermediaries (Goel, 1995), SoT guides models to produce concise, structured reasoning steps that capture essential logic while avoiding full-sentence elaboration. These representations are analogous to mathematical notation or expert shorthand, preserving semantic fidelity while minimizing redundancy.

Code: <https://www.github.com/SimonAytes/SoT>

To implement this framework, we define three cognitively motivated reasoning paradigms: *Conceptual Chaining*, based on associative memory; *Chunked Symbolism*, grounded in working memory theory; and *Expert Lexicons*, inspired by domain-specific schemas used by specialists. Each paradigm is designed for a distinct class of reasoning tasks and is implemented using training-free prompts. To support adaptive paradigm selection, we incorporate a lightweight routing model that analyzes query structure to determine the most suitable reasoning style at inference time.

We extensively evaluate SoT on 18 reasoning datasets spanning mathematical, commonsense, logical, multi-hop, scientific, and medical domains. Experimental results show that SoT reduces output token usage by up to 84% compared to traditional CoT prompting, with no significant loss in accuracy—and even improving performance in some domains. Additional multilingual and multimodal evaluations demonstrate SoT’s ability to generalize across both languages and input modalities.

Our key contributions are as follows:

- We introduce *Sketch-of-Thought* (SoT), a prompting framework that leverages cognitively inspired reasoning paradigms to produce concise and structured model outputs.
- We present a lightweight routing model that dynamically selects the optimal reasoning paradigm based on the input query’s structure and semantics.
- On a battery of tests, we show that SoT significantly reduces token usage while maintaining or improving accuracy across diverse datasets, models, languages, and modalities.

2 Method

This section outlines the technical implementation of *Sketch-of-Thought* (SoT), a framework designed to improve the efficiency of reasoning in large language models while preserving performance. The end-to-end framework is illustrated in Figure 2.

2.1 Preliminary

We first provide a brief overview of large language models (LLMs) and their application to reasoning.

Large Language Models Large language models (LLMs) are trained on massive text corpora to predict the next token in a sequence. Formally, an

LLM with parameters θ receives an input token sequence x and produces an output sequence y , denoted as follows: $y = \text{LLM}_\theta(x)$. While LLMs exhibit strong performance across a variety of domains, their computational cost increases with the number of tokens processed and generated, making efficiency a critical concern.

Reasoning via Prompting Recent LLMs exhibit emergent reasoning capabilities—the ability to solve multi-step problems without task-specific training—when guided by carefully designed prompts. Letting P denote a prompt with instructions, exemplars, and the input query q , the process can be formalized as $r = \text{LLM}_\theta(P(q))$, where r includes both the reasoning trace $[s_1, s_2, \dots, s_m]$ and final answer a , as follows: $r = [s_1, s_2, \dots, s_m, a]$. Chain-of-Thought (CoT) (Wei et al., 2023) encourages decomposition into natural language steps.

Traditional CoT

Q: If a train travels 60 miles per hour for 3 hours, how far does it go?

A: <think> 1. I understand we need to find the total distance traveled by: A train moving at 60 miles per hour for a duration of 3 hours.

2. To calculate the distance, I'll use the formula:
Distance = Speed × Time
Distance = 60 miles/hour × 3 hours

3. Now I'll perform the calculation:
Distance = 60 × 3 = 180 miles

4. Verification:
This makes sense because the train moves 60 miles each hour. After 3 hours, it will have covered 3 times that distance. </think>

Answer: 180 miles

More advanced variants, for example, Tree-of-Thoughts and Graph-of-Thoughts (Yao et al., 2023; Besta et al., 2024), extend this idea by exploring alternative structural formulations. Despite their stylistic differences (leading to performance gains), all CoT-style prompting methods share a common limitation: increased token usage and reduced efficiency, when compared to standard direct-answer prompting (Arora and Zanette, 2025).

2.2 Sketch-of-Thought (SoT)

Sketch-of-Thought (SoT) mitigates reasoning inefficiency by restructuring how models express intermediate steps. Unlike prior methods that reduce prompt length via input compression (Jiang et al., 2023; Huang et al., 2024), SoT compresses the reasoning process using cognitively inspired prompts that elicit concise, structured steps.

Formally, we define different paradigm-specific prompt templates P_{SoT} , which steer the model to produce sketched reasoning: $[\hat{s}_1, \hat{s}_2, \dots, \hat{s}_m, a] = \text{LLM}_\theta(P_{\text{SoT}}(q))$, where each \hat{s}_i conveys the same

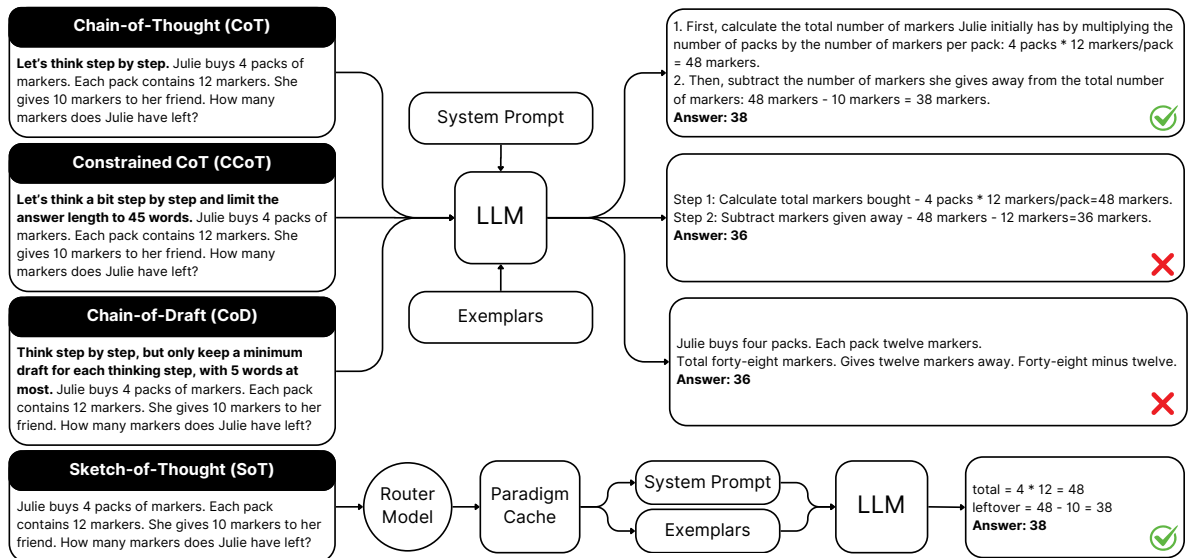
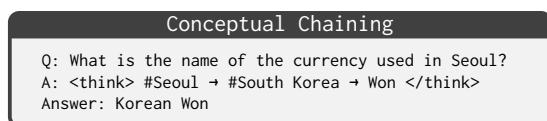


Figure 2: **Illustration of reasoning workflows**, including the input format, intermediate reasoning structure, and output style, across four prompting methods: Chain-of-Thought (CoT) (Wei et al., 2023), Constrained CoT (CCoT) (Nayab et al., 2025), Chain-of-Draft (CoD) (Xu et al., 2025), and Sketch-of-Thought (SoT). While CoT produces verbose natural language traces, CCoT and CoD apply explicit constraints on reasoning length. SoT introduces paradigm-guided sketching, yielding more compact yet structured intermediate steps via dynamic routing.

logical content as s_i (from CoT, for example), but using significantly fewer tokens, i.e., $|\hat{s}| < |s|$. These prompts enforce both linguistic constraints and cognitive structuring tailored to the task type.

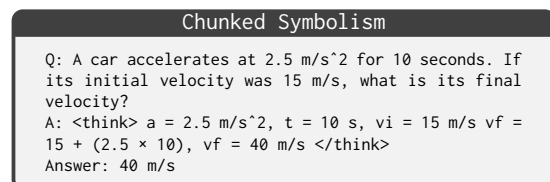
As an initial realization of SoT, we create three reasoning paradigms inspired by cognitive science, each designed to align with distinct patterns found across a range of reasoning tasks.

Conceptual Chaining. Rooted in cognitive science principles of how humans connect and retrieve related information, this paradigm creates concise logical sequences between key concepts. It draws from episodic buffer integration (Baddeley, 2000), the cognitive mechanism that temporarily holds and links information from different sources, and associative memory networks (Anderson, 1983), which describe how activating one concept automatically triggers related concepts in our minds (like how thinking of "rain" might immediately evoke "umbrella"). *Conceptual Chaining* extracts essential terms and presents reasoning as direct step-by-step pathways with minimal text.



Conceptual Chaining is particularly effective for commonsense, multi-hop, logical, and scientific reasoning tasks, where establishing structured relationships between ideas is critical.

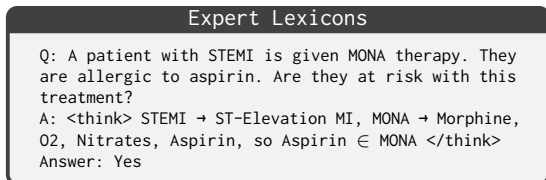
Chunked Symbolism. Based on working memory chunking theory (Miller, 1956), this paradigm organizes numerical and symbolic reasoning into compact, structured steps. This seminal cognitive science research showed that humans can only hold about 7 ± 2 (i.e., 5 to 9) distinct items in working memory at once, but we overcome this limitation by "chunking" related information into meaningful units—like remembering phone numbers as area code, prefix, and line number instead of 10 separate digits. *Chunked Symbolism* applies this principle by condensing mathematical reasoning into dense symbolic representations that pack more information into fewer tokens. It systematically extracts variables and performs operations while eliminating verbose explanations, using symbolic variables to transform natural language into a structured shorthand that preserves logical flow.



Chunked Symbolism excels in mathematical and arithmetic reasoning problems, where symbolic notation naturally compresses complex concepts.

Expert Lexicons. Inspired by expert schema research (Chi et al., 1981), this paradigm leverages domain-specific shorthand and specialized notation

to condense reasoning. This research demonstrated that experts in any field organize knowledge differently than novices—they develop mental frameworks (schemas) that allow them to quickly recognize patterns and use specialized terminology to communicate efficiently with peers. For example, a physician can convey complex medical conditions with a few acronyms that would require paragraphs of explanation for non-specialists. *Expert Lexicons* mimics this cognitive efficiency by employing domain-specific abbreviations, notation, and symbols that pack multiple concepts into single tokens. The example below demonstrates how domain-specialized reasoning can be compressed into concise notation while preserving the critical logical connections.



Expert Lexicons is particularly suited for technical disciplines, specialized reasoning tasks, and scenarios, where domain expertise enables significant information compression.

2.3 Adaptive Paradigm Selection

While manual selection among three paradigms is possible for each query based on heuristic rules, such an approach is impractical at scale. Instead, we introduce a lightweight routing model that selects the paradigm dynamically based on semantic and structural features of the input query.

Given a query q , the routing process is denoted as follows: $P_{\text{SoT}} = \text{ROUTER}(q)$, where P_{SoT} refers to the selected paradigm’s prompt-exemplar pair and ROUTER denotes the router model. We use DistilBERT (Sanh et al., 2020) as the base model due to its strong performance-efficiency trade-off and minimal inference overhead (see Appendix C.1).

Router Training We train the router model using 14,200 machine-labeled examples drawn from the training splits of the datasets outlined in Section 3.1. Each sample is labeled using GPT-4o (OpenAI, 2024), guided by a classification prompt derived from the paradigm definitions in Section 2.2. We provide this classification prompt in Appendix B.6. Additionally, we evaluate GPT-4o’s paradigm labeling performance in Appendix C.2.

To avoid overwhelming the router with irrelevant

input, we replace any long or non-textual context (e.g., images or documents) with a special placeholder token (e.g., [CONTEXT HERE]). This ensures that the model focuses solely on the question itself, which typically contains sufficient cues for determining the appropriate reasoning style.

3 Experimental Setup

3.1 Datasets

To ensure a comprehensive evaluation, we validate Sketch-of-Thought (SoT) across 15 datasets spanning six categories of reasoning, following the taxonomy introduced by Sun et al. (2024). The datasets are as follows: **Mathematical Reasoning** includes GSM8K, SVAMP, AQUA-RAT, and DROP (Cobbe et al., 2021; Patel et al., 2021; Ling et al., 2017; Dua et al., 2019); **Commonsense Reasoning** includes CommonsenseQA, OpenbookQA, and StrategyQA (Talmor et al., 2019; Mihaylov et al., 2018; Geva et al., 2021); **Logical Reasoning** includes LogiQA and ReClor (Liu et al., 2020; Yu et al., 2020); **Multi-Hop Reasoning** includes HotPotQA and MuSiQue-Ans (Yang et al., 2018; Trivedi et al., 2022); **Scientific Reasoning** includes QASC and Worldtree (Khot et al., 2020; Jansen et al., 2018); and **Medical Reasoning** includes PubMedQA and MedQA (Jin et al., 2019, 2020).

Beyond English textual reasoning, we include two additional evaluation tracks: a multilingual experiment using MMLU and its professionally translated variant MMMLU (Hendrycks et al., 2021), and a multimodal experiment using GQA (Hudson and Manning, 2019) and the image-based subset of ScienceQA (Lu et al., 2022). Further details regarding the datasets are provided in Appendix A.1.

3.2 Baselines

We mainly compare SoT against three established prompting-based reasoning strategies. Chain-of-Thought (CoT) (Wei et al., 2023) elicits step-by-step natural language reasoning. Constrained CoT (CCoT) (Nayab et al., 2025) introduces a global verbosity constraint, limiting the total reasoning chain to a fixed number of words—in our case, 45 words (CCoT-45). Chain-of-Draft (CoD) (Xu et al., 2025) adopts a similar compression strategy but imposes constraints at the step level, requiring each intermediate step be no longer than five words.

3.3 Implementation Details

A diverse set of instruction-tuned LLMs is selected, spanning both open-weight and proprietary offer-

Table 1: **Main Experimental Results.** Results are shown for Sketch-of-Thought (SoT), Chain-of-Thought (CoT) (Wei et al., 2023), Constrained Chain-of-Thought (CCoT) (Nayab et al., 2025), and Chain-of-Draft (CoD) (Xu et al., 2025). Results are grouped by reasoning type, with each entry representing the average over all associated datasets. "Acc" denotes accuracy and "Tkn" denotes the number of output tokens. In the Overall section, we report two additional metrics: the token reduction percentage (shown as "Red.") and the change in accuracy between CoT and the baseline (shown as " Δ "). The best results are in **bold** and the second-best are underlined.

		Reasoning Task															
		Mathematical		Commonsense		Logical		Multi-Hop		Scientific		Medical		Overall			
Method		Acc	Tkn	Acc	Tkn	Acc	Tkn	Acc	Tkn	Acc	Tkn	Acc	Tkn	Acc \uparrow	Tkn \downarrow	Red. \uparrow	$\Delta \uparrow$
Qwen 2.5-32B	CoT	84.17	222	91.48	177	71.23	298	79.44	155	92.89	213	67.66	292	82.24	222	-	-
	CoD	71.94	53	89.48	38	72.89	45	80.00	41	90.00	42	58.89	47	77.32	45	79.75	-4.92
	CCoT	80.50	76	88.82	49	72.78	60	80.11	54	88.89	49	57.66	65	79.16	61	72.56	-3.08
	SoT	86.94	88	92.00	34	71.00	66	81.89	43	91.34	31	61.11	63	82.30	57	<u>74.36</u>	0.06
Qwen 2.5-14B	CoT	83.00	190	91.41	150	67.00	248	77.67	149	90.89	164	65.11	234	80.50	187	-	-
	CoD	69.22	63	89.04	41	66.22	47	80.44	46	89.44	43	59.00	52	75.61	50	73.23	-4.89
	CCoT	81.33	115	90.52	58	70.00	89	78.89	91	89.44	55	61.44	86	79.76	85	54.49	-0.74
	SoT	82.72	78	90.89	37	67.44	63	79.89	45	90.89	37	62.56	63	80.34	56	<u>70.02</u>	-0.16
Qwen 2.5-7B	CoT	77.94	186	86.52	158	65.67	284	73.22	137	88.67	181	57.11	249	74.86	199	-	-
	CoD	66.83	57	84.74	37	64.33	49	76.11	43	87.00	39	55.89	48	72.48	46	76.88	-2.38
	CCoT	78.00	81	84.15	45	63.67	63	78.89	53	82.78	44	50.33	61	66.30	58	70.85	-8.56
	SoT	79.28	72	86.74	30	62.00	66	76.11	45	87.22	27	53.78	107	74.19	58	<u>70.85</u>	-0.67
Llama 3.1-8B	CoT	72.56	235	81.92	209	51.22	292	74.56	193	85.78	260	65.00	323	72.61	247	-	-
	CoD	55.28	73	80.67	45	47.22	58	73.22	49	81.00	47	66.22	55	66.56	56	77.31	-6.05
	CCoT	65.22	88	80.89	58	51.00	73	75.45	60	85.00	57	68.11	73	70.84	70	71.64	-1.77
	SoT	64.67	78	81.41	36	48.11	71	77.11	44	83.56	35	66.44	63	70.22	57	<u>76.91</u>	-2.39
Llama 3.2-11B	CoT	70.55	232	82.74	216	50.33	297	73.45	198	85.78	263	68.44	334	72.43	251	-	-
	CoD	56.17	67	80.89	43	48.22	51	74.00	46	79.44	44	65.00	50	66.71	52	79.25	-5.72
	CCoT	64.56	79	80.81	59	51.89	69	73.00	62	84.22	57	68.34	71	70.37	67	73.27	-2.06
	SoT	64.50	75	81.48	35	45.34	69	77.89	44	79.44	36	66.56	64	69.39	55	<u>78.06</u>	-3.04
GPT-4o	CoT	85.44	240	92.74	200	74.78	311	81.56	156	93.22	240	75.22	308	84.64	240	-	-
	CoD	83.17	71	87.11	50	71.56	62	82.56	53	90.67	55	46.33	63	78.41	60	74.95	-6.23
	CCoT	83.72	93	90.59	63	71.22	69	82.33	70	90.22	63	56.22	71	80.44	74	69.11	-4.20
	SoT	86.17	69	92.52	39	73.22	80	84.78	47	92.56	39	72.44	61	84.55	57	76.20	-0.09
Claude Sonnet 3.5	CoT	87.11	233	91.26	242	75.22	314	81.67	206	93.89	264	75.67	321	85.01	258	-	-
	CoD	82.00	78	91.33	61	75.78	96	82.00	63	91.33	67	76.22	105	83.51	77	70.16	-1.50
	CCoT	82.94	97	72.44	80	64.67	91	80.89	85	68.33	83	55.78	103	72.56	90	65.12	-12.45
	SoT	84.06	85	91.11	59	75.00	112	84.44	57	91.78	62	77.78	116	84.50	80	<u>68.99</u>	-0.51
All Models	CoT	80.11	220	88.30	193	65.06	292	77.37	171	90.16	226	67.75	294	78.12	233	-	-
	CoD	69.23	66	86.18	45	63.75	58	78.33	49	86.98	48	61.08	60	74.26	54	76.82	-3.86
	CCoT	76.61	90	81.17	59	63.60	73	78.51	68	81.27	58	59.70	76	73.48	71	69.53	-4.64
	SoT	78.33	78	88.02	39	63.16	75	80.30	46	88.11	38	65.81	77	77.29	59	<u>74.68</u>	-0.83

ings. These include Qwen-2.5 in 7B, 14B, and 32B variants (Team, 2024), LLaMA-3.1-8B (Meta, 2024a), LLaMA-3.2-11B (Meta, 2024b), GPT-4o (OpenAI, 2024), and Claude Sonnet 3.5 (Anthropic, 2024). For experiments involving multimodal inputs, we use Qwen-2.5-VL-7B (Team, 2025), which supports visual input processing. Unless otherwise specified, Qwen-2.5-32B serves as the default model for all other experiments. We use a temperature value of 0.5 for all models to balance output stability and diversity. For open-source models, inference is accelerated using FlashAttention2 (Dao, 2023). We sample 150 questions from each dataset for the sake of computational costs, and report the averaged performance over three independent runs per question. For the router model, we fine-tune DistilBERT with cross-entropy loss over 5 epochs, using a batch size of 64 and a learning rate of $2e^{-5}$. During inference, the router

processes the core input query. Following previous work, we use few-shot prompting to illustrate the required reasoning style, with exemplars being generated by prompting Qwen-2.5-32B with the method-specific prompt and selecting high-quality outputs. Further information regarding prompts and exemplars can be found in Appendix B.

3.4 Evaluation Protocol

We evaluate using two primary metrics: accuracy and output token count. For multiple-choice, yes/no, or numeric tasks, accuracy is computed via exact match with the ground truth. For open-ended generation, we follow the LLM-as-a-judge paradigm (Liu et al., 2023), using GPT-4o (OpenAI, 2024) to assess correctness. Answers are extracted according to the output format (see Appendix B.2). We analyze efficiency through the total number of generated tokens in the intermediate reasoning.

4 Results and Discussion

4.1 Overall Performance

As shown in Table 1, Sketch-of-Thought (SoT) consistently reduces output token count while minimizing the impact on reasoning accuracy across all evaluated models. On average, SoT achieves a token reduction of over 74% relative to CoT, with accuracy deviations typically within 1%. These trends hold across both open-weight models and proprietary models, confirming SoT’s generalizability across architectures and model families. SoT also demonstrates strong stability across reasoning tasks, consistently balancing token reduction with minimal accuracy variance, unlike other baselines which exhibit greater fluctuations. Notably, across all runs, we found that SoT consistently reduces token usage while having a statistically insignificant impact on accuracy ($p < 0.05$).

4.2 Model-wise Trends

Performance gains with SoT are especially notable in the Qwen family of models. On Qwen-2.5-32B, SoT achieves 82.30% accuracy—slightly above CoT’s 82.24%—while reducing output token count by 74.36%. Similar patterns hold at the 14B and 7B scales, where SoT maintains accuracy within 1% of CoT while reducing output length by over 70%. On GPT-4o, SoT achieves 84.55% accuracy—just 0.09% below CoT—while reducing token usage by 76%. Claude Sonnet 3.5 shows similar behavior, with SoT reaching 84.50% accuracy versus CoT’s 85.01%, alongside a 68% reduction in tokens. Results on LLaMA-3.1 and 3.2 indicate stronger compression (up to 78%) but slightly wider accuracy gaps (up to 3%). These findings confirm that SoT performs reliably across model families, consistently achieving strong token reductions with minimal accuracy degradation.

4.3 Paradigm-Task Alignment

Task-level results indicate that SoT’s effectiveness is most pronounced in reasoning settings with inherently compressible logic. In mathematical tasks, SoT closely matches the performance of CoT in the majority of settings. For example, in the Qwen-2.5-32B setting, SoT achieves 86.94% accuracy compared to 84.17% for CoT, while reducing average output length from 222 to 88 tokens. These gains are attributable to the effectiveness of the *Chunked Symbolism* paradigm in representing arithmetic reasoning concisely, which is the dominant paradigm

Table 2: **Results of Ensemble Methods.** Comparison of SoT and CoT in ensemble reasoning pipelines.

Approach	Method	Tkn	Acc	Red.	Δ
Self-Consistency	CoT	680	81.86	–	–
	SoT	176	81.90	74.1	0.04
Self-Refine	CoT	614	80.53	–	–
	SoT	244	80.80	60.3	0.27
Multi-Agent Debate	CoT	766	81.87	–	–
	SoT	238	82.44	68.9	0.57

for this category of reasoning (see Appendix C.3).

In commonsense and multi-hop reasoning, SoT maintains strong performance while achieving substantial compression. In the Qwen-2.5-32B setting, SoT reaches 92.00% accuracy on commonsense tasks using just 34 tokens on average, compared to 91.48% at 177 tokens under CoT. These improvements are driven by the *Conceptual Chaining* paradigm, which is the prevailing strategy for these reasoning categories and effectively captures structured relationships between ideas.

Domain-specialized tasks, such as PubMedQA and QASC, show more variability in accuracy across models, reflecting the inherent complexity of technical reasoning. Nevertheless, the *Expert Lexicons* paradigm remains effective at compressing domain-specific reasoning, often using half as many tokens as CoT while preserving competitive accuracy. Across all categories, SoT maintains competitive performance with far shorter outputs than CoT, underscoring its adaptive nature.

Further discussion on paradigm distribution across datasets can be found in Appendix C.

4.4 Token-Constrained Alternatives

Compared to other compression-focused prompting strategies such as Chain-of-Draft (CoD) and Constrained CoT (CCoT), SoT provides a more favorable trade-off between brevity and performance. Although CoD yields the most aggressive reductions in output length, it suffers notable accuracy degradation—for example, a 6.2% decline on GPT-4o despite a 75% token reduction. CCoT offers more balanced results, but still lags behind SoT in both efficiency and generalization across reasoning types. Although cases exist where these methods perform better in either accuracy or token reduction, there is no such case where these methods outperform SoT in both. In all observed settings, SoT achieves similar or better accuracy than these methods alongside competitive token reduction.

Table 3: **Multilingual Results.** Performance comparison of CoT and SoT across different languages.

Lang.	Method	Tkn	Acc	Red.	Δ
Korean	CoT	308	74.20	–	–
	SoT	49	73.40	84.09	-0.80
Italian	CoT	332	76.40	–	–
	SoT	57	75.07	82.83	-1.33
German	CoT	306	76.40	–	–
	SoT	48	76.07	84.31	-0.33

4.5 Ensemble Reasoning Methods

To examine SoT’s compatibility with ensemble-style reasoning methods, we integrate it into three established frameworks. Self-Consistency (Wang et al., 2023b) aggregates multiple reasoning paths by majority voting to improve answer stability. Self-Refine (Ranaldi and Freitas, 2024) enables iterative refinement of reasoning traces through reflection-based prompting. Multi-Agent Debate (Du et al., 2023) simulates deliberation among independent agents, each producing a rationale before converging on a final answer. In each case, we follow the original methodology but substitute SoT in place of CoT as the core reasoning strategy. Further implementation details, including prompts and hyperparameters, are provided in Appendix D.

Table 2 reports results from integrating SoT into three ensemble reasoning frameworks. In all cases, SoT improves performance relative to CoT, while substantially reducing output length. For instance, in the Self-Refine setting, SoT improves accuracy by 0.27% while generating 60% fewer tokens per response. In the Multi-Agent Debate framework, SoT yields a 0.57% accuracy increase alongside a 69% token reduction. These results indicate that SoT can be effectively substituted into more complex, multi-pass prompting pipelines, retaining its advantages in both efficiency and output quality.

4.6 Multilingual Reasoning

To evaluate SoT’s performance in non-English settings, we conduct a multilingual evaluation using Korean, Italian, and German subsets of MMMLU (Hendrycks et al., 2021). For each language, we select the same set of 500 questions from each language and generate three responses, for an effective sample size of 1,500. To maintain consistent paradigm selection across languages, each non-English query is paired with its English counterpart and routed using the same routing model. The se-

Table 4: **Multimodal Results.** Performance comparison of CoT and SoT for multimodal reasoning tasks.

Dataset	Method	Tkn	Acc	Red.	Δ
ScienceQA	CoT	136	79.60	–	–
	SoT	27	86.20	80.15	6.60
GQA	CoT	79	74.47	–	–
	SoT	19	71.93	75.95	-2.50

lected paradigm prompt and associated exemplars are then translated into the target language using GPT-4o (OpenAI, 2024), preserving both semantic fidelity and structural constraints.

As summarized in Table 3, SoT reduces output length by over 80% in all three languages while incurring a modest decrease in accuracy, ranging from -0.33% to -1.33%. These findings suggest that the sketching paradigms underlying SoT generalize across linguistic structures and preserve core reasoning logic beyond English.

4.7 Multimodal Reasoning

To assess SoT’s extensibility to multimodal scenarios, we evaluate its performance using Qwen-2.5-VL-7B (Team, 2025) on 500 multiple-choice samples from both GQA (Hudson and Manning, 2019) and the image-based subset of ScienceQA (Lu et al., 2022). Each sample is run three times for an effective sample size of 1,500. As in the unimodal setting, paradigm selection is handled by the router model. Images and supplementary materials are replaced with a placeholder token during routing (see Section 2.3), allowing the router to focus on the question text. We reuse the same text-only exemplars from the primary experiments.

Results from multimodal evaluations are shown in Table 4. On ScienceQA, SoT reduces output length by 80% while outperforming CoT by 6.60%. On GQA, however, we observed a 2.50% reduction in accuracy when using SoT while reducing output length by 75%. The accuracy degradation in GQA likely reflects the difficulty of applying abstract sketching methods to tasks requiring fine-grained visual grounding. Another possible explanation is that the text-only exemplars, while effective in general, may not sufficiently prime the model for vision-intensive reasoning.

4.8 Analysis on Routing

To investigate the efficacy of our router model for paradigm assignment, we evaluate its ability to select appropriate reasoning paradigms across the

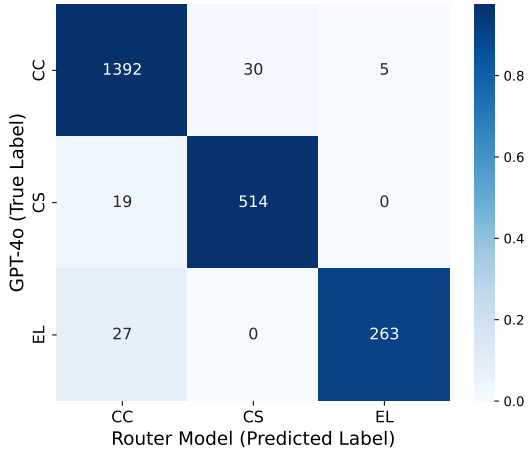


Figure 3: Confusion matrix illustrating the performance of the router model in selecting among the three SoT paradigms. Predictions are compared against GPT-4o-assigned ground truth labels.

2,250 samples used in our primary experiments (see Section 3.1). Ground-truth labels are produced by GPT-4o using the same labeling protocol as during training (see Section 2.3). As shown in Figure 3, the model achieves 96.4% overall accuracy, with high recall for the two most common paradigms, *Conceptual Chaining* (0.964) and *Chunked Symbolism* (0.975). Recall for *Expert Lexicons* is slightly lower at 0.907, largely due to class imbalance. However, this asymmetry is expected as *Expert Lexicons* is intentionally applied more conservatively given its specialized nature, and the router defaults to general paradigms in ambiguous cases to reduce risk of misapplication.

4.9 Paradigm-Task Alignment

To test if there is a significant difference between the performance of each paradigm in their respective tasks, we benchmark the performance of all three paradigms on datasets across different reasoning tasks. For any given dataset, we define the dominant paradigm as the paradigm which is assigned to the majority of samples. For example, from the paradigm definitions outlined in Section 2.2, we can assume that the expected-dominant paradigm of GSM8K is *Chunked Symbolism*. In Appendix C.3, we conduct an analysis of the expected versus actual dominant paradigms across all datasets to validate the router’s overall performance. In all cases the expected-dominant paradigm aligns with the actual-dominant paradigm.

However, this analysis of paradigm routing says nothing of the accuracy on these tasks. For this, we select one representative dataset per paradigm us-

Table 5: **Paradigm-Task Alignment.** Comparison of paradigm performance on various reasoning types. Expected-dominant paradigm is in **bold**.

Reasoning Type	Paradigm	Tkn	Acc
Mathematical <i>SVAMP</i>	Chunked Symbolism	30	93.70
	Expert Lexicons	33	92.60
	Conceptual Chaining	34	92.00
Medical <i>MedQA</i>	Chunked Symbolism	46	73.10
	Expert Lexicons	52	85.70
	Conceptual Chaining	73	81.10
Commonsense <i>CommonsenseQA</i>	Chunked Symbolism	8	40.00
	Expert Lexicons	18	82.90
	Conceptual Chaining	21	84.60

ing the previously-mentioned dominant-paradigm distribution in Table 10. We then run inference on these datasets using all paradigms and compare their performance in Table 5. Our findings indicate that, for all examined datasets, the expected-dominant paradigm outperforms all others in terms of task accuracy. Notably, a paradigm being dominant does not mean it will have the lowest token usage. For example, in medical reasoning, *Expert Lexicons* has the highest accuracy with 85.70% and, while *Chunked Symbolism* has the lowest token usage, its accuracy is far lower at 73.10%. These results demonstrate that different reasoning paradigms yield different performance levels depending on the task, and that selecting the optimal paradigm is critical for maximizing accuracy.

5 Related Work

Token-Efficient Reasoning A growing body of work targets the reduction of output length during language model reasoning. Concise Chain-of-Thought (Renze and Guven, 2024) and Constrained CoT (Nayab et al., 2025) apply fixed constraints on the number of steps or words in the reasoning trace. SCOTT (Wang et al., 2023a) uses a two-stage summarization pipeline that compresses verbose CoT outputs into shorter versions. While these methods reduce token usage, they rely on surface heuristics or summary-based rewriting, often reducing clarity. As an orthogonal direction, Coconut (Hao et al., 2024) bypasses token-based reasoning by operating entirely in the latent vector space, though this requires additional training procedures, making it inapplicable to frozen LLMs. In contrast, SoT rewrites reasoning steps using compact representations, yielding outputs that are both concise and interpretable.

Structured Reasoning Strategies Other approaches enhance reasoning by restructuring the generation process itself. Tree-of-Thoughts (Yao et al., 2023) and Graph-of-Thoughts (Besta et al., 2024) treat reasoning as a search over intermediate steps, producing graph-structured outputs. Self-Consistency (Wang et al., 2023b) improves stability by sampling multiple reasoning paths and selecting the majority answer. While these methods improve accuracy on certain tasks, they often incur significant increases in compute overhead. In contrast, SoT leverages a standard prompting interface to restructure internal reasoning, achieving efficiency gains without increasing inference complexity.

Prompt Compression and Adaptive Inference

Several techniques improve efficiency through prompt compression or selective computation. Chain-of-Draft (CoD) (Xu et al., 2025) uses densely packed natural language reasoning to reduce length, but this often comes at the cost of clarity and yields large performance drops on more complex reasoning tasks. CoT-Influx (Huang et al., 2024) and LLMingua (Jiang et al., 2023) prune or compress input exemplars to reduce prompt length. Cascaded inference (Yue et al., 2024) and compute-adaptive methods (Arora and Zanette, 2025) dynamically route examples to high-cost inference pipelines only when necessary. SoT differs by addressing compression as a representational design challenge: instead of relying on pruning or selection, it restructures how reasoning is expressed, guided by task-specific cognitive principles.

6 Conclusion

We present Sketch-of-Thought (SoT), a prompting framework that reduces token usage in language model reasoning by up to 84%, preserving accuracy in most tasks and incurring only minor trade-offs in others. SoT leverages cognitively inspired paradigms to generate compact yet semantically faithful reasoning traces, offering a practical alternative to verbose prompting. Extensive experiments across 18 reasoning datasets, multiple languages, and multimodal tasks demonstrate SoT’s broad applicability. Its compatibility with ensemble prompting strategies further reinforces its practical utility, particularly in resource-constrained settings. By reframing efficiency as a reasoning design challenge rather than a surface-level compression problem, SoT opens new directions for scalable, cognitively informed prompting.

Limitations and Future Work

Sketch-of-Thought (SoT) is designed for interpretable, efficient reasoning, and while our current approach performs well on a variety of tasks, there exist several interesting directions for future work.

Following prior work, our use of fixed exemplars per paradigm—intentionally chosen to preserve stylistic consistency and interpretability—may limit adaptability to subtle variations within a task type. Alternatively, a retrieval system could dynamically pull in-context exemplars from a larger pool based on the reasoning paradigm and question characteristics. These strategies could help to improve SoT’s flexibility across subtly different queries but also disparate tasks and domains.

Also, while we focused this work on evaluating cognitively grounded, prompt-based three paradigms, the framework is not limited to the three we present here. Future work may incorporate additional reasoning paradigms to better adapt SoT to downstream tasks such as code generation. These can be integrated by adding new paradigms, updating the sketching pool, and retraining the routing module accordingly.

Lastly, while our current multilingual experiments already demonstrate SoT’s stability across widely-spoken languages, evaluating its impact in low-resource languages is an exciting direction for future work.

Ethics Statement

This work builds on widely used public datasets and large language models (LLMs). All datasets used in our experiments are publicly available and cited accordingly. Where applicable, we follow dataset authors’ intended uses and licensing terms. All models are used in accordance with their respective licenses.

While Sketch-of-Thought (SoT) improves the efficiency of model reasoning, we acknowledge that compressing intermediate outputs may affect interpretability in certain high-stakes settings. We encourage caution when applying SoT in domains such as healthcare or legal analysis, where full transparency of reasoning steps may be essential.

Further, our router model was trained using annotations generated via GPT-4o, and as such may reflect biases present in the underlying model. We recommend further evaluation before deploying SoT in sensitive or high-stakes settings.

Acknowledgments

This work was supported by the Institute for Information & Communications Technology Planning & Evaluation (IITP) grant funded by the Korea government (MSIT) (RS-2019-II190075, Artificial Intelligence Graduate School Program (KAIST), and RS-2022-II220713, Meta-learning Applicable to Real-world Problems); the National Research Foundation of Korea (NRF) grant funded by the Korea government (MSIT) (RS-2023-00256259); the grant of the Korea Machine Learning Ledger Orchestration for Drug Discovery Project (K-MELLODDY) funded by the Ministry of Health & Welfare and the Ministry of Science and ICT, Republic of Korea (RS-2024-00460870); Institute of Information & Communications Technology Planning & Evaluation (IITP) with the grant funded by the Ministry of Science and ICT (MSIT) of the Republic of Korea in connection with the Global AI Frontier Lab International Collaborative Research (RS-2024-00469482 & RS-2024-00509279); and the Artificial Intelligence Industrial Convergence Cluster Development Project funded by the Ministry of Science and ICT (MSIT, Korea) & Gwangju Metropolitan City.

References

- John R. Anderson. 1983. *A spreading activation theory of memory*. *Journal of Verbal Learning and Verbal Behavior*, 22(3):261–295.
- Anthropic. 2024. *Claude 3.5 sonnet*.
- Daman Arora and Andrea Zanette. 2025. *Training language models to reason efficiently*. *Preprint*, arXiv:2502.04463.
- A. Baddeley. 2000. *The episodic buffer: a new component of working memory?* *Trends in Cognitive Sciences*, 4(11):417–423.
- Maciej Besta, Nils Blach, Ales Kubicek, Robert Gerstenberger, Michal Podstawski, Lukas Gianinazzi, Joanna Gajda, Tomasz Lehmann, 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.
- Sébastien Bubeck, Varun Chandrasekaran, Ronen Eldan, Johannes Gehrike, Eric Horvitz, Ece Kamar, Peter Lee, Yin Tat Lee, Yuanzhi Li, Scott Lundberg, Harsha Nori, Hamid Palangi, Marco Tulio Ribeiro, and Yi Zhang. 2023. *Sparks of artificial general intelligence: Early experiments with gpt-4*. *Preprint*, arXiv:2303.12712.
- Michelene T. H. Chi, Paul J. Feltovich, and Robert Glaser. 1981. *Categorization and Representation of Physics Problems by Experts and Novices*. *Cognitive Science*, 5(2):121–152.
- 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 [cs].
- Tri Dao. 2023. *Flashattention-2: Faster attention with better parallelism and work partitioning*. *Preprint*, arXiv:2307.08691.
- Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2018. *BERT: pre-training of deep bidirectional transformers for language understanding*. *CoRR*, abs/1810.04805.
- Yilun Du, Shuang Li, Antonio Torralba, Joshua B. Tenenbaum, and Igor Mordatch. 2023. *Improving factuality and reasoning in language models through multiagent debate*. *Preprint*, arXiv:2305.14325.
- Dheeru Dua, Yizhong Wang, Pradeep Dasigi, Gabriel Stanovsky, Sameer Singh, and Matt Gardner. 2019. *DROP: A Reading Comprehension Benchmark Requiring Discrete Reasoning Over Paragraphs*. *arXiv preprint*. ArXiv:1903.00161 [cs].
- Mor Geva, Daniel Khashabi, Elad Segal, Tushar Khot, Dan Roth, and Jonathan Berant. 2021. *Did Aristotle Use a Laptop? A Question Answering Benchmark with Implicit Reasoning Strategies*. *arXiv preprint*. ArXiv:2101.02235 [cs].
- Vinod Goel. 1995. *Sketches of Thought*. MIT Press.
- Shibo Hao, Sainbayar Sukhbaatar, DiJia Su, Xian Li, Zhiting Hu, Jason Weston, and Yuandong Tian. 2024. *Training large language models to reason in a continuous latent space*. *Preprint*, arXiv:2412.06769.
- Dan Hendrycks, Collin Burns, Steven Basart, Andy Zou, Mantas Mazeika, Dawn Song, and Jacob Steinhardt. 2021. *Measuring Massive Multitask Language Understanding*. *arXiv preprint*. ArXiv:2009.03300 [cs].
- Xijie Huang, Li Lina Zhang, Kwang-Ting Cheng, Fan Yang, and Mao Yang. 2024. *Fewer is more: Boosting llm reasoning with reinforced context pruning*. *Preprint*, arXiv:2312.08901.
- Drew A Hudson and Christopher D Manning. 2019. *Gqa: A new dataset for real-world visual reasoning and compositional question answering*. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 6700–6709.
- Peter A. Jansen, Elizabeth Wainwright, Steven Marmorstein, and Clayton T. Morrison. 2018. *WorldTree: A Corpus of Explanation Graphs for Elementary Science Questions supporting Multi-Hop Inference*. *arXiv preprint*. ArXiv:1802.03052 [cs].

- Huiqiang Jiang, Qianhui Wu, Chin-Yew Lin, Yuqing Yang, and Lili Qiu. 2023. [Lmlingua: Compressing prompts for accelerated inference of large language models](#). *Preprint*, arXiv:2310.05736.
- Di Jin, Eileen Pan, Nassim Oufattole, Wei-Hung Weng, Hanyi Fang, and Peter Szolovits. 2020. [What Disease does this Patient Have? A Large-scale Open Domain Question Answering Dataset from Medical Exams](#). *arXiv preprint*. ArXiv:2009.13081 [cs].
- Qiao Jin, Bhuwan Dhingra, Zhengping Liu, William Cohen, and Xinghua Lu. 2019. [PubMedQA: A Dataset for Biomedical Research Question Answering](#). In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, pages 2567–2577, Hong Kong, China. Association for Computational Linguistics.
- Jaap Jumelet, Willem Zuidema, and Arabella Sinclair. 2024. [Do language models exhibit human-like structural priming effects?](#) *Preprint*, arXiv:2406.04847.
- Tushar Khot, Peter Clark, Michal Guerin, Peter Jansen, and Ashish Sabharwal. 2020. [QASC: A Dataset for Question Answering via Sentence Composition](#). *arXiv preprint*. ArXiv:1910.11473 [cs].
- Wang Ling, Dani Yogatama, Chris Dyer, and Phil Blunsom. 2017. [Program Induction by Rationale Generation: Learning to Solve and Explain Algebraic Word Problems](#). *arXiv preprint*. ArXiv:1705.04146 [cs].
- Jian Liu, Leyang Cui, Hanmeng Liu, Dandan Huang, Yile Wang, and Yue Zhang. 2020. [LogiQA: A Challenge Dataset for Machine Reading Comprehension with Logical Reasoning](#). *arXiv preprint*. ArXiv:2007.08124 [cs].
- Yang Liu, Dan Iter, Yichong Xu, Shuohang Wang, Ruochen Xu, and Chenguang Zhu. 2023. [G-eval: NLG evaluation using gpt-4 with better human alignment](#). In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pages 2511–2522, Singapore. 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)*.
- Meta. 2024a. [Introducing llama 3.1: Our most capable models to date](#).
- Meta. 2024b. [Llama 3.2: Revolutionizing edge ai and vision with open, customizable models](#).
- 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](#). *arXiv preprint*. ArXiv:1809.02789 [cs].
- George A. Miller. 1956. [The magical number seven, plus or minus two: Some limits on our capacity for processing information](#). *Psychological Review*, 63(2):81–97. Place: US Publisher: American Psychological Association.
- Sania Nayab, Giulio Rossolini, Marco Simoni, Andrea Saracino, Giorgio Buttazzo, Nicolamaria Manes, and Fabrizio Giacomelli. 2025. [Concise thoughts: Impact of output length on llm reasoning and cost](#). *Preprint*, arXiv:2407.19825.
- OpenAI. 2024. [Gpt-4o system card](#). *Preprint*, arXiv:2410.21276.
- Arkil Patel, Satwik Bhattamishra, and Navin Goyal. 2021. [Are NLP Models really able to Solve Simple Math Word Problems?](#) In *Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 2080–2094, Online. Association for Computational Linguistics.
- Alec Radford, Jeff Wu, Rewon Child, David Luan, Dario Amodei, and Ilya Sutskever. 2019. [Language models are unsupervised multitask learners](#).
- Leonardo Ranaldi and Andre Freitas. 2024. [Self-refine instruction-tuning for aligning reasoning in language models](#). In *Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing*, page 2325–2347. Association for Computational Linguistics.
- Matthew Renze and Erhan Guven. 2024. [The benefits of a concise chain of thought on problem-solving in large language models](#). *Preprint*, arXiv:2401.05618.
- Victor Sanh, Lysandre Debut, Julien Chaumond, and Thomas Wolf. 2020. [Distilbert, a distilled version of bert: smaller, faster, cheaper and lighter](#). *Preprint*, arXiv:1910.01108.
- Zayne Sprague, Fangcong Yin, Juan Diego Rodriguez, Dongwei Jiang, Manya Wadhwa, Prasann Singhal, Xinyu Zhao, Xi Ye, Kyle Mahowald, and Greg Durrett. 2024. [To cot or not to cot? chain-of-thought helps mainly on math and symbolic reasoning](#). *Preprint*, arXiv:2409.12183.
- Jiankai Sun, Chuanyang Zheng, Enze Xie, Zhengying Liu, Ruihang Chu, Jianing Qiu, Jiaqi Xu, Mingyu Ding, Hongyang Li, Mengzhe Geng, Yue Wu, Wenhai Wang, Junsong Chen, Zhangyue Yin, Xiaozhe Ren, Jie Fu, Junxian He, Wu Yuan, Qi Liu, Xihui Liu, Yu Li, Hao Dong, Yu Cheng, Ming Zhang, Pheng Ann Heng, Jifeng Dai, Ping Luo, Jingdong Wang, Ji-Rong Wen, Xipeng Qiu, Yike Guo, Hui Xiong, Qun Liu, and Zhenguo Li. 2024. [A Survey of Reasoning with Foundation Models](#). *arXiv preprint*. ArXiv:2312.11562 [cs].
- Alon Talmor, Jonathan Herzig, Nicholas Lourie, and Jonathan Berant. 2019. [CommonsenseQA: A Question Answering Challenge Targeting Commonsense Knowledge](#). In *Proceedings of the 2019 Conference*

- of the North American Chapter of the Association for Computational Linguistics: *Human Language Technologies, Volume 1 (Long and Short Papers)*, pages 4149–4158, Minneapolis, Minnesota. Association for Computational Linguistics.
- Qwen Team. 2024. [Qwen2.5: A party of foundation models](#).
- Qwen Team. 2025. [Qwen2.5-vl](#).
- Harsh Trivedi, Niranjan Balasubramanian, Tushar Khot, and Ashish Sabharwal. 2022. [MuSiQue: Multi-hop Questions via Single-hop Question Composition](#). *Transactions of the Association for Computational Linguistics*, 10:539–554. Place: Cambridge, MA Publisher: MIT Press.
- Peifeng Wang, Zhengyang Wang, Zheng Li, Yifan Gao, Bing Yin, and Xiang Ren. 2023a. [Scott: Self-consistent chain-of-thought distillation](#). *Preprint*, arXiv:2305.01879.
- Xuezhi Wang, Jason Wei, Dale Schuurmans, Quoc Le, Ed Chi, Sharan Narang, Aakanksha Chowdhery, and Denny Zhou. 2023b. [Self-consistency improves chain of thought reasoning in language models](#). *Preprint*, arXiv:2203.11171.
- Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Brian Ichter, Fei Xia, Ed Chi, Quoc Le, and Denny Zhou. 2023. [Chain-of-thought prompting elicits reasoning in large language models](#). *Preprint*, arXiv:2201.11903.
- Thomas Wolf, Lysandre Debut, Victor Sanh, Julien Chaumond, Clement Delangue, Anthony Moi, Pierric Cistac, Tim Rault, Rémi Louf, Morgan Funtowicz, Joe Davison, Sam Shleifer, Patrick von Platen, Clara Ma, Yacine Jernite, Julien Plu, Canwen Xu, Teven Le Scao, Sylvain Gugger, Mariama Drame, Quentin Lhoest, and Alexander M. Rush. 2020. [Transformers: State-of-the-art natural language processing](#). In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing: System Demonstrations*, pages 38–45, Online. Association for Computational Linguistics.
- Silei Xu, Wenhao Xie, Lingxiao Zhao, and Pengcheng He. 2025. [Chain of draft: Thinking faster by writing less](#). *Preprint*, arXiv:2502.18600.
- Zhilin Yang, Peng Qi, Saizheng Zhang, Yoshua Bengio, William W. Cohen, Ruslan Salakhutdinov, and Christopher D. Manning. 2018. [HotpotQA: A Dataset for Diverse, Explainable Multi-hop Question Answering](#). *arXiv preprint*. ArXiv:1809.09600 [cs].
- 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.
- Weihao Yu, Zihang Jiang, Yanfei Dong, and Jiashi Feng. 2020. [ReClor: A Reading Comprehension Dataset Requiring Logical Reasoning](#). *arXiv preprint*. ArXiv:2002.04326 [cs].
- Murong Yue, Jie Zhao, Min Zhang, Liang Du, and Ziyu Yao. 2024. [Large language model cascades with mixture of thoughts representations for cost-efficient reasoning](#). *Preprint*, arXiv:2310.03094.
- Wayne Xin Zhao, Kun Zhou, Junyi Li, Tianyi Tang, Xiaolei Wang, Yupeng Hou, Yingqian Min, Beichen Zhang, Junjie Zhang, Zican Dong, Yifan Du, Chen Yang, Yushuo Chen, Zhipeng Chen, Jinhao Jiang, Ruiyang Ren, Yifan Li, Xinyu Tang, Zikang Liu, Peiyu Liu, Jian-Yun Nie, and Ji-Rong Wen. 2024. [A survey of large language models](#). *Preprint*, arXiv:2303.18223.

A Experimental Setup

A.1 Datasets

All datasets used in our experiments are publicly available and accessed via Hugging Face using the dataset IDs listed in Table 6. Where datasets included multiple subsets, we explicitly specified which subset was used in our experiments. All datasets are used in accordance with their respective licenses and terms of use.

A.2 Model Checkpoints

We use the following model checkpoints in our experiments:

Qwen 2.5 Family

- Qwen/Qwen2.5-7B-Instruct
- Qwen/Qwen2.5-14B-Instruct
- Qwen/Qwen2.5-32B-Instruct
- Qwen/Qwen2.5-VL-7B-Instruct

Llama 3 Family

- meta-llama/Meta-Llama-3-8B-Instruct
- meta-llama/Meta-Llama-3-11B-Instruct

Closed-source Models

- gpt-4o-2024-11-20
- claude-3-5-sonnet-20241022

All open-weight models were accessed through Hugging Face via the transformers library (Wolf et al., 2020) and evaluated in their instruction-tuned form. Closed-weight models such as GPT-4o and Claude Sonnet 3.5 were accessed through their respective Python API wrappers. All models are used in accordance with their licenses.

A.3 Inference Environment

All experiments were conducted on 2 x A5000 24GB GPUs on a Linux distribution running CUDA 12.1. For inference, we use FlashAttention2 (Dao, 2023) for acceleration. All models were run in bfloat16 precision where supported. No parameter fine-tuning or additional adaptation was applied to the LLMs during experimentation.

A.4 Reproducibility

All experiments were conducted with the same fixed random seed, 42, to ensure reproducibility across runs. We used a consistent temperature of 0.5 for all models across all methods and tasks. For few-shot setups, exemplars were selected prior to evaluation and held constant across all trials. Token counts were measured using the default tokenizer associated with each model’s checkpoint. For closed source models accessed via the API, token counts were obtained through the token logs found in the returned inference object.

B Prompting Framework

B.1 System Prompt Format

Each paradigm-specific system prompt follows a consistent structure composed of four sections. Among these sections, *Role & Objective* is the most extensive. This design choice is motivated by two factors. First, since our method is training-free and prompt-based, we found that smaller LLMs require more explicit and structured instructions to reliably follow the desired output format. In early trials, such models often ignored or deviated from intended behavior when given minimal guidance. Second, recent work on structural priming in LLMs shows that lexical and semantic cues in preceding context influence syntactic expectations and downstream predictions (Jumelet et al., 2024). We leverage this by “setting the stage” for the model to reason within the appropriate paradigm.

Role & Objective Provides background on the paradigm, including its cognitive basis and theoretical motivation. It outlines representative use cases and serves as a semantic primer to help align the model’s reasoning style with the paradigm.

Application Steps Describes a step-by-step procedure for applying the paradigm to solve a problem. This includes identifying relevant concepts, performing transformations, and following best practices for structuring the reasoning process.

Rules & Directives Specifies required tone, structure, and formatting constraints. It highlights common failure modes—such as verbosity, redundancy, or incorrect notation—and explicitly defines output style requirements (see Appendix B.2).

Closing Statement Ends with a reminder to adhere to the formatting guidelines, reinforcing the objective of concise, structured reasoning.

GitHub: github.com/SimonAytes/SoT
Router: huggingface.co/saytes/SoT_DistilBERT

Table 6: **Dataset Information.** Comprehensive details of datasets used for our experiments.

Dataset	Citation	HF ID	Train Split:Subset	Train Size	Test Split:Subset	Test Size
GSM8K	Cobbe et al. (2021)	gsm8k	main:train	1000	main:test	150
SVAMP	Patel et al. (2021)	ChilleD/SVAMP	train	700	test	150
AQUA-RAT	Ling et al. (2017)	aqua_rat	train	1000	test	150
DROP	Dua et al. (2019)	drop	train	1000	validation	150
OpenbookQA	Mihaylov et al. (2018)	openbookqa	train	1000	test	150
StrategyQA	Geva et al. (2021)	ChilleD/StrategyQA	train	1000	test	150
LogiQA	Liu et al. (2020)	lucasmccabe/logiqa	train	1000	test	150
Reclor	Yu et al. (2020)	metaeval/reclor	train	1000	validation	150
HotPotQA	Yang et al. (2018)	hotpot_qa	distractor:train	1000	distractor:validation	150
MuSiQue-Ans	Trivedi et al. (2022)	dgslibisey/MuSiQue	train	1000	validation	150
QASC	Khot et al. (2020)	allenai/qasc	train	1000	validation	150
Worldtree	Jansen et al. (2018)	nguyen-brat/worldtree	train (last 1000)	1000	train (rest)	150
PubMedQA	Jin et al. (2019)	qiaojin/PubMedQA	pqa_labeled (last 500)	500	pqa_labeled (first 150)	150
MedQA	Jin et al. (2020)	bigbio/med_qa	med_qa_en_source:train	1000	med_qa_en_source:validation	150
CommonsenseQA	Talmor et al. (2019)	tau/commonsense_qa	train	1000	validation	150
MMLU	Hendrycks et al. (2021)	cais/mmlu	—	—	test:all	500
MIMLU	Hendrycks et al. (2021)	openai/MIMLU	—	—	test:K0_KR, DE_DE, IT_IT	500
ScienceQA	Lu et al. (2022)	lms-lab/ScienceQA	—	—	test:ScienceQA-IMG	500
GQA	Hudson and Manning (2019)	lms-lab/GQA	—	—	val:val_all_images	500

B.2 Output Conventions

To ensure consistent evaluation and accurate token-level comparisons, all outputs follow a strict formatting protocol:

- **Answers** must be enclosed in `\boxed{. . .}`.
- **Reasoning traces** must appear within `<think>` and `</think>` tags.

This formatting allows for reliable programmatic parsing and segmentation of outputs into intermediate reasoning and final answers, supporting reproducibility and enabling fair evaluation across prompting methods and models.

Because all experiments are conducted using instruction-tuned LLMs with no additional fine-tuning to enforce output structure, we explicitly reserve space in both the *Rules & Directives* and *Closing Statement* sections of each system prompt to reinforce these formatting requirements.

In practice, we find that providing exemplars alone is insufficient for enforcing consistent formatting. In early experiments, models frequently omitted structural tags or deviated from the expected format when prompted using exemplars only. After incorporating explicit formatting instructions into the system prompt, the rate of malformed or non-compliant outputs dropped to near zero across all paradigms and model variants.

B.3 Paradigm Prompts

We provide reference versions of our paradigm system prompts for *Conceptual Chaining*, *Chunked Symbolism*, and *Expert Lexicons* in Figures 5, 6, and 7, respectively. Parts of the prompts have been adjusted to render correctly in this document. We direct the interested reader to our public code repositories for full, code-ready prompts.

B.4 In-Context Exemplars

We provide three in-context exemplars for each method evaluated in our study to guide model outputs during inference. For Sketch-of-Thought (SoT), a separate set of exemplars is constructed for each paradigm to match the distinct reasoning styles each paradigm is designed to elicit. Example questions are manually selected to reflect typical tasks associated with each paradigm’s target use cases. To construct exemplars, we first generate candidate responses using Qwen-2.5-32B with the corresponding system prompt, then manually select outputs that most faithfully demonstrate the paradigm’s intended structure, clarity, and conciseness. This results in three curated exemplars per paradigm.

For baseline methods—Chain-of-Thought (CoT), Constrained CoT (CCoT), and Chain-of-Draft (CoD)—we apply the same method. Each method is prompted using its respective strategy, and the most stylistically representative outputs are selected. Because these baselines do not dynamically adapt to the query type, we ensure fair coverage by drawing exemplars from the same three reasoning categories used for SoT (e.g., commonsense, mathematical, medical). One exemplar is selected per category, yielding a total of three per method. All exemplars are held fixed across all experiments.

B.5 Combined System Prompt

The question naturally arises as to whether the Sketch-of-Thought (SoT) framework would have benefited from a more streamlined approach in which the LLM is prompted with a single joint prompt describing all three paradigms, allowing the model to select one or apply them jointly. While promising in principle, our exploration of this ap-

Table 7: **Router Ablation Results.** Comparison of benchmark performance among candidate router model architectures. Accuracy and latency values are averaged across all test cases. VRAM denotes the peak memory usage of the model during the experiment in megabytes.

Model	Params	Acc	Latency (s)	VRAM
DistilBERT	67M	90.31	0.0118	283 MB
GPT-2	137M	91.11	0.0107	652 MB
BERT-base	110M	90.93	0.0139	445 MB
BERT-large	336M	88.93	0.0259	1309 MB

proach in early development resulted in high rates of formatting errors, especially for models with fewer than 70B parameters. These smaller models frequently failed to follow paradigm-specific instructions or ignored the structured reasoning techniques altogether, reverting to verbose or default reasoning styles. We believe this reflects a limitation of training-free prompting for smaller models. Furthermore, as discussed in our future work section, SoT is designed to scale to many more than three paradigms. In that setting, a joint prompt would incur significant token overhead from combining system prompts and exemplars, undercutting the efficiency gains SoT is designed to provide.

B.6 Classification Prompt

The router model used to assign paradigms was trained using GPT-4o-generated labels. The classification prompt presented each query and instructed the model to assign one of the three paradigms based on reasoning characteristics, following the heuristic definitions given in Section 2.2. A reference version of the classification prompt is shown in Figure 8. To conserve space we omit repetitive text in this version. We direct the interested reader to our public code repositories for the full, unabridged classification prompt.

C Paradigm Assignment Analysis

C.1 Router Architecture Ablation

To support our choice of architecture, we conducted an ablation comparing four base models for routing: DistilBERT (Sanh et al., 2020), GPT-2 (Radford et al., 2019), BERT-base, and BERT-large (Devlin et al., 2018). All models were trained on the same task with identical supervision and hyperparameters. We measured test accuracy, latency, and peak memory usage during inference for the classification of 1,250 samples from our test set. The results of this ablation are shown in Table 7.

While GPT-2 and BERT-base achieved slightly higher accuracy than DistilBERT, they also incurred significantly higher memory demand and higher latency. BERT-large had the highest parameter count but the lowest accuracy. Overall, DistilBERT achieved the best tradeoff across accuracy, latency, and memory overhead, confirming that it is the ideal architecture for the purpose of our experiments.

C.2 Analysis of Machine Labeling

As outlined in Section 2.3 and 4.8, our method relies on using GPT-4o (OpenAI, 2024) for machine-labeling tasks—both for the training and evaluation of our router. To evaluate the quality of GPT-4o’s paradigm assignment, we designed a human study to measure agreement between human annotators and GPT-4o. We randomly selected 30 samples—uniformly distributed across the three paradigms—and recruited nine annotators with no prior hands-on knowledge of the Sketch-of-Thought framework. Annotators were given the exact same system prompt used by GPT-4o and were allowed to reference it freely while annotating (see Appendix B.6).

We evaluated the agreement between the human annotators and machine labeling with Cohen’s Kappa (k). The results from this experiment show substantial agreement between GPT-4o and the human-majority label ($k = 0.95$), with strong inter-human agreement ($k = 0.72$) supporting the reliability of our human study and highlighting the quality of the machine-labeled training data.

C.3 Routing Distribution by Dataset

To better understand how SoT paradigms align with reasoning task types, we analyze the output of our router model across the datasets used in our primary experiments (Section 3.1). Table 10 reports the predicted paradigm distribution, dominant paradigm, and its agreement with an expected paradigm label defined based on the paradigm descriptions in Section 2.2. These counts reflect router predictions over the 150 samples used per dataset in our primary experiments. The dominant paradigm is defined as the one with the highest frequency within a dataset, and we compare this to the expected paradigm, which is assigned based on prior task categorizations and paradigm design goals.

The router’s predictions match expectations in all 15 datasets, with 100% agreement between dominant and expected paradigms. Most datasets are

routed to a single paradigm, reflecting high confidence and class purity. In a few edge cases (e.g., DROP, LogiQA, QASC), we observe minor cross-paradigm overlap, though these do not shift the dominant label. This behavior aligns with our router’s conservative design, which favors general-purpose paradigms (especially *Conceptual Chaining*) in ambiguous scenarios.

As expected, *Conceptual Chaining* dominates in commonsense, logical, and multi-hop datasets (e.g., StrategyQA, HotPotQA, Reclor), where relational inference is critical.

Chunked Symbolism is used exclusively in mathematical tasks (e.g., GSM8K, AQUA, SVAMP), where symbolic notation offers the clearest compression benefit. In DROP, which mixes symbolic and textual reasoning, some samples are routed to *Conceptual Chaining*, reflecting the complementary reasoning patterns among the paradigms.

Expert Lexicons is most common in domain-specific datasets like PubMedQA and MedQA. Occasional routing to *Conceptual Chaining* in these cases reflects the router’s conservative fallback behavior, favoring general-purpose paradigms when confidence is low—a design choice that reduces the risk of applying technical conventions in inappropriate contexts.

C.4 Paradigm Assignment Examples

To illustrate how the SoT router assigns paradigms to diverse questions, we present three representative examples below—one from each SoT paradigm. Each example includes the query as processed by the router and the assigned paradigm. Appendix C.5 showcases an edge case that demonstrates the system’s conservative fallback behavior.

Chunked Symbolism (GSM8K)

Query:
Darrell and Allen’s ages are in the ratio of 7:11. If their total age now is 162, calculate Allen’s age 10 years from now.

Assigned Paradigm:
Chunked Symbolism

Conceptual Chaining (OpenbookQA)

Query:
Polar bears require
Choices:
A. a tropical environment
B. a frigid environment
C. a tepid environment
D. a warm environment

Assigned Paradigm:
Conceptual Chaining

Expert Lexicons (PubMedQA)

Query:
Is the holmium:YAG laser the best intracorporeal lithotripter for the ureter?
Choices: Yes, No, Maybe

Assigned Paradigm:
Expert Lexicons

C.5 Router Misalignment

The results from Section 4.9 show that using the correct paradigm consistently yields the highest accuracy. While the drop in performance varies by task, the trend is consistent across datasets. Notably, output token usage remained relatively stable across paradigms due to fixed prompts and exemplar lengths, suggesting that routing errors primarily affect accuracy rather than efficiency.

To examine a case in which the predicted paradigm of a query differs from its parent dataset’s dominant paradigm, we consider the following.

Edge Case: Conceptual Chaining (PubMedQA)

Query:
[Context Here] Question: Birth characteristics and risk of low intellectual performance in early adulthood: are the associations confounded by socioeconomic factors in adolescence or familial effects?
Choices: Yes, No, Maybe

Assigned Paradigm:
Conceptual Chaining

For this query, the router assigns the query to *Conceptual Chaining* rather than *Expert Lexicons*, despite the latter being the dominant paradigm for that dataset (see Appendix C.3). However, this is not necessarily an incorrect prediction. The question centers on causal inference, which aligns with *Conceptual Chaining*. Although the context is medical, the query does not rely on the domain-specific jargon typical of *Expert Lexicons*. Because PubMedQA blends domain expertise with general causal reasoning, it contains a nontrivial number of such mixed-paradigm cases (approximately 66% *Expert Lexicons*, 33% *Conceptual Chaining*).

D Extended Results

D.1 Per-Dataset Results

We report per-dataset results from the primary experiments across all model families and prompting strategies in Tables 8 and 9. Results are shown for each individual dataset with accuracy and token counts averaged across three runs per dataset. For further information regarding the primary experiments, see Section 3 for the experimental design and Section 4 for the results and discussion.

D.2 Multi-Agent Debate

To evaluate whether Sketch-of-Thought (SoT) remains effective in ensemble-style deliberation, we incorporate it into the Multi-Agent Debate (MAD) framework (Du et al., 2023). This method simulates independent agents answering the same question and iteratively revising their answers through multi-round critique. For the Multi-Agent Debate setup, we preserve the paradigm-specific SoT system prompt and introduce a debate prompt that allows agents to revise their reasoning in response to other agents’ answers. The debate prompt is structured to request updated responses while retaining the specified output formatting conventions.

```
Multi-Agent Debate Prompt

You are participating in a multi-agent debate. Other agents have responded as follows:
#Agent 1:
<think> [agent 1’s reasoning] </think>
Answer: [agent 1’s answer]
Your previous answer was:
<think> [your previous reasoning] </think>
Answer: [your previous answer]
Would you like to revise your reasoning or stick with it? Please provide your updated reasoning inside <think>...</think> tags and your final answer inside \boxed{}
```

Each debate run involves three agents and a maximum of three rounds. In the first round, all agents independently generate answers using SoT prompts selected by the router model (see Section 2.3). In subsequent rounds, each agent receives the other agents’ reasoning and has the opportunity to revise its answer using the shared debate prompt above. Debates terminate early if all agents converge on the same answer. If consensus is not reached within three rounds, a majority vote determines the final answer. The rationale of the majority-aligned agent is retained as the final justification.

Notably, we find that using SoT does not have a notable impact on the number of rounds-per-query. For CoT we observed an average of 1.14 rounds-per-query, similar to SoT’s 1.11. Results are shown in Table 2 and discussed in Section 4.

D.3 Self-Consistency

The Self-Consistency framework (Wang et al., 2023b) is an ensemble reasoning method where multiple reasoning paths are generated from the same input query and the final answer is chosen via majority vote. This method improves answer stability and can, in some cases, mitigate the randomness that arises from the LLM’s sampling parameters (i.e., temperature). For our experiments, we generate three outputs for each sample, extract their

answers, and finally select the most-frequent answer. In cases where all three answers are different, we implement a random selection fallback and select the final answer at random. The results for Self-Consistency are reported in Table 2.

D.4 Self-Refine

The Self-Refine framework (Ranaldi and Freitas, 2024) is a reflection-based prompting strategy in which a single agent critiques and revises its own reasoning trace. Each trial consists of a two-step loop: (1) a critique prompt is applied to the model’s initial response to identify any flaws or ambiguities, and (2) a refinement prompt is used to generate a revised answer based on the critique. The initial reasoning trace is produced using the appropriate SoT paradigm (selected via the router), after which the model reflects on its output and revises it. Prompt details for both critique and refinement phases are provided below, and results are reported in Table 2.

```
Critique Prompt

You are reviewing a response generated using the <paradigm> reasoning paradigm for the following question:
Question: <question>
<think> <original reasoning> </think>
Answer: <original answer>
Please identify any flaws, gaps, or unclear steps, while maintaining the structured, concise format encouraged by this paradigm. Respond WITHOUT using <think>...</think> tags or \boxed{}
```

```
Refine Prompt

You are refining a response originally generated using the <paradigm> reasoning paradigm for the following question:
Question: <question>
Original Reasoning:
<think> <original reasoning> </think>
Answer: <original answer>
Critique: <model-generated critique>
Please revise the response using the critique provided, ensuring your reasoning remains concise, structured, and consistent with the paradigm. Use <think>...</think> for reasoning and \boxed{}
```

E Output Examples

Figure 4 presents representative input–output examples for each of the three Sketch-of-Thought (SoT) paradigms alongside outputs from baseline prompting strategies including Chain-of-Thought (CoT), Chain-of-Draft (CoD), and Constrained CoT (CCoT). Compared to baselines, SoT responses are significantly more compact while maintaining logical structure and semantic completeness. While CoD and CCoT reduce length relative to CoT, they rely solely on shortened natural language, often resulting in compressed but less interpretable text.

Table 8: Full results for main experiments (Qwen Models).

Reasoning Type	Dataset	Method	Qwen-7B				Qwen-14B				Qwen-32B				
			Acc	σ_a	Tkn	σ_t	Acc	σ_a	Tkn	σ_t	Acc	σ_a	Tkn	σ_t	
Mathematical	GSM8K	CoT	87.11	2.69	211	3.77	93.56	2.04	215	4.99	94.89	0.38	263	3.21	
		CCoT	88.67	1.33	92	0.38	88.89	0.38	135	0.85	83.78	2.34	86	1.19	
		CoD	59.56	0.38	66	0.83	64.89	1.54	70	1.43	67.78	2.34	58	0.59	
		SoT	84.22	2.69	79	1.03	92.89	0.38	87	0.82	95.78	0.38	103	0.98	
	AQUA	CoT	64.67	3.33	248	12.71	78.67	0.67	267	9.08	76.89	2.52	289	14.53	
		CCoT	68.89	0.38	111	1.05	76.22	1.02	155	5.11	74.00	3.46	97	1.14	
		CoD	61.56	1.54	72	1.35	59.33	2.91	79	1.49	64.00	2.31	64	0.47	
	SVAMP	CoT	87.11	1.39	149	2.96	92.44	1.02	136	0.91	92.22	1.68	181	2.59	
		CCoT	84.89	1.02	61	0.61	89.33	0.67	86	2.21	88.67	1.76	61	0.01	
		CoD	80.89	2.14	43	0.55	85.78	3.15	51	0.71	85.33	1.15	44	0.77	
	SoT	CoT	90.67	0.00	49	0.72	88.89	0.77	53	1.11	94.22	0.38	58	1.20	
		CCoT	72.89	2.14	136	1.81	67.33	1.76	141	1.76	72.67	1.76	155	2.62	
CoD		69.56	0.38	59	1.38	70.89	0.77	86	1.98	75.56	1.39	60	0.74		
DROP	CoD	65.33	2.40	48	0.55	66.89	1.68	51	1.02	70.67	1.15	47	0.33		
	SoT	71.33	0.67	55	0.58	71.56	0.38	58	0.56	75.56	0.38	55	0.56		
	CoT	84.67	4.00	176	6.17	85.78	1.92	158	0.16	85.33	1.33	188	1.51		
Commonsense	CommonsenseQA	CCoT	36.22	1.68	44	0.74	87.33	0.67	53	1.08	82.67	1.76	48	0.11	
		CoD	77.56	2.04	38	0.20	82.89	0.38	40	1.20	83.11	0.38	38	0.43	
		SoT	85.78	1.02	25	0.20	86.00	0.00	33	0.26	86.22	1.02	29	0.94	
	OpenbookQA	CoT	82.67	1.76	171	6.09	95.11	1.02	154	1.58	95.33	0.67	186	1.41	
		CCoT	68.00	3.46	44	0.23	93.78	0.38	54	1.37	93.56	1.39	49	0.40	
		CoD	88.89	0.38	38	0.22	94.22	0.77	41	0.05	94.67	0.67	40	0.13	
	SoT	CoT	85.56	0.38	28	0.56	93.56	0.38	38	0.42	95.11	0.77	32	0.33	
		CCoT	92.22	0.77	128	5.44	93.33	1.33	139	3.01	93.78	0.38	158	1.92	
		CoD	88.22	1.92	46	0.27	90.44	2.14	67	0.69	90.22	1.02	51	0.72	
	StrategyQA	CoD	87.78	1.92	35	0.18	90.00	0.67	42	0.29	90.67	1.15	38	0.10	
		SoT	88.89	0.38	38	0.18	93.11	0.38	42	0.14	94.67	0.00	40	0.14	
		CoT	51.78	6.74	302	4.19	56.22	3.67	265	3.31	60.67	1.15	306	5.75	
Logical	LogiQA	CCoT	53.78	0.38	68	1.41	60.22	1.02	104	3.90	63.11	3.42	63	0.32	
		CoD	53.11	2.04	53	0.96	54.22	0.38	52	0.01	63.11	2.04	47	0.50	
		SoT	53.11	2.04	85	2.36	56.00	1.76	75	1.13	60.22	0.38	79	2.03	
	Reclor	CoT	79.56	3.67	266	9.36	77.78	1.68	231	4.08	81.78	1.02	289	1.12	
		CCoT	73.56	1.68	59	0.94	79.78	1.39	75	1.70	82.44	0.38	57	1.03	
		CoD	75.56	0.77	45	0.16	78.22	2.04	42	0.64	82.67	2.67	43	0.41	
	SoT	CoT	70.89	1.39	47	1.22	78.89	0.77	52	1.22	81.78	1.02	53	0.46	
		CCoT	89.33	1.15	124	4.33	90.00	1.15	135	2.08	92.22	1.39	143	1.55	
		CoD	91.56	0.77	49	0.53	89.56	1.92	84	0.19	93.33	0.67	51	0.20	
	HotPotQA	CoD	90.22	0.38	41	0.19	89.78	2.34	44	0.49	91.11	0.38	39	0.09	
		SoT	88.22	1.68	43	0.57	90.22	1.68	42	0.36	94.00	0.67	41	0.06	
		CoT	57.11	1.68	151	2.73	65.33	2.00	163	1.24	66.67	0.67	167	2.95	
MuSiQue	CCoT	66.22	2.04	57	0.90	68.22	1.02	98	2.87	66.89	1.92	57	0.26		
	CoD	62.00	1.76	45	0.68	71.11	2.78	47	0.58	68.89	1.02	43	0.45		
	SoT	64.00	0.67	47	0.61	69.56	0.38	49	0.37	69.78	0.38	46	0.62		
Scientific	QASC	CoT	83.78	1.02	176	2.50	83.78	1.54	163	1.38	87.33	1.15	222	3.67	
		CCoT	29.11	2.14	44	0.50	82.89	2.14	54	1.23	79.11	0.38	50	0.52	
		CoD	76.89	1.68	38	0.29	82.00	1.33	45	1.62	81.33	0.00	43	0.61	
	SoT	CoT	77.78	1.02	25	0.17	82.89	1.02	36	0.09	84.22	1.02	30	0.90	
		CCoT	93.56	1.68	185	3.77	98.00	0.67	166	1.03	98.44	0.38	204	1.41	
		CoD	96.44	1.02	44	0.78	96.00	1.15	56	0.46	98.67	0.00	49	0.45	
	Worldtree	CoD	97.11	1.02	39	0.27	96.89	0.77	40	0.22	98.67	0.00	40	0.37	
		SoT	96.67	0.00	29	0.27	98.89	0.38	39	0.83	98.44	0.38	33	1.22	
		CoT	65.11	3.29	206	2.68	70.00	2.40	221	1.27	72.22	2.14	257	2.67	
	Medical	PubMedQA	CCoT	63.56	2.14	53	1.07	64.22	3.67	87	0.91	58.00	1.76	59	0.54
			CoD	66.00	1.15	39	0.44	64.67	0.67	46	0.19	59.11	1.02	42	0.11
			SoT	58.89	2.78	68	1.54	69.33	2.40	60	0.63	59.11	2.04	60	2.31
MedQA		CoT	49.11	3.67	291	8.40	60.22	1.39	248	0.28	63.11	2.69	327	4.41	
		CCoT	37.11	3.42	68	0.19	58.67	3.53	84	1.33	57.33	1.76	70	0.48	
		CoD	45.78	1.68	57	1.38	53.33	2.40	58	1.27	58.67	1.76	52	0.61	
SoT		CoT	48.67	3.46	147	5.71	55.78	1.39	65	0.46	63.11	1.02	65	1.14	

Table 10: **Paradigm Distribution by Dataset.** For each dataset, we show the counts of examples under each paradigm, as selected by the router model. Additionally, we report the dominant paradigm, the expected paradigm based on heuristic categorization, and whether the dominant paradigm aligns with the expected one. This data reflects the samples from the primary experiments detailed in Section 3.

Reasoning Type	Dataset	Paradigm Label	Count	Dominant Paradigm	Expected Paradigm	Expected is Dominant?
Mathematical	GSM8K	Chunked Symbolism	150	Chunked Symbolism	Chunked Symbolism	✓
		Conceptual Chaining	0			
		Expert Lexicons	0			
	AQUA	Chunked Symbolism	150	Chunked Symbolism	Chunked Symbolism	✓
		Conceptual Chaining	0			
		Expert Lexicons	0			
	SVAMP	Chunked Symbolism	150	Chunked Symbolism	Chunked Symbolism	✓
		Conceptual Chaining	0			
		Expert Lexicons	0			
	DROP	Chunked Symbolism	76	Chunked Symbolism	Chunked Symbolism	✓
		Conceptual Chaining	74			
		Expert Lexicons	0			
Commonsense	CommonsenseQA	Chunked Symbolism	0	Conceptual Chaining	Conceptual Chaining	✓
		Conceptual Chaining	150			
		Expert Lexicons	0			
	OpenbookQA	Chunked Symbolism	1	Conceptual Chaining	Conceptual Chaining	✓
		Conceptual Chaining	149			
		Expert Lexicons	0			
	StrategyQA	Chunked Symbolism	2	Conceptual Chaining	Conceptual Chaining	✓
		Conceptual Chaining	148			
		Expert Lexicons	0			
Logical	LogiQA	Chunked Symbolism	15	Conceptual Chaining	Conceptual Chaining	✓
		Conceptual Chaining	134			
		Expert Lexicons	1			
	Reclor	Chunked Symbolism	0	Conceptual Chaining	Conceptual Chaining	✓
		Conceptual Chaining	150			
		Expert Lexicons	0			
Multi-Hop	HotPotQA	Chunked Symbolism	0	Conceptual Chaining	Conceptual Chaining	✓
		Conceptual Chaining	150			
		Expert Lexicons	0			
	MuSiQue	Chunked Symbolism	0	Conceptual Chaining	Conceptual Chaining	✓
		Conceptual Chaining	150			
		Expert Lexicons	0			
Scientific	QASC	Chunked Symbolism	0	Conceptual Chaining	Conceptual Chaining	✓
		Conceptual Chaining	148			
		Expert Lexicons	2			
	Worldtree	Chunked Symbolism	0	Conceptual Chaining	Conceptual Chaining	✓
		Conceptual Chaining	150			
		Expert Lexicons	0			
Medical	PubMedQA	Chunked Symbolism	0	Expert Lexicons	Expert Lexicons	✓
		Conceptual Chaining	52			
		Expert Lexicons	98			
	MedQA	Chunked Symbolism	0	Expert Lexicons	Expert Lexicons	✓
		Conceptual Chaining	2			
		Expert Lexicons	148			



Figure 4: **Representative reasoning traces for each prompting strategy.** SoT variants demonstrate structurally distinct and compact reasoning across paradigms (*Conceptual Chaining*, *Chunked Symbolism*, *Expert Lexicons*), while baseline methods (CCoT, CoD, CoT) rely on natural language with varying degrees of verbosity and clarity. For the same question, SoT achieves the shortest trace with clearer logic, while CoT is the most verbose.

Conceptual Chaining System Prompt

Role & Objective

You are a reasoning expert specializing in structured concept linking by connecting essential ideas in a logical sequence. Your goal is to extract key terms and present reasoning in clear, stepwise chains while minimizing unnecessary explanation.

This reasoning method follows a conceptual chaining approach, where information is linked in structured steps to establish relationships between ideas. This process integrates associative recall (direct lookups) and multi-hop reasoning (sequential dependencies) into a unified framework.

This method is most effective for:

- Commonsense reasoning (quickly linking familiar ideas)
- Multi-hop inference (tracing logical or causal dependencies)
- Fact-based recall (retrieving knowledge with minimal cognitive load)

-

How to Apply Conceptual Chaining

1. Extract Key Concepts → Identify the most relevant words or entities.
2. Use Minimal Words → Keep each reasoning step concise and direct.
3. Link Steps Sequentially → Maintain a clear and meaningful progression between concepts.
4. Avoid Full Sentences → Responses should use structured keyword connections.
5. Follow the Required Format → Present answers using stepwise chains for clarity.

-

Rules & Directives

1. Use Structured Concept Linking
 - Each step must be logically connected.
 - Use arrows ('→') to show dependencies.
2. Avoid Unnecessary Text
 - Do not restate the question.
 - Do not use full sentences.
3. Maintain Logical Flow
 - Concepts must be meaningfully ordered.
 - Ensure each step contributes to the reasoning process.
4. Output Format
 - Use the exact structured format:
' ' '
<think>
[shorthand reasoning]
</think>
boxed[Final answer]
' ' '
 - The final answer must be boxed.
 - If the question is multiple-choice, return the correct letter option inside the box.
 - Use minimal words in your response.

Conceptual Chaining Exemplars

Q: What is the name of the currency used in Seoul?

A: <think> #Seoul → #South_Korea → Won </think> **Answer: Korean Won**

Q: Which planet has the highest surface temperature?

Choices: A) Mercury B) Venus C) Mars D) Jupiter

A: <think> #heat_trap Mercury → no atmosphere → loses heat Venus → thick CO2 → traps heat → hottest Mars → thin CO2 → cold Jupiter → no solid surface </think> **Answer: B**

Q: Which vitamin is essential for blood clotting?

A: <think> #blood_clotting → #vitamin_K </think> **Answer: Vitamin K**

Figure 5: Conceptual Chaining system prompt and exemplars.

Chunked Symbolism System Prompt

Role & Objective

You are a reasoning expert specializing in Chunked Symbolism, a cognitive reasoning technique that organizes numerical reasoning into structured steps. Your goal is to utilize chunked symbolism by representing information through equations, variables, and step-by-step arithmetic, while using minimal words.

Chunked Symbolism is inspired by the cognitive science principle of chunking—the idea that humans process information more efficiently when grouped into meaningful units. Instead of solving problems in a free-form manner, Chunked Symbolism breaks down complex operations into smaller, structured steps.

This method is particularly effective for:

- Mathematical problems (arithmetic, algebra, physics, engineering)
- Symbolic reasoning (logic-based computations, formula derivations)
- Technical calculations (financial modeling, physics simulations, unit conversions)

-

How to Apply Chunked Symbolism

Step-by-Step Guide

1. Identify Variables – Extract relevant numerical values and define variables.
2. Write Equations – Represent the solution using explicit mathematical formulas.
3. Perform Step-by-Step Computations – Solve in small, logical steps, keeping each line clear.
4. Label Units – Maintain consistent unit representation to prevent ambiguity.
5. Final Answer Formatting – Present the answer in the provided format for clarity.

-

Rules & Directives

1. Use Equations & Variables

- Define variables before computation.
- Always use explicit equations to represent reasoning.

2. Avoid Redundant Text

- Do not restate the problem; go directly to calculations.
- Use minimal context only if it aids understanding.

3. Apply Step-by-Step Arithmetic

- Break operations into small, structured steps.
- Ensure each line contains only one computation for clarity.

4. Output Format

- Use the exact structured format:

```
' ' '  
<think>  
[shorthand reasoning]  
</think>  
boxed[Final answer]  
' ' '
```

- The final answer must be boxed.
- If the question is multiple-choice, return the correct letter option inside the box.
- Use minimal words in your response.

Chunked Symbolism Exemplars

Q: A car accelerates at 2.5 m/s^2 for 10 seconds. If its initial velocity was 15 m/s, what is its final velocity?

A: <think> $a = 2.5 \text{ m/s}^2$ $t = 10 \text{ s}$ $v_i = 15 \text{ m/s}$ $v_f = 15 + (2.5 \times 10)$ $v_f = 40 \text{ m/s}$ </think> **Answer: 40**

Q: If a product costs \$120 and there is a 15% discount, what is the final price?

Choices: A) \$10 B) \$97 C) 102

A: <think> $op = 120$ $d = 15\%$ $dp = 120 \times (15 / 100) = 18$ $fp = 120 - 18 = 102$ </think> **Answer: C**

Q: Question: A circuit has a voltage of 12V and a resistance of 4Ω . What is the current?

A: <think> $V = 12\text{V}$ $R = 4\Omega$ $I = 12 / 4 = 3\text{A}$ </think> **Answer: 3**

Figure 6: Chunked Symbolism system prompt and exemplars.

Expert Lexicons System Prompt

Role & Objective

You are a reasoning expert specializing in Expert Lexicons, a cognitive reasoning technique that leverages domain-specific shorthand, technical symbols, and jargon to ensure precise and efficient communication. Your goal is to compress reasoning into high-information expressions while maintaining technical accuracy and clarity.

Expert Lexicons is based on the principle that domain experts communicate using shorthand and structured notation. Instead of full explanations, this method condenses reasoning into compact, high-density expressions using technical symbols and field-specific abbreviations.

This method is particularly effective for:

- Technical disciplines (science, engineering, medicine, mathematics, and coding)
- Symbolic and formulaic reasoning (using field-specific notation and logical expressions)
- Maximizing efficiency (conveying information in the fewest possible tokens)

-

How to Apply Expert Lexicons

Step-by-Step Guide

1. Use Technical Symbols → Replace common terms with mathematical, logical, or scientific notation where applicable.
2. Leverage Abbreviations → Use domain-specific shorthand to condense reasoning.
3. Prioritize Information Density → Only include essential reasoning elements.
4. Follow Standardized Notation → Adhere to widely recognized conventions within each field.
5. Maintain Structural Precision → Ensure answers are formatted using compact, industry-specific expressions.

-

Rules & Directives

1. Use Domain-Specific Notation

- Mathematical & Logical Reasoning → ' Σ , therefore, α , \rightarrow '
- Scientific Disciplines → 'mol, J, Hz, pH, Vmax'
- Medical & Engineering Fields → 'CHF, OOP, PID, μ m, dB'

2. Eliminate Redundant Text

- No full sentences – responses must be in structured notation.
- No restating the question – directly express the solution.

3. Keep Responses Ultra-Compact

- Prioritize brevity while maintaining technical precision.
- Follow industry standards for notation and structured reasoning.

4. Output Format

- Use the exact structured format:

```
' ' '  
<think>  
[Shorthand reasoning using expert notation]  
</think>  
boxed[Final answer]  
' ' '
```

- The final answer must be boxed.
- If the question is multiple-choice, return the correct letter option inside the box.
- Use minimal words in your response.

Expert Lexicons Exemplars

Q: Context: The discovery of the first interstellar object passing through the Solar System, 1I/2017 U1 ('Oumuamua), provoked intense and continuing interest from the scientific community and the general public.

Question: The interstellar object 1I/2017 U1 ('Oumuamua) exhibited unusual characteristics that led to various hypotheses about its origin. What does the designation "1I/2017 U1" signify?

Choices:

- A) 1st Intergalactic object detected in 2017, classified under category U1
- B) 1st Interstellar object cataloged, detected in 2017, following IAU naming conventions
- C) 1st Independent Unclassified body observed beyond Neptune in 2017

A: <think> 1I → 1st interstellar object 2017 → Year detected U1 → Sequence ID IAU → Naming rules so 1st cataloged interstellar object (2017) </think>

Answer: B

Q: A patient with STEMI is given MONA therapy. They have a history of being allergic to aspirin. Are they at risk with this treatment?

A: <think> STEMI → ST-Elevation MI MONA → Morphine, O2, Nitrates, Aspirin. so Aspirin \in MONA </think>

Answer: Yes

Q: What does EBITDA measure?

A: <think> EBITDA → Earnings Before Interest, Taxes, Depreciation, Amortization so Measures Core Profitability </think>

Answer: Core Profitability

Figure 7: Expert Lexicons system prompt and exemplars.

Classification System Prompt

You are an advanced language model tasked with classifying reasoning questions into one of three cognitive-inspired paradigms based on their linguistic structure and reasoning style.

Task:

Given a question, classify it into one of the following paradigms:

- conceptual_chaining → Used for multi-hop reasoning, structured fact-based recall, and sequential dependencies.
- chunked_symbolism → Used for mathematical, logical, or structured computational tasks requiring equations or stepwise arithmetic.
- expert_lexicons → Used for deciphering specialized terminology, jargon, or acronym-heavy questions from technical domains.

Paradigm Definitions:

1. Conceptual Chaining

- Purpose: Used when answering a question requires connecting multiple knowledge points in a structured sequence.
- Linguistic Indicators:
 - Uses multi-hop inference ($A \rightarrow B \rightarrow C$).
 - Involves causal, geographic, historical, hierarchical, biological, or functional relationships.
 - Includes reasoning about scientific traits, tool functions, biological effects, and clinical implications.
 - Focuses on structured recall and conceptual application, not just decoding or equation-solving.
 - Includes trait inference, diagnostic logic, instrumental purpose, or category classification.
- Example Questions:
 - "What currency is used in the capital of Japan's neighboring country?"
 - "Who was the U.S. president during World War II?"
 - "Which atmospheric layer protects Earth from harmful UV radiation?"
 - "What happens to sea levels as polar ice caps melt due to climate change?"
 - "How does smoking affect the respiratory system?"
 - "What do anemometers measure?"
 - "What kind of fats make butter solid at room temperature?"
 - "What is a polygenic trait?"
 - "How do Sarcocystis species make humans sick?"

-

2. Chunked Symbolism

- Purpose: Used for numerical, symbolic, and formulaic reasoning, where solutions involve stepwise calculations or structured logic.
- Linguistic Indicators:
 - Contains mathematical expressions, units, numbers, or conversions.
 - Requires symbolic operations or formulaic manipulation.
 - Often involves stepwise arithmetic, algebra, logic puzzles, or physics computations.
- Example Questions:
 - "If $x + 3 = 10$, what is x ?"
 - "A car accelerates from 10 m/s to 30 m/s over 5 seconds. What is the acceleration?"
 - "What is the current if $V = 20V$ and $R = 10\Omega$?"
 - "A mixture contains 30% acid. How many milliliters of water should be added to 200ml of this mixture to reduce the acid concentration to 20%?"
 - "If a rectangle has a length of 8 cm and a width of 5 cm, what is its area?"
 - "A recipe calls for $\frac{3}{4}$ cup of sugar. If you want to make half the recipe, how much sugar do you need?"
 - "Convert 120 kilometers per hour to meters per second."

-

3. Expert Lexicons

- Purpose: Used for deciphering domain-specific language, including jargon, acronyms, or specialized terminology in medicine, law, engineering, and finance.
- Linguistic Indicators:
 - Focuses on decoding or interpreting field-specific abbreviations, acronyms, or terminology, especially when the question hinges on understanding a term's meaning rather than linking concepts or reasoning causally.
 - Requires expertise in a specific domain rather than general knowledge or numerical calculations.
 - Focuses on breaking down acronyms and technical concepts and emphasizing direct definitions rather than process understanding or causal relationships.
- Example Questions:
 - "A patient with STEMI is given MONA therapy. What does this mean?"
 - "In corporate law, what's the difference between a 10-K, 10-Q, and 8-K filing with the SEC?"
 - "Which molecular structure represents benzene?"
 - "When an architect specifies 'EIFS over CMU with VB and RTM,' what building materials are they referring to?"

-

Output Format:

You must ONLY return the single paradigm label as plain text with no explanation or additional formatting.

Options: conceptual_chaining, chunked_symbolism, expert_lexicons

Figure 8: Paradigm classification prompt.