REASONING TO LEARN FROM LATENT THOUGHTS

Anonymous authors

000

001 002 003

004

006

008

009

010

011

012

013

014

016

017

018

019

021

024

025

026

028 029

031

033

034

037

040

041

042

043

044

046

047

048

051

052

Paper under double-blind review

ABSTRACT

Compute scaling for language model (LM) pretraining has outpaced the growth of human-written texts, leading to concerns that data will become the bottleneck to LM scaling. To continue scaling pretraining in this data-constrained regime, we propose that explicitly modeling and inferring the *latent thoughts* that underlie the text generation process can significantly improve pretraining data efficiency. Intuitively, our approach views web text as the compressed final outcome of a verbose human thought process and that the latent thoughts contain important contextual knowledge and reasoning steps that are critical to data-efficient learning. We empirically demonstrate the effectiveness of our approach through data-constrained continued pretraining for math. We first show that synthetic data approaches to inferring latent thoughts significantly improve data efficiency, outperforming training on the same amount of raw data. Furthermore, we demonstrate latent thought inference without a strong teacher, where an LM bootstraps its own performance by using an EM algorithm to iteratively improve the capability of the trained LM and the quality of thought-augmented pretraining data. We show that a 1B LM can bootstrap its performance across at least three iterations and significantly outperform baselines trained on raw data, with increasing gains from additional inference compute when performing the E-step. The gains from inference scaling and EM iterations suggest new opportunities for scaling data-constrained pretraining.

1 Introduction

Human-written text is the culmination of an underlying thought process—when we write, there is often an internal dialogue that clarifies or even determines the written word. However, modern language models (LMs) (Radford et al., 2019; Brown et al., 2020; OpenAI, 2023; Dubey et al., 2024) are pretrained directly on the final results of this process in a highly *compressed* form (e.g., research papers). This may explain why LMs struggle with data efficiency and require almost the entire human-written web to learn (Kaplan et al., 2020; Hoffmann et al., 2022). Since the rate of growth in pretraining compute is far greater than that of the web itself (Villalobos et al., 2022; Muennighoff et al., 2024), we may soon enter a data-constrained regime, motivating data efficiency approaches to extract more capabilities from limited web data.

In contrast to LMs, humans learn very efficiently from the same compressed text, which suggests the possibility of significantly improving data-efficient pretraining. In this work, we focus on *how* we learn as one potential cause for this gap. For example, when we read a research paper, we analyze specific claims, integrate them with prior knowledge, and attempt to "decompress" the author's original thought process. In other words, we use reasoning in service of learning, to infer the internal dialogue that undergirds the observed text. We refer to this procedure—augmenting the observed data with inferred thoughts to enable more efficient learning—as **reasoning to learn**.

Inspired by this, we introduce an LM pretraining approach that implements this reasoning-to-learn paradigm to improve data efficiency (Fig. 1). Specifically, we approach language modeling from a latent variable perspective, where the observed data X depends on underlying latent thoughts Z. We train our LMs to learn from observed data X augmented with the latents Z by modeling the joint distribution p(Z,X). The main challenge is synthesizing (and learning to synthesize) Z with a latent generator $q(Z\mid X)$ (Fig. 2a). One key insight of our work is that for a natural language latent thought Z, the LM itself provides a strong prior for producing latent thoughts (via its reasoning and theory-of-mind abilities (Wei et al., 2022a)). This observation turns latent thought inference into a synthetic data generation problem and has significant practical benefits—it allows us to leverage the strong capabilities of existing LMs, share weights between the LM and the latent thought generator, and simplify training into a small modification to the standard pretraining pipeline (Fig. 2b).

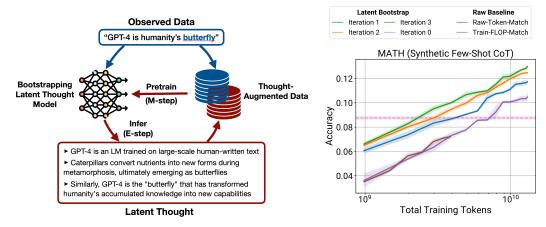


Figure 1: **Reasoning to learn**. (Left) Motivated by how humans apply deliberate thinking to learn from limited data, we train an LM to infer (or "decompress") latent thoughts underlying the highly compressed observed data. These synthesized latent thoughts augment the raw observed data during pretraining, improving the LM's data efficiency. This procedure can be iteratively applied through an EM algorithm (Fig. 4) and form a model self-improvement loop where increasingly capable LMs synthesize more effective latent thoughts, which in turn train more capable models. (Right) Our results demonstrate consistent improvement in model performance across bootstrap iterations.

We show that training a model with latent thoughts enables it to produce higher-quality latent thoughts, allowing a model to bootstrap its "reasoning to learn" ability with only a small amount of initial supervision. We demonstrate this through a simple Expectation-Maximization based approach which we refer to as **Boo**tstrapping Latent Thoughts (**BoLT**) that enables an iterative improvement of the latent thought generator (Fig. 4). Importantly, we show that BoLT can take advantage of additional inference compute to further improve data efficiency. In particular, the E-step in BoLT makes use of a Monte-Carlo estimator that serves as a non-parametric "policy improvement operator", where the approximate posterior $q(Z \mid X)$ approaches the true posterior as the number of samples increases. We find in our experiments that BoLT is able to take advantage of additional samples (at least four) to improve its data efficiency and bootstrap its performance for at least three iterations, opening the possibility of new ways of scaling pretraining data efficiency.

We validate the effectiveness of our approach in improving model capabilities in data-constrained setups. As a testbed, our experiments continually pretrain a TinyLlama (Zhang et al., 2024) model on a limited amount of data from a reasoning-intensive corpus FineMath (Lozhkov et al., 2024).

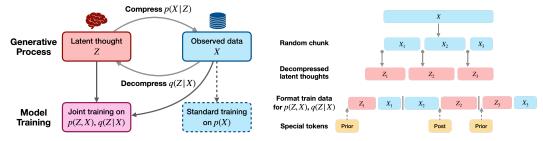
- First, we demonstrate the potential of training LMs with data augmented with latent thoughts: when using GPT-4o-mini (Hurst et al., 2024) to synthesize latent thoughts, the trained LM achieves 25.4% on MATH (Hendrycks et al., 2021a), significantly outperforming LMs trained on raw data (5.74%) or synthetic Chain-of-Thought (Wei et al., 2022a) style paraphrases (19.4%) (Fig. 5).
- Furthermore, we show that our BoLT algorithm enables an LM to bootstrap its performance on limited data. When trained on a fixed amount of data (Fig. 6), we find LMs trained with self-generated latents improve across iterations on both likelihood metrics and MATH evaluations (Fig. 7). Crucially, these gains require no task-specific data and stem purely from improved latent thought quality across bootstrap iterations.

Altogether, our results suggest that the powerful reasoning primitives of LMs may be leveraged to extract more capabilities from limited, task-agnostic data during pretraining.

2 Related Work

Here we briefly review the most relevant works and include an extended discussion in Appx. B.

Synthetic data for training LMs Recent work has demonstrated the benefits of synthetic data to improve training efficiency, typically by prompting frontier models to generate data for model training. These approaches have yielded competitive performance at small data and model scales (Eldan and Li, 2023; Gunasekar et al., 2023; Abdin et al., 2024a;b). Another paradigm closer to our work is to augment raw data by paraphrasing (Maini et al., 2024) or describing relationships among extracted entities (Yang et al., 2025), which uses real data to ground the synthetic data generation. Our work differs in two key aspects. First, we take a latent variable perspective that augments (rather than replaces) each text chunk X with latent thoughts Z, which we show is crucial to downstream



(a) Learning by decompression

(b) Training with next-token predictions

Figure 2: **Reasoning to learn with latent thought models** (a) The latent thought model is trained to "decompress" plausible human thoughts underlying the observed data (i.e., q(Z|X)) and to utilize the latent thoughts in learning more efficiently from the data (i.e., p(Z,X)). (b) The latent thought is modeled for each chunk of text in an autoregressive manner and in the same discrete text space. Given paired data $\{(Z_n, X_n)\}_{n=1}^N$, we use standard next-token prediction to train a single LM as both the p(Z,X) and q(Z|X), by randomly placing Z_n before or after X_n in the sequence. This strategy allows for minimal modifications to the standard LM pretraining pipeline.

performance in Sec. 5.2. Additionally, our work studies *bootstrapping* of the data generator rather than the distillation settings considered in most previous works.

Learning to reason using pretraining data Closer to our setting, some recent works attempt to enhance LM reasoning capabilities using pretraining data without explicit external rewards. Geiping et al. (2025) pretrain a looped transformer (Dehghani et al., 2019; Giannou et al., 2023) on general web text, using continuous hidden states to model a thought process. Zelikman et al. (2024) propose a domain-agnostic post-training method that uses reinforcement learning on pretraining data to learn "thought tokens" that improve reasoning. Our work is distinct in the goal of leveraging reasoning to improve pretraining data efficiency, which leads us to consider scalable, synthetic data—based approaches. This leads to a number of benefits, including a simple training method and embarassingly parallel latent generation. Our work provides a complementary perspective on how latent variable and synthetic data perspectives can enable both "learning to reason" and "reasoning to learn".

3 REASONING TO LEARN WITH LATENT THOUGHT MODELS

In this section, we introduce our "reasoning to learn" approach to pretraining LMs. Our key idea is to model the latent thoughts underlying the human data generation process, and train a latent thought model both to *reason* about the latent thoughts underlying pretraining data and to *learn* from the data based on the synthesized latent thoughts.

3.1 Latent Thought Models

Learning by decompression Human-written data on the web is a *compressed* representation of an underlying thought process. For example, when Geoffrey Hinton wrote "GPT-4 is humanity's butterfly", he drew upon his knowledge of how GPT-4 was trained on a large corpus of human-generated data and distilled human intelligence, and analogized this process to a butterfly's metamorphosis. Understanding the context and reasoning steps behind the observed data can facilitate deeper understanding of the text they generate, but such reasoning traces are usually not explicitly presented in our pretraining data. Our hypothesis is twofold: augmenting pretraining data with underlying human thoughts can significantly improve its learning efficiency, and that this ability to infer latent thoughts can be learned and improved after being trained with these thoughts during pretraining.

Formalizing latent thought models We formalize this from the perspective of latent variable models, as illustrated in Fig. 2a. We model the human thoughts underlying the observed data X as the latent Z, and the generative process as a joint distribution $p(Z,X) = p(Z)p(X \mid Z)$. Instead of training an LM to directly model the observed data p(X), we train it both to "decompress" the latent thoughts from the observed data (i.e., approximate posterior $q(Z \mid X)$) and to learn from the data using the synthesized latent thoughts (i.e., p(Z,X)); we call this LM a latent thought model. The latent thought Z is modeled in the same discrete text space as the observed data X (just as human thoughts can often be expressed in natural language). By augmenting the data with latent thoughts, the LM's learning process resembles the human process of reasoning to learn, where we deliberately think through the data to better absorb it. Importantly, we view latent thoughts as being encoded in natural language. This allows us to initialize $q(Z \mid X)$ using supervision from an existing model prompted to infer latent reasoning and background context, as well as to jointly model all conditional

You are provided with a pair of web document prefix and suffix. Your task is to insert latent thoughts between them underlying the creation of the suffix conditioned on the prefix. The latent thoughts should include: the missing background knowledge and the reasoning traces underlying each claim (especially, step-by-step derivations or logical reasoning). [..omit..]

(a) Prompt for GPT-4o-mini to synthesize latent thoughts.

```
<StartOfLatent><Prior> [..omit..] If the data is not centered, the computed variances
along the axes will be skewed, leading to misleading results in the identification of PCs.
Uncentered data would cause the variance to reflect the mean the data rather than
[..omit..] <EndOfLatent>
## The first step, data decentralization
 ..omit..] We directly decentralize the data (that is, the mean value of the data is at
the far point). If the data is not decentralized, we cannot find the optimal
dimensionality reduction.
<StartOfLatent><Prior> To find the optimal axis for PCA, we utilize the covariance matrix
of the decentralized data. The eigenvalues of this matrix indicate the amount of variance
captured by each PC, while the corresponding eigenvectors provide the directions. PC1 is
the eigenvector associated with the largest eigenvalue, representing the direction of
maximum variance. [..omit..] <EndOfLatent>
## The second step is to find the new most marked axis
How do we find the most standard axis to achieve principal component analysis? That is,
the greater the distance between the projected point and the origin of the coordinate axis
, the better [..omit..
```

(b) Example synthetic latent thoughts for a document about PCA's mathematical derivation.

Figure 3: We use GPT-4o-mini to synthesize latent thoughts to train the initial latent thought model. The synthetic latents as shown in (b) typically contain the background knowledge and reasoning not explicitly stated in the raw data, presented in a consistent and clean form. The prompt and example are simplified for clarity; see Prompt F.1.1 and Appx. G.2 for the full prompt and additional examples.

distributions for Z and X using a single, autoregressive language model. We describe the training and inference processes for latent thought models in the following.

3.2 Training with Synthetic Latent Thoughts

We adopt an autoregressive generative model of latent thoughts that is compatible with standard language modeling. Given a document X, we randomly chunk it into segments $\{X_n\}_{n=1}^N$ at the sentence boundaries and aim to infer the latent thought Z_n underlying each chunk X_n conditioned on the previous context (see Fig. 2b, top).

Synthesizing latent thoughts for training Human-generated data like internet text does not naturally come with underlying latent thoughts. Therefore, we need to synthesize the latent thought Z from some surrogate $\tilde{q}(Z\mid X)$ of the true posterior to augment the observed data X for training the latent thought model. $\tilde{q}(Z\mid X)$ can either be instantiated as a frontier model, or the approximate posterior model $q(Z\mid X)$ itself for bootstrapping as in our EM algorithm. For example, we can prompt GPT-40-mini (Hurst et al., 2024) to synthesize the latent thoughts by inferring missing reasoning steps or background knowledge (see Fig. 3 for the prompt and examples).

Training latent thought models We develop a simple method to train both a joint model p(Z,X) and an approximate posterior model $q(Z \mid X)$ with only minor modifications to the standard LM training pipeline (Fig. 2b). Since both Z and X are presented in the same discrete text space, we train models with standard next-token predictions. Given the synthetic latent thoughts paired with the observed data $\{(Z_n, X_n)\}_{n=1}^N$, we train both the joint and the posterior as the same model by formatting the data as conditional maximum likelihood estimation: we place Z_n as the suffix after X_n in the sequence to train the approximate posterior $q(Z_n \mid X_n)$, and place Z_n as the prefix before X_n to train the joint $p(Z_n, X_n)$ (see the bottom of Fig. 2b). We format the data in these two modes with a random coin flip and use two special tokens $\{Prior\}$ and $\{Post\}$ to differentiate them. All latents are wrapped within the special $\{StartofLatent\}$ and $\{EndofLatent\}$ tokens to differentiate them from the raw observed data (Fig. 3b). The formatted data can be directly fed into the standard LM pretraining pipeline to train both the joint and the posterior with next-token predictions.

3.3 PREDICTION WITH LATENT THOUGHTS

Since models have been pretrained to utilize latent thoughts to predict the following text, they can also use them to perform CoT reasoning for problem solving during downstream evaluation.

Figure 4: **Bootstrapping latent thoughts (BoLT)** in an iterative Expectation-Maximization algorithm. In the E-step, we use Monte Carlo sampling as a "policy improvement operator" to obtain higher-quality latent thoughts. This boosts learning efficiency in the M-step, enabling the training of more capable LMs that synthesize better latent thoughts.

Reasoning with CoT in the latent space During CoT prompting (Nye et al., 2021; Wei et al., 2022a), it is common to augment few-shot examples with explicit reasoning chains. For standard LMs, these thoughts are simply part of the text prompt. However, in the case of latent thought models, we now have the more natural option of putting the CoT prompts into the latent space Z. We implement this by changing the formatting of the few-shot examples: given a question Q and its answer A, we wrap the CoT within the special tokens (<StartofLatent> and <EndofLatent>) and add the <Prior> prefix to indicate that this reasoning should occur in the latent space for producing the following answer (see Appx. F.2 for few-shot examples). We find that latent thought models are highly effective latent CoT reasoners at inference time (see Appx. G.1 for examples), achieving better downstream performance than when reasoning in the observed text space X.

4 BOOTSTRAPPING LATENT THOUGHT MODELS

The data efficiency of training latent thought models relies on the surrogate posterior $\tilde{q}(Z\mid X)$ used to synthesize the latent thoughts, which limits the potential of our approach in advancing frontier model capabilities. To overcome this limitation, we introduce an Expectation-Maximization (EM) algorithm called **Bo**otstrapping Latent Thoughts (**BoLT**), illustrated in Fig. 4. The key idea is to instantiate $\tilde{q}(Z\mid X)$ with a non-parametric, enhanced version of the model's approximate posterior $q(X\mid Z)$ via Monte Carlo sampling. This enables a self-improvement loop, where more effective latents are synthesized at the E-step, leading to more capable models with improved learning efficiency at the M-step, which in turn improve latent quality at the next E-step.

4.1 EXPECTATION-MAXIMIZATION WITH MONTE CARLO SAMPLING

At each iteration t in EM, we have access to the current latent thought model \mathcal{M}_t . This model parameterizes both the approximate posterior $q(Z \mid X; \mathcal{M}_t)$ and the joint $p(Z, X; \mathcal{M}_t)$ through the same data formatting and special token usage used during training (Sec. 3.2). A naive instantiation of EM would be to alternate between sampling latents Z from the current posterior $q(Z \mid X; \mathcal{M}_t)$ in the E-step and training a model with the sampled latents Z in the M-step. However, it is unclear that training a model on its self-generated latents can improve it. We address this gap with Monte Carlo sampling, which induces a surrogate posterior $\tilde{q}(Z \mid X; \mathcal{M}_t)$ that is provably better than the current one to serve as a "policy improvement operator". We detail the procedure below.

E-step: Synthesizing better latents with Monte Carlo sampling At the E-step, we sample K latents $\{Z^{(k)}\}_{k=1}^K$ from the current posterior $q(Z \mid X; \mathcal{M}_t)$ and compute their importance weights $w^{(k)} = \frac{p(Z^{(k)}, X; \mathcal{M}_t)}{q(Z^{(k)} \mid X; \mathcal{M}_t)}$. Intuitively, $w^{(k)}$ upweights a latent that is both simple (high $p(Z^{(k)}; \mathcal{M}_t)$) and predictive of the data (high $p(X \mid Z^{(k)}; \mathcal{M}_t)$), and downweights too obvious ones $(q(Z^{(k)} \mid X; \mathcal{M}_t))$ in the denominator). Then, we resample one latent from the categorical distribution proportional to the importance weights, i.e., $j \sim \operatorname{Cat}(w^{(k)})$, and use $Z^{(j)}$ as the final latent. This procedure is known to sample from a posterior $\tilde{q}(Z \mid X; \mathcal{M}_t)$ that achieves the importance-weighted evidence lower bound (ELBO), which is tighter than the naive ELBO (Burda et al., 2015; Cremer et al., 2017):

$$\log p(X) \ge \underset{Z \sim \tilde{q}}{\mathbb{E}} \left[\log \frac{p(Z, X)}{\tilde{q}(Z \mid X)} \right] \ge \underset{Z^{(k)} \sim q}{\mathbb{E}} \left[\log \frac{1}{K} \sum_{k=1}^{K} \frac{p(Z^{(k)}, X)}{q(Z^{(k)} \mid X)} \right] \ge \underset{Z \sim q}{\mathbb{E}} \left[\log \frac{p(Z, X)}{q(Z \mid X)} \right] \tag{1}$$

The bound is tighter with more samples K and \tilde{q} approximates the true posterior when $K \to \infty$. Crucially, this procedure offers a potential way to improve the training efficiency of LMs by scaling up the inference compute to select a better latent from more samples.

M-step: Training the model with bootstrapped latents At the M-step, we use the latents synthesized by the current model \mathcal{M}_t and pair them with the raw data corpus \mathcal{X} to train the next model \mathcal{M}_{t+1} , following the same training procedure in Sec. 3.2. Note that the new posterior $q(Z \mid X; \mathcal{M}_{t+1})$

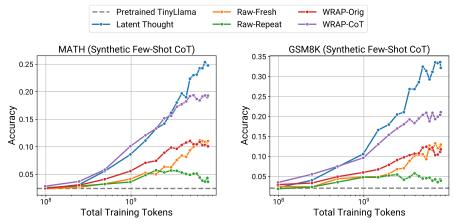


Figure 5: Training on raw data augmented with synthetic latent thoughts significantly boosts LM data efficiency. We study a data-constrained setup where a TinyLlama-1.1B model is continually pretrained on 480M FineMath tokens with an 8B training token budget. GPT-4o-mini generates synthetic latent thoughts (as a proxy for human ground-truth thoughts) for the corpus. This augmented data ("Latent Thought") leads to substantial gains over both training on 8B unique raw tokens ("Raw-Fresh") and paraphrased data with synthetic reasoning traces ("WRAP-CoT"). The synthetic token budget for all synthetic data methods is fixed to \sim 1.1B tokens. The advantages of our method are even more pronounced under prompt variations (see Fig. D.2) and when normalized by the effective raw tokens seen by each method (Fig. D.3).

is trained to approximate the improved posterior $\tilde{q}(Z \mid X; \mathcal{M}_t)$, and the new joint $p(Z, X; \mathcal{M}_{t+1})$ is trained by augmenting the data with bootstrapped latents from $\tilde{q}(Z \mid X; \mathcal{M}_t)$.

Full algorithm Our full approach is described in Algorithm 1. Since the initial base LM has not been trained as a latent thought model, we first train it on a small corpus of warmstart data with synthetic latent thoughts (synthesized by, e.g., GPT-40-mini), and then apply our EM algorithm on a much larger corpus. The warmstart \rightarrow iterative EM loop resembles the supervised finetuning \rightarrow reinforcement learning pipeline (Ouyang et al., 2022) that is widely adopted in LM post-training, though our approach operates purely on task-agnostic pretraining data instead of supervised data.

5 SYNTHETIC LATENT THOUGHTS IMPROVE LM DATA EFFICIENCY

In this section, we first demonstrate that the data efficiency of LMs can be significantly improved through joint training on observed data and latent thoughts. We show this by using a capable LM (GPT-4o-mini) as the surrogate posterior $\tilde{q}(X\,|\,Z)$ to synthesize latent thoughts to augment a fixed corpus, and compare its performance against several baseline approaches.

5.1 EXPERIMENTAL SETUP

Continued pretraining on a reasoning-intensive corpus We conduct continued pretraining (CPT) of TinyLlama-1.1B (Zhang et al., 2024) on a reasoning-intensive corpus FineMath-4+ (Lozhkov et al., 2024). While reasoning to learn is domain-agnostic, we choose this setup to obtain meaningful downstream accuracy readouts under an academic budget. We choose TinyLlama-1.1B because it has not been specifically trained on mathematical reasoning data, enabling clean comparisons. To assess the data efficiency of various methods, we adopt a *data-constrained* setup by fixing the total number of unique "raw" tokens for CPT to 480M. We fix a total CPT compute budget of 8B tokens, with methods training on additional synthetic data and/or for multiple epochs on raw data.

Synthetic generation of latent thoughts We use GPT-40-mini as the surrogate posterior $\tilde{q}(Z \mid X)$ and generate latent thoughts over the CPT corpus. We prompt the model with Prompt F.1.1 and temperature 0.7. We split each document into chunks of sentences, where each chunk contains a random number of sentences following a Poisson distribution with a mean of 8, truncated between 1 and 20. We generate one latent for each chunk using the previous L=3 chunks as the prefix context. Overall, the token ratio between synthetic and raw tokens is 2.3:1 (\sim 1.1B synthetic tokens), and under the 8B token compute budget, we CPT on the thought-augmented data for \sim 5 epochs.

Baselines We compare our approach with several natural baselines in the data-constrained regime using a combination of data repetition and synthetic data generation: 1) **Raw-Repeat**: We train on the raw corpus for \sim 16 epochs with early stopping; 2) **Raw-Fresh**: As an anticipated upper bound, we train on 8B *unique* tokens without repetition; 3) **WRAP-Orig**: WRAP (Maini et al., 2024) rephrases the data in diverse styles, and we prompt GPT-40-mini to rephrase each document with their prompts.

Table 1: **Downstream performance of different methods and ablation studies**. For our method, it is crucial to embed the synthetic thoughts in a latent space Z separate from the observed raw text X, and to utilize Z for CoT reasoning during downstream evaluation.

Data	MATH	GSM8K	MMLU-STEM
Raw-Repeat	5.74	5.76	27.31
Raw-Fresh	11.18	13.27	30.63
WRAP-Orig	11.06	12.43	31.40
WRAP-CoT	19.36	21.08	34.51
Latent Thought (Ours)	25.38	33.59	35.87
- mixing latents in raw text space during training	22.38	20.17	33.33
- using CoT in raw text space during eval	20.34	22.97	31.78

4) **WRAP-CoT**: To probe whether the gains of our approach arise from simply including synthetic reasoning traces in the training data, we develop a WRAP variant that prompts GPT-4o-mini to rephrase docs with interspersed reasoning steps. For both WRAP baselines, we generated multiple paraphrases per document based on their synthetic-raw token ratios to match the total synthetic tokens of our approach. We also tuned the synthetic-raw mixture coefficient and found entirely synthetic works best (see Figs. D.1a and D.1b). We include more baseline details in Appx. D.1.

Training & Evaluation We CPT using AdamW (Loshchilov and Hutter, 2019) with $\beta_1=0.9$, $\beta_2=0.95$, weight decay of 0.01, and a tuned learning rate of 1e-4 (see Appx. D for details). We evaluate trained LMs on the popular reasoning benchmarks MATH (Hendrycks et al., 2021a), GSM8K (Cobbe et al., 2021), and MMLU-STEM (Hendrycks et al., 2021b) with few-shot CoT prompting (Wei et al., 2022a). We use both the default CoT prompts and prompts with CoTs synthesized by GPT-40-mini, and report the latter by default as it uniformly performed better across methods (see Fig. D.2). Few-shot CoTs are placed in the latent space Z for our method, and in the raw text space X for baselines. We refer the reader to Appx. D.1 for further experimental details.

5.2 RESULTS

Training with synthetic latent thoughts improves data efficiency over baselines Fig. 5 and Table 1 show downstream performance during CPT for various methods. Training with synthetic latent thoughts substantially outperforms all baselines, even outperforming training on an equivalent amount of unique raw tokens ("Raw-Fresh"). While rephrasing-based synthetic data generation methods do improve over raw data baselines—particularly the variant incorporating reasoning steps ("WRAP-CoT")—they still considerably underperform our approach. Moreover, our method is more robust to prompt variations compared to baselines that suffer considerable degradation with standard CoT prompts (see Fig. D.2). Notably, since our method trains jointly on synthetic latent thoughts and raw data (at a 2.3:1 token ratio), it achieves better performance while seeing the raw data fewer times (3.3x) than the baselines. We include a comparison normalized by the effective raw tokens seen by each method in Fig. D.3, where our method demonstrates even more significant gains.

Learning and utilizing latent thoughts in a separate latent space is critical A key distinction between our approach and baselines (such as WRAP-CoT) is that we model thoughts in a latent space Z separate from the observed text X, which we hypothesize to improve performance by explicitly handling the language modeling task from the thought process. We ablate this design in Table 1. First, we demonstrate the importance of jointly modeling thought and document chunk pairs (Z, X). We test a variant of our approach that does not separate the latent thoughts from the raw text X and mixes them with randomly sampled (unpaired) text chunks at a 1:1 ratio. This variant (second-to-last row) performed comparably to WRAP-CoT but significantly worse than our full method, demonstrating that our gains primarily stem from the latent model design rather than merely from the quality of the synthetic thought data. Furthermore, we assess the effectiveness of reasoning with CoT in the latent space during downstream evaluation. For our trained latent thought models, we instead provide all few-shot CoTs in the raw text space X similar to the baselines and suppress the generation of all special latent tokens during evaluation. This test-time intervention degrades performance (last row), demonstrating the benefit of explicitly using latent thoughts during downstream evaluation.

6 LMs Can Self-Improve by Bootstrapping Latent Thoughts

We have demonstrated that training jointly on raw texts and accompanying latent thoughts synthesized by a powerful LM significantly improves data efficiency. In this section, we take a step towards LMs that can *self-improve* on limited pretraining data, by investigating whether the BoLT algorithm enables monotonic improvement across iterations. We focus on a scientific setup for understanding whether BoLT can enable model self-bootstrapping by iteratively generating higher quality latents (see Fig. 6).

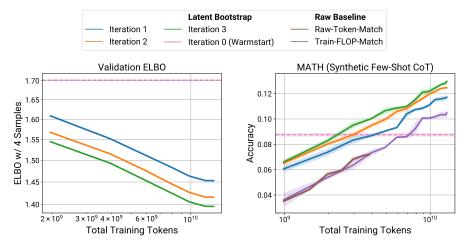
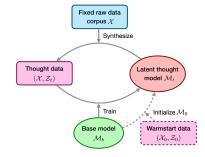


Figure 7: **BoLT bootstraps performance on a fixed corpus.** We bootstrap on a fixed corpus of 1.92B raw tokens and warmstart with 240M raw tokens, observing consistent gains in ELBO (left) and downstream MATH evaluations (see Fig. E.2 for detailed results across prompt variations). Results are over 5 training runs for bootstrapping and 3 runs for baselines to reduce variability.

Specifically, we fix the training corpus \mathcal{X} and retrain the model \mathcal{M}_t from scratch at each iteration. With a limited amount of fixed data, the final model performance at each iteration is determined by the quality of synthesized latents \mathcal{Z}_t , enabling a controlled study. In Appx. E.2, we also study a more practical scenario where LMs are continually bootstrapped from previous models (without retraining from scratch) and observe similar bootstrapping effects.



6.1 BOOTSTRAPPING ON FIXED DATA

6.1.1 EXPERIMENTAL SETUP

Figure 6: Bootstrapping on fixed data.

Bootstrapping configuration We use TinyLlama-1.1B as the base model \mathcal{M}_b and fix the training corpus as 1.92B raw tokens from FineMath-4+. For the warmstart data \mathcal{X}_0 , we use a separate corpus of 240M raw tokens with GPT-4o-mini synthetic latents as \mathcal{Z}_0 (\sim 550M synthesized tokens), following Sec. 5.1. We train the warmstart model \mathcal{M}_0 for 1 epoch on $(\mathcal{X}_0, \mathcal{Z}_0)$. At each bootstrap iteration t, we use the current trained model \mathcal{M}_{t-1} to instantiate the posterior $q(Z \mid X; \mathcal{M}_{t-1})$ at temperature 1.0, from which we sample latents for each text chunk (with chunks split following Sec. 5.1). We sample K=4 latents by default and select one by sampling proportional to their importance weights (Sec. 4), which we found to provide decent performance gains while maintaining a reasonable cost for synthetic data generation. The synthetic:raw token ratio is \sim 2.4:1, and we train LMs on the augmented data for 2 epochs at each iteration, corresponding to \sim 13B total training tokens.

Evaluation Due to high run-to-run variability, we require high signal-to-noise ratio evaluation metrics that differentiate model quality across bootstrap iterations. We first use **likelihood measures**: we use a holdout set of 1,000 documents from FineMath-4+ to measure test likelihoods. For latent thought models, we measure the 4-sample ELBO (Burda et al., 2015) as in Eq. (1) for a tighter bound. While the ELBO could be a loose bound on the true negative log-likelihood (NLL), it serves as a smooth diagnostic metric to demonstrate that our model self-improves predictably with respect to our chosen training metric. Second, we use **downstream metrics**: we evaluate performance on MATH and GSM8K at each bootstrap iteration with few-shot CoT prompting, following Sec. 5.1. We exclude MMLU-STEM as performance fell within the noise floor (<28%). Because few-shot prompting performance may be confounded by the LM's in-context learning ability (Dominguez-Olmedo et al., 2025), we also perform finetuning-based evaluations on each benchmark with synthetic CoT traces that measure the model performance after being finetuned on training splits. For clarity, we plot the best performance achieved during training for each run. Exact performance curves are in Fig. E.3 with similar relative comparisons. We refer readers to Appx. E.1 for further details.

6.1.2 RESULTS

BoLT enables monotonic self-improvement across multiple iterations Fig. 7 shows the model performance after each bootstrap iteration. Our results demonstrate that BoLT consistently improves performance over multiple iterations in both ELBO and MATH accuracy, where more capable models

generate higher quality latents that further enhance model capabilities. Performance improvements remain clearly observable on the smoother ELBO and NLL (Fig. E.1) metrics through the fourth iteration. Downstream MATH performance shows consistent improvement through the third iteration, with diminishing returns by the fourth. This plateau may be due to the discontinuous nature of few-shot evaluations (Wei et al., 2022b; Schaeffer et al., 2023) that can mask smaller improvements; after task-specific finetuning, gains persist through the fourth iteration (Fig. 8). We speculate that the saturation point could correlate with scale, as our preliminary smaller-scale experiments showed few-shot performance plateauing after the second iteration. Qualitatively, we observe instances where the latents at earlier iterations make mistakes that are corrected in later ones (see Example G.3.3).

Training on self-generated latents outperforms training on raw data We compare the downstream MATH performance of BoLT-trained models with those trained on raw data (Fig. 7). We include a FLOP-matched baseline trained for the same total training tokens as ours (3.4× more passes over the corpus), and a raw token-matched baseline trained on the same count of raw tokens as ours (2 epochs over the corpus). BoLT significantly outperforms both baselines, demonstrating the effectiveness of reasoning to learn with self-generated latents for improved data efficiency. Gains persist in finetuning evaluations on MATH & GSM8K (Fig. 8). Latent examples are provided in Appx. G.3.

Synthesizing latents with more Monte Carlo samples improves data quality A key component of our approach is the use of Monte Carlo sampling: we draw multiple samples and reweight them to select one, inducing an improved posterior and serving as a "policy improvement operator". We study the impact of the number of MC samples K on latent quality, synthesizing latents with K=1,2,4,8 at the first iteration and training models on data augmented with them (Fig. 9). We find that the

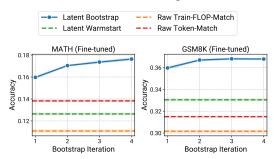


Figure 8: The performance gains of BoLT across multiple iterations and over raw data baselines remain robust in finetuning evaluations.

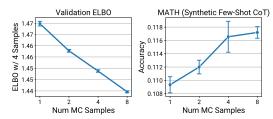


Figure 9: More Monte Carlo samples lead to improved latent quality and better trained models (further results in Fig. E.4).

ELBO and MATH accuracy improve monotonically in the number of MC samples, demonstrating a potential avenue for scaling inference compute in improving pretraining efficiency.

7 Conclusion

We have introduced reasoning to learn—a new approach to data-efficient LM pretraining by deliberately thinking through the observed data. We show that training LMs with synthetic latent thoughts significantly improves learning efficiency and downstream performance in data-constrained setups. Moreover, we instantiate an EM algorithm—Bootstrapping Latent Thoughts (BoLT)—which enables model self-improvement, where more capable models synthesize higher-quality latents that in turn enable greater learning efficiency. We extensively validate BoLT with continued pretraining on math data, and show that BoLT demonstrates steeper scaling and consistent downstream gains across multiple bootstrap iterations. Our work demonstrates the promise of explicit inference of underlying latent thoughts in improving data efficiency on task-agnostic pretraining data.

Limitations and future work Our work represents a proof-of-concept under a constrained compute budget, using a 1B parameter model and a few billion tokens of math-specific data. Future work should validate the approach at larger scales, with general-domain or even multimodal data, to further explore the potential for unlocking general reasoning capabilities beyond domain-specific reinforcement learning approaches (Jaech et al., 2024; Guo et al., 2025). We focused on a particular instantiation of our approach and did not exhaust design choices such as alternative latent structures or hyperparameter configurations. Moreover, bootstrapping on self-generated data may amplify specific biases; we observed a case where few-shot performance on GSM8K deteriorates over multiple BoLT iterations (see Fig. E.5)—such effects warrant more extensive investigation. We provide a detailed discussion in Appx. C.

REFERENCES

- Alec Radford, Jeff Wu, Rewon Child, David Luan, Dario Amodei, and Ilya Sutskever. Language models are unsupervised multitask learners. 2019.
- Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. Language models are few-shot learners. *Advances in neural information processing systems*, 33:1877–1901, 2020.
 - OpenAI. Gpt-4 technical report, 2023.
- Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha Letman, Akhil Mathur, Alan Schelten, Amy Yang, Angela Fan, et al. The llama 3 herd of models. arXiv preprint arXiv:2407.21783, 2024.
 - Jared Kaplan, Sam McCandlish, Tom Henighan, Tom B Brown, Benjamin Chess, Rewon Child, Scott Gray, Alec Radford, Jeffrey Wu, and Dario Amodei. Scaling laws for neural language models. *arXiv* preprint arXiv:2001.08361, 2020.
 - Jordan Hoffmann, Sebastian Borgeaud, Arthur Mensch, Elena Buchatskaya, Trevor Cai, Eliza Rutherford, Diego de Las Casas, Lisa Anne Hendricks, Johannes Welbl, Aidan Clark, et al. Training compute-optimal large language models. *arXiv preprint arXiv:2203.15556*, 2022.
 - Pablo Villalobos, Anson Ho, Jaime Sevilla, Tamay Besiroglu, Lennart Heim, and Marius Hobbhahn. Will we run out of data? limits of llm scaling based on human-generated data. *arXiv preprint arXiv:2211.04325*, 2022.
 - Niklas Muennighoff, Alexander Rush, Boaz Barak, Teven Le Scao, Nouamane Tazi, Aleksandra Piktus, Sampo Pyysalo, Thomas Wolf, and Colin A Raffel. Scaling data-constrained language models. *Advances in Neural Information Processing Systems*, 36, 2024.
 - Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Fei Xia, Ed Chi, Quoc V Le, Denny Zhou, et al. Chain-of-thought prompting elicits reasoning in large language models. *Advances in neural information processing systems*, 2022a.
 - Peiyuan Zhang, Guangtao Zeng, Tianduo Wang, and Wei Lu. Tinyllama: An open-source small language model. *arXiv preprint arXiv:2401.02385*, 2024.
 - Anton Lozhkov, Loubna Ben Allal, Elie Bakouch, Leandro von Werra, and Thomas Wolf. Finemath: the finest collection of mathematical content, 2024. URL https://huggingface.co/datasets/HuggingFaceTB/finemath.
 - Aaron Hurst, Adam Lerer, Adam P Goucher, Adam Perelman, Aditya Ramesh, Aidan Clark, AJ Ostrow, Akila Welihinda, Alan Hayes, Alec Radford, et al. Gpt-4o system card. *arXiv preprint arXiv:2410.21276*, 2024.
 - Dan Hendrycks, Collin Burns, Saurav Kadavath, Akul Arora, Steven Basart, Eric Tang, Dawn Song, and Jacob Steinhardt. Measuring mathematical problem solving with the MATH dataset. In *Thirty-fifth Conference on Neural Information Processing Systems Datasets and Benchmarks Track (Round 2)*, 2021a.
 - Ronen Eldan and Yuanzhi Li. Tinystories: How small can language models be and still speak coherent english? *arXiv preprint arXiv:2305.07759*, 2023.
- Suriya Gunasekar, Yi Zhang, Jyoti Aneja, Caio César Teodoro Mendes, Allie Del Giorno, Sivakanth Gopi, Mojan Javaheripi, Piero Kauffmann, Gustavo de Rosa, Olli Saarikivi, et al. Textbooks are all you need. *arXiv preprint arXiv:2306.11644*, 2023.
- Yuanzhi Li, Sébastien Bubeck, Ronen Eldan, Allie Del Giorno, Suriya Gunasekar, and Yin Tat Lee. Textbooks are all you need ii: phi-1.5 technical report. *arXiv preprint arXiv:2309.05463*, 2023.
- Marah Abdin, Jyoti Aneja, Hany Awadalla, Ahmed Awadallah, Ammar Ahmad Awan, Nguyen Bach, Amit Bahree, Arash Bakhtiari, Jianmin Bao, Harkirat Behl, et al. Phi-3 technical report: A highly capable language model locally on your phone. *arXiv preprint arXiv:2404.14219*, 2024a.

543

544

546

547

548

549 550

551

552

553

554

555

556

558

559

560

561

562

563

565

566

567

568

569 570

571 572

573

574

575

576

577

578

579

580

581 582

583

584

585

586

587

588

589

590

591

- 540 Marah Abdin, Jyoti Aneja, Harkirat Behl, Sébastien Bubeck, Ronen Eldan, Suriya Gunasekar, Michael Harrison, Russell J Hewett, Mojan Javaheripi, Piero Kauffmann, et al. Phi-4 technical 542 report. arXiv preprint arXiv:2412.08905, 2024b.
 - Pratyush Maini, Skyler Seto, He Bai, David Grangier, Yizhe Zhang, and Navdeep Jaitly. Rephrasing the web: A recipe for compute & data-efficient language modeling. In *Proceedings of the 62nd* Annual Meeting of the Association for Computational Linguistics, 2024.
 - Zitong Yang, Neil Band, Shuangping Li, Emmanuel Candes, and Tatsunori Hashimoto. Synthetic continued pretraining. In The Thirteenth International Conference on Learning Representations, 2025. URL https://openreview.net/forum?id=07yvxWDSla.
 - Jonas Geiping, Sean McLeish, Neel Jain, John Kirchenbauer, Siddharth Singh, Brian R Bartoldson, Bhavya Kailkhura, Abhinav Bhatele, and Tom Goldstein. Scaling up test-time compute with latent reasoning: A recurrent depth approach. arXiv preprint arXiv:2502.05171, 2025.
 - Mostafa Dehghani, Stephan Gouws, Oriol Vinyals, Jakob Uszkoreit, and Lukasz Kaiser. Universal transformers. In International Conference on Learning Representations, 2019. URL https: //openreview.net/forum?id=HyzdRiR9Y7.
 - Angeliki Giannou, Shashank Rajput, Jy-Yong Sohn, Kangwook Lee, Jason D. Lee, and Dimitris Papailiopoulos. Looped transformers as programmable computers. In Andreas Krause, Emma Brunskill, Kyunghyun Cho, Barbara Engelhardt, Sivan Sabato, and Jonathan Scarlett, editors, Proceedings of the 40th International Conference on Machine Learning, volume 202 of Proceedings of Machine Learning Research, pages 11398–11442. PMLR, 23–29 Jul 2023. URL https://proceedings.mlr.press/v202/giannou23a.html.
 - Eric Zelikman, Georges Raif Harik, Yijia Shao, Varuna Jayasiri, Nick Haber, and Noah Goodman. Quiet-STar: Language models can teach themselves to think before speaking. In First Conference on Language Modeling, 2024. URL https://openreview.net/forum?id= oRXPiSOGH9.
 - Maxwell Nye, Anders Johan Andreassen, Guy Gur-Ari, Henryk Michalewski, Jacob Austin, David Bieber, David Dohan, Aitor Lewkowycz, Maarten Bosma, David Luan, et al. Show your work: Scratchpads for intermediate computation with language models. 2021.
 - Yuri Burda, Roger Grosse, and Ruslan Salakhutdinov. Importance weighted autoencoders. arXiv preprint arXiv:1509.00519, 2015.
 - Chris Cremer, Quaid Morris, and David Duvenaud. Reinterpreting importance-weighted autoencoders. arXiv preprint arXiv:1704.02916, 2017.
 - Long Ouyang, Jeffrey Wu, Xu Jiang, Diogo Almeida, Carroll Wainwright, Pamela Mishkin, Chong Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, et al. Training language models to follow instructions with human feedback. Advances in neural information processing systems, 35:27730– 27744, 2022.
 - Ilya Loshchilov and Frank Hutter. Decoupled weight decay regularization. In International Conference on Learning Representations, 2019. URL https://openreview.net/forum?id= Bkg6RiCqY7.
 - Karl Cobbe, Vineet Kosaraju, Mohammad Bavarian, Mark Chen, Heewoo Jun, Lukasz Kaiser, Matthias Plappert, Jerry Tworek, Jacob Hilton, Reiichiro Nakano, et al. Training verifiers to solve math word problems. arXiv preprint arXiv:2110.14168, 2021.
 - Dan Hendrycks, Collin Burns, Steven Basart, Andy Zou, Mantas Mazeika, Dawn Song, and Jacob Steinhardt. Measuring massive multitask language understanding. In International Conference on Learning Representations, 2021b. URL https://openreview.net/forum?id= d7KBjmI3GmQ.
 - Ricardo Dominguez-Olmedo, Florian E. Dorner, and Moritz Hardt. Training on the test task confounds evaluation and emergence. In The Thirteenth International Conference on Learning Representations, 2025.
 - Jason Wei, Yi Tay, Rishi Bommasani, Colin Raffel, Barret Zoph, Sebastian Borgeaud, Dani Yogatama, Maarten Bosma, Denny Zhou, Donald Metzler, Ed H. Chi, Tatsunori Hashimoto, Oriol

- Vinyals, Percy Liang, Jeff Dean, and William Fedus. Emergent abilities of large language models. *Transactions on Machine Learning Research*, 2022b. ISSN 2835-8856.
 - Rylan Schaeffer, Brando Miranda, and Sanmi Koyejo. Are emergent abilities of large language models a mirage? In *Thirty-seventh Conference on Neural Information Processing Systems*, 2023.
 - Aaron Jaech, Adam Kalai, Adam Lerer, Adam Richardson, Ahmed El-Kishky, Aiden Low, Alec Helyar, Aleksander Madry, Alex Beutel, Alex Carney, et al. Openai o1 system card. *arXiv preprint arXiv:2412.16720*, 2024.
 - Daya Guo, Dejian Yang, Haowei Zhang, Junxiao Song, Ruoyu Zhang, Runxin Xu, Qihao Zhu, Shirong Ma, Peiyi Wang, Xiao Bi, et al. Deepseek-r1: Incentivizing reasoning capability in Ilms via reinforcement learning. *arXiv preprint arXiv:2501.12948*, 2025.
 - Bingbin Liu, Sebastien Bubeck, Ronen Eldan, Janardhan Kulkarni, Yuanzhi Li, Anh Nguyen, Rachel Ward, and Yi Zhang. Tinygsm: achieving> 80% on gsm8k with small language models. *arXiv* preprint arXiv:2312.09241, 2023.
 - Ziyang Luo, Can Xu, Pu Zhao, Qingfeng Sun, Xiubo Geng, Wenxiang Hu, Chongyang Tao, Jing Ma, Qingwei Lin, and Daxin Jiang. Wizardcoder: Empowering code large language models with evol-instruct. *arXiv preprint arXiv:2306.08568*, 2023.
 - Yuxiang Wei, Zhe Wang, Jiawei Liu, Yifeng Ding, and Lingming Zhang. Magicoder: Source code is all you need. *arXiv preprint arXiv:2312.02120*, 2023.
 - Rohan Taori, Ishaan Gulrajani, Tianyi Zhang, Yann Dubois, Xuechen Li, Carlos Guestrin, Percy Liang, and Tatsunori B. Hashimoto. Stanford alpaca: An instruction-following llama model. https://github.com/tatsu-lab/stanford_alpaca, 2023.
 - Can Xu, Qingfeng Sun, Kai Zheng, Xiubo Geng, Pu Zhao, Jiazhan Feng, Chongyang Tao, and Daxin Jiang. Wizardlm: Empowering large language models to follow complex instructions. *arXiv* preprint arXiv:2304.12244, 2023.
 - Tao Ge, Xin Chan, Xiaoyang Wang, Dian Yu, Haitao Mi, and Dong Yu. Scaling synthetic data creation with 1,000,000,000 personas. *arXiv preprint arXiv:2406.20094*, 2024.
 - Gonzalo Martínez, Lauren Watson, Pedro Reviriego, José Alberto Hernández, Marc Juarez, and Rik Sarkar. Towards understanding the interplay of generative artificial intelligence and the internet. In *International Workshop on Epistemic Uncertainty in Artificial Intelligence*, pages 59–73. Springer, 2023.
 - Rohan Taori and Tatsunori Hashimoto. Data feedback loops: Model-driven amplification of dataset biases. In *International Conference on Machine Learning*, pages 33883–33920. PMLR, 2023.
 - Sina Alemohammad, Josue Casco-Rodriguez, Lorenzo Luzi, Ahmed Imtiaz Humayun, Hossein Babaei, Daniel LeJeune, Ali Siahkoohi, and Richard G Baraniuk. Self-consuming generative models go mad. *arXiv preprint arXiv:2307.01850*, 2023.
 - Elvis Dohmatob, Yunzhen Feng, Pu Yang, Francois Charton, and Julia Kempe. A tale of tails: Model collapse as a change of scaling laws. In *Forty-first International Conference on Machine Learning*, 2024. URL https://openreview.net/forum?id=KVvku47shW.
 - Joshua Kazdan, Rylan Schaeffer, Apratim Dey, Matthias Gerstgrasser, Rafael Rafailov, David L Donoho, and Sanmi Koyejo. Collapse or thrive? perils and promises of synthetic data in a self-generating world. *arXiv preprint arXiv:2410.16713*, 2024.
 - David Silver, Aja Huang, Chris J. Maddison, Arthur Guez, L. Sifre, George van den Driessche, Julian Schrittwieser, Ioannis Antonoglou, Vedavyas Panneershelvam, Marc Lanctot, Sander Dieleman, Dominik Grewe, John Nham, Nal Kalchbrenner, Ilya Sutskever, Timothy P. Lillicrap, Madeleine Leach, Koray Kavukcuoglu, Thore Graepel, and Demis Hassabis. Mastering the game of go with deep neural networks and tree search. *Nature*, 529:484–489, 2016.
- David Silver, Thomas Hubert, Julian Schrittwieser, Ioannis Antonoglou, Matthew Lai, Arthur Guez, Marc Lanctot, L. Sifre, Dharshan Kumaran, Thore Graepel, Timothy P. Lillicrap, Karen Simonyan, and Demis Hassabis. Mastering chess and shogi by self-play with a general reinforcement learning algorithm. *ArXiv*, abs/1712.01815, 2017a.

- Trieu H. Trinh, Yuhuai Wu, Quoc V. Le, He He, and Thang Luong. Solving olympiad geometry without human demonstrations. *Nature*, 625:476 482, 2024.
- Eric Zelikman, Yuhuai Wu, Jesse Mu, and Noah Goodman. STar: Bootstrapping reasoning with reasoning. In *Advances in Neural Information Processing Systems*, 2022. URL https://openreview.net/forum?id=_3ELRdg2sgI.
 - Jiaxin Huang, Shixiang Shane Gu, Le Hou, Yuexin Wu, Xuezhi Wang, Hongkun Yu, and Jiawei Han. Large language models can self-improve, 2022. URL https://arxiv.org/abs/2210.11610.
 - Avi Singh, John D Co-Reyes, Rishabh Agarwal, Ankesh Anand, Piyush Patil, Xavier Garcia, Peter J Liu, James Harrison, Jaehoon Lee, Kelvin Xu, et al. Beyond human data: Scaling self-training for problem-solving with language models. *arXiv* preprint arXiv:2312.06585, 2023.
 - Yuntian Deng, Yejin Choi, and Stuart Shieber. From explicit cot to implicit cot: Learning to internalize cot step by step. *arXiv preprint arXiv:2405.14838*, 2024.
 - Shibo Hao, Sainbayar Sukhbaatar, DiJia Su, Xian Li, Zhiting Hu, Jason Weston, and Yuandong Tian. Training large language models to reason in a continuous latent space. *arXiv* preprint *arXiv*:2412.06769, 2024.
 - Deqian Kong, Minglu Zhao, Dehong Xu, Bo Pang, Shu Wang, Edouardo Honig, Zhangzhang Si, Chuan Li, Jianwen Xie, Sirui Xie, et al. Scalable language models with posterior inference of latent thought vectors. *arXiv* preprint arXiv:2502.01567, 2025.
 - Matthew Douglas Hoffman, Du Phan, David Dohan, Sholto Douglas, Tuan Anh Le, Aaron Parisi, Pavel Sountsov, Charles Sutton, Sharad Vikram, and Rif A Saurous. Training chain-of-thought via latent-variable inference. *Advances in Neural Information Processing Systems*, 36, 2024.
 - Edward J Hu, Moksh Jain, Eric Elmoznino, Younesse Kaddar, Guillaume Lajoie, Yoshua Bengio, and Nikolay Malkin. Amortizing intractable inference in large language models. In *The Twelfth International Conference on Learning Representations*, 2024a.
 - Haolin Chen, Yihao Feng, Zuxin Liu, Weiran Yao, Akshara Prabhakar, Shelby Heinecke, Ricky Ho, Phil Mui, Silvio Savarese, Caiming Xiong, et al. Language models are hidden reasoners: Unlocking latent reasoning capabilities via self-rewarding. *arXiv preprint arXiv:2411.04282*, 2024.
 - Han Zhong, Yutong Yin, Shenao Zhang, Xiaojun Xu, Yuanxin Liu, Yifei Zuo, Zhihan Liu, Boyi Liu, Sirui Zheng, Hongyi Guo, et al. Brite: Bootstrapping reinforced thinking process to enhance language model reasoning. *arXiv preprint arXiv:2501.18858*, 2025.
 - Dongwei Jiang, Guoxuan Wang, Yining Lu, Andrew Wang, Jingyu Zhang, Chuyu Liu, Benjamin Van Durme, and Daniel Khashabi. Rationalyst: Pre-training process-supervision for improving reasoning. *arXiv preprint arXiv:2410.01044*, 2024.
 - Mojan Javaheripi, Sébastien Bubeck, Marah Abdin, Jyoti Aneja, Sebastien Bubeck, Caio César Teodoro Mendes, Weizhu Chen, Allie Del Giorno, Ronen Eldan, Sivakanth Gopi, et al. Phi-2: The surprising power of small language models. *Microsoft Research Blog*, 1(3):3, 2023.
 - David Silver, Julian Schrittwieser, Karen Simonyan, Ioannis Antonoglou, Aja Huang, Arthur Guez, Thomas Hubert, Lucas Baker, Matthew Lai, Adrian Bolton, et al. Mastering the game of Go without human knowledge. *Nature*, 550(7676):354–359, 2017b.
 - Arnaud Doucet et al. Sequential Monte Carlo methods in practice, volume 1. Springer, 2001.
 - Tim Brooks, Bill Peebles, Connor Holmes, Will DePue, Yufei Guo, Li Jing, David Schnurr, Joe Taylor, Troy Luhman, Eric Luhman, Clarence Ng, Ricky Wang, and Aditya Ramesh. Video generation models as world simulators. 2024. URL https://openai.com/research/video-generation-models-as-world-simulators.
 - Aitor Lewkowycz, Anders Andreassen, David Dohan, Ethan Dyer, Henryk Michalewski, Vinay Ramasesh, Ambrose Slone, Cem Anil, Imanol Schlag, Theo Gutman-Solo, et al. Solving quantitative reasoning problems with language models. *Advances in Neural Information Processing Systems*, 35:3843–3857, 2022.

- Jason Wei, Maarten Bosma, Vincent Y Zhao, Kelvin Guu, Adams Wei Yu, Brian Lester, Nan Du, Andrew M Dai, and Quoc V Le. Finetuned language models are zero-shot learners. *arXiv preprint arXiv:2109.01652*, 2021.
- Leo Gao, Jonathan Tow, Baber Abbasi, Stella Biderman, Sid Black, Anthony DiPofi, Charles Foster, Laurence Golding, Jeffrey Hsu, Alain Le Noac'h, Haonan Li, Kyle McDonell, Niklas Muennighoff, Chris Ociepa, Jason Phang, Laria Reynolds, Hailey Schoelkopf, Aviya Skowron, Lintang Sutawika, Eric Tang, Anish Thite, Ben Wang, Kevin Wang, and Andy Zou. A framework for few-shot language model evaluation, 07 2024. URL https://zenodo.org/records/12608602.
- HuggingFace. Math-verify: A robust mathematical expression evaluation system, 2025. URL https://github.com/huggingface/Math-Verify. Accessed: 2025-02-23.
- Jeffrey Li, Alex Fang, Georgios Smyrnis, Maor Ivgi, Matt Jordan, Samir Yitzhak Gadre, Hritik Bansal, Etash Guha, Sedrick Scott Keh, Kushal Arora, et al. Datacomp-lm: In search of the next generation of training sets for language models. *Advances in Neural Information Processing Systems*, 37:14200–14282, 2024.
- Shengding Hu, Yuge Tu, Xu Han, Ganqu Cui, Chaoqun He, Weilin Zhao, Xiang Long, Zhi Zheng, Yewei Fang, Yuxiang Huang, et al. Minicpm: Unveiling the potential of small language models with scalable training strategies. In *First Conference on Language Modeling*, 2024b.

Appendix

Table of Contents

A	A Synthetic Data Generation Experiments		
	A.1	Experimental Details	25
	A.2	Additional Results	25
B	Boo	otstrapping Experiments	28
	B.1	Bootstrapping on Fixed Data	28
	B.2	Continual Bootstrapping	31
C	C Prompts		
	C.1	Prompts for Synthetic Data Generation	33
	C.2	Prompts for Few-Show CoT Evaluation	34
		C.2.1 Synthetic Few-Shot CoT Prompts	34
		C.2.2 Standard Few-Shot CoT Prompts	42
D	Exa	amples of Latent Thoughts	46
	D.1	Examples of Downstream Transfer of Latent Thoughts	46
	D.2	Examples of GPT-4o-mini Synthesized Latent Thoughts	48
	D.3	Examples of TinyLlama-1B Self-Generated Latent Thoughts	57
	D.4	Failure Examples of Bootstrapped Models on GSM8K	64

FULL ALGORITHM

 The full algorithm of our approach is described in Algorithm 1. Since the initial base LM has not been trained as a latent thought model, we first train it on a small corpus of warmstart data with synthetic latent thoughts (e.g., by GPT-4o-mini). The warmstart → iterative EM loop resembles the supervised finetuning \rightarrow reinforcement learning pipeline (Ouyang et al., 2022) that is widely adopted in LM post-training, though our approach operates purely on task-agnostic pretraining data instead of supervised data. Our algorithm can be instantiated in different setups depending on the sources of data for bootstrapping: we may either have a new corpus of data \mathcal{X}_t at each iteration where we can initialize the training of \mathcal{M}_t from \mathcal{M}_{t-1} , or we may have a fixed corpus (with all \mathcal{X}_t being the same) and retrain the model from scratch on better latents at each iteration; see Algorithm 1 for details.

Algorithm 1: Bootstrapping Latent Thought Models (BoLT)

```
Input: Base model \mathcal{M}_b, warmstart data (\mathcal{X}_0, \mathcal{Z}_0), raw data corpora \{\mathcal{X}_t\}_{t=1}^T
Parameters: EM iterations T, Monte Carlo samples K, context window size L
Output: Bootstrapped model \mathcal{M}_T
/* Initialize with warmstart data
\mathcal{M}_0 = \text{TrainLM} \left( \mathcal{M}_b, (\mathcal{X}_0, \mathcal{Z}_0) \right)
                                                                                    > Train with next-token predictions as in Fig. 2b
/* Iteratively train the latent thought model with EM, as illustrated in Fig. 4
for t = 1, 2, \cdots, T do
      /* E-step: synthesize latent thoughts with the current model
                                                                                                                                                                         */
      for \forall X \in \mathcal{X}_t do
             \{X_n\}_{n=1}^N = \operatorname{ChunkData}(X) for n=1,2,\cdots,N do
                                                                       \triangleright Randomly chunk the raw text into N chunks
              | C<sub>n</sub> = X<sub>[n-L:n-1]</sub> | Set context window | \left\{Z_n^{(k)}\right\}_{k=1}^K \sim q(Z \mid X_n, C_n; \mathcal{M}_{t-1}) | Sample latents from the model posterior | w_n^{(k)} = \frac{p(Z_n^{(k)}, X_n \mid C_n; \mathcal{M}_{t-1})}{q(Z_n^{(k)} \mid X_n, C_n; \mathcal{M}_{t-1})} | Set context window | Sample latents from the model posterior | Weight by the model likelihood | Z_n = Z_n^{(j)}, j \sim \operatorname{Cat}\left(\left\{w_n^{(k)}\right\}_{k=1}^K\right) | Resample with importance weights
            (\mathcal{X}_t, \mathcal{Z}_t).append (\{X_n\}_{n=1}^N, \{Z_n\}_{n=1}^N)
                                                                                                  */
      /* M-step: train the model with bootstrapped thought data
      \mathcal{M}_{\text{init}} = \left\{ egin{array}{ll} \mathcal{M}_b & 	ext{if retrain from scratch} \ \mathcal{M}_{t-1} & 	ext{otherwise} \end{array} 
ight.
      \mathcal{M}_t = \text{TrainLM}\left(\mathcal{M}_{\text{init}}, (\mathcal{X}_t, \mathcal{Z}_t)\right)
                                                                                    > Train with next-token predictions as in Fig. 2b
end
return \mathcal{M}_T
```

B EXTENDED RELATED WORK

864

865

866

867

868

870

871

872

873

874

875

876

877

878

879

880

881

882

883

884

885

886

887

889

890

891

892

893

894

895

896

897

898

899

900

901

902

903

904

905

906

907

914 915 916

Synthetic data for training LMs Recent work demonstrates the benefits of synthetic data in improving training efficiency, obtaining strong models at remarkably small data or model scales. Eldan and Li (2023) synthesize a dataset of short stories with simple words and train very small LMs to generate coherent stories. Gunasekar et al. (2023) synthesize textbooks and exercises to train a small LM with strong reasoning performance. Similar approaches have been adopted in follow-up works pretraining competitive small LMs (Li et al., 2023; Abdin et al., 2024a;b) or post-training LMs with synthetic supervised data to enhance their reasoning (Liu et al., 2023; Luo et al., 2023; Wei et al., 2023) or instruction following capabilities (Taori et al., 2023; Xu et al., 2023; Ge et al., 2024). These approaches carefully curate prompts and/or a seed corpus to promote diversity in the generated corpus; failure to do so may lead to mode collapse (Martínez et al., 2023; Taori and Hashimoto, 2023; Alemohammad et al., 2023; Dohmatob et al., 2024; Kazdan et al., 2024). Another paradigm closer to our work is to augment raw pretraining data by paraphrasing (Maini et al., 2024) or describing relationships among extracted entities (Yang et al., 2025), which may avoid mode collapse due to the use of real data as a seed corpus. Our work differs in two key aspects. First, we take a latent variable perspective that augments (rather than replaces) each document chunk X with latent thoughts Z. We show this is crucial to downstream performance, outperforming strong synthetic data generation baselines in Sec. 5.2. Additionally, our work studies *bootstrapping* of the data generator rather than the teacher-student distillation settings considered in most synthetic data approaches.

Learning to reason using external supervision An increasingly popular area of work improves the reasoning capabilities of LMs using external supervision. Most paradigms rely on a verifiable reward signal, enabling the use of reinforcement learning and/or self-play (Silver et al., 2016; 2017a; Trinh et al., 2024). This approach has been successfully applied to bootstrap reasoning capabilities in math and coding tasks with verifiable rewards (Zelikman et al., 2022; Huang et al., 2022; Singh et al., 2023; Jaech et al., 2024; Guo et al., 2025). An alternative approach uses supervised finetuning data to internalize verbalized chains-of-thought into continuous hidden states (Deng et al., 2024; Hao et al., 2024; Kong et al., 2025). Some recent works take a latent variable inference perspective similar to ours, treating reasoning traces as latent variables and deriving variational inference (Hoffman et al., 2024; Hu et al., 2024a; Chen et al., 2024; Zhong et al., 2025) or Expectation-Maximization (Singh et al., 2023) approaches to optimization. Unlike our approach, these methods are domain-specific and rely on verifiable rewards or other external supervision.

Learning to reason using pretraining data Closer to our setting, other works enhance LM reasoning capabilities using pretraining data and no explicit external rewards. Geiping et al. (2025) pretrain a looped transformer (Dehghani et al., 2019; Giannou et al., 2023) on general web text and find that the continuous hidden states tend to converge with more applications of the backbone transformer, in a thought-like process. Closest to our work among these is Zelikman et al. (2024) who propose a domain-agnostic, post-training method that uses reinforcement learning to learn "thought tokens" to improve reasoning with pretraining data. While our work is also domain-agnostic and learns latent thoughts, our goal of leveraging reasoning to improve pretraining data efficiency is distinct and leads us to consider scalable, synthetic data-based approaches rather than reinforcement learning ones. This leads to a number of benefits, including a simple training method and embarassingly parallel latent generation. Our work provides a complementary perspective on how latent variable and synthetic data perspectives can enable both "learning to reason" and "reasoning to learn". Lastly, Jiang et al. (2024) also adopt a synthetic data approach, but focus on generating rationales for pretraining a process reward model. Our work differs in our application of reasoning as a primitive to enable data-efficient learning during pretraining, as well as our focus on enabling a model to bootstrap its own latent thought synthesis abilities.

C DISCUSSION

C.1 Broader Implications

Learning to reason and reasoning to learn Recent RL-based approaches such as o1 (Jaech et al., 2024) and R1 (Guo et al., 2025) *learn to reason* with LMs, using supervised data and verified rewards to unlock the reasoning capabilities of a strong pretrained model. In contrast, our work demonstrates the potential of *reasoning to learn*, by training a model to learn more from observed data by reasoning through it. We show promising results on one reasoning-heavy domain, but the likelihood-based supervision makes it applicable to more general domains, and we believe the same approach can enable more data-efficient learning of capabilities at scale. These two paradigms are complementary, and the base models trained with our approach may serve as a better starting point for task-specific RL by transferring their thoughts trained at scale.

Using asynchronous synthetic data generation to improve synchronous training efficiency Our and recent works (Javaheripi et al., 2023; Maini et al., 2024; Yang et al., 2025) demonstrate that synthetic data generation can enable improved scaling trends for synchronous pretraining on centralized compute. Because synthetic data generation can be distributed across disparate resources (Silver et al., 2017b), this approach shifts a portion of the overall synchronous pretraining compute to an asynchronous workload. This trend may inspire changes in infrastructure design for LM pretraining, enabling the effective use of distributed resources with low-bandwidth interconnects. Additionally, developing more scalable synthetic data generation techniques, such as increasing Monte Carlo samples in BoLT, could enable another axis for scaling the asynchronous compute used in pretraining.

C.2 LIMITATIONS

Constrained experimental setup under compute budget Due to our compute budget, our experimental setup was constrained to a 1B parameter LM and continued pretraining on a few billion tokens of mostly mathematical text. These choices were made to demonstrate a proof-of-concept; our use of a small model enables faster inference to synthesize billions of latent tokens, while continued pretraining on reasoning-intensive data enables measurable differences among methods in downstream evaluations. We are hopeful that future work will test reasoning to learn at larger scales and on general-domain pretraining data.

Limited exploration of design choices We focused on a particular instantiation of reasoning to learn without extensive testing of design choices such as: the generative structure of latent thoughts (currently latents are modeled autoregressively for each text chunk), the initial warmstart data generated by different models or prompts, more efficient Monte Carlo sampling techniques such as Sequential Monte Carlo (Doucet et al., 2001), and various hyperparameters like the chunk size, etc.

Side-effects of bootstrapping Bootstrapping LMs on their own synthetic data at scale may amplify specific biases and ultimately lead to unintended consequences in model behavior. We observed one possible example of such side effects: the few-shot CoT performance on GSM8K degraded with BoLT iterations, though an alternative explanation is that more optimization on mathematical tokens may degrade the language understanding necessary in GSM8K (see Fig. E.5). Future work should more extensively investigate the side effects of bootstrapping during pretraining.

C.3 FUTURE DIRECTIONS

Bootstrapping on general-domain pretraining data We believe the most exciting application of our approach could be in enabling models to self-improve on general-domain pretraining-scale data without task-specific supervision. This is in stark contrast to the recent approaches which improve pretrained LMs with reinforcement learning on task-specific labeled data (Jaech et al., 2024; Guo et al., 2025). By bootstrapping LMs with latent thoughts on pretraining-scale data, we believe it may be possible for models to acquire more general-domain reasoning capabilities that are useful beyond specific domains.

Application to general data modalities Our approach is not limited to text data and may be applied to general data modalities. This is because every piece of human-generated data is a result of an underlying latent generative process that is typically not observable. For example, when working with video data, one could train a multimodal latent model to extract the creative intent or emotional pacing behind scene transitions, rather than focusing solely on pixel-level features. For non-textual data, these underlying latent structures may be even more obscure, and models trained on such data

typically demonstrates much worse data efficiency than LMs (Brooks et al., 2024). Consequently, augmenting these data modalities with learned latents could possibly yield more pronounced gains in data efficiency than observed in the text domain.

Hierarchical latent structures Our current instantiation models the latent generation process in an autoregressive manner, where each latent is generated for a single chunk of text conditioned on the previous context. While we have demonstrated the effectiveness of this approach on reasoning-intensive data, this local latent thought structure remains inherently "myopic". It may not be sufficiently expressive to capture the hierarchical planning processes that humans employ when creating complex, long-form content such as research papers, novels, and large-scale codebases. Future work could explore more sophisticated latent structures that mirror human planning hierarchies, potentially incorporating both high-level planning and low-level reasoning.

D SYNTHETIC DATA GENERATION EXPERIMENTS

D.1 EXPERIMENTAL DETAILS

Baselines We compare our approach with several natural baselines in the data-constrained regime (Muennighoff et al., 2024) using a combination of data repetition and synthetic data generation:

- Raw-Repeat: We train on the raw CPT corpus for \sim 16 epochs with early stopping.
- Raw-Fresh: As an anticipated upper bound, we CPT on 8B unique tokens from FineMath-4+, without any repetition.
- WRAP-Orig: WRAP (Maini et al., 2024) rephrases the data in four diverse styles: easy (with simple language), hard (with complex language), Wikipedia (high-quality), and question-answer. We prompt GPT-40-mini (see Prompt F.1.2) with temperature 0.7 to rephrase each document in these styles. The average token ratio between synthetic paraphrases and raw tokens is 0.48:1, so we generate 5 paraphrases per document to approximately match the total synthetic tokens of our approach. While Maini et al. (2024) found that mixing synthetic and raw data works best, we tuned the mixture coefficient and found that entirely synthetic works best (see Fig. D.1a).
- WRAP-CoT: To probe whether the gains of our approach arise from simply including synthetic reasoning traces in the training data, we develop a WRAP variant that prompts GPT-4o-mini to rephrase documents with interspersed reasoning steps (Prompt F.1.3). This strong WRAP baseline allows us to assess whether it is key to maintain thoughts in a separate latent space to explain the corresponding text, rather than directly in the raw text space. The synthetic-raw token ratio is 0.7:1, so we generate 4 paraphrases per document for a total of ~1.3B synthetic tokens. As above, we tuned the synthetic-raw mixture coefficient and found entirely synthetic works best (see Fig. D.1b).

Training We use a cosine learning rate schedule with a 1000 step warmup and peak learning rate of 1e-4 (tuned over {1e-5, 3e-5, 1e-4, 3e-4, 1e-3} in initial experiments). All models are trained with sequence length 2048 and batch size 96 on 4 x H200 GPUs.

Evaluation To ensure robust evaluation, we evaluate with two distinct sets of few-shot CoT prompts (see Appx. F.2): (1) Standard prompts from previous works, i.e., Minerva CoT (Lewkowycz et al., 2022) for MATH, the default CoT from (Wei et al., 2022a) for GSM8K, and the FLAN CoT (Wei et al., 2021) for MMLU-STEM; (2) Synthetic CoT prompts, where the CoT traces are synthesized by GPT-40-mini using Prompt F.1.1, using the question as the prefix and answer as the suffix. We report the results with synthetic CoT prompts by default, as they performed better uniformly across methods (see Fig. D.2). All models are evaluated with a temperature 0.0. We conducted all evaluations using the LM Evaluation Harness (Gao et al., 2024), and the Math-Verify evaluator (HuggingFace, 2025) for scoring MATH final answers against the ground-truth.

Tuning the mixing ratio of WRAP baselines We have tuned the ratio of mixing raw data with paraphrased data for the WRAP baselines, similar to Maini et al. (2024). For each document during training, we randomly select either the raw data or the paraphrased data according to a fixed mixing ratio drawn from {0.0, 0.25, 0.5}, where 0.0 means using only paraphrased data without mixing raw data. The results are shown in Fig. D.1. For both WRAP-Orig and WRAP-CoT, applying no mixing (ratio=0.0) leads to the best performance, which we used in our experiments.

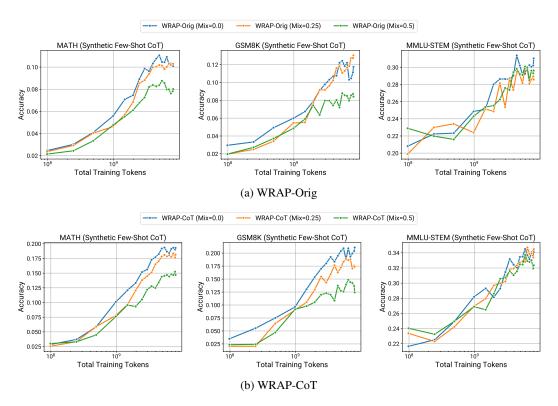


Figure D.1: Training purely on paraphrased data without mixing raw data generally leads to the best performance for WRAP baselines.

D.2 ADDITIONAL RESULTS

Evaluation results with different few-shot CoT prompts To facilitate a robust evaluation of model downstream performance, we have tested each benchmark using two sets of few-shot CoT prompts. Besides the synthetic CoT prompts that we used by default, we also tested the standard CoT prompts from previous works – specifically, the Minerva CoT (Lewkowycz et al., 2022) for MATH, the default CoT from (Wei et al., 2022a) for GSM8K, and the FLAN CoT (Wei et al., 2021) for MMLU-STEM. See Appx. F.2 for detailed prompts. The evaluation results are shown in Fig. D.2. Our method demonstrates robust performance across different prompt variations, maintaining consistent gains over baselines. In contrast, several baselines, most notably WRAP-CoT, exhibit substantial performance degradation when evaluated with standard prompts, leading to an even wider performance gap between our method and the baselines.

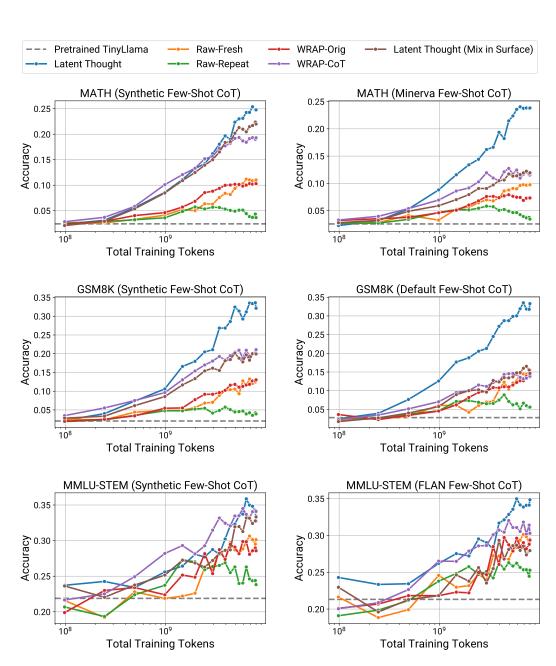


Figure D.2: Evaluation results with different few-shot CoT prompts: synthetic CoT prompts (left) vs. standard CoT prompts used in previous works (right). Our method demonstrates robust performance across different prompt variations and consistent gains over baselines. The performance gap is more pronounced when using standard CoT prompts due to degradation of the baseline performance.

 Comparison normalized by the effective raw tokens seen by each method Fig. 5 compares different methods with the same amount of training tokens. However, since different methods produce different amounts of training tokens per raw document, this means each method processes a different number of raw documents during training. In particular, our method was trained on 3.3 times more raw tokens (due to the latent to raw token ratio of 2.3) than the raw data baselines. To provide a complementary perspective, Fig. D.3 shows performance when methods are normalized by the effective number of raw tokens seen during training. For WRAP baselines, we computed the effective raw tokens based on the number of raw tokens that were paraphrased (i.e., training tokens divided by synthetic-to-raw token ratio). We find that the gains of our method over baselines are more significant, highlighting the data efficiency of our method.

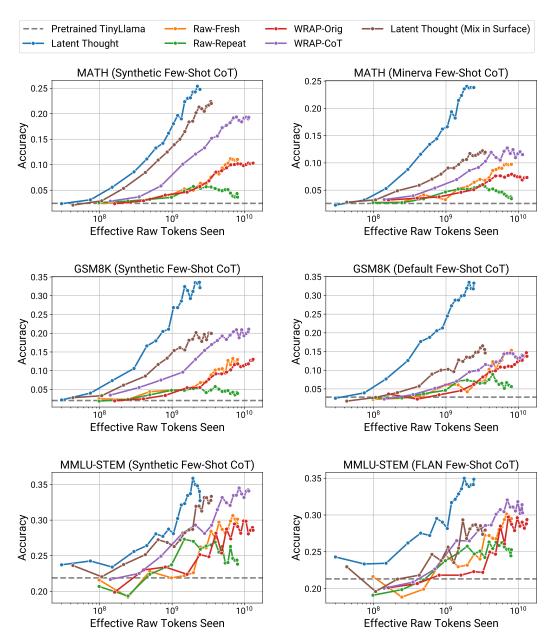


Figure D.3: Comparison normalized by the effective raw tokens seen by each method. Our method achieves more significant gains over baselines, highlighting its data efficiency.

E BOOTSTRAPPING EXPERIMENTS

E.1 BOOTSTRAPPING ON FIXED DATA

E.1.1 EXPERIMENTAL DETAILS

Bootstrapping configuration We distribute the latent generation workload over a cluster of H100/200, A100/6000/5000/40, L40, and RTX3090 GPUs. We follow the training configuration in Sec. 5.1, except we now use batch size 192 on 8 x H200 GPUs for bootstrapping training runs.

Downstream metrics and finetuned evaluation setup We take the training set for MATH & GSM8K and finetune on this set. We synthetically generate CoT traces for the finetuning set using our prompt in Prompt F.1.1, by treating the question as the prefix and the answer as the suffix. We finetune using AdamW with standard hyperparameters $\beta_1 = 0.9$, $\beta_2 = 0.95$, weight decay 0, and a cosine learning rate schedule with a warmup over the first 5% of steps. We split out 10% of the train set as a validation set and tune hyperparameters to obtain a learning rate of 1e-4, batch size 64, and 5 epochs. We report test accuracy of the final model checkpoint with CoT, using standard sampling hyperparameters top_k= 50, top_p= 0.9, and temperature 0.6. Due to the small test set size of GSM8K, we follow Guo et al. (2025) and calculate the mean accuracy over 16 random samples for each test question to reduce variance.

Reducing run-to-run variability We report means and standard errors over 5 train seeds for bootstrapping and 3 train seeds for baselines. For finetuning evaluations, standard error is over 3 upstream CPT \times 5 downstream finetuning = 15 total seeds.

E.1.2 ADDITIONAL RESULTS

Log-likelihood evaluation In Fig. E.1, we measure the negative log-likelihood (NLL) of latent thought models at each bootstrapping iteration on a holdout validation set. Note that this is not an ideal evaluation metric for latent thought models as they are not trained to directly optimize NLL and the evaluation does not utilize the latent thoughts. Nevertheless, we report it for reference and to provide a relative comparison across bootstrap iterations as a smooth evaluation metric. Our results show that latent thought models demonstrate lower NLL over multiple iterations with smooth and clear gains up to the fourth iteration, which is consistent with the ELBO results in Fig. 7. For additional context, we also include raw data baselines, which directly optimize NLL during training and therefore achieve substantially better NLL performance.

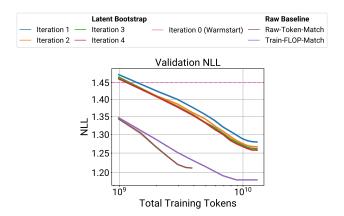


Figure E.1: Evaluation results of negative log-likelihood (NLL) on a holdout validation set. Our latent thought models demonstrate lower NLL over multiple iterations with clear gains, despite not being directly trained to optimize NLL (unlike raw data baselines).

Detailed evaluation results on MATH In Fig. E.2, we include detailed evaluation results on MATH on both our synthetic CoT prompt and the standard Minerva CoT prompt. We find that the gains of latent thought models across multiple bootstrap iterations are robust across prompt variations. The gains seem to plateau after the third iteration, which might partially be due to the discrete nature of downstream evaluations. The gains over the baseline of training on raw data is even more significant on the Minerva CoT prompt, indicating the effectiveness of training models with self-generated latent thoughts.

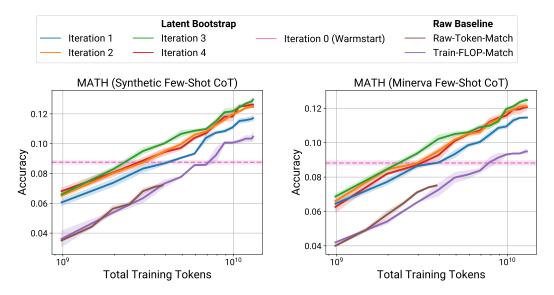


Figure E.2: Detailed evaluation results on MATH. The gains of latent thought models across multiple bootstrap iterations are robust across prompt variations, and the gains over the baseline of training on raw data is more significant on the Minerva CoT prompt.

Exact model performance curves In our main results (Fig. 7), we show the performance curves of the best model during training at each bootstrap iteration. In Fig. E.3, we show the exact performance curves of all models during training, which demonstrates a cosnistent relative comparison but with a slightly larger variance.

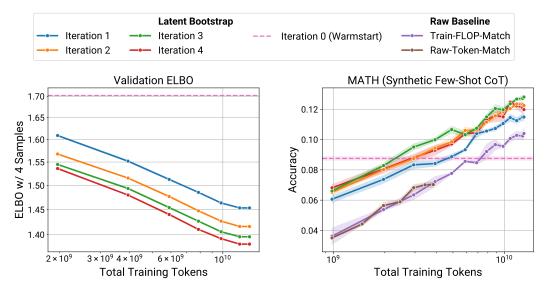


Figure E.3: Exact model performance curves during training. We observe a consistent relative comparison with Fig. 7 but with a slightly larger variance.

 Additional results on scaling MC samples In Fig. E.4, we include additional results of scaling MC samples, following the same setup as Fig. 9. The performance gains remain robust across different prompt variations and when evaluated GSM8K.

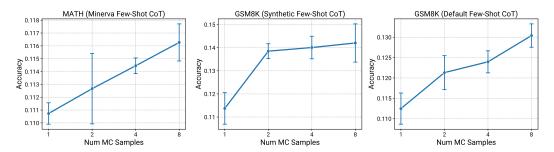


Figure E.4: Additional results on scaling MC samples. The performance gains of increasing MC samples persist across prompt variations (left) and when evaluated on GSM8K (middle and right).

"Negative" bootstrapping results In Fig. E.5, we include some negative bootstrapping results where the model performance on GSM8K with few-shot prompting deteriorates over multiple bootstrap iterations (but still outperforms the raw data baselines), which differs from our results on MATH (Fig. 7) or on fine-tuning evaluations (Fig. 8, right). A potential explanation for the degradation could be that as the model optimizes better on the mathematical-heavy FineMath data, it becomes worse at natural language understanding. We tested this hypothesis on a holdout validation set from the general-domain DCLM (Li et al., 2024) data, shown in Fig. E.5 (right). We find that models trained more on FineMath typically get worse on DCLM NLL (as evidenced by bootstrapped models performing worse than the warmstart model, and the train-FLOP-matched baseline performance on GSM8K might be attributed to their decreased natural language understanding capabilities due to increased optimization on the FineMath training data. We have also included some failure examples of bootstrapped models on GSM8K in Appx. G.4, where the bootstrapped models made mistake seemingly due to misinterpretation of the math word problems.

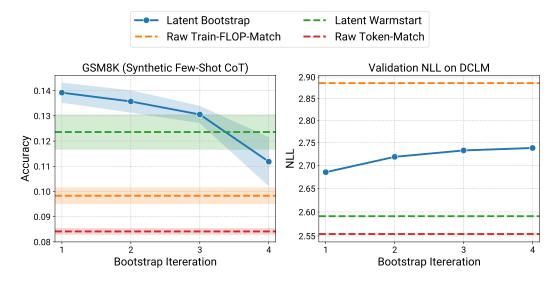


Figure E.5: "Negative" bootstrapping results. The bootstrapped models' performance on GSM8K deteriorates over multiple bootstrap iterations (left), which might be attributed to their decreased natural language understanding capabilities due to increased optimization on the FineMath training data, as evidenced by their worse NLL on DCLM (right).

Decomposed ELBO across bootstrap iterations To provide more insights on how the latent thoughts evolve across bootstrap iterations, we decompose the ELBO objective as follows and measure each term respectively across bootstrap iterations:

$$ELBO = \mathbb{E}_{q(Z|X)}[\log p(X|Z) + \log p(Z)] + H[q(Z|X)]$$

= $\mathbb{E}_{q(Z|X)}[\log p(X|Z)] - D_{KL}[q(Z|X)||p(Z)]$

Here, the three terms can be interpreted as:

- Utility: $\log p(X|Z)$ how well Z helps predict X,
- Simplicity: $\log p(Z)$ how well the prior fits Z,
- Diversity: H[q(Z|X)] the entropy of the thought generator.

In Fig. E.6, we plot each term (measured in nats per sample) over four bootstrapping iterations and observe the following trends:

- Utility ($\log P(X|Z)$) decreased slightly, from about -1290 to -1320 nats/sample,
- Simplicity ($\log P(Z)$) improved substantially, from about -3090 to -2920 nats/sample,
- Diversity (H[Q(Z|X)]) also decreased slightly, from about 2790 to 2730 nats/sample.

These changes contributed to an overall ELBO improvement from about -1590 to -1510 nats/sample. This suggests that, in our regime, there remains a large gap between the posterior and the prior, and the bootstrapping process optimizes for narrowing this gap—i.e., making the latent thoughts more learnable—while sacrificing a small amount of predictability and diversity. Ultimately, this improves the overall ELBO and downstream accuracy, and therefore we consider it to be desirable.

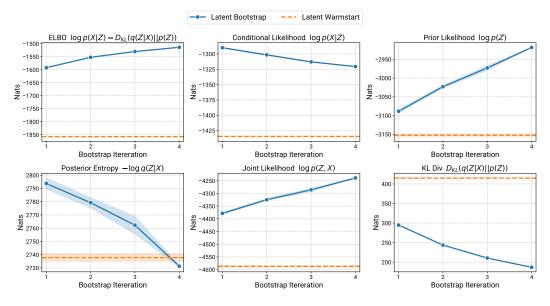


Figure E.6: Decomposed ELBO across bootstrap iterations. Latent thought models prioritize making the latents simpler and more learnale (higher $\log p(Z)$) while sacrificing a small amount of predictability (lower $\log p(X|Z)$) and diversity (lower H[q(Z|X)]).

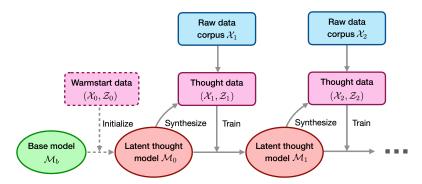


Figure E.7: **Continual bootstrapping**: We study a more practical scenario in which an LM is continually pretrained on a series of corpora, using the model at each iteration to generate latent thoughts for the subsequent corpus to train itself.

E.2 CONTINUAL BOOTSTRAPPING

Our previous experiment confirms that BoLT leads to iterative improvement of the latent thought model. However, re-training the latent thought model from scratch at each iteration could be wasteful, and a more realistic setting would be to perform those experiments as part of continual learning, where the model updates its latent thought generator as it processes more and more data. We show that BoLT continues to work in this setting.

E.2.1 EXPERIMENTAL SETUP

Bootstrapping configuration We use TinyLlama-1.1B as the base model \mathcal{M}_b and 240M raw tokens from the FineMath-4+ dataset as the warmstart data \mathcal{X}_0 . We conduct T=4 bootstrap iterations, where each iteration t uses a distinct subset of the FineMath-4+ dataset with logarithmically increasing size as \mathcal{X}_t —960M, 960M, 1.92B, and 3.84B raw tokens respectively. Models are trained for 1 epoch on the warmstart data and 2 epochs on the bootstrap data at each iteration. To mitigate forgetting and model degradation during continual training, we use the warmup-stable-decay (WSD, Hu et al., 2024b) schedule. In particular, the learning rate is warmed up for 1000 steps only at the warmstart stage and then maintained at a constant value until the decay phase. For each training stage (including both warmstart and bootstrap iterations), we linearly decay the learning rate during the final 15% of training steps, where the final checkpoint is obtained as \mathcal{M}_t for evaluation and synthetic latent generation. When training \mathcal{M}_{t+1} at next iteration, we initialize both the model and the optimizer states from the pre-decay checkpoint, and continue training with the same constant learning rate without re-warming up. The learning rate is set to 3e-5, which was tuned in our preliminary experiments to mitigate forgetting and achieve stable transitions across iterations. All other configurations follow Sec. 6.1.1, such as the use of 4 MC samples for synthetic latent generation.

Baseline comparison The key question to answer in this setup is whether iterative improvement of the latent thought model leads to performance benefits. To understand this, we compare to a baseline where we stop iterative improvement at time t', fixing the latent thought generation model. We investigate this by comparing our bootstrapped models (that use the most recent and capable model \mathcal{M}_t to synthesize latents for \mathcal{M}_{t+1} at each iteration) against the alternatives of fixing the latent generator at a previous iteration $\mathcal{M}_{t'}$ for training subsequent models \mathcal{M}_t , $\forall t > t'$. We conduct 3 training runs for each training setup to reduce the run-to-run variability and measure the models' performance at each iteration following the same evaluation protocol described in Sec. 6.1.1.

E.2.2 RESULTS

Continual, iterative improvement of latent thought models Fig. E.8 shows the best model performance at each bootstrap iteration, comparing the use of the bootstrapped model to generate latents for the next training corpus, versus the use of models from previous iterations to generate latents. Our results demonstrate that using the more capable bootstrapped models lead to a steeper scaling trend in likelihood-based metrics (see Fig. E.8 left for ELBO and Fig. E.9 left for NLL results), which demonstrates the bootstrapping effects of our approach in improving data scaling with higher-quality latents. The gains are also reflected in the downstream MATH performance (Fig. E.8 right), where the bootstrapped models consistently outperform the fixed-latent-generator baselines, with an increasingly pronounced gap in performance at later iterations. We include the detailed model

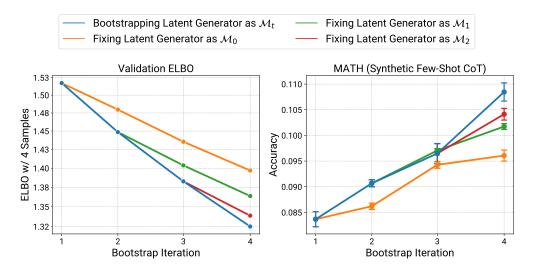


Figure E.8: **BoLT can bootstrap in continual learning settings.** We instantiate the setup illustrated in Fig. E.7 by training on a series of corpora and compare two options: bootstrapping the latent model across all four iterations (blue), or fixing the latent generator at earlier iterations. Our bootstrapped models lead to a steeper scaling trend in likelihood-based metrics (left) and consistently outperform the fixed-latent-generator baselines on downstream MATH performance (right), with gains becoming more pronounced at later iterations. Results over 3 training runs are reported. Due to the log-sized scaling of the x-axis, the scaling trends here are comparable to a traditional data-scaling law. See Fig. E.9 for additional results and Fig. E.10 for detailed evaluations over the training runs.

performance during training runs in Fig. E.10. Collectively, these results demonstrate the potential of our approach in forming a model self-improvement loop—where the more capable latent thought models produce higher quality latents that lead to better learning efficiency.

Additional evaluation results In Fig. E.9, we include additional evaluation results of NLL on the holdout validation set and the MATH performance on the Minerva CoT prompt. We find that NLL evaluation demonstrates a similar trend as the ELBO evaluation in Fig. E.8 (left), where the bootstrapped models lead to a steeper scaling trend than the fixed-latent-generator baselines. For the MATH performance on the Minerva CoT prompt, the bootstrapped models also demonstrate consistent gains, even though the performance gap at the forth iteration is a bit less pronounced than using the synthetic CoT prompt (Fig. E.8 right).

Detailed evaluation results during training runs In Fig. E.10, we include the evaluation results of each model during the training runs. We plot the best model performance during each training run to denoise the evaluation results and report the average over 3 runs, following the same practice in Fig. 7. From the plots, we can observe a clear difference of scaling trend between the bootstrapped models and the fixed-latent-generator baselines, especially for the likelihood-based metrics (top row).

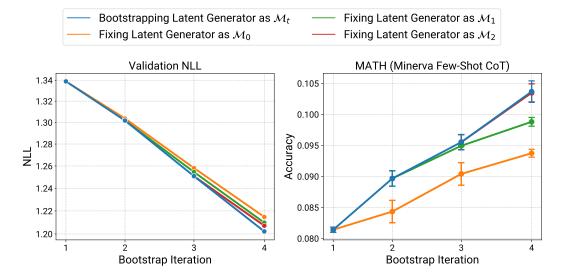


Figure E.9: Additional evaluation results of NLL on the holdout validation set (left) and MATH performance on the Minerva CoT prompt (right). Similar to our main results (Fig. E.8), the bootstrapped models lead to a steeper scaling trend in likelihood-based metrics and consistent improvement on downstream MATH performance.

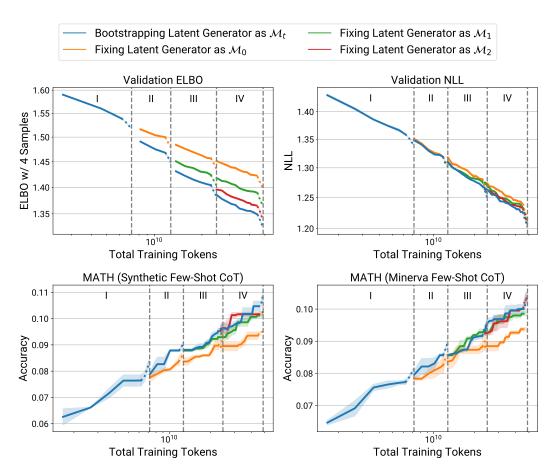


Figure E.10: Detailed evaluation results during training runs. We report the best model performance during each training run to denoise the results and report the average over 3 runs. The dashed curves denote the learning rate decay phase of each training run, where the final models are used for evaluation and latent generation. Models right before the decay phase are used for continual training at the next iteration.

F PROMPTS

F.1 Prompts for Synthetic Data Generation

Prompt F.1.1: GPT-4o-mini to generate latent thoughts on warmstart data

System

You are an advanced AI system, highly knowledgeable and capable of deeply understanding and reasoning through any web document

User

You are provided with a pair of web document prefix and suffix. Your task is to insert latent thoughts between them underlying the creation of the suffix conditioned on the prefix. The latent thoughts should include: the missing background knowledge and the reasoning traces underlying each claim (especially, step-by-step derivations or logical reasoning).

Prefix {prefix}

Suffix {suffix}

Your turn

Now provide the latent thoughts. Use concise, simple, and declarative language. Do not give any supporting remarks or references to the terms 'prefix' and 'suffix', as this output will go directly into a computer program. Do not apply any markdown formatting or text embellishments. Optimize the content to ensure every word is informative, avoid vague language like 'xxx is essential'. Emphasize on the suffix without repeating the content in the prefix. Focus on implicit reasoning and background knowledge that is not explicitly stated in the suffix, and use concrete logical reasoning or mathematical derivations when applicable.

Prompt F.1.2: GPT-4o-mini to rephrase the data with WRAP (Maini et al., 2024) prompts

Prompt 1

For the following website, give me a paraphrase of it using a very small vocabulary and extremely simple sentences that a toddler will understand.

Website

1654 {text}

Prompt II

For the following website, give me a paraphrase of it using very terse and abstruse language that only an erudite scholar will understand. Replace simple words and phrases with rare and complex ones.

Website

1660 {text}

1661 # Prompt III

For the following website, give me a diverse paraphrase of it in high quality English language, as in sentences on Wikipedia.

Website

{text}

1667 # Prompt IV

Convert the following website into a conversational format with multiple tags of 'Question:' followed by 'Answer:'.

Website {text}

1684 1685

1686

1687

1688

1693

1694

1695

1697

1699

1700

1701

1702

1703

```
1675
                   Prompt F.1.3: GPT-4o-mini to rephrase the data with explicit reasoning steps
1676
1677
           ## User
1678
           You are provided with the following document. Give me a diverse paraphrase in high
           quality English language that preserves all information in the original document. Also,
1679
           insert detailed reasoning thoughts (especially, step-by-step mathematical derivations or
1680
           logical reasoning). Do not output '## Paraphrased Document'.
1681
           ## Document
1682
           {text}
1683
```

F.2 PROMPTS FOR FEW-SHOW COT EVALUATION

Here we include the detailed prompts used in our few-shot CoT evaluation. In particular, we used distinct sets of few-shot CoT prompts (see Appx. F.2): (1) Synthetic CoT prompts (Appx. F.2.1) – where the CoT traces are synthetically generated by GPT-4o-mini using Prompt F.1.1 given the question as prefix and answer as suffix; (2) Standard prompts from previous works (Appx. F.2.2) – specifically, the Minerva CoT (Lewkowycz et al., 2022) for MATH, the default CoT from (Wei et al., 2022a) for GSM8K, and the FLAN CoT (Wei et al., 2021) for MMLU-STEM. We include the prompts used for our latent thought models for illustration – for baselines we used the exact same prompts with special latent tokens (i.e., <StartofLatent>, <Prior>, and <EndofLatent>) removed.

For each prompt, we include the few-shot CoT examples and a randomly sampled test question. For MATH, we used the same few-shot questions as in the Minerva prompts. For GSM8K, we used a randomly selected 8 problems from the training set as the few-shot examples in our synthetic CoT prompt, which aligns better with the actual test problem distribution than the easier problems in the default CoT prompt. For MMLU-STEM, we randomly selected 4 problems from 4 subjects in our synthetic CoT prompt, while FLAN CoT prompts use subject-specific few-shot examples for each subject (we include the one for high school statistics for an example).

F.2.1 SYNTHETIC FEW-SHOT COT PROMPTS

```
1704
                                 Prompt F.2.1: Synthetic Few-Shot CoT for MATH
1705
            Problem:
1706
            Find the domain of the expression \frac{\sqrt{x-2}}{\sqrt{5-x}}.
1707
1708
            <StartOfLatent><Prior>
1709
            To find the domain of the expression, we need to ensure the numerator and denominator
            are both defined and the denominator is not zero.
1710
1711
            Start with the numerator: \sqrt{x-2}. This square root is defined when the expression inside
1712
             is non-negative. Thus, we require:
            x - 2 \ge 0
1713
            x \ge 2.
1714
1715
            Next, consider the denominator: \sqrt{5-x}. This square root is defined when the expression
            inside is non-negative as well, and it must also be strictly positive to avoid division
1716
            by zero. Therefore, we need:
1717
            5 - x > 0
1718
            x < 5.
1719
            Now, we combine the two inequalities. From the first inequality, we have x \geq 2, and
1720
            from the second inequality, we have x < 5.
1721
            The combined conditions are:
1722
            2 < x < 5.
1723
            In interval notation, this is expressed as [2, 5].
1724
1725
            Thus, the final answer is [2, 5).
1726
            The final answer is [2,5)
1727
```

```
1728
1729
            Final Answer: The final answer is [2,5). I hope it is correct.
1730
1731
            If \det \mathbf{A} = 2 and \det \mathbf{B} = 12, then find \det(\mathbf{AB}).
1732
1733
             <StartOfLatent><Prior>
1734
            To find the determinant of the product of two matrices, we use the property that states
            \det(AB) = \det(A) \cdot \det(B). This property holds for any square matrices A and B of the
1735
            same size.
1736
            Given that \det \mathbf{A} = 2 and \det \mathbf{B} = 12, we can directly apply this property.
1737
1738
1739
            1. Start with the known determinants: \det \mathbf{A}=2 and \det \mathbf{B}=12.
1740
            2. Multiply these two values together: 2 \cdot 12.
1741
            3. Perform the multiplication: 2 \cdot 12 = 24.
1742
            Thus, based on these calculations, we conclude that \det(\mathbf{AB}) = 24.
1743
            The final answer is 24
1744
1745
            Final Answer: The final answer is 24. I hope it is correct.
1746
1747
            Terrell usually lifts two 20-pound weights 12 times. If he uses two 15-pound weights
1748
            instead, how many times must Terrell lift them in order to lift the same total weight?
1749
1750
            <StartOfLatent><Prior>
            To determine how many times Terrell must lift the two 15-pound weights to equal the
1751
            total weight lifted with the two 20-pound weights, we first calculate the total weight
1752
            lifted with the 20-pound weights.
1753
            Each 20-pound weight contributes 20 pounds, and since he lifts two weights, the total
1754
            weight per lift is 20 pounds + 20 pounds = 40 pounds. He lifts this total 12 times,
1755
            resulting in a total weight of 40 pounds * 12 = 480 pounds.
1756
            Next, we compute how many times he needs to lift the two 15-pound weights to match this
1757
            total weight. Each 15-pound weight contributes 15 pounds, so the total weight per lift
1758
            with the 15-pound weights is 15 pounds + 15 pounds = 30 pounds.
1759
            To find the number of lifts required to reach the same total weight of 480 pounds, we
1760
            set up the equation: 30 pounds \star x lifts = 480 pounds, where x represents the number of
1761
1762
            Solving for x, we divide both sides by 30 pounds:
1763
            x = 480 pounds / 30 pounds = 16 lifts.
1764
            Thus, Terrell must lift the two 15-pound weights 16 times to equal the total weight he
1765
             lifts with the two 20-pound weights.
1766
            The final answer is 16
1767
             <EndOfLatent>
1768
            Final Answer: The final answer is 16. I hope it is correct.
1769
1770
            If the system of equations
1771
1772
                                                      6x - 4y = a,
1773
                                                      6y - 9x = b.
1774
1775
            has a solution (x,y) where x and y are both nonzero,
1776
            find \frac{a}{b}, assuming b is nonzero.
1777
             <StartOfLatent><Prior>
1778
            To find the ratio \frac{a}{h}, we start with the system of equations given. The first equation is
1779
             6x - 4y = a, and the second is 6y - 9x = b.
1780
            We can express a in terms of x and y from the first equation:
1781
```

1835

```
1782
              a = 6x - 4y.
1783
1784
              Next, we rearrange the second equation to express b:
1785
1786
              To find the ratio \frac{a}{b}, we substitute the expressions we derived:
1787
              \frac{a}{b} = \frac{6x - 4y}{6y - 9x}.
1788
              Next, we need to simplify this expression. We can factor out a common factor in the
1789
              numerator and the denominator. First, observe that both a and b can be rewritten in a
1790
              way that may reveal their relationship:
1791
              In the numerator, we can rearrange it as 6x - 4y = 2(3x - 2y).
              In the denominator, we rearrange 6y - 9x = 3(2y - 3x).
1792
1793
              Now, substituting these factorizations back into our ratio gives us:
1794
              \frac{a}{b} = \frac{2(3x-2y)}{3(2y-3x)}.
1795
              Next, we can simplify further. Notice that 2y-3x can be rewritten as -(3x-2y):
1796
              \frac{a}{b} = \frac{2(3x-2y)}{3(-(3x-2y))}
1797
1798
              This simplifies to:
1799
               \frac{a}{b} = \frac{2}{-3} = -\frac{2}{3}.
1800
              Since we assumed b is nonzero, this leads us to conclude that the system of equations
1801
              has a specified relationship between \boldsymbol{a} and \boldsymbol{b} when both \boldsymbol{x} and \boldsymbol{y} are nonzero.
1802
              The final answer is -\frac{2}{3}
1803
              <EndOfLatent>
1804
              Final Answer: The final answer is -\frac{2}{3}. I hope it is correct.
1805
              Problem:
1806
              Rationalize the denominator: \frac{1}{2\sqrt{7}}.
1807
1808
              <StartOfLatent><Prior>
1809
```

Prompt F.2.2: Synthetic Few-Shot CoT for GSM8K

```
1812
1813
            ##Problem
            Question: Ben works 8-hour shifts in a furniture shop. It takes him 5 hours to build 1
1814
            rocking chair. How many chairs can he build in 10 days?
1815
            Final Answer:
1816
            <StartOfLatent><Prior>
            To find out how many rocking chairs Ben can build in 10 days, first determine the total
1817
            number of hours he works in that period. Since he works 8-hour shifts, over 10 days, he
1818
            works a total of 8 hours/day multiplied by 10 days, which equals 80 hours.
1819
            Next, calculate how many rocking chairs Ben can build in those 80 hours. Each rocking
1820
            chair takes him 5 hours to complete. To find the total number of chairs, divide the
1821
            total hours worked by the time taken to build one chair: 80 hours divided by 5 hours/
            chair equals 16 chairs.
1822
1823
            Thus, in 10 days, Ben can build 16 rocking chairs. Therefore, the answer is 16.
1824
            <EndOfLatent>
            The final answer is 16. I hope it is correct.
1825
1826
1827
            ##Problem
1828
            Question: A family has three adults and children, both girls and boys. They went out for
1829
             a family trip and prepared 3 dozen boiled eggs for the trip. Every adult got 3 eggs,
            and the boys each received 1 more egg than each girl since the girls had some snacks.
1830
            How many boys went on the trip if the number of girls was 7 and each girl received an
1831
            egg?
1832
            Final Answer:
            <StartOfLatent><Prior>
1833
1834
```

```
1836
            To determine the number of boys on the trip, we first analyze the information given. We
1837
            know there are 7 girls. Each girl received 1 egg. Therefore, the total number of eggs
1838
            given to the girls is:
1839
            7 \text{ qirls} * 1 \text{ eqq/qirl} = 7 \text{ eqqs}.
1840
1841
            Next, we find out how many eggs were allocated to the adults. There are 3 adults, and
1842
            each received 3 eggs, which results in:
1843
            3 adults * 3 eggs/adult = 9 eggs.
1844
            Now we calculate the total number of eggs distributed so far:
1845
1846
            7 eggs (to girls) + 9 eggs (to adults) = 16 eggs.
1847
            Since the family prepared 3 dozen boiled eggs, we find the total number of eggs:
1848
1849
            3 \text{ dozen} = 3 * 12 = 36 \text{ eggs}.
1850
            Now we determine how many eggs remain for the boys:
1851
1852
            36 total eggs - 16 eggs already distributed = 20 eggs available for the boys.
1853
            According to the problem, each boy received 1 more egg than each girl, which means each
1854
            boy received:
1855
            1 egg (for the girl) + 1 egg = 2 eggs.
1856
1857
            Let B represent the number of boys. The total number of eggs given to the boys can be
1858
            expressed as:
1859
            B boys * 2 eggs/boy = 2B eggs.
1860
1861
            Setting the available eggs for the boys equal to the eggs they received, we have:
1862
1863
1864
            To find B, we solve the equation:
1865
            B = 20 / 2 = 10.
1866
            Thus, the number of boys who went on the trip is 10. Therefore, the answer is 10.
1867
            <EndOfLatent>
1868
            The final answer is 10. I hope it is correct.
1869
1870
1871
            ##Problem
1872
            Question: A deep-sea monster rises from the waters once every hundred years to feast on
            a ship and sate its hunger. Over three hundred years, it has consumed 847 people. Ships
1873
            have been built larger over time, so each new ship has twice as many people as the last
1874
            ship. How many people were on the ship the monster ate in the first hundred years?
1875
            Final Answer:
            <StartOfLatent><Prior>
1876
            To find the number of people on the ship the monster ate in the first hundred years, we
1877
            need to understand the pattern of ship sizes and the total consumption over three
1878
            hundred vears.
1879
            The monster consumes a ship every hundred years, meaning it consumes a total of three
1880
            ships over three hundred years. The problem states that ships have been built larger
            over time, specifically that each new ship has twice as many people as the last ship.
1881
1882
            Let's denote the number of people on the ship in the first hundred years as \boldsymbol{x}. According
1883
             to the problem, the second ship, which the monster eats in the second hundred years,
            has twice as many people, or 2x. The third ship has twice the number of people as the
1884
            second ship, or 4x.
1885
            Now we can express the total number of people consumed over the three hundred years as:
1886
            x (first ship) + 2x (second ship) + 4x (third ship) = 847 people.
1887
1888
            Combining these terms gives us:
```

```
1890
            x + 2x + 4x = 7x.
1891
1892
            Setting this equal to the total consumption:
1893
1894
            To find x, we divide both sides by 7:
1895
            x = 847 / 7 = 121.
1896
            Thus, the number of people on the ship the monster ate in the first hundred years is 121.
1897
             Therefore, the answer is 121.
1898
            <EndOfLatent>
            The final answer is 121. I hope it is correct.
1899
1900
1901
            ##Problem
1902
            Question: Nancy is filling an aquarium for her fish. She fills it halfway and goes to
1903
            answer the door. While she's gone, her cat knocks the aquarium over and spills half the
1904
            water in it. Then Nancy comes back and triples the amount of water in the aquarium. If
            the aquarium is 4 feet long, 6 feet wide, and 3 feet high, how many cubic feet of water
1905
            are in the aquarium?
1906
            Final Answer:
1907
            <StartOfLatent><Prior>
            To find the total volume of water in the aquarium, calculate the aquarium's total volume
1908
             first. The formula for the volume of a rectangular prism is length multiplied by width
1909
            multiplied by height. Here, the dimensions are 4 feet long, 6 feet wide, and 3 feet high.
1910
1911
            Calculate the total volume:
1912
            Volume = length x width x height
            Volume = 4 feet x 6 feet x 3 feet = 72 cubic feet.
1913
1914
            Initially, Nancy fills the aquarium halfway. Therefore, the amount of water after she
1915
            fills it is:
            Half of 72 cubic feet = 72 / 2 = 36 cubic feet.
1916
1917
            While Nancy is away, the cat spills half of this amount:
1918
            Amount spilled = 36 cubic feet / 2 = 18 cubic feet.
            Thus, the remaining water in the aquarium after the spill is:
1919
            36 cubic feet - 18 cubic feet = 18 cubic feet.
1920
1921
            When Nancy returns, she triples the remaining amount of water. Calculate the new amount
            of water:
1922
            New amount = 3 \times 18 cubic feet = 54 cubic feet.
1923
1924
            The final amount of water in the aquarium in cubic feet is 54. Therefore, the answer is
1925
1926
            The final answer is 54. I hope it is correct.
1927
1928
1929
            ##Problem
            Question: There are 30 different nuts in a bowl. If 5/6 of the nuts were eaten, how many
1930
             nuts were left?
1931
            Final Answer:
1932
            To determine how many nuts are left after 5/6 of them are eaten, first, calculate how
1933
            many nuts were consumed. Start with the total number of nuts, which is 30. To find the
1934
            number of nuts eaten, multiply 30 by 5/6.
1935
            30 * (5/6) = 25 \text{ nuts eaten.}
1936
1937
            Next, subtract the number of nuts eaten from the total number of nuts to find how many
            are left.
1938
1939
            30 - 25 = 5 nuts remaining.
1940
            Thus, after consuming 5/6 of the nuts, 5 nuts are left in the bowl. Therefore, the
1941
            answer is 5.
1942
            <EndOfLatent>
1943
```

```
1944
            The final answer is 5. I hope it is correct.
1945
1946
1947
            ##Problem
1948
            Question: Yuan is 7 years older than his brother David and twice his age. How old is
1949
            David?
1950
            Final Answer:
            <StartOfLatent><Prior>
1951
            Let Y represent Yuan's age and D represent David's age.
1952
1953
            From the information given, we can establish two equations based on the relationships
            described.
1954
1955
            First, the statement "Yuan is 7 years older than his brother David" can be
            mathematically expressed as:
1956
            Y = D + 7.
1957
1958
            Second, the statement "Yuan is twice his age" translates to:
            Y = 2D.
1959
1960
            Now, we can set the two equations equal to each other since both represent Yuan's age:
1961
1962
            To isolate D, we need to rearrange the equation:
1963
            7 = 2D - D,
            7 = D
1964
1965
            Thus, David's age is determined to be 7 years old. Therefore, the answer is 7.
1966
            The final answer is 7. I hope it is correct.
1967
1968
1969
            ##Problem
1970
            Question: Joey wants to buy the latest released pair of designer High Jump basketball
1971
            sneakers. He plans to mow 3 neighbors' lawns for $8 a lawn, sell 2 collectible figures
1972
            to his friends for $9 each, and work an after-school job for 10 hours at $5 per hour. If
             his earnings just cover the price of the High Jump sneakers, how much do the shoes cost
1973
1974
            Final Answer:
1975
            <StartOfLatent><Prior>
            To find the total earnings Joey accumulates, we need to calculate each source of income
1976
            step-by-step.
1977
            First, calculate the income from mowing lawns. Joey mows 3 lawns at a rate of $8 per
1978
            lawn. Therefore, the total income from mowing is:
1979
            3 \text{ lawns} \star \$8/\text{lawn} = \$24.
1980
            Next, calculate the income from selling collectible figures. Joey sells 2 figures at a
1981
            price of $9 each. Thus, the total income from selling figures is:
1982
            2 figures * $9/figure = $18.
1983
            Now, calculate the income from his after-school job. Joey works for 10 hours at a rate
1984
            of $5 per hour. Therefore, his income from the job is:
1985
            10 hours * $5/hour = $50.
1986
            Now, sum all these earnings to find the total amount Joey has:
1987
            $24 (mowing) + $18 (selling figures) + $50 (job) = $92.
1988
            Since Joey's total earnings of $92 exactly cover the price of the High Jump basketball
1989
            sneakers, it can be concluded that the cost of the sneakers is $92. Therefore, the
1990
            answer is 92.
1991
            <EndOfLatent>
            The final answer is 92. I hope it is correct.
1992
1993
1994
            ##Problem
1995
            Question: Before getting to work, Jess has a few errands to run. Jess has to walk 11
1996
            blocks to the store and 6 blocks to the gallery, before walking the final 8 blocks to
1997
```

```
1998
            arrive at work. If Jess has already walked 5 blocks, how many more blocks must she walk
1999
            before arriving at work?
2000
            Final Answer:
2001
            <StartOfLatent><Prior>
            To determine the total distance Jess needs to walk to work, first calculate the total
2002
            number of blocks from the store and the gallery to work. The distance to the store is 11
             blocks, and the distance to the gallery is 6 blocks. Thus, the total distance to work
2004
            is the sum of these two distances plus the final stretch to work, which is 8 blocks.
            Total distance = 11 blocks (to store) + 6 blocks (to gallery) + 8 blocks (to work) = 25
2006
            blocks.
2007
            Next, subtract the distance Jess has already walked, which is 5 blocks, from the total
2008
            distance calculated.
2009
            Remaining distance = Total distance - Distance already walked = 25 blocks - 5 blocks =
2010
            20 blocks.
2011
2012
            Thus, Jess must walk 20 more blocks to reach work. Therefore, the answer is 20.
2013
            The final answer is 20. I hope it is correct.
2014
2015
2016
            ##Problem
2017
            Question: Marcell and Beatrice are having a contest to see who can eat the most fruit
            roll-ups, so they unroll as many as they can find. Unfortunately, someone makes a
2018
            mistake and Beatrice's was two roll-ups wide and 24 rolls up long while Marcell's was 3
2019
            roll-ups wide and 14 roll-ups long. If they both ate their entire amount, how many did
2020
            they eat on average?
            Final Answer:
2021
            <StartOfLatent><Prior>
2022
```

Prompt F.2.3: Synthetic Few-Shot CoT for MMLU-STEM

```
2025
2026
2027
            Question: A microwave oven is connected to an outlet, 120 V, and draws a current of 2
2028
            amps. At what rate is energy being used by the microwave oven?
             (A) 10 W
2029
             (B) 30 W
2030
             (C) 60 W
             (D) 240 W
2031
            Answer:
2032
             <StartOfLatent><Prior>
2033
            To determine the rate of energy usage, we apply the formula for electrical power, which
            is given by P = V \times I, where P is power in watts, V is voltage in volts, and I is
2034
            current in amps.
2035
2036
            In this scenario, the voltage V is 120 V and the current I is 2 A.
2037
            Calculating the power:
2038
            P = 120 V \times 2 A = 240 W.
2039
            This calculation shows that the microwave oven consumes energy at a rate of 240 watts.
2040
2041
            Thus, the answer is (D) 240 W.
             <EndOfLatent>
2042
            The answer is (D) 240W.
2043
2044
2045
2046
            Question: The variable x varies directly as the square of y, and y varies directly as
2047
            the cube of z. If x equals -16 when z equals 2, what is the value of x when z equals \frac{1}{2}?
             (A) -1
2048
             (B) 16
2049
            (C) -\frac{1}{256} (D) \frac{1}{16}
2050
            Answer:
```

```
2052
             <StartOfLatent><Prior>
2053
             To solve the problem, we need to establish the relationships between the variables x, y,
2054
             and z based on the given conditions.
2055
            Firstly, since x varies directly as the square of y, we can express this relationship
2056
            mathematically as:
2057
            x = k1 * y^2, where k1 is a constant.
2058
            Next, since y varies directly as the cube of z, we can express this relationship as:
2059
             y = k2 * z^3, where k2 is another constant.
2060
            Now we can substitute y in the first equation with its expression in terms of z:
2061
             x = k1 * (k2 * z^3)^2
2062
             x = k1 * k2^2 * z^6.
2063
             This shows that x varies directly as the sixth power of z, leading us to express this as
2064
2065
             x = k * z^6, where k = k1 * k2^2.
2066
            Given that x = -16 when z = 2, we can substitute these values into the equation to find
2067
             k:
2068
             -16 = k * (2^6)
2069
             -16 = k \star 64
             k = -16 / 64
2070
            k = -1/4.
2071
            Now we have the relationship for x in terms of z:
2072
             x = -1/4 * z^6.
2073
2074
            Next, we need to find the value of x when z = 1/2. We substitute z = 1/2 into the
2075
            equation:
             x = -1/4 * (1/2)^6
2076
            x = -1/4 * (1/64)
2077
             x = -1/256
2078
             Therefore, the value of x when z = 1/2 is -1/256. Thus, the answer is (C) -\frac{1}{256}.
2079
2080
             <EndOfLatent>
2081
2082
2083
             ##Problem
            Ouestion: Which expression is equivalent to 5 x 9?
2084
             (A) (5 \times 4) \times (6 \times 5)
2085
             (B) (5 \times 5) + (5 \times 4)
             (C) (5 \times 5) + (5 \times 9)
2086
             (D) (5 x 9) x (6 x 9)
2087
             Answer:
2088
             <StartOfLatent><Prior>
             To determine which expression is equivalent to 5 \times 9, we need to evaluate each option
2089
             step by step.
2090
2091
            Option (A): (5 \times 4) \times (6 \times 5)
             This expression simplifies to 20 x 30, which equals 600. This is not equal to 5 x 9 = 45.
2092
2093
2094
            Option (B): (5 \times 5) + (5 \times 4)
            This expression simplifies to 25 + 20, which equals 45. This matches 5 \times 9.
2095
2096
             Option (C): (5 \times 5) + (5 \times 9)
            This expression simplifies to 25 + 45, which equals 70. This is not equal to 5 \times 9.
2097
2098
             Option (D): (5 \times 9) \times (6 \times 9)
2099
            This expression simplifies to 45 \times 54, which equals 2430. This is not equal to 5 \times 9.
2100
             Thus, the answer is (B) (5 \times 5) + (5 \times 4).
2101
             <EndOfLatent>
            The answer is (B) (5 \times 5) + (5 \times 4).
2102
2103
2104
2105
```

```
2106
            ##Problem
2107
            Question: A new smartwatch is manufactured in one part of a factory, then secured for
2108
            shipping in another, independent part of the factory. The weight of the smartwatch has a
2109
             mean of 62 grams and a standard deviation of 1.0 grams. The weight of the packaging (
            box, user's quide, bubble wrap, etc.) has a mean of 456 grams and a standard deviation
2110
            of 6 grams. Together, the distribution of the weight of the smartwatch and its packaging
2111
             would have the following mean and standard deviation:
2112
            (A) Mean 518 grams; standard deviation 7.0 grams
            (B) Mean 518 grams; standard deviation 3.5 grams
2113
            (C) Mean 518 grams; standard deviation 6.1 grams
2114
            (D) Mean 394 grams; standard deviation 6.1 grams
2115
            Answer:
            <StartOfLatent><Prior>
2116
            To find the combined mean and standard deviation of the smartwatch and packaging, we
2117
            start by calculating the mean of the total weight. The mean weight of the smartwatch is
            62 grams, and the mean weight of the packaging is 456 grams. The total mean weight is
2118
            calculated by adding these two means:
2119
2120
            Mean total weight = Mean smartwatch + Mean packaging
            Mean total weight = 62 grams + 456 grams = 518 grams.
2121
2122
            Next, we need to calculate the standard deviation of the total weight. Since the weights
2123
            of the smartwatch and packaging are independent, we can use the formula for the sum of
            two independent random variables:
2124
2125
            Standard deviation of total weight = sqrt((Standard deviation of smartwatch)^2 + (
            Standard deviation of packaging) ^2).
2126
2127
            The standard deviation of the smartwatch is 1.0 grams, and the standard deviation of the
2128
             packaging is 6 grams. Now, we apply the values to the formula:
2129
            Standard deviation of total weight = sqrt((1.0 \text{ grams})^2 + (6 \text{ grams})^2)
2130
            = sqrt(1 + 36)
2131
            = sqrt(37)
            \approx 6.08 grams.
2132
2133
            Rounding this value to one decimal place gives us approximately 6.1 grams.
2134
            Therefore, the combined mean weight is 518 grams, and the combined standard deviation is
2135
             approximately 6.1 grams, which corresponds to option (C). Thus, the answer is (C) Mean
2136
            518 grams; standard deviation 6.1 grams.
2137
            <EndOfLatent:
            The answer is (C) Mean 518 grams; standard deviation 6.1 grams.
2138
2139
2140
            ##Problem
2141
            Question: Mr. Bee's statistics class had a standard deviation of 11.2 on a standardized
2142
            test, while Mr. Em's class had a standard deviation of 5.6 on the same test. Which of
            the following is the most reasonable conclusion concerning the two classes' performance
2143
            on the test?
2144
            (A) Mr. Bee's class is less heterogeneous than Mr. Em's.
2145
            (B) Mr. Em's class is more homogeneous than Mr. Bee's.
            (C) Mr. Bee's class performed twice as well as Mr. Em's.
2146
            (D) Mr. Em's class did not do as well as Mr. Bee's.
2147
            Answer:
            <StartOfLatent><Prior>
2148
```

F.2.2 STANDARD FEW-SHOT COT PROMPTS

21492150

2151

21522153

2154

2155 2156

2157

2158

2159

Prompt F.2.4: Minerva Few-Shot CoT for MATH (Lewkowycz et al., 2022)

```
Problem: Find the domain of the expression \frac{\sqrt{x-2}}{\sqrt{5-x}}. 
 <StartOfLatent><Prior> The expressions inside each square root must be non-negative. Therefore, x-2\geq 0, so x\geq 2, and 5-x\geq 0, so x\leq 5. Also, the denominator cannot be equal to zero, so 5-x>0, which gives x<5. Therefore, the domain of the expression is (2,5).
```

```
2160
              The final answer is [2,5)
2161
2162
2163
             Final Answer: The final answer is [2,5). I hope it is correct.
2164
             Problem:
2165
              If \det \mathbf{A} = 2 and \det \mathbf{B} = 12, then find \det(\mathbf{AB}).
2166
              <StartOfLatent><Prior>
2167
             We have that \det(\mathbf{AB}) = (\det \mathbf{A})(\det \mathbf{B}) = (2)(12) = \boxed{24}
2168
              The final answer is 24
2169
2170
              <EndOfLatent>
             Final Answer: The final answer is 24. I hope it is correct.
2171
2172
             Terrell usually lifts two 20-pound weights 12 times. If he uses two 15-pound weights
2173
              instead, how many times must Terrell lift them in order to lift the same total weight?
2174
2175
              <StartOfLatent><Prior>
              If Terrell lifts two 20-pound weights 12 times, he lifts a total of 2\cdot 12\cdot 20=480 pounds
2176
              of weight. If he lifts two 15-pound weights instead for n times, he will lift a total
2177
             of 2\cdot 15\cdot n=30n pounds of weight. Equating this to 480 pounds, we can solve for n\colon
2178
2179
                                                        30n = 480
2180
                                                        n = 480/30 = \boxed{16}
2181
2182
              The final answer is 16
2183
2184
              <EndOfLatent>
             Final Answer: The final answer is 16. I hope it is correct.
2185
2186
             Problem:
2187
             If the system of equations
2188
2189
                                                          6x - 4y = a,
2190
                                                          6y - 9x = b.
2191
2192
             has a solution (x,y) where x and y are both nonzero,
2193
              find \frac{a}{b}, assuming b is nonzero.
2194
2195
              <StartOfLatent><Prior>
              If we multiply the first equation by -\frac{3}{2}, we obtain
2196
2197
                                                        6y - 9x = -\frac{3}{2}a.
2198
2199
2200
             Since we also know that 6y-9x=b, we have
2201
2202
                                                    -\frac{3}{2}a = b \Rightarrow \frac{a}{b} = \boxed{-\frac{2}{3}}
2203
2204
2205
              The final answer is -\frac{2}{3}
2206
2207
              Final Answer: The final answer is -\frac{2}{3}. I hope it is correct.
2208
2209
             Problem:
             Rationalize the denominator: \frac{1}{2\sqrt{7}}.
2210
2211
              <StartOfLatent><Prior>
2212
```

```
2214
                    Prompt F.2.5: Default Few-Shot CoT for GSM8K from (Wei et al., 2022a)
2215
2216
            ##Problem
2217
            Q: Olivia has 23.Sheboughtfive bagels for 3 each. How much money does she have left?
2218
            <StartOfLatent><Prior>
2219
            Olivia had 23 dollars. 5 bagels for 3 dollars each will be 5 \times 3 = 15 dollars. So she
2220
            has 23 - 15 dollars left. 23 - 15 is 8.
            <EndOfLatent>
2221
            The answer is 8.
2222
2223
2224
            ##Problem
2225
            Q: Jason had 20 lollipops. He gave Denny some lollipops. Now Jason has 12 lollipops. How
             many lollipops did Jason give to Denny?
2226
2227
            <StartOfLatent><Prior>
2228
            Jason started with 20 lollipops. Then he had 12 after giving some to Denny. So he gave
            Denny 20 - 12 = 8.
2229
            <EndOfLatent>
2230
            The answer is 8.
2231
2232
2233
            Q: There are 15 trees in the grove. Grove workers will plant trees in the grove today.
2234
            After they are done, there will be 21 trees. How many trees did the grove workers plant
2235
            today?
2236
            A:
2237
            <StartOfLatent><Prior>
            There are 15 trees originally. Then there were 21 trees after some more were planted. So
2238
             there must have been 21 - 15 = 6.
2239
            <EndOfLatent>
            The answer is 6.
2240
2241
2242
            ##Problem
2243
            Q: If there are 3 cars in the parking lot and 2 more cars arrive, how many cars are in
2244
            the parking lot?
2245
            A:
            <StartOfLatent><Prior>
2246
            There are originally 3 cars. 2 more cars arrive. 3 + 2 = 5.
2247
            <EndOfLatent>
2248
            The answer is 5.
2249
2250
            ##Problem
2251
            Q: There were nine computers in the server room. Five more computers were installed each
2252
             day, from monday to thursday. How many computers are now in the server room?
2253
2254
            There were originally 9 computers. For each of 4 days, 5 more computers were added. So 5
2255
             * 4 = 20 computers were added. 9 + 20 is 29.
            <EndOfLatent>
2256
            The answer is 29.
2257
2258
2259
            ##Problem
2260
            Q: Michael had 58 golf balls. On tuesday, he lost 23 golf balls. On wednesday, he lost 2
2261
            more. How many golf balls did he have at the end of wednesday?
            A:
2262
2263
            Michael started with 58 golf balls. After losing 23 on tuesday, he had 58 - 23 = 35.
            After losing 2 more, he had 35 - 2 = 33 golf balls.
2264
            <EndOfLatent>
2265
            The answer is 33.
2266
2267
```

2320

2321

experiment?

```
2268
2269
2270
            ##Problem
2271
            Q: Leah had 32 chocolates and her sister had 42. If they ate 35, how many pieces do they
             have left in total?
2272
2273
2274
            Originally, Leah had 32 chocolates. Her sister had 42. So in total they had 32 + 42 = 74.
            After eating 35, they had 74 - 35 = 39.
2275
            <EndOfLatent>
2276
            The answer is 39.
2277
2278
2279
            ##Problem
            Q: Shawn has five toys. For Christmas, he got two toys each from his mom and dad. How
2280
            many toys does he have now?
2281
            A:
2282
            Shawn started with 5 toys. If he got 2 toys each from his mom and dad, then that is 4
2283
            more tovs. 5 + 4 = 9.
2284
            <EndOfLatent>
2285
            The answer is 9.
2286
2287
            ##Problem
2288
            Q: Marcell and Beatrice are having a contest to see who can eat the most fruit roll-ups,
2289
             so they unroll as many as they can find. Unfortunately, someone makes a mistake and
2290
            Beatrice's was two roll-ups wide and 24 rolls up long while Marcell's was 3 roll-ups
2291
            wide and 14 roll-ups long. If they both ate their entire amount, how many did they eat
            on average?
            A:
2293
            <StartOfLatent><Prior>
2294
```

2296 Prompt F.2.6: FLAN Few-Shot CoT for MMLU-STEM from (Wei et al., 2021) 2297 2298 The following are multiple choice questions (with answers) about high school statistics. 2299 ##Problem 2300 Question: A new smartwatch is manufactured in one part of a factory, then secured for shipping in another, independent part of the factory. The weight of the smartwatch has a 2301 mean of 62 grams and a standard deviation of 1.0 grams. The weight of the packaging (2302 box, user's guide, bubble wrap, etc.) has a mean of 456 grams and a standard deviation 2303 of 6 grams. Together, the distribution of the weight of the smartwatch and its packaging would have the following mean and standard deviation: (A) Mean 518 grams; standard deviation 7.0 grams (B) Mean 518 grams; standard deviation 2305 3.5 grams (C) Mean 518 grams; standard deviation 6.1 grams (D) Mean 394 grams; standard 2306 deviation 6.1 grams Answer: <StartOfLatent><Prior> 2308 Let's think step by step. Since the weight of the watch and the weight of the packaging 2309 are independent random variables, the mean and variance of their sum is equal to the sum of their individual means and variances. So the mean is 62 + 456 = 518 grams, and the 2310 variances is $1.0^2 + 6.0^2 = 37$, leading to a standard deviation of 6.1 grams. The 2311 answer is (C). <EndOfLatent 2312 The answer is (C). 2313 2314 2315 ##Problem 2316 Question: After a frost warning was issued, the owner of a large orange grove asked his 2317 workers to spray all his trees with water. The water was supposed to freeze and form a protective covering of ice around the orange blossom. Nevertheless, the owner suspected 2318 that some trees suffered considerable damage due to the frost. To estimate the 2319 proportion of trees that suffered more than 50 percent damage due to the frost, he took

a random sample of 100 trees from his grove. What is the response variable in this

```
2322
            (A) The proportion of trees that suffered more than 50 percent damage due to frost. (B)
            The number of trees affected by the frost. (C) The number of trees sampled from the
2324
            grove. (D) For each sampled tree, whether it suffered more than 50 percent damage or at
2325
            most 50 percent damage.
            Answer:
2326
            <StartOfLatent><Prior>
2327
            Let's think step by step. In this experiment, the response variable is what is measured.
2328
             For each tree, what is measured is whether or not it suffered more than 50 percent
            damage due to the frost. The answer is (D).
2329
            <EndOfLatent>
2330
            The answer is (D).
2331
2332
2333
            ##Problem
            Question: Suppose X and Y are random variables with E(X) = 37, var(X) = 5, E(Y) = 62,
2334
            and var(Y) = 12. What are the expected value and variance of the random variable X + Y?
2335
            (A) E(X + Y) = 99, Var(X + Y) = 8.5 (B) E(X + Y) = 99, Var(X + Y) = 13 (C) E(X + Y) = 13
2336
            99, var(X + Y) = 17 (D) There is insufficient information to answer this question.
            Answer:
2337
            <StartOfLatent><Prior>
2338
            Let's think step by step. While means of sums of random variables add (regardless of
2339
            whether the variables are independent) in order to determine the variance of a sum of
            random variables, we need to know not just their individual variances but the covariance
2340
             of the two variables, which is not given in this problem. The answer is (D).
2341
            The answer is (D).
2342
2343
2344
2345
            ##Problem
            Question: Which of the following sets has the smallest standard deviation? Which has the
2346
             largest?
2347
            I: \{1,2,3\}
            II: {-10,10}
2348
            III: {100}
2349
            (A) I, II (B) II, III (C) III, I (D) III, II
2350
            Answer:
            <StartOfLatent><Prior>
2351
            Let's think step by step. The variance of distribution I is the expected squared
2352
            deviation from its mean (which is 2), so the variance is 2/3 . The variance of
2353
            distribution II is 10^2 (because both elements are 10 away from the mean of zero). The
            variance of distribution III is 0, since it has a single entry. So distribution III has
2354
            the smallest standard deviation and distribution II has the largest. The answer is (D).
2355
2356
            The answer is (D).
2357
2358
            ##Problem
2359
            Question: Mr. Bee's statistics class had a standard deviation of 11.2 on a standardized
2360
            test, while Mr. Em's class had a standard deviation of 5.6 on the same test. Which of
2361
            the following is the most reasonable conclusion concerning the two classes' performance
2362
            (A) Mr. Bee's class is less heterogeneous than Mr. Em's. (B) Mr. Em's class is more
2363
            homogeneous than Mr. Bee's. (C) Mr. Bee's class performed twice as well as Mr. Em's. (D)
             Mr. Em's class did not do as well as Mr. Bee's.
2364
            Answer:
2365
            <StartOfLatent><Prior>
2366
```

G Examples of Latent Thoughts

23672368

2369 2370

2371

23722373

2374

G.1 Examples of Downstream Transfer of Latent Thoughts

Example G.1.1: Solving algebra problems on MATH

```
Problem: What is the distance between the two intersections of y=x^2 and x+y=1? 
 <StartOfLatent><Prior>
```

```
2376
             To find the distance between the two intersections of the curves y=x^2 and x+y=1, we
2377
             first need to determine the points of intersection. Setting the equations equal to each
2378
             other gives us x^2 = x + 1. Rearranging this leads to x^2 - x - 1 = 0. We can solve this
2379
             quadratic equation using the quadratic formula, x = [1 \pm sqrt(1 + 4)] / 2, which
             simplifies to x = [1 \pm sqrt(5)] / 2. This results in two intersection points: (1 + sqrt
2380
             (5))/2 and (1 - sqrt(5))/2.
2382
            Next, we need to find the distance between these two points. The distance formula
            between two points (x1, y1) and (x2, y2) is given by d = sqrt((x2 - x1)^2 + (y2 - y1)^2).
             Substituting our points into this formula, we have d = sqrt(((1 + sqrt(5))/2 - (1 - sqrt(5))/2))
2384
             sqrt(5))/2)^2 + ((1 + sqrt(5))/2 - (1 - sqrt(5))/2)^2).
2385
             Calculating the differences in the x-coordinates and y-coordinates separately, we find:
2386
2387
             1. For the x-coordinates: (1 + sqrt(5))/2 - (1 - sqrt(5))/2 = (1 + sqrt(5) - 1 + sqrt(5))
             )/2 = (2 \operatorname{sqrt}(5))/2 = \operatorname{sqrt}(5).
2388
             2. For the y-coordinates: (1 + sqrt(5))/2 - (1 - sqrt(5))/2 = (1 + sqrt(5) - 1 + sqrt(5))
2389
             )/2 = (2 \operatorname{sqrt}(5))/2 = \operatorname{sqrt}(5).
2390
             Now, substituting these values back into the distance formula gives us d = sqrt((sqrt(5)
2391
             )^2 + (sqrt(5))^2 = sqrt(5 + 5) = sqrt(10).
2392
2393
             Thus, the final answer is sqrt(10).
             <EndOfLatent>
2394
2395
             Final Answer: The final answer is \sqrt{10}. I hope it is correct.
2396
```

Example G.1.2: Solving combinatorial problems on MATH

```
2399
2400
            Problem: Pat is to select six cookies from a tray containing only chocolate chip,
            oatmeal, and peanut butter cookies. There are at least six of each of these three kinds
2401
            of cookies on the tray. How many
2402
            different assortments of six cookies can be selected? (Note that cookies of the same
2403
            type are not distinguishable.)
2404
            <StartOfLatent><Prior>
2405
            To determine the number of different assortments of six cookies from a tray containing
2406
            chocolate chip, oatmeal, and peanut butter cookies, we need to consider the constraints
            imposed by the requirement of at least six of each type of cookie.
2407
2408
            First, we recognize that the total number of cookies is fixed at six. The problem can be
             approached using combinatorial methods, specifically the "stars and bars" theorem. This
2409
             theorem helps us find the number of ways to distribute indistinguishable objects (the
2410
            cookies) into distinguishable boxes (the types of cookies).
2411
            Let:
2412
            - x = number of chocolate chip cookies
2413
            - v = number of oatmeal cookies
2414
            -z = number of peanut butter cookies
2415
            The total number of cookies is given by the equation:
2416
            x + y + z = 6
2417
            The stars and bars theorem states that the number of ways to distribute \boldsymbol{n}
2418
            indistinguishable objects into k distinguishable boxes is given by the formula:
2419
2420
            (n + k - 1) choose (k - 1)
2421
            In our case, n = 6 (the total number of cookies) and k = 3 (the types of cookies:
2422
            chocolate chip, oatmeal, and peanut butter).
2423
            Applying the formula:
2424
2425
            (6 + 3 - 1) choose (3 - 1) = 8 choose 2
2426
            Calculating this gives:
2427
2428
            8 choose 2 = 8! / (2!(8-2)!) = 8! / (2! * 6!) = (8 * 7) / (2 * 1) = 28
2429
```

```
Thus, there are 28 different assortments of six cookies that can be selected from the
tray, ensuring that at least six of each type are included.
<EndOfLatent>
Final Answer: The final answer is 28. I hope it is correct.
```

2430

2431

2432

2433

2434

Example G.1.3: Solving number theory problems on MATH

```
2438
            Problem: How many positive factors does 48 have?
2439
2440
            <StartOfLatent><Prior>
            To find the number of positive factors of 48, we start by determining its prime
2441
            factorization. The prime factorization of 48 is 2^4 * 3^1.
2442
2443
            Next, we apply the formula for finding the number of positive factors from the prime
            factorization. For a number expressed as pl^e1 * p2^e2 * ... * pk^ek, the total number
2444
            of factors is given by (e1 + 1)(e2 + 1)...(ek + 1).
2445
2446
            In the case of 48, we have:
            - For the prime factor 2, the exponent is 4.
2447
            - For the prime factor 3, the exponent is 1.
2448
            Applying the formula: (4 + 1)(1 + 1) = 5 * 2 = 10. Thus, 48 has 10 positive factors.
2449
2450
            The final answer is 10.
2451
            <EndOfLatent>
2452
            Final Answer: The final answer is 10. I hope it is correct.
2453
```

2454 2455

Example G.1.4: Solving math word problems on GSM8K

```
2456
2457
            Question: Farmer Brown has 60 animals on his farm, all either chickens or cows. He has
2458
            twice as many chickens as cows. How many legs do the animals have, all together?
2459
            Final Answer:
2460
            <StartOfLatent><Prior>
2461
            To find the total number of legs, we need to establish the relationship between the
2462
            number of chickens and cows. Let the number of cows be represented as C. According to
            the problem, the number of
2463
            chickens is twice that of cows, so we can express the number of chickens as 2C.
2464
            Next, we know that the total number of animals is 60. Therefore, we can set up the
2465
            equation:
2466
2467
            C + 2C = 60
2468
            This simplifies to:
2469
2470
2471
            To find C, we divide both sides by 3:
2472
            C = 20
2473
2474
            Now that we have the number of cows, we can find the number of chickens:
2475
            Number of chickens = 2C = 2 \times 20 = 40
2476
2477
            Now we can calculate the total number of legs. Each chicken has 2 legs, and each cow has
2478
             4 legs. The total number of legs can be calculated as follows:
2479
            Total legs = Number of chickens * Legs per chicken + Number of cows * Legs per cow
2480
2481
            Substituting the known values:
2482
            Total legs = 40 * 2 + 20 * 4
2483
```

```
Calculating this gives:

Total legs = 80 + 80 = 160

Thus, the total number of legs is 160.

<EndOfLatent>
The final answer is 160.
```

G.2 Examples of GPT-40-mini Synthesized Latent Thoughts

Example G.2.1: Latent thoughts recover missing background knowledge and reasoning

<StartOfLatent>

Principal Component Analysis (PCA) is a statistical technique used for dimensionality reduction, which is crucial in data analysis and machine learning. The objective of PCA is to identify the directions (principal components) in which the variance of the data is maximized. This process involves transforming the original features into a new set of uncorrelated features, ordered by the amount of variance they capture.

To begin with, we consider a dataset with multiple dimensions, where each dimension represents a feature. The first step in PCA is to center the data by subtracting the mean of each feature from the dataset. This step is vital because it ensures that the new coordinate system is based around the origin, allowing for a more accurate analysis of variance.

Next, we compute the covariance matrix of the decentralized data. The covariance matrix captures the relationships between the different features, indicating how they vary together. The covariance values reveal the spread and the direction of the data in the feature space.

Following the covariance matrix calculation, we need to determine the eigenvalues and eigenvectors of this matrix. The eigenvectors represent the directions of the principal components, while the eigenvalues indicate the magnitude of variance in those directions. The larger the eigenvalue, the more significant the corresponding eigenvector is in explaining the variance in the dataset.

To select the principal components, we rank the eigenvectors based on their eigenvalues in descending order. The top k eigenvectors, where k is the number of dimensions we wish to retain, form a new feature space that captures the most variance.

Finally, we project the original, decentralized data onto this new feature space by multiplying the original data matrix by the matrix of the top k eigenvectors. This projection results in a reduced dataset that retains the essential characteristics of the original data while minimizing information loss. The effectiveness of PCA lies in its ability to simplify complex datasets, making subsequent analyses more manageable and interpretable.

The basic idea of PCA

Principal component analysis is to find out the most important aspects in the data and replace the original data with the most important aspects in the data. To put it bluntly, it is to reduce the data from n-dimensional to n'-dimensional, and hope that the data set of n'-dimensional features retains most of the information as much as possible.

PCA mathematical derivation (maximum variance method)

For example, if the data is a two-dimensional plum, we project the original data (blue dots) onto the new most marked axis (yellow and blue crosshairs). The way to find this new coordinate axis is to find the point (red point) where the data is projected on the new coordinate axis and the distance from the origin of the new coordinate axis is the largest, which is the maximum variance method.

<StartOfLatent>

Decentralization adjusts the data such that its mean is zero, which is crucial for PCA. If the data is not centered, the computed variances along the axes will be skewed, leading to misleading results in the identification of principal components. The principal component should ideally capture the largest variance, and uncentered data

2538 would cause the variance to reflect the mean of the data rather than the true 2539 distribution of the data points. 2540 2541 Without decentralization, the projection of the original data onto the new axis would not accurately reflect the inherent structure of the data, undermining the effectiveness 2542 of dimensionality reduction. Thus, the step of data decentralization is not merely a 2543 preprocessing step; it fundamentally ensures that PCA captures the true variance of the 2544 data, allowing for an accurate reduction in dimensionality while retaining as much information as possible. 2545 <EndOfLatent> 2546 2547 ## The first step, data decentralization 2548 The discussion of the new coordinate axis and the old coordinate axis is very 2549 troublesome. We directly decentralize the data (that is, the mean value of the data is 2550 at the far point). 2551 If the data is not decentralized, we cannot find the optimal dimensionality reduction. 2552 This step is necessary 2553 <StartOfLatent> 2554 To find the optimal axis for PCA, we utilize the covariance matrix of the decentralized 2555 data. The eigenvalues of this matrix indicate the amount of variance captured by each principal component, while the corresponding eigenvectors provide the direction of these 2556 components. The first principal component (PC1) is the eigenvector associated with the 2557 largest eigenvalue, representing the direction of maximum variance. 2558 The distance from the projected data points to the origin of the new axis (PC1) is 2559 calculated as the projection of the original data points onto this axis. The objective 2560 is to maximize the sum of squared distances $(\sum d_i^2)$, which mathematically quantifies the 2561 variance captured by the principal component. This optimization problem can be framed as maximizing the Rayleigh quotient for the covariance matrix, leading to the derivation of the eigenvalues and eigenvectors. 2563 <EndOfLatent> 2564 ## The second step is to find the new most marked axis 2565 2566 How do we find the best most standard axis to achieve principal component analysis? That is, the greater the distance between the projected point and the origin of the 2567 coordinate axis, the better (this is the maximum variance) 2568 2569 As shown in the figure, the red dotted line is the new coordinate axis, we call it PC1; the green dot is the initial data sample point; the green cross is the point projected 2570 on the new coordinate; d1, d2, d3...d6 are the projected ones The distance from the 2571 point to the far point. 2572 All we need to do is find the largest sum of squares (that's $\sum d_i^2$). 2573 2574 Here comes the math! !!! 2575 <StartOfLatent> 2576 To find the optimal new coordinate axis for dimensionality reduction, we must analyze 2577 the relationships among the data points. This involves calculating the correlation coefficient matrix, which quantifies how much each feature varies with others. 2578 2579 Once we have the correlation matrix, we compute its eigenvalues and eigenvectors. The 2580 eigenvalues indicate the amount of variance captured by each corresponding eigenvector. The eigenvector with the highest eigenvalue gives us the direction of the axis that 2581 maximizes variance among the data points. 2582 In principal component analysis (PCA), we focus on the principal components (PCs) that 2583 represent the highest variance first. This ensures that we retain the most important 2584 information while reducing dimensions. We select a percentage of these components, based 2585 on the cumulative variance they explain, to maintain a balance between data fidelity and dimensionality reduction. 2586 2587 By projecting the original data onto these new axes (PC1, PC2, etc.), we can effectively 2588 reduce the number of features while preserving the essential structure and relationships within the data. The choice of how many components to keep is guided by 2589 examining the explained variance ratios. The goal is to achieve a compact representation 2590

that still captures the underlying patterns in the dataset.

2592 | <EndOfLatent

As can be seen from the figure above, when we are looking for the best newest axis, we are actually finding eigenvalues and eigenvectors by finding the data correlation coefficient matrix

The third step is to choose a few percent of the data you need

When we perform PCA dimensionality reduction to find new coordinates, the number of coordinates is the same as the number of data features. But we project the data on the new coordinate axis in order to express the information of the entire data with as few features as possible.

In the figure below, we can see that we have obtained two coordinate axes, PC1 and PC2, respectively.

<StartOfLatent>

To understand why the coordinate axes PC1 and PC2 are chosen, we first recognize that PCA aims to maximize the variance captured by the new axes. The eigenvalues obtained from the eigenvalue decomposition of the data correlation coefficient matrix indicate the variance explained by each principal component.

The eigenvector corresponding to the largest eigenvalue defines the direction of PC1, capturing the most significant variance in the data. The second eigenvector, corresponding to the second largest eigenvalue, defines PC2 and is orthogonal to PC1, ensuring that these axes are independent and do not introduce redundancy in the representation of the data.

The information captured by each principal component can be quantified by the proportion of the total variance they explain. In this case, PC1 accounts for 83% of the variance, indicating that it captures the most critical features of the data. PC2, while still informative, only captures 17% of the variance.

When reducing dimensions, it is logical to prioritize PC1 over PC2 due to its higher variance contribution. This prioritization allows for a more efficient representation of the data with fewer dimensions, simplifying analysis while retaining essential information

In three dimensions, using only PC1 and PC2 to represent the data can effectively capture the underlying structure without the noise, as the dimensionality reduction is inherently aimed at compressing the data while minimizing loss. This approach demonstrates PCA's effectiveness in unsupervised learning, particularly for data compression and denoising, making it a widely adopted technique in various practical applications.

The two coordinate axes are perpendicular to each other and do not interfere with each other.

Among them, the information on the data on PC1 accounts for 83%, and the information on the data on PC2 accounts for 17%. Anyone can see which coordinate axis should be chosen to represent the information after dimensionality reduction.

When we look at three dimensions, we can see that if we use PC1 and PC2 to reduce the dimension

As a dimensionality reduction method for unsupervised learning, it only needs eigenvalue decomposition to compress and denoise data . Therefore, it is widely used in practical scenarios.

<StartOfLatent>

When performing PCA, we seek to maximize the variance captured by each principal component. Since variance is central to determining the principal components, squaring the differences ensures that components capturing larger variances are prioritized. This aligns with the goal of PCA, which is to identify directions in the data that account for the greatest amount of variability.

The orthogonality of principal components guarantees that the information captured by each component is independent of the others. This independence is crucial for effective dimensionality reduction because it ensures that the chosen components do not overlap in the information they represent. The simplification of calculations through eigenvalue

decomposition further supports this process, as it systematically identifies the components associated with the highest variance.

In terms of interpretation, while PCA effectively reduces dimensionality and retains the most significant variance, it also introduces ambiguity regarding the meaning of the principal components. Each principal component is a linear combination of the original features, making it challenging to directly interpret their significance in relation to the original data.

The potential loss of small variance components is another critical consideration. Although these components may seem less important, they could contain subtle but meaningful differences in the dataset. Discarding them may lead to overlooking valuable insights during subsequent analysis.

Lastly, the effectiveness of PCA is contingent on the assumptions of data distribution. PCA assumes that the data is normally distributed; deviations from this assumption can lead to suboptimal results. Therefore, understanding the underlying data characteristics is essential for ensuring PCA's effectiveness in dimensionality reduction and information retention.

The main advantages of the PCA algorithm are:

- 1) Only the variance needs to be used to measure the amount of information, and it is not affected by factors other than the data set.
- 2) The principal components are orthogonal to each other, which can eliminate the mutual influence factors between the original data components.
- 3) The calculation method is simple, and the main operation is eigenvalue decomposition, which is easy to implement.

The main disadvantages of the PCA algorithm are:

- 1) The meaning of each feature dimension of the principal component has a certain degree of ambiguity, which is not as strong as the interpretation of the original sample features.
- 2) Non-principal components with small variance may also contain important information on sample differences, and discarding due to dimensionality reduction may have an impact on subsequent data processing.
- 3) When the data distribution is not a normal distribution, the effect is not very good

<StartOfLatent>

To understand why variance is squared in PCA, we must first recognize that variance is a measure of the spread of data points in a dataset. Squaring the differences from the mean when calculating variance emphasizes larger deviations. This is important because larger deviations contribute more significantly to the overall variability. By squaring these differences, we ensure that both positive and negative deviations contribute positively to the variance, which allows us to assess the total spread of the data accurately. Without squaring, positive and negative distances could nullify each other, leading to misleading interpretations of data relationships. Squaring transforms all distances into non-negative values, ensuring that the overall variance remains a measure of spread without cancellation effects.

Using the correlation coefficient matrix instead of the covariance matrix is fundamental when attributes have different units. The correlation coefficient standardizes the data, allowing for comparison across variables that are measured on different scales. This standardization is achieved by dividing the covariance by the product of the standard deviations of the two variables involved. Consequently, it enables a clear interpretation of the strength and direction of relationships without the influence of differing units.

Dimensioning in PCA is necessary when attributes have different units. If attributes are in the same unit, covariance can directly reflect the degree of correlation. However, when units differ, covariance alone is insufficient as it does not convey the actual strength of the relationship. The correlation coefficient, derived by dividing the covariance by the standard deviations, effectively normalizes the data, ensuring that the relationships are comparable. This normalization is vital for accurately assessing how similar the attributes are in their variations.

2700 2701 # question 2702 2703 1. Why square it ? 2704 It is because if the distance is not squared, there will be positive and negative 2705 distances, which will cancel each other out. 2706 2.**Why is the correlation coefficient matrix instead of the covariance matrix? ** You will know this when you look down 2707 2708 # Does PCA need to be dimensioned? 2709 (1) When: when the unit of each attribute is the same (for example, both are kg, both 2710 are meters), each attribute is comparable. Therefore, it is enough to directly calculate 2711 the covariance between attributes. The size of the original covariance does not indicate the degree of correlation (covariance only indicates positive or negative 2712 correlation), But when the units are the same, we can think that the greater the 2713 covariance, the greater the correlation 2714 (2) When the units of each attribute are different (for example, one is kg and the other 2715 is meter), at this time, due to the different units, the covariance does not indicate 2716 the degree of correlation. At this time, we need to use the correlation coefficient to 2717 2718 The formula of the correlation coefficient (that is, the correlation coefficient matrix 2719 is divided by two standard deviations, where dividing by the standard deviation is a way of dimensioning). It eliminates the influence of the range of change of two variables, 2720 but simply reflects the degree of similarity between the two variables per unit change. 2721 2722 ### Guess you like 2723 Origin blog.csdn.net/CSTGYinZong/article/details/127097464 2724 Recommended 2725 Ranking Daily 2726 2727

Example G.2.2: Latent thoughts elaborate physical knowledge

<StartOfLatent>

2728

272927302731

2732

2733

2734

2735

2736

2737 2738

2739

2740

2741

27422743

2744

2745

2746

2747

27482749

2750

2751 2752

2753

The draining of water creates a vortex due to fluid dynamics principles. In fluid mechanics, when water flows towards a drain, it encounters resistance and varying velocity, resulting in a circular motion. The initial swirl can be caused by a small disturbance, which can be a random fluctuation in water movement. As this disturbance propagates, it influences adjacent water particles, leading to a self-reinforcing loop of motion. This phenomenon can be described by the Navier-Stokes equations, which govern fluid flow and demonstrate how momentum and viscosity interact to create vortices.

In physics, time and space are interwoven in the fabric of the universe, as described by the theory of relativity. Just as water spirals down a drain, objects in spacetime can exhibit similar swirling patterns due to gravitational forces. Massive objects, like planets or stars, warp spacetime around them, creating paths that can be likened to the water's spiral. The curvature of spacetime can be mathematically represented by Einstein 's field equations, which illustrate how mass influences the geometry of the universe.

The analogy suggests that just as water spirals towards a drain through self-reinforcing interactions, objects in spacetime move along geodesics-paths determined by the curvature of spacetime. The motion of celestial bodies can be influenced by the gravitational pull of nearby masses, creating a dynamic interplay that mirrors the swirling water. This interconnectedness in physical phenomena reveals deeper insights into the nature of motion and force, both in fluids and in the cosmos.

Circling the Drain

In Relativity by Brian Koberlein4 Comments

If you've ever watched water drain from a bathtub, you've seen that it doesn't flow into the drain in a straight line, but rather swirls around the drain. Contrary to popular

belief, this is not due to the rotation of the Earth but rather random currents in the water. The reason the spiral forms is because it is self-reinforcing. Water near the drain spirals a bit due to a random current, which drags the water behind it slightly in the same direction. Over time, the motion builds up until you have a rapid spiral around the drain.

There is a similar effect with space and time.

<StartOfLatent>

Gravity is a fundamental force that can be understood through the concepts of space and time. The relationship between mass and the curvature of space-time is a key aspect of Einstein's General Relativity. When mass is present, it causes a distortion in the fabric of space-time, resulting in a gravitational field. This distortion leads objects to follow geodesics, which are the curved paths determined by the curvature of space-time, rather than straight lines.

The concept of frame dragging is significant in understanding how rotating masses affect nearby space-time. When a massive object, like Earth, rotates, it imparts motion to the surrounding space-time. This effect can be quantified through mathematical formulations derived from the Einstein field equations, which relate the mass-energy content of an object to the curvature of space-time.

The Lense-Thirring effect specifically describes how the rotation of a massive body influences the orbits of nearby objects. It can be expressed in terms of the angular momentum of the rotating mass and its distance from the orbiting object. The mathematical representation allows for predictions of the precession of orbits in the vicinity of rotating masses.

While the analogy to water spiraling down a drain illustrates the concept of self-reinforcement in motion, it is important to clarify that space-time does not exhibit fluid dynamics. The spiral effect in space-time due to frame dragging is a result of geometric properties rather than a physical flow. This distinction is crucial for accurately understanding the implications of relativity in astrophysical contexts.

Through careful measurements and experiments, such as those conducted with satellites and gyroscopes, the frame dragging effect can be observed, confirming the predictions made by General Relativity. These observations provide evidence for the intricate relationship between mass, rotation, and the curvature of space-time, enhancing our understanding of gravitational phenomena in the universe.

I've talked before about how gravity is due to a curvature of space and time. The presence of mass bends space around it, and the resulting curvature means that objects follow curved paths rather than straight ones. As a result, an object's motion near a large mass like the Earth looks as if it is due to a force, which we call gravity.

It turns out that the rotation of a mass also distorts space and time. For example, as the Earth rotates, it drags the nearby space along with it (an effect known as frame dragging). Just like the drain spiral, this effect builds up, and as a result, space spirals a bit around the Earth. You have to be a bit careful with this comparison. Spacetime doesn't "flow" the way water does, but the spiral effect is somewhat similar.

Near the Earth, this frame dragging is very small, but it can be measured through an effect known as the Lense-Thirring effect.

<StartOfLatent>

The gyroscope experiment measures the effect of frame dragging caused by Earth's rotation. A gyroscope maintains its axis of rotation due to the conservation of angular momentum. If frame dragging occurs, the space around the gyroscope changes, altering its path of motion. This results in precession, where the gyroscope's axis shifts over time.

The measurement intervals in the experiment are likely influenced by the orbital period of the gyroscope, which can be approximately 1000 days. This period corresponds to the time it takes for the gyroscope to complete multiple orbits around the Earth, allowing sufficient time for the frame dragging effects to accumulate.

The Lense-Thirring effect is mathematically modeled using general relativity principles. The precession rate can be calculated with the formula for frame dragging, which involves the mass of the rotating body, the distance from the center of mass, and the angular momentum.

The agreement between the observed and predicted precession validates the theoretical framework of general relativity. The consistent measurements reinforce the understanding of how rotating masses influence spacetime geometry.

The analogy of water spiraling down a drain illustrates the concept of frame dragging but simplifies complex spacetime dynamics. Understanding these dynamics requires knowledge of tensor calculus and the geodesic equations governing motion in curved spacetime.

Overall, the experiment demonstrates a fundamental aspect of general relativity, linking mass, rotation, and the structure of spacetime, while providing empirical evidence for theoretical predictions.

<EndOfLatent>

Basically, you put a gyroscope in orbit and see if its axis of rotation changes. If there is no frame dragging, then the orientation of the gyroscope shouldn't change. If there is frame dragging, then the spiral twist of space and time will cause the gyroscope to precess, and its orientation will slowly change over time.

We've actually done this experiment, and you can see the results in the figure below. The black line represents the change in orientation over time, and the red line is the predicted change via the Lense-Thirring effect. As you can see, they agree very well.

So the next time you watch water circling the drain, you will know that a similar effect occurs with space and time itself.

1. Do you know why there is max in measurements roughly every 1000 days?

<StartOfLatent>

The gyroscope's behavior in orbit is influenced by gravitational interactions and the geometry of spacetime. The concept of frame dragging arises from general relativity, where massive objects like Earth distort the curvature of spacetime around them. This distortion affects objects in orbit.

The maximum in measurements every 1000 days can be attributed to the elliptical nature of the orbit. An elliptical orbit means that the distance between the gyroscope and Earth varies periodically. The gravitational force experienced by the gyroscope is strongest when it is closest to Earth, at the perigee of its orbit.

The Lense-Thirring effect predicts that frame dragging will vary depending on the distance to the massive body. As the gyroscope approaches the perigee, the gravitational influence increases, leading to a greater precession rate. Thus, the maximum frame dragging occurs at this closest point, resulting in a peak in the gyroscope's orientation change.

The period of 1000 days corresponds to the orbital characteristics of the gyroscope, defined by Kepler's laws of planetary motion. These laws describe how the orbiting body travels faster when closer to the massive body and slower when farther away. The relationship between the orbital period and the semi-major axis can be expressed using Kepler's third law, where the square of the orbital period is proportional to the cube of the semi-major axis of the orbit.

In summary, the gyroscope's orbit experiences maximum frame dragging as it reaches its closest point to Earth every 1000 days due to the predictable nature of its elliptical orbit, influenced by gravitational forces and the curvature of spacetime.

<EndOfLatent>

- 1. Space has to wind to some local maximum of curvature, and then snap back. Since Earth 's mass is constant, this maximum winding should be predictable.
- 2. What i think is that , by the time you put the gyroscope in orbit there is a max distance and a min distance from earth (supposing that orbit is an elipse) . So if the period is almost 1000 days then every 1000 days it will reach the min distance from earth where the gravity field is max so the frame dragging should be max .

<StartOfLatent>

The phenomenon of water spiraling down a drain is influenced by various factors beyond mere random currents. The Coriolis effect, resulting from Earth's rotation, plays a significant role in determining the direction of fluid motion. The Earth's rotation creates a force that affects moving objects, including fluids, causing them to follow a

curved path rather than a straight line. This effect is more pronounced over larger 2863 distances and timescales, but it can still impact smaller systems like water draining in 2864 a sink. 2865 If the motion were entirely random, one would expect a uniform distribution of spiral 2866 directions. However, the consistent direction observed suggests a systematic influence, which can be explained by the Coriolis effect. The interaction between the water's 2868 initial velocity, the shape of the drain, and the rotation of the Earth results in a dominant rotational direction. 2870 Mathematically, the Coriolis force can be expressed as $F_c=2m(v)(\Omega)\sin(\phi)$, where m is the mass of the fluid element, v is its velocity, Ω is the angular velocity of the Earth 2871 , and ϕ is the latitude. This force acts perpendicular to the velocity of the fluid, 2872 leading to a predictable rotation direction. 2873 Moreover, factors such as the geometry of the sink and the initial conditions of the 2874 water's movement can create a bias towards a particular direction of spin. The 2875 conservation of angular momentum further supports this, as any slight initial spin will 2876 be amplified due to the closed system of the drain. Thus, while randomness may contribute to local turbulence, the underlying forces impart a consistent rotational 2877 behavior that aligns with the Earth's rotation. 2878 <EndOfLatent> 2879 2. If the water spiralling down the drain is just due to "random currents in the water" 2880 and not the Earths rotation, why is the spiral always the same direction of rotation? If 2881 it was truly random, surely it would be 50/50 on direction of rotation? 2882

Example G.2.3: Latent thoughts elaborate mathematical derivations

```
2885
            <StartOfLatent>
            To solve the integral of the product of tangent functions, we start by recognizing that
2887
            the tangent identity can simplify the problem. The tangent function can be expressed in
2888
            terms of sine and cosine. We can also leverage the angle addition formula for tangent.
2889
            We define t_n = \tan(nx). The product of tangents can be expressed using the identity for
2890
            tangent of sums:
2891
2892
            \tan A \tan B = \frac{\tan A + \tan B}{1 \tan A \tan B}
2893
            This allows us to express \tan(5x) in terms of \tan(3x) and \tan(2x):
2894
2895
            \tan(5x) = \frac{\tan(3x) + \tan(2x)}{1 - \tan(3x)\tan(2x)}
2896
            Rearranging this gives us:
2897
            \tan(5x) - \tan(3x) - \tan(2x) = \tan(5x)\tan(3x)\tan(2x)
2899
            This leads to the equation:
2900
2901
            t_5 - t_3 - t_2 = t_5 t_3 t_2
2902
            To integrate, we substitute this result back into the integral:
2903
2904
            \int \tan(5x)\tan(3x)\tan(2x)\,dx = \int (\tan(5x) - \tan(3x) - \tan(2x))\,dx
2905
            Now, we can integrate each term separately. The integral of \tan(nx) is known:
2906
2907
             \int \tan(nx) \, dx = -\frac{1}{n} \log|\cos(nx)| + C
2908
            Applying this to our terms yields:
2909
2910
            1. For \tan(5x): \int \tan(5x) dx = -\frac{1}{5} \log|\cos(5x)|
2911
            2. For \tan(3x): \int \tan(3x) dx = -\frac{1}{3} \log|\cos(3x)|
2912
            3. For \tan(2x): \int \tan(2x) \, dx = -\frac{1}{2} \log|\cos(2x)|
2913
2914
            Substituting these results back into the integral, we compile the contributions:
2915
```

 $\int \tan(5x)\tan(3x)\tan(2x)\,dx = -\frac{1}{5}\log|\cos(5x)| + \frac{1}{3}\log|\cos(3x)| - \frac{1}{2}\log|\cos(2x)| + C$ Finally, we can rearrange this into a single expression if desired, but it is already presented clearly as separate contributions from each tangent function. The solution is complete with the integration constants summed together at the end. <EndOfLatent> # How to find the integral $\int \tan(5x) \tan(3x) \tan(2x) dx$? Find the integral $\int \tan(5x) \tan(3x) \tan(2x) dx$. This question is posted in a maths group in Facebook. What way should we use to solve integral like this? Thanks in advance #### Solutions Collecting From Web of "How to find the integral $\int \tan(5x)\tan(3x)\tan(2x) dx$?" Let $t_n = \tan(nx)$, we have $t_5 = \frac{t_3 + t_2}{1 - t_3 t_2} \iff t_5 - t_5 t_3 t_2 = t_3 + t_2 \implies t_5 t_3 t_2 = t_5 - t_3 - t_2$ $\int \tan(5x)\tan(3x)\tan(2x)dx = \int (\tan(5x) - \tan(3x) - \tan(2x)) dx$ $=\frac{1}{2}\log\cos(2x)+\frac{1}{3}\log\cos(3x)-\frac{1}{5}\log\cos(5x)+$ const.

Example G.2.4: Latent thoughts demystify complex Mathematica logs

<StartOfLatent>

 To evaluate the integral $\int e^{3\coth^{-1}(ax)}(c-\frac{c}{ax})^{3/2}\,dx$, we begin by analyzing the term $e^{3\coth^{-1}(ax)}$. The hyperbolic cotangent inverse function, $\coth^{-1}(u)$, can be expressed in terms of logarithmic functions: $\coth^{-1}(u)=\frac{1}{2}\ln\left(\frac{u+1}{u-1}\right)$. This implies that $e^{3\coth^{-1}(ax)}$ transforms into an exponential function involving ax and logarithmic components.

Next, we need to explore the term $(c-\frac{c}{ax})^{3/2}$. This indicates a dependency on x that can simplify under certain substitutions. Notably, as x approaches infinity, this expression simplifies to $c^{3/2}$, suggesting that limits and behavior at infinity might be important in evaluating the integral.

Rearranging the term yields $(c(1-\frac{1}{ax}))^{3/2}$. This opens up a potential substitution where we let u=ax. Consequently, dx transforms into $\frac{du}{a}$, leading to a new form of the integral in terms of u.

The expression requires careful treatment of the limits of integration and potential singularities, particularly considering the behavior of $\frac{1}{ax}$ as x approaches zero.

To derive the result, we will apply integration techniques that may include integration by parts or special functions. The presence of $e^{3\coth^{-1}(ax)}$ suggests that the integral is non-elementary, potentially requiring numerical methods or special functions for evaluation.

The final expression involves several terms combined from the integration steps:

- 1. The first term appears to arise from the integration of the polynomial-like expression after applying the substitution and integration techniques.

 2. The second term involves derivatives of the logarithmic form that comes from differentiating the exponentiated terms.
- 3. The third term stems from integrating the product of the hyperbolic function and the resultant polynomial.

Each of these components is derived by applying the fundamental theorem of calculus, respecting the chain rule where appropriate, and ensuring all terms are accounted for in the final assembly of the integral's evaluation.

This leads to the concluding expression, which captures the complexity of the integral and reflects the interplay between exponential, logarithmic, and polynomial forms. The combination of these terms indicates a sophisticated integration process that likely necessitates verification through differentiation to confirm correctness. <EndOfLatent>

3.458 $\int e^{3 \coth^{-1}(ax)} (c - \frac{c}{ax})^{3/2} dx$

Optimal. Leaf size=118
$$\frac{c^3 x \left(1 - \frac{1}{a^2 x^2}\right)^{3/2}}{\left(c - \frac{c}{ax}\right)^{3/2}} - \frac{3c^2 \sqrt{1 - \frac{1}{a^2 x^2}}}{a \sqrt{c - \frac{c}{ax}}} + \frac{3c^{3/2} \tanh^{-1} \left(\frac{\sqrt{c} \sqrt{1 - \frac{1}{a^2 x^2}}}{\sqrt{c - \frac{c}{ax}}}\right)^{3/2}}{a}$$

[Out]

 $(-3*c^2*Sqrt[1 - 1/(a^2*x^2)])/(a*Sqrt[c - c/(a*x)]) + (c^3*(1 - 1/(a^2*x^2))^(3/2)*x)/(c - c/(a*x))^(3/2) + (3*c^3/2)*ArcTanh[(Sqrt[c]*Sqrt[1 - 1/(a^2*x^2)])/Sqrt[c - c/(a*x)])/a$

<StartOfLatent>

The integral involves the function $e^{3\coth^{-1}(ax)}$, where the inverse hyperbolic cotangent function, ArcCoth, is transformed to ArcTanh for simplification in the calculations. The transformation is made using the identity that relates the two functions, leveraging properties of hyperbolic functions.

The expression $(c-\frac{c}{ax})^{3/2}$ in the integrand indicates that the variable dependencies of c and a must be handled carefully. This leads to the necessity of simplifying the expression to isolate the variable x and expressing it in a form amenable to integration.

The optimal antiderivative size of 118 signifies the complexity of the expression being handled and the efficiency of the algorithm used. The normalized size of 1 indicates that the final expression is effectively simplified to its most concise form, making it easier to handle computationally.

The number of steps used, 5, denotes the sequential approach taken to derive the antiderivative, indicating that multiple integration techniques may have been employed, such as substitution or integration by parts.

The rules employed, specifically 6177, 863, 865, 875, and 208, correspond to integration techniques or transformations utilized in the process. Each rule applies specific mathematical properties or identities that facilitate the evaluation of the integral.

The final expression includes three distinct terms, each derived from different components of the original integrand. The term involving ArcTanh suggests the utilization of hyperbolic identities, while the other terms reflect polynomial behavior modulated by square root functions.

The successful verification of the antiderivative confirms that the computed result matches the expected behavior of the integral, adhering to the properties of continuous functions over the specified domain of integration, further supporting the correctness of the applied rules and techniques.
<EndOfLatent>

Rubi [A] time = 0.214601, antiderivative size = 118, normalized size of antiderivative = 1., number of steps used = 5, number of rules used = 5, integrand size = 24, $\frac{\text{number of rules}}{\text{integrand size}}$ = 0.208, Rules used = {6177, 863, 865, 875, 208}

$$\frac{c^3x\left(1-\frac{1}{a^2x^2}\right)^{3/2}}{\left(c-\frac{c}{ax}\right)^{3/2}}-\frac{3c^2\sqrt{1-\frac{1}{a^2x^2}}}{a\sqrt{c-\frac{c}{ax}}}+\frac{3c^{3/2}\tanh^{-1}\left(\frac{\sqrt{c}\sqrt{1-\frac{1}{a^2x^2}}}{\sqrt{c-\frac{c}{ax}}}\right)}{a}$$

Antiderivative was successfully verified.

[In]

 $Int[E^{(3*ArcCoth[a*x])*(c - c/(a*x))^{(3/2),x}]$

'

3025 3026

3027

3028 3029

3030 3031

3032 3033

```
[Out]

(-3*c^2*Sqrt[1 - 1/(a^2*x^2)])/(a*Sqrt[c - c/(a*x)]) + (c^3*(1 - 1/(a^2*x^2))^(3/2)*x)/(c - c/(a*x))^(3/2) + (
3*c^(3/2)*ArcTanh[(Sqrt[c]*Sqrt[1 - 1/(a^2*x^2)])/Sqrt[c - c/(a*x)]])/a

Rule 6177
```

G.3 EXAMPLES OF TINYLLAMA-1B SELF-GENERATED LATENT THOUGHTS

Example G.3.1: Models think through the data for learning

```
3034
3035
             <StartOfLatent>
            Understanding the concept of factorial requires knowledge of basic arithmetic operations
3036
             and the rules governing multiplication. A factorial is defined for a non-negative
3037
             integer n, expressed as n!. This definition implies that the result of n! is calculated
3038
            by taking n and multiplying it by the next integer down from it, thereby accumulating
            all integers from 1 to n.
3039
3040
            To derive n! explicitly, one can compute it using the formula:
3041
            n! = n \cdot (n-1) \cdot (n-2) \cdot \dots \cdot 2 \cdot 1.
3042
3043
            This recursive relationship can be visualized as multiplying the number by the number
            immediately preceding it, illustrating how each additional integer reduces the result.
3044
3045
            The factorial function grows extremely fast, with the value of 1! being 1 and 2! being 2.
3046
             This fact showcases the exponential growth of factorials as n increases.
3047
            In contexts such as probability, combinatorics, and calculus, the factorial is pivotal
3048
            in calculating permutations and combinations, where arrangements and selections from a
            set are analyzed. For example, the number of ways to arrange n distinct objects is given
3049
             by n!, which explains its widespread application in these fields.
3050
3051
            Thus, the factorial concept signifies a foundational element in mathematics, bridging
            basic arithmetic with more complex applications in various scientific and mathematical
3052
            disciplines. Understanding its calculation and implications is crucial for further
3053
            mathematical exploration and real-world problem-solving.
3054
            <EndOfLatent>
3055
            # Factorial number
3056
3057
            In mathematics symbolized by placing the "!" (known as the exclamation mark or bang)
            after a number, it represents multiplying a number by all whole numbers smaller than it.
3058
3059
            ## Definition
3060
            A factorial is defined by the product
3061
3062
            n! = 1 \cdot 2 \cdot 3 \cdot \cdot \cdot (n{-}2) \cdot (n{-}1) \cdot n
3063
            for n > 1.
3064
3065
            <StartOfLatent>
3066
            The factorial notation represents the product of descending positive integers, which
3067
            means for any integer n, the factorial is calculated as n multiplied by the product of
3068
            all positive integers from (n-1) down to 1. This can be expressed mathematically as n! = n
             n \times (n-1) \times (n-2) \times ... \times 2 \times 1.
3069
3070
            For n = 0, the factorial is defined as 0! = 1 because the product of no numbers is 1.
3071
            This is a foundational rule in combinatorics, as it ensures that there is one way to
            choose zero objects.
3072
3073
            To evaluate larger factorials, the recurrence relation is useful. This relation is
            derived from the definition of factorial, indicating that n! can be computed directly or
3074
             recursively based on smaller factorials. Specifically, it states that n! is equal to n
3075
            multiplied by (n-1)!, which means to find the factorial of n, one can take the factorial
3076
             of the factorial of (n-1).
3077
```

```
3078
            The computational implementation in computer programming languages like Python utilizes
3079
            this definition directly for standard input, while Java and C# provide corresponding
3080
            classes to facilitate easy access and manipulation of factorial values.
3081
            For example, to calculate 5!, one applies the defined multiplication: 5! = 5 \times 4 \times 3 \times 2
3082
             \mathbf{x} 1, ultimately simplifying to 120. This demonstrates the practical application of the
3083
             factorial definition in programming contexts, where efficient computation of large
3084
            factorial values is necessary.
            <EndOfLatent>
3085
3086
            The same written as mathmatical product
3087
            n! = \prod_{i=1}^{n} i.
3088
3089
            and as recurrence relation
3090
            n! = n \cdot (n-1)!
3091
3092
            ## Examples
3093
            5! = 5 * 4 * 3 * 2 * 1 = 120
3094
            10!
3095
            <StartOfLatent>
3096
            To derive the value 10!, we need to multiply 10 distinct positive integers: 1, 2, 3, 4,
3097
            5, 6, 7, 8, 9, and 10. Each integer contributes to the factorial, which is a product of
            all integers from 1 to n. The multiplication of these integers provides the result:
3098
3099
            1! = 1
3100
            2! = 2
            3! = 6
3101
            4! = 24
3102
            5! = 120
3103
            6! = 720
            7! = 5040
3104
            8! = 40320
3105
            9! = 362880
3106
            10! = 3628800
3107
            This demonstrates the factorial's recursive nature, where each factorial is defined from
3108
             the product of the integers preceding it. The first factorial starts with 1! equal to
3109
            1, establishing the base case. The subsequent calculations follow the established
            pattern of multiplication for increasing integers.
3110
3111
            The base cases for factorial calculations include:
3112
            - n! = n for n = 1, where the product of all integers equals 1.
            - n! = 1 for n = 0, as it is defined that 0! = 1.
3113
            - n! = 0 for n < 0, as established conventionally in factorial definitions.
3114
            In the case of higher-order factorials, the pattern continues with:
3115
            - n! = n for odd n, and 1 for n = 0.
3116
            - n! = n! div n for n > 0, denoted as n / 1!.
3117
            This provides a comprehensive understanding of how factorials are computed across
3118
            various integers, including the calculation of 10!, confirming the established sequence
3119
            of factorial products.
            <EndOfLatent
3120
3121
            = 10 * 9 * 8 * 7 * 6 * 5 * 4 * 3 * 2 * 1 = 3628800
3122
```

Example G.3.2: Models learn math problem solving with latent thoughts

3123 3124

3125 3126

3127

3128

3129 3130

```
<StartOfLatent>
To determine the percent increase and decrease, we first need to understand that
percentages represent a portion of a whole, hence percent increase or decrease is a way
to express this portion in terms of 100.

For an increase, the formula is:
```

```
3132
            Percent Increase = (New Value - Original Value) / Original Value * 100.
3133
3134
            For a decrease, the formula is similar:
3135
            Percent Decrease = (Original Value - New Value) / Original Value * 100.
3136
3137
            In the calculation of an increase in pay from $10.00 to $12.00, we observe that the
3138
            increase is $2.00. Thus, the percent increase can be calculated as follows:
3139
            1. Identify the original value (PV = $10.00) and the new value (NPV = $12.00).
3140
            2. Calculate the difference: NPV - PV = $12.00 - $10.00 = $2.00.
            3. Divide the difference by the original value: $2.00 / $10.00 = 0.2.
3141
            4. Multiply by 100 to convert to a percentage: 0.2 \times 100 = 20\%.
3142
3143
            This indicates that Ann's pay increased by 20%.
3144
            For the example involving the sale price of $50.00 after a 20% discount, we follow a
3145
            similar process. The discount represents a reduction in price, equivalent to a decrease
3146
            of $20 from the original price. Therefore, we can calculate:
3147
            1. Identify the original price (\$50.00) and the discounted price (\$50.00 - \$20.00 =
3148
            2. Calculate the percentage decrease: (\$30.00 / \$50.00) * 100 = 60\%.
3149
3150
            This shows that the sale price is 60% of the original price, indicating a substantial
3151
            reduction from the original value.
            <EndOfLatent>
3152
3153
            # Calculating Percentage Increase And Decrease
3154
3155
            ## Learn About Calculating Percent Increase And Decrease With The Following Examples And
             Interactive Exercises.
3156
3157
            Example 1: Ann works in a supermarket for $10.00 per hour. If her pay is increased to
            $12.00, then what is her percent increase in pay?
3158
3159
            Analysis: When finding the percent increase, we take the absolute value of the
3160
            difference and divide it by the original value. The resulting decimal is then converted
            to a percent.
3161
3162
            Solution:
3163
            Answer: The percent increase in Ann's pay is 20%.
3164
3165
3166
            <StartOfLatent>
3167
            To calculate a percent decrease, it's important to understand the concept of a
3168
            percentage as a way to express a number as a fraction of 100. A percent decrease is
            essentially finding the remaining value after a certain amount is taken away from a
3169
            total.
3170
            In the case of the staff at a company, the number of employees decreased from 40 to 29.
3171
            To find the percent decrease, we first determine the difference between the new and old
3172
            values. This difference is 29 - 40 = -11. The next step involves finding the absolute
3173
            value of this difference, which is 11.
3174
            To convert this absolute value into a decimal, we divide 11 by 40. This calculation
3175
            results in 0.275. To express this decimal as a percent, we multiply 0.275 by 100,
3176
            converting it entirely into a percent form. This results in a percent decrease of 27.5%,
             indicating that the number of employees has decreased by 27.5% from the original count.
3177
3178
            This process illustrates the mechanics of both increasing and decreasing values in terms
3179
             of percentages, which is crucial for interpreting life scenarios like salary or
            employment changes. Understanding these calculations enables accurate financial planning
3180
             and assessment of trends in employment numbers.
3181
            <EndOfLatent>
3182
            Let's look at an example of percent decrease.
3183
3184
```

3186 Example 2: The staff at a company went from 40 to 29 employees. What is the percent 3187 decrease in staff? 3188 3189 Analysis: When finding the percent decrease, we take the absolute value of the difference and divide it by the original value. The resulting decimal is then converted 3190 to a percent. 3191 3192 Solution: 3193 3194 <StartOfLatent> 3195 To determine the percent increase, first identify the original value and the new value. 3196 The original pay is \$10.00, and the new pay is \$12.00. Calculate the difference between 3197 these amounts: \$12.00 - \$10.00 = \$2.00. Next, divide this difference by the original pay : \$2.00 / \$10.00 = 0.2. Multiply the result by 100 to convert it to a percent: $0.2 \, \star \, 100$ 3198 = 20%. Thus, Ann's pay increased by 20% during her hourly wage increase. 3199 3200 In the second example, we analyze the staff reduction from 40 to 29. The new value is 29 represent the reduced employees. Subtract the original count from the new count: 29 -3201 40^{-} -11. This negative result indicates an increase in the number of staff; therefore, 3202 the percent decrease is calculated as follows: (negative) / (original) = 11 / 40 = 0.275. 3203 Convert this decimal into a percent by multiplying by 100: 0.275 * 100 = 27.5%. Thus, there was a 27.5% decrease in staff. 3204 <EndOfLatent> 3205 Answer: There was a 27.5% decrease in staff. 3206 3207 Percent increase and percent decrease are measures of percent change, which is the 3208 extent to which something gains or loses value. Percent changes are useful to help 3209 people understand changes in a value over time. Let's look at some more examples of percent increase and decrease. 3210 3211 In Example 1, we divided by 10, which was the lower number. In Example 2, we divided by 40, which was the higher number. Students often get confused by this. Remember that the 3212 procedure above asked us to divide by the original value. Another way to remember the 3213 procedure is to subtract the old value from the new value and then divide by the old 3214 value. Convert the resulting decimal to a percent. The formula is shown below. 3215 3216 3217 <StartOfLatent> To find the percent increase in a cost, we first need to understand the concept of 3218 percentage change. The formula for percent increase is: 3219 3220 Percent Increase = (New Value - Original Value) / Original Value * 100. 3221 In the first example, the original quantity of the item is 75 cents per pound, and it increases to 81 cents per pound. The absolute increase in cost is calculated as: 3222 3223 Absolute Increase = New Value - Original Value = 81 - 75 = 6 cents. 3224 3225 Next, we determine the decimal representation of the increase: 3226 Absolute Increase = 6 cents. 3227 3228 Since we are looking for the percent increase, we divide this absolute increase by the 3229 3230 Percent Increase = (6 cents / 75 cents) * 100 = 0.08 or 8%. 3231 This calculation shows that the cost of the item increased by 8% percentage. 3232 3233 In the second example, when four feet are removed from a 12-foot board, the original length is 12 feet. We need to find the decrease in length. The length removed is 4 feet. 3234 Therefore, the decrease in length is: 3235 Decrease in Length = Original Length - Removed Length = 12 - 4 = 8 feet. 3236 3237 Next, we calculate the percent decrease in length: 3238 3239

```
3240
            Percent Decrease = (Decrease in Length / Original Length) * 100 = (8 / 12) * 100 = 66%.
3241
3242
            Both calculations show how to calculate percent decrease based on original and new
3243
            values. Understanding these steps is crucial for effectively applying percentage
            calculations in various contexts.
3244
            <EndOfLatent>
3245
3246
            Example 3: At a supermarket, a certain item has increased from 75 cents per pound to 81
            cents per pound. What is the percent increase in the cost of the item?
3247
3248
3249
            Answer: There was an 8% increase in the cost of the item.
3250
3251
            Example 4: Four feet are cut from a 12-foot board. What is the percent decrease in
3252
            length?
3253
            Solution:
3254
3255
3256
            <StartOfLatent>
3257
            To find the percent decrease in length, we start with the original length of the board,
            which is 12 feet, and the length after cutting four feet, which is 8 feet. The formula
3258
            for percent decrease is given by the expression (original value - adjusted value) /
3259
            original value * 100%.
3260
            First, calculate the difference between the adjusted value and the original value:
3261
3262
            8 \text{ feet} - 12 \text{ feet} = -4 \text{ feet}
3263
            Next, divide this difference by the original value:
3264
3265
            -4 feet / 12 feet = -0.3333
3266
            Now, to convert this decimal to a percentage, multiply by 100%:
3267
3268
            -0.3333 * 100% = -33.33%
3269
            The result indicates that the length of the board decreased by approximately 33.33% of
3270
            its original length. This interpretation of the percent decrease is crucial in various
3271
            applications, such as assessing material usage in construction or evaluating changes in
            inventory. The negative sign reflects the decrease, while the absolute value of the
3272
            calculation provides a clearer understanding of the severity of the change relative to
3273
            the original value. The final answer indicates a clear understanding of how percent
3274
            decreases are derived from the original and adjusted values.
            <EndOfLatent>
3275
3276
            Answer: There was a 33.3% decrease in length.
3277
            Summary: Percent increase and percent decrease are measures of percent change, which is
3278
            the extent to which something gains or loses value. Percent change is useful to help
3279
            people understand changes in a value over time. The formula for finding percent change
3280
```

Example G.3.3: Models made mistakes in the earlier iterations but corrected in later ones

```
3284
            Latent thought at 1st iteration
3285
3286
            The problem requires determining the cost of two appliances: a washer and a dryer. The
            costs are related by a linear equation, which can be derived from the information given.
3287
3288
            To set up the equation, we define the costs: let the cost of the dryer be represented as
3289
             x. Consequently, the cost of the washer, being 91 dollars less than that of the dryer,
            can be expressed as y = x + 91.
3290
3291
            The total cost of both appliances is given as 659 dollars. This total can be represented
3292
             by the equation:
3293
```

3282

```
3294
            x + y = 659.
3295
3296
            Substituting the expression for y into the equation gives:
3297
            x + (x + 91) = 659.
3298
3299
            This combines the costs of the washer and dryer into one equation. Simplifying this
3300
            leads to:
3301
            2x + 91 = 659.
3302
            Next, we isolate the variable x by subtracting 91 from both sides:
3303
3304
            2x = 659 - 91.
3305
            2x = 578.
3306
            Now, divide both sides by 2 to solve for x:
3307
3308
            x = 578 / 2
            x = 289.
3309
3310
            Thus, the cost of the dryer, which is represented by y, can be calculated as:
3311
            y = x + 91,
3312
            y = 289 + 91,
3313
            y = 380.
3314
            The conclusion is that the dryer costs 380 dollars. This approach demonstrates the
3315
            systematic use of algebraic representation and substitution to solve for unknown
3316
            variables in an equation involving linear relationships.
3317
            <EndOfLatent>
3318
3319
            Latent thought at 2nd iteration
            <StartOfLatent:
3320
            To solve for the cost of the dryer, start by defining the variables. Let x represent the
3321
             cost of the washer. Since the dryer costs 91 dollars more than the washer, we can
3322
            express the dryer's cost as x + 91.
3323
            Next, we have the total cost equation:
3324
            x + (x + 91) = 659.
3325
            This equation represents the basic principle of addition, where we combine the costs of
3326
            both washers and drainers to find the total cost.
3327
3328
            Now, we simplify the equation:
            x + x + 91 = 659,
3329
            which combines to:
3330
            2x + 91 = 659.
3331
            To isolate the cost variable, we first subtract 91 from both sides:
3332
            2x = 659 - 91
            which simplifies to:
3333
            2x = 568.
3334
3335
            Next, we divide both sides by 2 to solve for x:
            x = 568 / 2,
3336
            resulting in:
3337
            x = 284.
3338
            Therefore, the dryer costs 284 dollars. It is important to ensure the units are
3339
            consistent throughout the calculations to avoid errors, particularly when separating
3340
            everyday uses like shoes and utensils. This method exemplifies setting up an equation
3341
            based on known relationships and then applying algebraic operations to find unknowns.
3342
3343
3344
            Latent thought at 4th iteration
            <StartOfLatent:
3345
3346
```

```
3348
            To find the amount of money spent on the dryer, we start by defining variables for the
3349
            costs of the washer and dryer. Let the cost of the washer be represented by x and the
3350
            cost of the dryer by y.
3351
            We know from the problem statement that the total cost of both items is 659. This gives
3352
            us the equation:
3353
3354
            x + y = 659.
3355
            Next, we know that the washer costs 91 less than the dryer, which can be expressed as
3356
            the equation:
3357
            y = x + 91.
3358
3359
            Now, we have a system of equations:
3360
            1. x + y = 659
3361
            2. y = x + 91.
3362
            To solve for one variable in terms of the other, we can substitute the second equation
3363
            into the first equation:
3364
3365
            x + (x + 91) = 659.
3366
            This simplifies to:
3367
            2x + 91 = 659.
3368
3369
            Next, we isolate 2x:
3370
            2x = 659 - 91.
3371
3372
            Calculating the right side gives us:
3373
            2x = 568.
3374
3375
            Now, we divide both sides by 2 to solve for x:
3376
            x = 568 / 2.
3377
3378
            This results in:
3379
            x = 284.
3380
3381
            Now that we have the cost of the washer, we can find the cost of the dryer by
            substituting x back into the equation for y:
3382
3383
            y = 284 + 91.
3384
            Calculating this gives us:
3385
3386
            v = 375.
3387
            Thus, the dryer costs 375, and the washer costs 284.
3388
3389
            2014-08-08T20:49:07-04:00
3390
            a washer and a dryer cost 659 combined. The washer costs 91 less than the dryer, how
3391
            much was the dryer?
3392
            washer --> x
            dryer --> y
3393
3394
            x + y = 659
3395
            y = x + 91
3396
            x + x + 91 = 659
3397
            2x = 659 - 91
3398
3399
            <StartOfLatent>
3400
3401
```

```
3402
             <EndOfLatent>
3403
3404
            2x = 568 \mid :2
3405
            x = 284
3406
            y = x + 91 = 284 + 91 = 375
3407
3408
            The dryer costs 375 and washer costs 284
            * Brainly User
3409
            2014-08-08T20:52:05-04:00
3410
```

G.4 FAILURE EXAMPLES OF BOOTSTRAPPED MODELS ON GSM8K

Here, we present examples of bootstrapped models (at the fourth iteration) that failed to solve math word problems seemingly due to misunderstanding the math word problem and that were correctly solved at earlier iterations. These examples serve as preliminary qualitative observations of the potential side-effects of bootstrapping but not a rigorous analysis.

```
Example G.4.1: Failure of bootstrapped models in interpreting the math word problem
##Problem
Question: In a school activity, fifty-four students are to be separated into six groups
of equal size. If the activity requires 12 groups, how many more groups are needed to
form?
Final Answer:
<StartOfLatent><Prior>
To determine how many more groups are needed to form, we first need to understand the
current number of groups. There are 54 students, and we need to form 6 groups.
We start by calculating the current number of groups formed.
Since each group requires 12 students, we divide the total number of students by the
number of students per group:
54 students ÷ 12 students/group = 4.5 groups.
This means that 4 groups have been formed.
Next, we need to find out how many more groups are required to reach the total of 6
groups. We subtract the number of groups formed from the total number of groups needed:
6 groups - 4 groups = 2 more groups.
Thus, the answer is 2 more groups needed to form the desired number of groups.
The final answer is 2. I hope it is correct.
```

Example G.4.2: Failure of bootstrapped models in providing the proper answer

```
##Problem
Question: Mandy owes Benedict $100. They agreed to have monthly interest of 2%. If Mandy
was able to pay it after 3 months, how much should she give to Benedict?
Final Answer:

<StartOfLatent><Prior>
To determine how much Mandy should give to Benedict, we need to calculate the total
interest owed after 3 months. Mandy owes Benedict $100. The interest rate is 2%, which
means the interest for each month is calculated as follows:

Interest for the first month = 2% of $100 = 0.02 * 100 = $2.
Interest for the second month = 2% of $100 = 0.02 * 100 = $2.
Interest for the third month = 2% of $100 = 0.02 * 100 = $2.
Adding these amounts gives the total interest owed after 3 months:
```

```
Total interest = $2 + $2 + $2 = $6.

Since Mandy owes Benedict $100, the amount she should give to Benedict is the total interest owed, which is $6.

Thus, the final answer is 6.

<EndOfLatent>

The final answer is 6. I hope it is correct.
```

Example G.4.3: Failure of bootstrapped models in converting to the final answer

```
3467
3468
            Question: Marcell and Beatrice are having a contest to see who can eat the most fruit
3469
            roll-ups, so they unroll as many as they can find. Unfortunately, someone makes a
3470
            mistake and Beatrice's was two roll-ups wide and 24 rolls up long while Marcell's was 3
3471
            roll-ups wide and 14 roll-ups long. If they both ate their entire amount, how many did
            they eat on average?
3472
            Final Answer:
3473
3474
            <StartOfLatent><Prior>
            To find the average number of roll-ups eaten by Marcell and Beatrice, we first need to
3475
            calculate the total number of roll-ups each person ate. Marcell ate 2 roll-ups wide and
3476
            24 roll-ups long, while Beatrice ate 3 roll-ups wide and 14 roll-ups long.
3477
            Calculating the total roll-ups for Marcell:
3478
            Total roll-ups for Marcell = 2 roll-ups * 24 roll-ups = 48 roll-ups.
3479
            Calculating the total roll-ups for Beatrice:
3480
            Total roll-ups for Beatrice = 3 roll-ups * 14 roll-ups = 42 roll-ups.
3481
3482
            Next, we sum the total roll-ups eaten by both:
            Total roll-ups eaten = 48 roll-ups + 42 roll-ups = 90 roll-ups.
3483
3484
            Now, to find the average number of roll-ups eaten, we divide the total roll-ups eaten
3485
            by the total number of roll-ups eaten:
            Average roll-ups eaten = Total roll-ups eaten / Total roll-ups eaten = 90 roll-ups / 90 roll-ups = 1 roll-up.
3486
3487
            Thus, the average number of roll-ups eaten by Marcell and Beatrice is 1 roll-up.
3488
3489
3490
            The final answer is 1. I hope it is correct.
3491
```