

TRAINING LARGE LANGUAGE MODELS TO REASON IN A CONTINUOUS LATENT SPACE

Anonymous authors

Paper under double-blind review

ABSTRACT

Large language models are restricted to reason in the “language space”, where they typically express the reasoning process with a chain-of-thoughts (CoT) to solve a complex reasoning problem. However, we argue that language space may not be the optimal reasoning space. For example, most word tokens are primarily for textual coherence and not essential for reasoning, while some critical tokens require complex planning and pose huge challenges to LLMs. To explore the potential of LLM reasoning in an unrestricted latent space instead of using human language, we introduce a new paradigm COCONUT (Chain of **C**ontinuous **T**hought). We utilize the last hidden state of the LLM as a representation of the reasoning state (termed “continuous thought”). Rather than decoding this into a word token, we feed it back to the LLM as the subsequent input embedding directly in the continuous space. Experiments show that COCONUT can effectively augment the LLM on several reasoning tasks. It even outperforms CoT in certain logical reasoning tasks that require substantial planning, despite generating fewer tokens during inference. More interestingly, we observe an advanced reasoning patterns emerging from latent reasoning: the continuous thought can encode multiple potential next reasoning steps, allowing the model to perform a breadth-first search (BFS) to solve the problem, rather than prematurely committing to a single deterministic path like CoT. These findings demonstrate the promise of latent reasoning and offer valuable insights for future research on latent reasoning methods.

1 INTRODUCTION

Large language models (LLMs) have demonstrated remarkable reasoning abilities, emerging from extensive pretraining on human language (Dubey et al., 2024; Achiam et al., 2023). While the next token prediction is an effective training objective, it imposes a fundamental constraint on the LLM as a reasoning machine: the reasoning process of LLMs must be generated in word tokens. For example, a prevalent approach, known as chain-of-thought (CoT) reasoning (Wei et al., 2022), involves prompting or training LLMs to generate solutions step-by-step using natural language. However, this stands in stark contrast to human cognition. Neuroimaging studies have consistently shown that the language network – a set of brain regions responsible for language comprehension and production – remains largely inactive during various reasoning tasks (Amalric & Dehaene, 2019; Monti et al., 2012; 2007; 2009; Fedorenko et al., 2011). More evidence has indicated that human language is optimized for communication rather than reasoning (Fedorenko et al., 2024).

A significant problem arises when LLMs are required to output language during reasoning: the “reasoning amount” behind each token varies greatly, yet current LLM architectures allocate nearly the same computing budget for predicting every token. Most tokens in a reasoning chain are generated solely for fluency, contributing little to the actual reasoning process. On the contrary, some critical tokens require complex planning and pose huge challenges to LLMs. While previous work has attempted to fix these problems by prompting LLMs to generate succinct reasoning chains (Madaan & Yazdanbakhsh, 2022), or performing additional reasoning before generating some critical tokens (Zelikman et al., 2024), these solutions remain constrained within the language space and do not solve the problems fundamentally. Ideally, LLMs should be allowed to reason freely in an unconstrained latent space and only translate the outcomes into language once the reasoning process is complete.

054
055
056
057
058
059
060
061
062
063
064
065
066
067
068
069
070
071
072
073
074
075
076
077
078
079
080
081
082
083
084
085
086
087
088
089
090
091
092
093
094
095
096
097
098
099
100
101
102
103
104
105
106
107

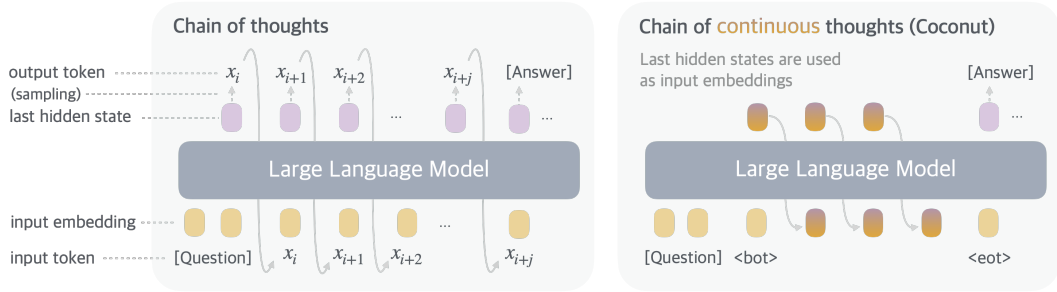


Figure 1: A comparison of CoT and COCONUT. In CoT, the model generates the reasoning process as a word token sequence (e.g., $[x_i, x_{i+1}, \dots, x_{i+j}]$ in the figure). COCONUT (Chain of Continuous Thoughts) regards the last hidden state as a representation of reasoning state (termed “continuous thought”), and directly uses it as the next input embedding. This allows the LLM to reason in an unrestricted latent space instead of language space.

We aim to explore LLM reasoning in the latent space by introducing a novel paradigm, COCONUT (Chain of Continuous Thought). It involves a simple modification to the traditional CoT process. Instead of mapping between hidden states and language tokens using the language model head and embedding layer, COCONUT directly feeds the last hidden state (a continuous thought) as the input embedding for the next token (Figure 1). This modification frees the reasoning from language space, and the architecture can be optimized end-to-end by gradient descent, as continuous thoughts are fully differentiable. To enhance the training of these continuous thoughts, we employ a multi-stage training strategy inspired by Deng et al. (2024), which effectively utilizes language reasoning chains to guide the training process.

The experiments demonstrate that COCONUT successfully enhances the reasoning capabilities of LLMs. Specifically, on math reasoning problems (GSM8k, Cobbe et al., 2021), using more continuous thoughts is shown to be beneficial to reasoning accuracy, mirroring the effects of language reasoning chains. This indicates the potential to scale and solve increasingly challenging problems by chaining more continuous thoughts. On logical reasoning problems including ProntoQA (Saparov & He, 2022), and our newly proposed ProsQA (Section 4.1) which requires stronger planning ability, COCONUT and some of its variants even surpasses language-based CoT methods, while generating significantly fewer tokens during inference.

Interestingly, the removal of language space constraints has led to a novel reasoning pattern. By manipulating the COCONUT model to switch between latent reasoning and language reasoning, we are able to unveil the latent reasoning process. Unlike language-based reasoning, continuous thoughts in COCONUT can encode multiple potential next steps simultaneously, allowing for a reasoning process akin to breadth-first search (BFS). While the model may not initially make the correct decision, it can maintain all possible options within the continuous thoughts and progressively eliminate incorrect paths through reasoning, guided by some implicit value functions. This advanced reasoning mechanism surpasses traditional CoT approaches, even though the model is not explicitly trained or instructed to operate in this manner, as seen in previous works (Yao et al., 2023; Hao et al., 2023). We believe that these findings underscore the potential of latent reasoning and could provide valuable insights for future research.

2 RELATED WORK

Chain-of-thought (CoT) reasoning. We use the term chain-of-thought broadly to refer to methods that generate an intermediate reasoning process in language before outputting the final answer. This includes prompting LLMs (Wei et al., 2022; Khot et al., 2022; Zhou et al., 2022), or training LLMs to generate reasoning chains, either with supervised fine-tuning (Yue et al., 2023; Yu et al., 2023) or reinforcement learning (Wang et al., 2024; Havrilla et al., 2024; Shao et al., 2024; Yu et al., 2024a). Madaan & Yazdanbakhsh (2022) classified the tokens in CoT into symbols, patterns, and text, and proposed to guide the LLM to generate concise CoT based on analysis of their roles. Recent theoretical analyses have demonstrated the usefulness of CoT from the perspective of model

expressivity (Feng et al., 2023; Merrill & Sabharwal, 2023; Li et al., 2024). By employing CoT, the effective depth of the transformer increases because the generated outputs are looped back to the input (Feng et al., 2023). These analyses, combined with the established effectiveness of CoT, motivated our exploration of continuous thoughts, in contrast to other latent reasoning methods. While CoT has proven effective for certain tasks, its autoregressive generation nature makes it challenging to mimic human reasoning on more complex problems (LeCun, 2022; Hao et al., 2023), which typically require planning and search. There are works that equip LLMs with explicit tree search algorithms (Xie et al., 2023; Yao et al., 2023; Hao et al., 2023), or train the LLM on search dynamics and trajectories (Lehnert et al., 2024; Gandhi et al., 2024). In our analysis, we find that after removing the constraint of language space, a new reasoning pattern similar to BFS emerges, even though the model is not explicitly trained in this way.

Latent reasoning of LLM. Previous works mostly define latent reasoning of LLM as the hidden computing in transformers (Yang et al., 2024; Biran et al., 2024). Yang et al. (2024) constructed a dataset of two-hop reasoning problems and discovered that it is possible to recover the intermediate variable from the hidden representation of LLMs. Biran et al. (2024) further proposed to intervene the latent reasoning by “back-patching” the hidden representation. Another line of work has discovered that, even if the model generates a CoT to reason, the model may actually utilize a different latent reasoning process. This phenomenon is known as the unfaithfulness of CoT reasoning (Wang et al., 2022; Turpin et al., 2024). To enhance the latent reasoning of LLM, previous research proposed to augment it with additional tokens. Goyal et al. (2023) pretrained model by randomly inserting a learnable `<pause>` tokens to the corpus. This improves LLM’s performance on a variety of tasks, especially when followed by supervised finetuning with `<pause>` tokens. On the other hand, Pfau et al. (2024) further explored the usage of filler tokens, e.g., “. . .”, and concluded that they work well for highly parallelizable problems. However, these methods do not extend the expressivity of the LLM like CoT (Pfau et al., 2024); hence, they may not scale to more general and complex reasoning problems. Recently, it has also been found that one can “internalize” the chain of thought reasoning into latent reasoning with knowledge distillation (Deng et al., 2023) or a special training curriculum which gradually shortens CoT (Deng et al., 2024). Yu et al. (2024b) also proposed to distill a model that can reason latently from data generated with complex reasoning algorithms. These training methods can be combined to our framework, and specifically, we find that breaking down the learning of continuous thoughts into multiple stages, inspired by iCoT (Deng et al., 2024), is very beneficial for the training.

3 COCONUT: CHAIN OF CONTINUOUS THOUGHTS

In this section, we introduce our new paradigm COCONUT (Chain of Continuous Thoughts) for reasoning outside the language space. We begin by introducing the background and notations of language models. For an input sequence $x = (x_1, \dots, x_T)$, the standard large language model \mathcal{M} can be described as:

$$H_t = \text{Transformer}(E_t + P_t)$$

$$\mathcal{M}(x_{t+1} | x_{\leq t}) = \text{softmax}(Wh_t)$$

where $E_t = [e(x_1), e(x_2), \dots, e(x_t)]$ is the sequence of token embeddings up to position t ; $P_t = [p(1), p(2), \dots, p(t)]$ is the sequence of positional embeddings up to position t ; $H_t \in \mathbb{R}^{t \times d}$ is the matrix of the last hidden states for all tokens up to position t ; h_t is the last hidden state of position t , i.e., $h_t = H_t[t, :]$; $e(\cdot)$ is the token embedding function; $p(\cdot)$ is the positional embedding function; W is the parameter of the language model head.

Method Overview. In the proposed COCONUT method, the LLM switches between the “language mode” and “latent mode” (Figure 1). In language mode, the model operates as a standard language model, autoregressively generating the next token. In latent mode, it directly utilizes the last hidden state as the next input embedding. This last hidden state represents the current reasoning state, termed as a “continuous thought”.

Special tokens `<bot>` and `<eot>` are employed to mark the beginning and end of the latent mode, respectively. As an example, we assume latent reasoning occurs between positions i and j , i.e., $x_i = \text{<bot>}$ and $x_j = \text{<eot>}$. When the model is in the latent mode

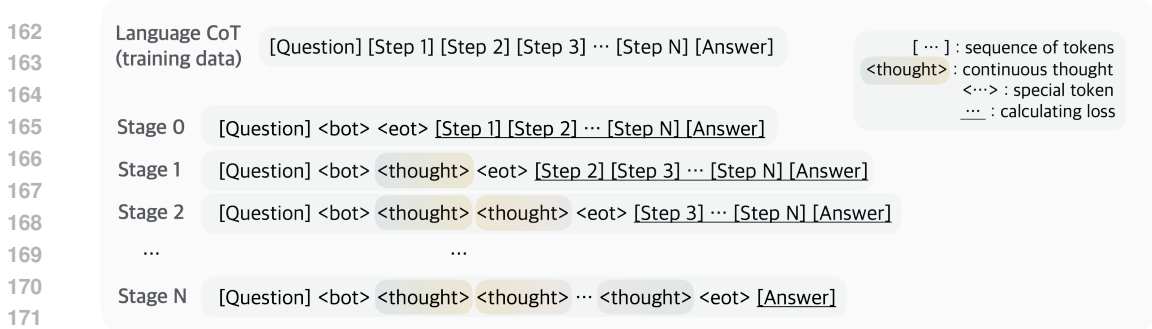


Figure 2: The training procedure of COCONUT. At each stage, we integrate c additional continuous thought ($c = 1$ in this example), and remove one reasoning step in the training data. The cross-entropy loss is then calculated on the remaining tokens after continuous thoughts.

($i < t < j$), we use the last hidden state from the previous token to replace the input embedding, i.e., $E_t = [e(x_1), e(x_2), \dots, e(x_i), h_i, h_{i+1}, \dots, h_{t-1}]$. After the latent mode finishes, ($t \geq j$), the input after position reverts to using the token embedding, i.e., $E_t = [e(x_1), e(x_2), \dots, e(x_i), h_i, h_{i+1}, \dots, h_{j-1}, e(x_j), \dots, e(x_t)]$. It is noteworthy that $\mathcal{M}(x_{t+1} | x_{\leq t})$ is not defined when $i < t < j$, since the latent thought is not intended to be mapped back to language space. However, $\text{softmax}(Wh_t)$ can still be calculated for probing purposes (see Section 4).

Training Procedure. In this work, we focus on a problem-solving setting where the model receives a question as input and is expected to generate an answer through a reasoning process. We leverage language CoT data to supervise continuous thought by implementing a multi-stage training curriculum inspired by Deng et al. (2024). As shown in Figure 2, in the initial stage, the model is trained on regular CoT instances. In the subsequent stages, at the k -th stage, the first k reasoning steps in the CoT are replaced with $k \times c$ continuous thoughts¹, where c is a hyperparameter controlling the number of latent thoughts replacing a single language reasoning step. Following Deng et al. (2024), we also reset the optimizer state when training stages switch. We insert `<bot>` and `<eot>` tokens to encapsulate the continuous thoughts.

During the training process, we mask the loss on questions and latent thoughts. It is important to note that the objective does not encourage the continuous thought to *compress the removed language thought*, but rather to *facilitate the prediction of future reasoning*. Therefore, it’s possible for the LLM to learn a more effective representation compared to language reasoning steps.

Training Details. Our proposed continuous thoughts are fully differentiable, allowing backpropagation. We perform $n + 1$ forward passes when n latent thoughts are scheduled in the current training stage, computing a new latent thought with each pass and then conducting an additional forward pass to obtain a loss on the remaining text sequence. While we can save any repetitive computing by using KV cache, the sequential nature of the multiple forward passes poses challenges for parallelism. Further optimizing the training efficiency of COCONUT remains an important direction for future research.

Inference Process. The inference process for COCONUT is analogous to standard language model decoding, except that in latent mode, we directly feed the last hidden state as the next input embedding. A challenge lies in determining when to switch between latent and language modes. As we focus on the problem-solving setting, we insert a `<bot>` token immediately following the question tokens. For `<eot>`, we consider two potential strategies: a) train a binary classifier on latent thoughts to enable the model to autonomously decide when to terminate the latent thoughts, or b) always pad the latent thoughts to a constant length. We found that both approaches work comparably well. Therefore, we use the second option in our experiment for simplicity, unless specified otherwise.

¹If a reasoning chain is shorter than k steps, then all the language thoughts will be removed.

4 EXPERIMENTS

In this section, we validate the feasibility of LLM reasoning in latent space through experiments on three datasets. We mainly evaluate the accuracy by comparing the model-generated answers with the ground truth. The number of newly generated tokens per question is also listed, as a measure of reasoning efficiency.²

4.1 DATASETS

Math Reasoning. We use GSM8k (Cobbe et al., 2021) as the dataset for math reasoning. It consists of grade school-level math problems. Compared to other datasets of our experiments, the problems are more diverse and open-domain, closely resembling the real-world use cases. Through this task, we explore the potential of latent reasoning in practical applications. To train the model, we use a synthetic dataset generated by Deng et al. (2023).

Logical Reasoning. Logical reasoning involves the proper application of known conditions to prove or disprove a conclusion using logical rules. This requires the model to choose from multiple possible reasoning paths, where the correct decision often relies on exploration and planning ahead. This serves as a simplified simulation of more advanced reasoning tasks, such as automated theorem proving (Chen et al., 2023; DeepMind, 2024). We use 5-hop ProntoQA (Saparov & He, 2022) questions, with fictional concept names. For each problem, an tree-structured ontology is randomly generated and described in natural language as a set of known conditions. The model is asked to judge whether a given statement is correct based on these conditions.

We found that the generation process of ProntoQA was overly simplified, especially since the size of distracting branches in the ontology is always small, reducing the need for complex planning. To fix that, we apply a new dataset construction pipeline using randomly generated DAGs to structure the known conditions. The resulting dataset requires the model to perform substantial planning and searching over the graph to find the correct reasoning chain. We refer to this new dataset as the ProsQA (**P**roof with **S**earch **Q**uestion-**A**nswering). A visualized example is shown in Figure 6. More details of datasets can be found in Appendix A.

4.2 EXPERIMENTAL SETUP

We pre-trained GPT-2 (Radford et al., 2019) as the base model for all experiments. The learning rate is set to 1×10^{-4} while the effective batch size is 128. Following Deng et al. (2024), we also reset the optimizer when training stages switch.

Math Reasoning. By default, we use 2 latent thoughts (i.e., $c = 2$) for each reasoning step. We analyze the correlation between performance and c in Section 4.4. The model goes through 3 stages besides the initial stage. Then, we will have an additional stage, where we still use $3 \times c$ continuous thoughts as in the last stage, but remove all the remaining language reasoning chain. This handles the long-tail distribution of reasoning chains longer than 3 steps. We train the model for 6 epochs in the initial stage, and 3 epochs in each remaining stage.

Logical Reasoning. We use one continuous thought for every reasoning step (i.e., $c = 1$). The model goes through 6 training stages in addition to the initial stage, because the maximum number of reasoning steps is 6 in these two datasets, and the model fully reasons with continuous thoughts to solve the problems in the last stage. We train the model for 5 epochs per stage.

For all datasets, after the standard schedule, the model stays in the final training stage, until the 50th epoch. We select the checkpoint based on the accuracy on the validation set. For inference, we manually set the number of continuous thoughts to be consistent with their final training stage. We use greedy decoding for all experiments.

4.3 BASELINES AND ABLATIONS

We consider the following baselines: (1) *CoT*: We use the complete reasoning chains to train the language model with supervised finetuning, and during inference, the model generates a reasoning

²One continuous thought is counted as one token since the computational cost is essentially the same.

Method	GSM8k		ProntoQA		ProsQA	
	Acc. (%)	# Tokens	Acc. (%)	# Tokens	Acc. (%)	# Tokens
CoT	42.9 ±0.2	25.0	98.8 ±0.8	92.5	77.5 ±1.9	49.4
No-CoT	16.5 ±0.5	2.2	93.8 ±0.7	3.0	76.7 ±1.0	8.2
iCoT	30.0*	2.2	99.8 ±0.3	3.0	98.2 ±0.3	8.2
Pause Token	16.4 ±1.8	2.2	77.7 ±21.0	3.0	75.9 ±0.7	8.2
COCONUT (Ours)	34.1 ±1.5	8.2	99.8 ±0.2	9.0	97.0 ±0.3	14.2
- w/o curriculum	14.4 ±0.8	8.2	52.4 ±0.4	9.0	76.1 ±0.2	14.2
- w/o thought	21.6 ±0.5	2.3	99.9 ±0.1	3.0	95.5 ±1.1	8.2
- pause as thought	24.1 ±0.7	2.2	100.0 ±0.1	3.0	96.6 ±0.8	8.2

Table 1: Results on three datasets. Higher accuracy indicates stronger reasoning ability, while generating fewer tokens indicates better efficiency. *The result of *iCoT* is from Deng et al. (2024).

chain before outputting an answer. (2) *No-CoT*: The LLM is trained to directly generate the answer without using a reasoning chain. (3) *iCoT* (Deng et al., 2024): The model is trained with language reasoning chains and follows a carefully designed schedule that “internalizes” CoT. As the training goes on, tokens at the beginning of the reasoning chain are gradually removed until only the answer remains. During inference, the model directly predicts the answer. (4) *Pause token* (Goyal et al., 2023): The model is trained using only the question and answer, without a reasoning chain. However, different from *No-CoT*, special `<pause>` tokens are inserted between the question and answer, which are believed to provide the model with additional computational capacity to derive the answer. For a fair comparison, the number of `<pause>` tokens is set the same as continuous thoughts in COCONUT.

We also evaluate some variants of our method: (1) *w/o curriculum*: Instead of the multi-stage training, we directly use the data from the last stage which only includes questions and answers to train COCONUT. The model uses continuous thoughts to solve the whole problem. (2) *w/o thought*: We keep the multi-stage training which removes initial reasoning steps gradually, but don’t use any continuous latent thought. While this is similar to *iCoT* in the high-level idea, the exact training schedule is set to be consistent with COCONUT, instead of *iCoT*. This ensures a more strict comparison. (3) *Pause as thought*: We use special `<pause>` tokens to replace the continuous thought, and apply the same multi-stage training scheme as COCONUT.

4.4 RESULTS AND DISCUSSION

We show the overall results on all datasets in Table 1. Continuous thoughts effectively enhance LLM reasoning, as shown by the consistent improvement over *no-CoT*. It even shows better performance than *CoT* on ProsQA. We describe several key conclusions from the experiments as follows.

“Chaining” continuous thoughts enhances reasoning. In conventional CoT, the output token serves as the next input, which is believed to increase the effective depth of LLMs and enhance their expressiveness (Feng et al., 2023). We explore whether latent space reasoning retains this property, as it would suggest that this method could scale to solve increasingly complex problems by chaining multiple latent thoughts.

In our experiments with GSM8k, we found that COCONUT outperformed other architectures trained with similar strategies, particularly surpassing the latest baseline, *iCoT* (Deng et al., 2024). The performance is significantly better than COCONUT (*pause as thought*) which also enables more computation in the LLMs. While Pfau et al. (2024) empirically shows that filler tokens, such as the special `<pause>` tokens, can benefit highly parallelizable problems, our results show that COCONUT architecture is more effective for general problems, e.g., math word problems, where a reasoning step often heavily depends on previous steps. Additionally,

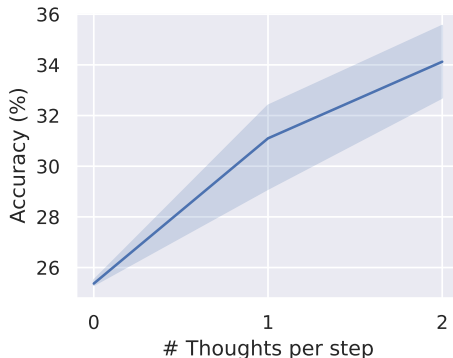


Figure 3: Accuracy on GSM8k with different number of continuous thoughts.

we experimented with adjusting the hyperparameter c , which controls the number of latent thoughts corresponding to one language reasoning step. As we increased c from 0 to 1 to 2, the model’s performance steadily improved (Figure 3). These results strongly suggest that a chaining effect similar to CoT can be observed in the latent space.

In two other synthetic tasks, we found that the variants of COCONUT (*w/o thoughts* or *pause as thought*), and *iCoT* also achieve impressive accuracy. This indicates that in these tasks, the model’s computational capacity may not be the bottleneck. In contrast, GSM8k, being an open-domain question-answering task, likely involves more complex contextual understanding and modeling, placing higher demands on computational capability.

Latent Reasoning Excels Language Reasoning in Planning. Some complex reasoning tasks require the model to “look ahead” to assess whether a particular step is the right choice. Among the datasets used in our experiments, GSM8k consists of grade-school-level math word problems, allowing for intuitive judgment of the next reasoning step; ProntoQA has distracting branches of small sizes, which makes it relatively easy to determine the next step too. In contrast, ProsQA is based on a randomly generated DAG structure, posing a significant challenge to the model’s planning abilities. Reasoning in language space cannot effectively solve the problem. As shown in the table, *CoT* doesn’t show significant improvement over *No-CoT*. On the contrary, COCONUT, some of its variants and *iCoT* significantly improve the reasoning on ProsQA. This suggests an advantage in using latent space over language tokens for tasks requiring extensive planning. We conduct in-depth analysis of the latent reasoning process in Section 5.

The LLM still needs guidance to learn continuous thoughts. In the ideal case, the model should learn the most effective continuous thoughts automatically through gradient descent on questions and answers (i.e., COCONUT *w/o curriculum*). However, from the experimental results, we found the models trained this way do not perform any better than no-CoT.

With the multi-stage curriculum which decomposes the training into easier objectives, COCONUT is able to achieve top performance across various tasks. The multi-stage training also integrates well with pause tokens (COCONUT-*pause as thought*). Despite using the same architecture and similar multi-stage training objectives, we observed a small gap between the performance of *iCoT* and COCONUT (*w/o thoughts*). The finer-grained removal schedule (token by token) and a few other tricks in *iCoT* may ease the training process. We leave combining *iCoT* and COCONUT as a future work. While the multi-stage training used for COCONUT has proven effective, further research is definitely needed to develop better and more general strategies for learning reasoning in latent space, especially without the supervision from language reasoning chains.

Continuous thoughts are efficient representations of reasoning. Though the continuous

thoughts are not intended to be decoded to language tokens, we can still use it as an intuitive interpretation of the latent reasoning. We show a case study in Figure 4 of a math word problem solved by COCONUT ($c = 1$). The first continuous thought can be decoded into tokens like “180”, “180” (with a space), and “9”. Note that, the reasoning trace for this problem should be $3 \times 3 \times 60 = 9 \times 60 = 540$, or $3 \times 3 \times 60 = 3 \times 180 = 540$. The interpretations of the first thought happen to be the first intermediate variables in the calculation. Moreover, it encodes a distribution of different traces into the continuous thoughts. As shown in Section 5.3, this feature enables a more advanced reasoning pattern for planning-intense reasoning tasks.

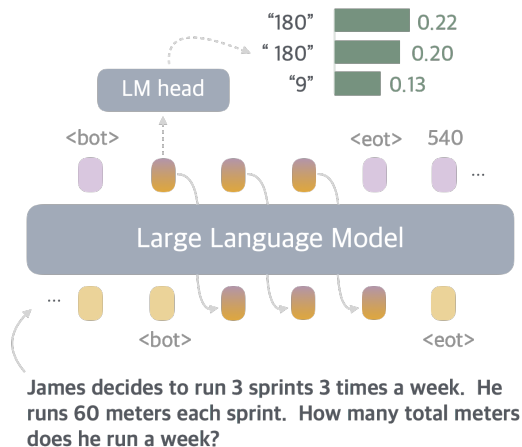


Figure 4: A case study where we decode the continuous thought into language tokens

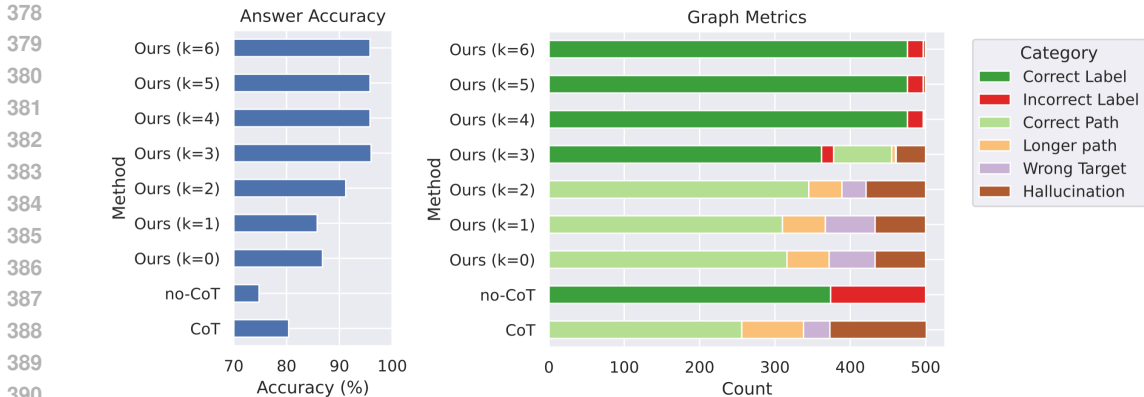


Figure 5: The answer accuracy and graph metrics (see Section 5.1) of multiple variants of COCONUT and baselines on ProsQA.

5 UNDERSTANDING THE LATENT REASONING IN COCONUT

Since COCONUT enables *the switch between language and continuous space reasoning*, we are able to manipulate the model to output language at a certain point, and then infer the previous latent reasoning process from it. This is especially helpful to understand why COCONUT and some other latent reasoning methods can outperform CoT on ProsQA while generating much fewer tokens. Our analysis surprisingly indicates that COCONUT allows the LLM to develop a fundamentally different reasoning pattern than CoT. The continuous thought not only encodes multiple partial reasoning paths, but also enables a latent search process similar to BFS.

5.1 EXPERIMENTAL SETUP

Method. We slightly modify the original training curriculum, so that at any training stage, data from other stages is mixed in with a certain probability ($p = 0.3$). This prevents the model from forgetting earlier stages after the complete training schedule. Therefore, it allows us to control the number of latent thoughts during inference by manually setting the `<eot>` token. When we enforce COCONUT to use k continuous thoughts during inference, the model should output the remaining reasoning chain in language, starting from the $k + 1$ step. In our experiments, we test variants with $k \in \{0, 1, 2, 3, 4, 5, 6\}$. Note that all these variants only differ in inference time while sharing the same model weights. Besides, we report the performance of *CoT* and *no-CoT* as references.

Metrics. We apply two sets of evaluation metrics. One of them is based on the correctness of the *final answer*, regardless of the reasoning process. It is the metric used in the main results (Section 4.4). To enable fine-grained analysis, we define another metric on the *reasoning process*. Assuming we have a complete language reasoning chain which specifies a path in the graph, we can classify it into (1) **Correct Path**: The output is one of the shortest paths to the correct answer. (2) **Longer Path**: A valid path that correctly answers the question but is longer than the shortest path. (3) **Hallucination**: The path includes nonexistent edges or is disconnected. (4) **Wrong Target**: A valid path in the graph, but the destination node is not the one being asked. These four categories naturally apply to the output from COCONUT ($k = 0$) and *CoT*, which generate the full path. For COCONUT with $k > 0$ that outputs only partial paths in language (with the initial steps in continuous reasoning), we classify the reasoning as a Correct Path *if a valid explanation can complete it*. Also, we define Longer Path and Wrong Target for partial paths similarly. If no valid explanation completes the path, it's classified as hallucination. In *no-CoT* and COCONUT with larger k , the model may only outputs the final answer without any partial path, it falls into (5) **Correct Label** or (6) **Incorrect Label**. These six categories cover all cases without overlap.

432
433
434
435
436
437
438
439
440
441
442
443
444
445
446
447
448
449
450
451
452
453
454
455
456
457
458
459
460
461
462
463
464
465
466
467
468
469
470
471
472
473
474
475
476
477
478
479
480
481
482
483
484
485

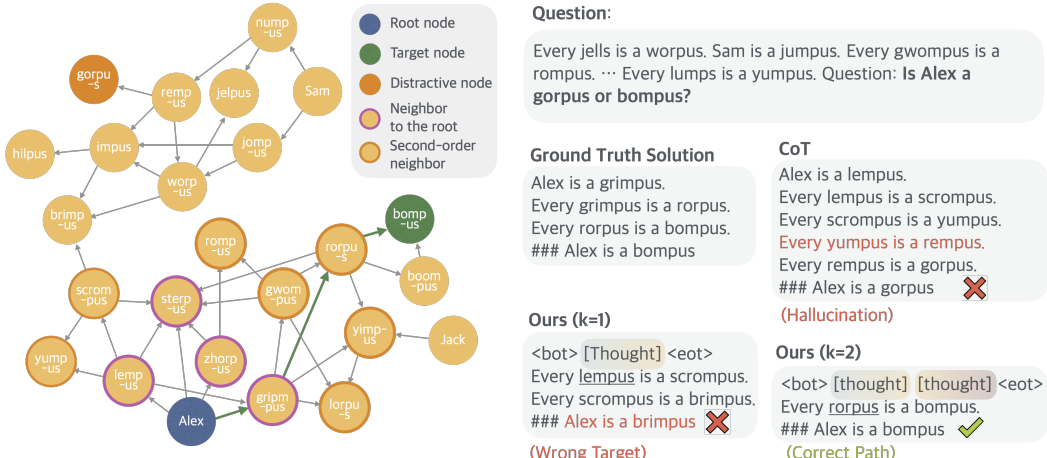


Figure 6: A case study of ProsQA. The model trained with *CoT* hallucinates an edge (*Every yumpus is a rempus*) after getting stuck in a dead end. COCONUT ($k=1$) outputs a path that ends with an irrelevant node. COCONUT ($k=2$) answers the question correctly.

5.2 RESULTS

Figure 5 shows a comparative analysis of different reasoning method on ProsQA. As more continuous thoughts are used, both the answer accuracy and the probability of predicting the correct path gradually increase. The rate of hallucination, which often occurs when the model makes a wrong move and gets stuck in a dead end, also decreased. A case study is shown in Figure 6, where *CoT* hallucinates an inexistent edge, COCONUT ($k = 1$) leads to a wrong target, but COCONUT ($k = 2$) successfully solves the problem. In this example, the model cannot accurately determine which edge to choose at the earlier step. However, as latent reasoning can avoid making a hard choice upfront, the model can progressively eliminate incorrect options in subsequent steps and achieves higher accuracy at the end of reasoning. We show more evidence and details of this reasoning process in Section 5.3.

The comparison between *CoT* and COCONUT ($k = 0$) reveals another interesting fact: even when COCONUT is forced to generate a complete reasoning chain, the accuracy of the answers is still higher than *CoT*. The generated reasoning paths are also more accurate with less hallucination. From this, we can infer that the training method of mixing different stages improves the model’s ability to plan ahead. The training objective of *CoT* always concentrates on the generation of the immediate next step, making the model “shortsighted”. In later stages of COCONUT training, the first few steps are excluded, allowing the model to focus more on future steps. This is similar to the principle of multi-token prediction pretraining (Gloeckle et al., 2024), which also helps improve the LLM’s ability to plan ahead. We leave more detailed analysis on this phenomenon to future work.

5.3 INTERPRETING THE LATENT TREE SEARCH

Actually, we can infer the reasoning encoded in latent thoughts from the model’s subsequent outputs (Figure 7). For instance, if we force the model to switch back to the language space after one latent thought (by placing <eot>), under greedy decoding, the model predicts “every lempus is a scrompus” as the next step. In this case, we have to assume that the latent thought encodes “Alex is a lempus”. Furthermore, if we do not use greedy decoding, and instead output the probability distribution of the token predicted by the model at the position of “lempus”, we can obtain a distribution of the reasoning step encoded by the latent thought. Similarly, we can get the probability of nodes in the second reasoning steps (Figure 7, right). This probability distribution can be viewed as the model’s implicit value function, that is, the estimated potential of a node to lead to the correct target. As shown in the figure, “lempus”, “zhorpus”, “grimpus”, and “sterpus” have a probability of 0.33, 0.16, 0.32, and 0.01, respectively. This indicates that in the first continuous thought, the model has mostly ruled out “sterpus” as an option but is still unable to determine which of the remaining three is the correct choice.

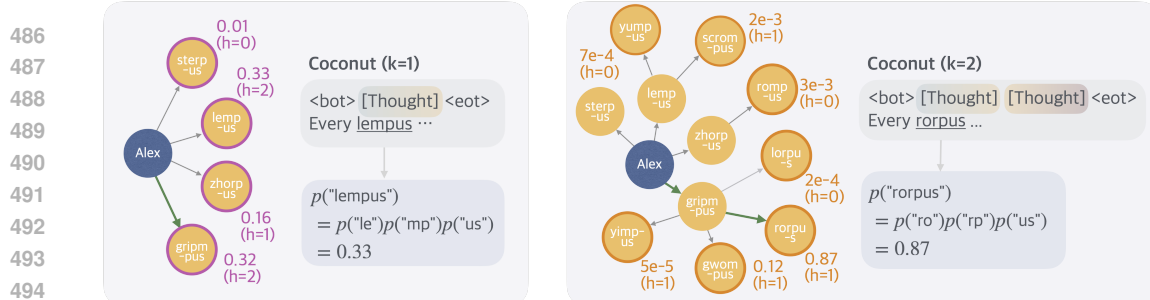


Figure 7: An illustration of the latent search trees. The example is consistent with Figure 6. The height of a node (denoted as h in the figure) is defined as the longest distance to any leaf nodes in the graph. We calculate the probability of the first concept predicted by the model following latent thoughts (e.g., “lempus” in the left figure). It is calculated as the multiplication of the probability of all tokens within the concept conditioned on previous context (omitted in the figure for brevity). This metric can be interpreted as an implicit value function estimated by the model, assessing the potential of each node leading to the correct answer.

A significant difference between “sterpus” and the other three options is that it is a leaf node (see Figure 6). In contrast, the other three nodes can still be further explored, which makes them harder to evaluate. We measure the height of each node, i.e., the shortest distance to any leaf nodes in the graph, as a proxy for the room of exploration. Based on the case discussed above, a natural hypothesis is that the lower a node is, the easier it is to estimate its value accurately. Indeed, in this case, the model is confused between “griimpus” and “lempus”, both of which has a height of 2.

To validate this hypothesis, we analyze the first and second latent steps on the whole test set. We can clearly see a trend in Figure 8. Generally, the model can effectively differentiate between correct and incorrect nodes (defined by whether they lead to the correct target node) when their heights are small, i.e., assigning a small value for incorrect nodes. However, it tends to become less accurate as the node heights increase. We can conclude that the model is not capable of doing an exhaustive search to evaluate the potential of a node, but relies more on heuristic features like the heights.

Therefore, it’s intuitive to understand why more latent thoughts makes reasoning easier: *as the search tree expands, the nodes under consideration are expected to have smaller heights*. When the height is small, the model is better at distinguishing correct nodes from incorrect nodes, and is more likely to output a correct answer.

6 CONCLUSION

In this paper, we presented COCONUT, a novel paradigm for reasoning in continuous latent space, aimed to address the inherent inefficiencies associated with traditional language-based reasoning in large language models. Through extensive experimentation on various datasets, we demonstrated that COCONUT significantly enhances LLM reasoning capabilities. Notably, our detailed analysis highlighted how an unconstrained latent space allows the model to develop an effective reasoning pattern similar to BFS. We anticipate that our findings will inspire further research into latent reasoning methods, contributing to the development of more intelligent machine reasoning system.

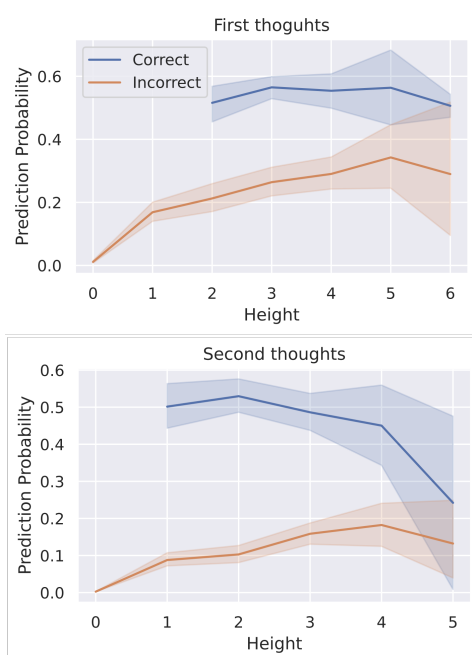


Figure 8: The correlation between prediction probability of concepts and their heights.

REFERENCES

- 540
541
542 Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Ale-
543 man, Diogo Almeida, Janko Altenschmidt, Sam Altman, Shyamal Anadkat, et al. Gpt-4 technical
544 report. *arXiv preprint arXiv:2303.08774*, 2023.
- 545 Marie Amalric and Stanislas Dehaene. A distinct cortical network for mathematical knowledge in
546 the human brain. *NeuroImage*, 189:19–31, 2019.
- 547
548 Eden Biran, Daniela Gottesman, Sohee Yang, Mor Geva, and Amir Globerson. Hopping too
549 late: Exploring the limitations of large language models on multi-hop queries. *arXiv preprint*
550 *arXiv:2406.12775*, 2024.
- 551 Wenhui Chen, Ming Yin, Max Ku, Pan Lu, Yixin Wan, Xueguang Ma, Jianyu Xu, Xinyi Wang, and
552 Tony Xia. Theoremqa: A theorem-driven question answering dataset. In *Proceedings of the 2023*
553 *Conference on Empirical Methods in Natural Language Processing*, pp. 7889–7901, 2023.
- 554
555 Karl Cobbe, Vineet Kosaraju, Mohammad Bavarian, Mark Chen, Heewoo Jun, Lukasz Kaiser,
556 Matthias Plappert, Jerry Tworek, Jacob Hilton, Reiichiro Nakano, et al. Training verifiers to
557 solve math word problems. *arXiv preprint arXiv:2110.14168*, 2021.
- 558 Google DeepMind. Ai achieves silver-medal standard solving international mathematical
559 olympiad problems, 2024. URL [https://deepmind.google/discover/blog/](https://deepmind.google/discover/blog/ai-solves-imo-problems-at-silver-medal-level/)
560 [ai-solves-imo-problems-at-silver-medal-level/](https://deepmind.google/discover/blog/ai-solves-imo-problems-at-silver-medal-level/).
- 561
562 Yuntian Deng, Kiran Prasad, Roland Fernandez, Paul Smolensky, Vishrav Chaudhary, and Stu-
563 art Shieber. Implicit chain of thought reasoning via knowledge distillation. *arXiv preprint*
564 *arXiv:2311.01460*, 2023.
- 565
566 Yuntian Deng, Yejin Choi, and Stuart Shieber. From explicit cot to implicit cot: Learning to inter-
567 nalize cot step by step. *arXiv preprint arXiv:2405.14838*, 2024.
- 568 Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha
569 Letman, Akhil Mathur, Alan Schelten, Amy Yang, Angela Fan, et al. The llama 3 herd of models.
570 *arXiv preprint arXiv:2407.21783*, 2024.
- 571
572 Evelina Fedorenko, Michael K Behr, and Nancy Kanwisher. Functional specificity for high-level
573 linguistic processing in the human brain. *Proceedings of the National Academy of Sciences*, 108
574 (39):16428–16433, 2011.
- 575
576 Evelina Fedorenko, Steven T Piantadosi, and Edward AF Gibson. Language is primarily a tool for
577 communication rather than thought. *Nature*, 630(8017):575–586, 2024.
- 578 Guhao Feng, Bohang Zhang, Yuntian Gu, Haotian Ye, Di He, and Liwei Wang. Towards revealing
579 the mystery behind chain of thought: a theoretical perspective. *Advances in Neural Information*
580 *Processing Systems*, 36, 2023.
- 581
582 Kanishk Gandhi, Denise Lee, Gabriel Grand, Muxin Liu, Winson Cheng, Archit Sharma, and
583 Noah D Goodman. Stream of search (sos): Learning to search in language. *arXiv preprint*
584 *arXiv:2404.03683*, 2024.
- 585 Fabian Gloeckle, Badr Youbi Idrissi, Baptiste Rozière, David Lopez-Paz, and Gabriel Syn-
586 naeve. Better & faster large language models via multi-token prediction. *arXiv preprint*
587 *arXiv:2404.19737*, 2024.
- 588
589 Sachin Goyal, Ziwei Ji, Ankit Singh Rawat, Aditya Krishna Menon, Sanjiv Kumar, and Vaishnavh
590 Nagarajan. Think before you speak: Training language models with pause tokens. *arXiv preprint*
591 *arXiv:2310.02226*, 2023.
- 592
593 Shibo Hao, Yi Gu, Haodi Ma, Joshua Jiahua Hong, Zhen Wang, Daisy Zhe Wang, and Zhiting Hu.
Reasoning with language model is planning with world model. *arXiv preprint arXiv:2305.14992*,
2023.

- 594 Alex Havrilla, Yuqing Du, Sharath Chandra Raparthy, Christoforos Nalmpantis, Jane Dwivedi-Yu,
595 Maksym Zhuravinskyi, Eric Hambro, Sainbayar Sukhbaatar, and Roberta Raileanu. Teaching
596 large language models to reason with reinforcement learning. *arXiv preprint arXiv:2403.04642*,
597 2024.
- 598 Tushar Khot, Harsh Trivedi, Matthew Finlayson, Yao Fu, Kyle Richardson, Peter Clark, and Ashish
599 Sabharwal. Decomposed prompting: A modular approach for solving complex tasks. *arXiv*
600 *preprint arXiv:2210.02406*, 2022.
- 601 Yann LeCun. A path towards autonomous machine intelligence version 0.9. 2, 2022-06-27. *Open*
602 *Review*, 62(1):1–62, 2022.
- 603 Lucas Lehnert, Sainbayar Sukhbaatar, Paul Mcvay, Michael Rabbat, and Yuandong Tian. Be-
604 yond a*: Better planning with transformers via search dynamics bootstrapping. *arXiv preprint*
605 *arXiv:2402.14083*, 2024.
- 606 Zhiyuan Li, Hong Liu, Denny Zhou, and Tengyu Ma. Chain of thought empowers transformers to
607 solve inherently serial problems. *arXiv preprint arXiv:2402.12875*, 2024.
- 608 Aman Madaan and Amir Yazdanbakhsh. Text and patterns: For effective chain of thought, it takes
609 two to tango. *arXiv preprint arXiv:2209.07686*, 2022.
- 610 William Merrill and Ashish Sabharwal. The expressive power of transformers with chain of
611 thought. *arXiv preprint arXiv:2310.07923*, 2023.
- 612 Martin M Monti, Daniel N Osherson, Michael J Martinez, and Lawrence M Parsons. Functional
613 neuroanatomy of deductive inference: a language-independent distributed network. *Neuroimage*,
614 37(3):1005–1016, 2007.
- 615 Martin M Monti, Lawrence M Parsons, and Daniel N Osherson. The boundaries of language and
616 thought in deductive inference. *Proceedings of the National Academy of Sciences*, 106(30):
617 12554–12559, 2009.
- 618 Martin M Monti, Lawrence M Parsons, and Daniel N Osherson. Thought beyond language: neural
619 dissociation of algebra and natural language. *Psychological science*, 23(8):914–922, 2012.
- 620 Jacob Pfau, William Merrill, and Samuel R Bowman. Let’s think dot by dot: Hidden computation
621 in transformer language models. *arXiv preprint arXiv:2404.15758*, 2024.
- 622 Alec Radford, Jeffrey Wu, Rewon Child, David Luan, Dario Amodei, Ilya Sutskever, et al. Language
623 models are unsupervised multitask learners. *OpenAI blog*, 1(8):9, 2019.
- 624 Abulhair Saparov and He He. Language models are greedy reasoners: A systematic formal analysis
625 of chain-of-thought. *arXiv preprint arXiv:2210.01240*, 2022.
- 626 Zhihong Shao, Peiyi Wang, Qihao Zhu, Runxin Xu, Junxiao Song, Mingchuan Zhang, YK Li,
627 Yu Wu, and Daya Guo. Deepseekmath: Pushing the limits of mathematical reasoning in open
628 language models. *arXiv preprint arXiv:2402.03300*, 2024.
- 629 Miles Turpin, Julian Michael, Ethan Perez, and Samuel Bowman. Language models don’t always
630 say what they think: unfaithful explanations in chain-of-thought prompting. *Advances in Neural*
631 *Information Processing Systems*, 36, 2024.
- 632 Boshi Wang, Sewon Min, Xiang Deng, Jiaming Shen, You Wu, Luke Zettlemoyer, and Huan Sun.
633 Towards understanding chain-of-thought prompting: An empirical study of what matters. *arXiv*
634 *preprint arXiv:2212.10001*, 2022.
- 635 Peiyi Wang, Lei Li, Zhihong Shao, Runxin Xu, Damai Dai, Yifei Li, Deli Chen, Yu Wu, and Zhifang
636 Sui. Math-shepherd: Verify and reinforce llms step-by-step without human annotations. In *Pro-*
637 *ceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume*
638 *1: Long Papers)*, pp. 9426–9439, 2024.
- 639 Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Fei Xia, Ed Chi, Quoc V Le, Denny
640 Zhou, et al. Chain-of-thought prompting elicits reasoning in large language models. *Advances in*
641 *neural information processing systems*, 35:24824–24837, 2022.

648 Yuxi Xie, Kenji Kawaguchi, Yiran Zhao, James Xu Zhao, Min-Yen Kan, Junxian He, and Michael
649 Xie. Self-evaluation guided beam search for reasoning. *Advances in Neural Information Process-*
650 *ing Systems*, 36, 2023.

651 Sohee Yang, Elena Gribovskaya, Nora Kassner, Mor Geva, and Sebastian Riedel. Do large language
652 models latently perform multi-hop reasoning? *arXiv preprint arXiv:2402.16837*, 2024.

653 Shunyu Yao, Dian Yu, Jeffrey Zhao, Izhak Shafran, Tom Griffiths, Yuan Cao, and Karthik
654 Narasimhan. Tree of thoughts: Deliberate problem solving with large language models. *Ad-*
655 *vances in Neural Information Processing Systems*, 36, 2023.

656 Fangxu Yu, Lai Jiang, Haoqiang Kang, Shibo Hao, and Lianhui Qin. Flow of reasoning: Efficient
657 training of llm policy with divergent thinking. *arXiv preprint arXiv:2406.05673*, 2024a.

658 Longhui Yu, Weisen Jiang, Han Shi, Jincheng Yu, Zhengying Liu, Yu Zhang, James T Kwok, Zhen-
659 guo Li, Adrian Weller, and Weiyang Liu. Metamath: Bootstrap your own mathematical questions
660 for large language models. *arXiv preprint arXiv:2309.12284*, 2023.

661 Ping Yu, Jing Xu, Jason Weston, and Ilia Kulikov. Distilling system 2 into system 1. *arXiv preprint*
662 *arXiv:2407.06023*, 2024b.

663 Xiang Yue, Xingwei Qu, Ge Zhang, Yao Fu, Wenhao Huang, Huan Sun, Yu Su, and Wenhui Chen.
664 Mammoth: Building math generalist models through hybrid instruction tuning. *arXiv preprint*
665 *arXiv:2309.05653*, 2023.

666 Eric Zelikman, Georges Harik, Yijia Shao, Varuna Jayasiri, Nick Haber, and Noah D Goodman.
667 Quiet-star: Language models can teach themselves to think before speaking. *arXiv preprint*
668 *arXiv:2403.09629*, 2024.

669 Denny Zhou, Nathanael Schärli, Le Hou, Jason Wei, Nathan Scales, Xuezhi Wang, Dale Schuur-
670 mans, Claire Cui, Olivier Bousquet, Quoc Le, et al. Least-to-most prompting enables complex
671 reasoning in large language models. *arXiv preprint arXiv:2205.10625*, 2022.

672
673
674
675
676
677
678
679
680
681
682
683
684
685
686
687
688
689
690
691
692
693
694
695
696
697
698
699
700
701

# Nodes	# Edges	Len of Shortest Path	# Shortest Paths
23.0	36.0	3.8	1.6

Table 2: Statistics of the graph structure in ProsQA.

Dataset	Training	Validation	Test
GSM8k	385,620	500	1319
ProntoQA	9,000	200	800
GSM8k	17,886	300	500

Table 3: Statistics of the datasets.

A DATASETS

A.1 CONSTRUCTION OF PROSQA

To construct the dataset, we need to define a set of entities (typical names like “Alex”, “Jack”, etc.) and a set of concepts (fictional words like “lorpus”, “rorpus”, etc., following Saparov & He (2022)).

The desired problem form is “Is [Entity] a [Concept A] or [Concept B]?”. Assume the correct answer is [Concept A], we will need to construct a graph, so that we can find a path between [Entity] and [Concept A], and make sure [Entity] and [Concept B] are not connected.

The overall idea to build the DAG is to gradually add more nodes. Every time a new node comes in, we randomly add edges from existing nodes to the new node. We first sample the in-degree following a Poisson distribution with a mean equal to 1.5, then sample the parents for this node. In this process, we need to make sure that any entity or concept cannot be the ancestor of both [Concept A] and [Concept B], in order to make a valid binary choice problem. Besides, we want to keep the family of [Concept A] and [Concept B] of similar sizes, otherwise the model may learn shortcuts.

Therefore, we implement a graph construction pipeline as follows: First, we initialize two nodes with labels 1 and 2. Then, for each new node, there is a probability p ($p = 0.35$) that it can only accept edges from nodes with label 1; and another probability p ($p = 0.35$) that it can only accept edges from nodes with label 2; otherwise the node can accept edges from any nodes. After sampling the incoming edges for the node, it will be assigned a label: 1 if all the parent nodes have label 1; 2 if all the parent nodes have label 2; 3 if there are both parent nodes with label 1 and 2; 0 if there are no parent nodes or all parent nodes are labeled 0.

All nodes without parents will be assigned an entity name, while others are given a concept names. These form the known conditions. To get the question, we use the first node as the [Entity], a node labeled with 1 as [Concept A], a node labeled with 2 as [Concept B]. The construction will ensure there is always a path from [Entity] to [Concept A] but not [Concept B]. We will find the [Concept A] and [Concept B] that makes the reasoning chain relatively long. Note that after rendering the graph into natural language, we will permute the position of [Concept A] and [Concept B] randomly. Given the symmetry of label 1 and 2, there is no risk for shortcut in the position of choice.

The statistics of the resulting dataset is listed in Table 2

A.2 STATISTICS

We show the size of all datasets in Table 3.