



BOOSTING ACCURACY AND EFFICIENCY OF BUDGET FORCING IN LLMs VIA REINFORCEMENT LEARNING FOR MATHEMATICAL REASONING

Bachelor's Project Thesis

Ravindra A. Tarunokusumo, s5142822, r.tarunokusumo@student.rug.nl

Supervisors: Rafael F. Cunha, r.f.cunha@rug.nl, Tsegaye M. Tashu t.m.tashu@rug.nl

Abstract: Test-time scaling methods have seen a rapid increase in popularity for its computational efficiency and parameter-independent training to improve reasoning performance on Large Language Models. One such method is called budget forcing, a decoding intervention strategy which allocates extra compute budget for thinking and elicits the inherent self-correcting behavior of the model. However, this relies on supervised fine-tuning (SFT) on long-context reasoning traces which causes performance degradation on smaller models due to verbose responses. For this reason, we offer a framework integrating reinforcement learning (RL) to improve token efficiency and boost the performance of a 1.5B model for mathematical reasoning. We demonstrate this using only 1.5K training samples and found that our SFT+RL model performed better on the GSM8K dataset with varying compute budgets. Our main findings showed an overall higher accuracy while significantly reducing its token usage by over 40% compared to the SFT model, revealing how RL can recover the losses due to long-context training and altogether improving performance in mathematical reasoning.

1 Introduction

Recent advances in Large Language Models (LLMs) have shown significant improvement in mathematical problem solving skills (Shao et al., 2024; Guo et al., 2025; Yang et al., 2024). Training LLMs to solve math problems typically require post-training methods such as supervised fine-tuning (SFT) to initialize reasoning behavior and reinforcement learning (RL) to encourage exploration over possible solutions (Luo et al., 2024; Yu et al., 2023; P. Wang et al., 2023). Notably, these methods scale up train-time compute which relies on updating the model parameters through back-propagation (Hoffmann et al., 2022; Kaplan et al., 2020).

However, a new scaling paradigm has emerged called test-time scaling whereby model responses are improved by allocating more compute during inference, allowing the model to “think” before giving its final output. There have been promising research both on its trade-off between pretraining—introducing a new scaling law for test-time

compute (Snell et al., 2024; Wu et al., 2024; Brown et al., 2024)—and as an alternative to SFT and RL (Cobbe et al., 2021a; X. Wang et al., 2022; Shinn et al., 2023; Madaan et al., 2023; Brown et al., 2024; C. Huang et al., 2025; J. Chen et al., 2025; E. Zhao et al., 2025). Recent releases have shown its viability in combination with large scale RL, with frontier models such as OpenAI’s o-series models (OpenAI, 2024, 2025), the DeepSeek-R1 model (Guo et al., 2025), the k1.5 model (Team et al., 2025) and Google’s Gemini models (DeepMind, 2025). These models exhibited state-of-the-art performances on notoriously difficult math benchmarks, performing multi-step reasoning process using Chain-of-Thought (CoT) prompting in their training phase, which elicits intermediate reasoning steps before outputting their final answer (Wei et al., 2022). However, these models required vast amounts of high-quality training data and millions of samples to train (Shao et al., 2024; Guo et al., 2025; Team et al., 2025), which makes it intractable for smaller compute budget and difficult to scale up.

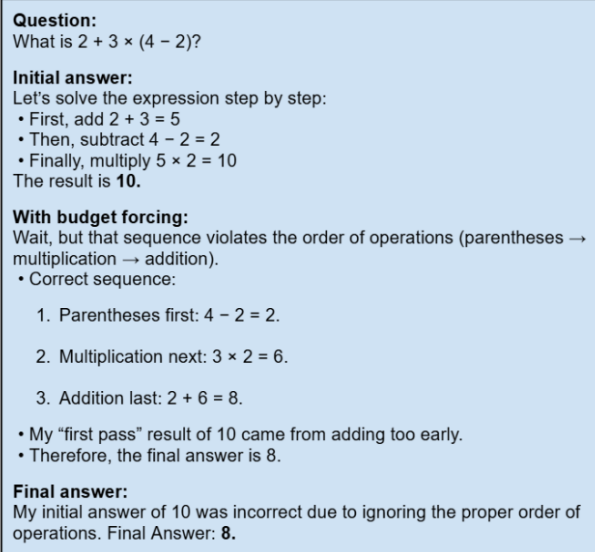


Figure 1.1: Example generation using budget forcing. At first, the model outputs an incorrect answer due to its faulty reasoning step. Then, the appended “Wait” token prompted the model to revise its reasoning and corrects itself.

A simple approach by Muennighoff et al. (2025) is called budget forcing. Budget forcing explicitly controls the number of “thinking tokens” generated by the model during inference. In practice, when the model is about to stop its CoT too early, a special token—often “Wait”—is injected to compel additional reasoning steps, thereby allowing the model to self-correct and self-verify its answers. Conversely, if the model exceeds the pre-determined token budget, an end-of-thinking delimiter is enforced to truncate the reasoning process and prompt the final answer. Figure 1.1 shows an example generation with budget forcing. This novel approach forgoes the need for data-intensive post-training, requiring only a thousand specially curated reasoning data during SFT, and instead relies on the model’s ability to self-correct. And while initializing self-correcting behavior during SFT has been shown to compliment RL training (Ma et al., 2025; Kumar et al., 2024; J. Zhang & Zuo, 2025), the current implementation of budget forcing by Muennighoff et al. (2025) relies solely on SFT.

Previous studies have shown that employing SFT and RL is crucial for optimal test-time scaling for mathematical reasoning (Setlur et al., 2025; Hou

et al., 2025; Havrilla et al., 2024; X. Zhang et al., 2025). Moreover, budget forcing with SFT alone often performs suboptimally and even degrades accuracy as compared to its base model for smaller LLMs (X. Zhang et al., 2025). Indeed, studies have shown that increasing the token budget often causes models to “overthink” (Sui et al., 2025), increasing computational overhead due to persistent back-and-forth between formulating answers and self-correcting itself, even on relatively simple question, and oftentimes ending up backtracking from the correct answer (Zeng et al., 2025; X. Chen et al., 2025). Moreover, we found that this sometimes leads the model to use up all of its token budget and outputs an incomplete generation where the model is cut-off from giving its final answer.

To address these problems, this study offers a framework which combines budget forcing with reinforcement learning for mathematical problem solving. In this framework, the model is first initialized using a modestly-sized reasoning SFT dataset after which RL is employed to refine its reasoning strategies in efficiently utilizing its token budget. We hypothesize that this phase will not require sample-intensive training, echoing the findings of Ye et al. (2025) that good performance is still achievable with less compute. During inference, budget forcing is applied to dynamically control the generation of “thinking tokens”, ensuring that the model effectively utilizes its latent reasoning capabilities without excessive compute. At the end, the model will be evaluated on a range of math benchmarks across several baselines and metrics to assess the effectiveness of incorporating RL training in employing budget forcing.

2 Methodology

In this section, we will first briefly outline each phase in the three-stage pipeline of the framework, starting with the SFT, RL and finally to the inference phase where we employ budget forcing during evaluation. We will discuss the core concepts behind these methods, the issues and setbacks with regards to implementing SFT and RL for mathematical reasoning and the solutions which we adopt to mitigate them. At the end, we provide several metrics to measure the effectiveness of employing these improvements.

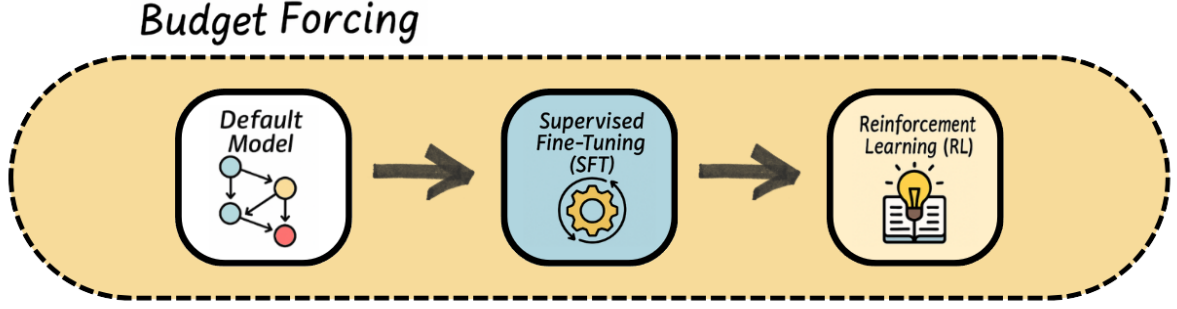


Figure 1.2: Pipeline of the framework. We start with the default model and apply supervised fine-tuning on a reasoning dataset to get our SFT model, and finally we use reinforcement learning to get our SFT+RL model. All three models will be tested using the same budget forcing inference scheme.

2.1 Supervised Fine-Tuning

Initializing reasoning behavior into the model requires a curated dataset containing reasoning traces in the form of intermediate steps for math problems towards the final solution (Yu et al., 2023; Shao et al., 2024; Guo et al., 2025; Ye et al., 2025; Muennighoff et al., 2025; Shen et al., 2025). Other than this, self-correction must also be present in the SFT dataset as budget forcing relies on self-correction. Though LLMs already possess intrinsic capacity for self-correction, they nevertheless fall short on accurately detecting errors and correcting them (J. Huang et al., 2023; Zeng et al., 2025) without additional fine-tuning (Gao et al., 2024; Kamoi et al., 2024).

It has been shown that a small and highly curated dataset containing no more than 1000 training samples (Ye et al., 2025; Sun et al., 2025; Muennighoff et al., 2025), and even as little as 500 is sufficient for this task (Sun et al., 2025). Moreover, Shen et al. (2025) showed that the length of the reasoning trace also plays an important role in enhancing performance. Therefore, a modestly-sized, specialized reasoning dataset with a sufficiently long sample length is enough to initialize reasoning behavior.

2.2 Reinforcement Learning

RL has proven to enhance mathematical reasoning on top of SFT (Ma et al., 2025; Kumar et al., 2024; J. Zhang & Zuo, 2025; X. Zhang et al., 2025), as well as in scaling test-time compute (Havrilla

et al., 2024; Hou et al., 2025). However, naive implementations of RL are often resource intensive and sample-inefficient with poor reward assignment, which risks unstable training, reward hacking and poor generalizability (Kumar et al., 2024; Havrilla et al., 2024; Hou et al., 2025; Setlur et al., 2024; P. Wang et al., 2023; Lightman et al., 2023; Shao et al., 2024). The core challenges with RL are therefore in the learning algorithm and designing robust reward models, both of which we will discuss in this section.

2.2.1 Overview of Algorithm

A novel algorithm called Group-Relative Policy Optimization (GRPO) addresses these issues by reformulating reward optimization as a group-relative comparison task (Shao et al., 2024; Guo et al., 2025). By calculating advantages as deviations from a response group’s mean reward, GRPO implicitly rewards reasoning patterns that correlate with success. This approach not only mitigates reward hacking—as models cannot game static rules—but also reduces sample complexity through variance-normalized training.

Unlike traditional PPO methods (Schulman et al., 2017) that rely on training a separate value function, GRPO estimates advantages by comparing each response within a group to the mean performance of that group. This design eliminates the need for a value network, significantly reducing training overhead.

2.2.2 Formulation

Concretely, given a question q , a group of G sampled outputs $\{o_1, o_2, \dots, o_G\}$ is drawn from the old policy π_{old} . Each output is evaluated by a reward model to obtain a scalar reward r_i , and the normalized advantage is computed as:

$$A_i = \frac{r_i - \bar{r}}{\sigma_r} \quad (2.1)$$

where \bar{r} is the mean reward of the group and σ_r is the standard deviation. This variance normalization helps mitigate reward variance across samples, thus improving sample efficiency. Formally, the GRPO algorithm is defined as:

$$\begin{aligned} \mathcal{J}_{GRPO}(\theta) = & \mathbb{E}[q \sim P(Q), o_i \sim \pi_{old}(o|q)] \\ & \frac{1}{G} \sum_{i=1}^G \left(\min(r_i A_i, \text{clip}(r_i, 1 - \epsilon, \right. \\ & \left. 1 + \epsilon) A_i) - \beta D_{KL}(\pi_{\theta} \parallel \pi_{ref}) \right) \end{aligned} \quad (2.2)$$

Where $P(Q)$ is the probability distribution over questions Q and ϵ and β are hyperparameters controlling clipping and the KL divergence. π_{θ} and π_{ref} are the current and reference policy respectively. GRPO closely resembles the original formulation of PPO in that it preserves the clipped surrogate objective and the adaptive KL divergence penalty term, ensuring policy updates are stable.

2.2.3 Reward Models

Two distinct reward models guide the learning process, based on traditional implementations (Shao et al., 2024; Guo et al., 2025):

1. **Accuracy Reward Model:** Scores the final answer’s correctness by comparing it with the ground-truth labels. However, this is preconditioned by how well the model can follow the prompt instruction in how to format its answer for extraction.
2. **Format Reward Model:** Evaluates whether the model outputs within a specified format. Usually, the outputs are structured in a template format with `<think>` and `<answer>` tags (Guo et al., 2025), but in the case of the SFT dataset, it uses the model’s chat template and a `\boxed{\}` format to enclose the final answer.

For this reason, we chose not to implement format reward models in our experiment since the model would have already been trained to output in that format.

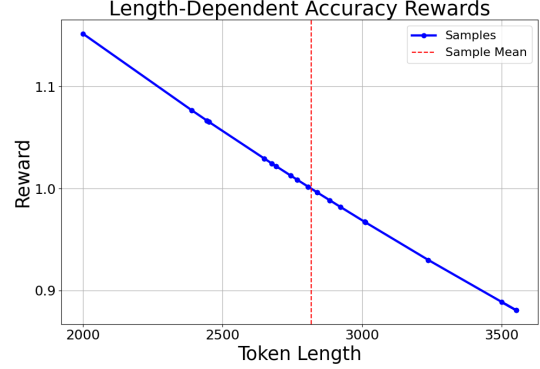


Figure 2.1: Accuracy rewards of 16 randomly generated (correct) dummy samples, following Equation 2.4 Samples with token length exceeding the sample mean receives a reward of < 1.0 and vice versa.

GRPO then integrates these signals in an outcome supervision framework, where the computed advantage is broadcasted across all tokens of the respective output. Though it is possible to incorporate outcome-level supervision, it leaves the model susceptible to reward hacking; moreover, Guo et al. (2025) reported failures on implementing it in their experiments.

However, challenges still remain with regards to binary reward signals, lacking differentiation and weak learning gradients. Furthermore, given the long CoT trajectories during the SFT stage, the model will have been trained to output verbose responses. Following J. Zhang & Zuo (2025), we modify the reward model to include a length-dependent accuracy rewards and a penalty term for incorrect solutions. Given a correct response o , the standardized length deviation is defined as:

$$z = \frac{|o| - \mu}{\sigma - \epsilon} \quad (2.3)$$

Where $\epsilon > 0$ is a small constant for numerical stability. Hence, the final reward accuracy is:

$$R_{\text{accuracy}}(o | q) = \begin{cases} \exp(-\alpha z), & \text{if } o \text{ is correct,} \\ -1.0, & \text{if } o \text{ is incorrect.} \end{cases} \quad (2.4)$$

Where $\alpha > 0$ is a tunable hyperparameter controlling the strength of length penalization.

Figure 2.1 visualizes how the rewards are distributed for correct responses of varying length. Shorter responses are encouraged to steer the model towards conciseness and avoids verbosity.

We also implemented an additional penalty function to address the issue of incomplete generations where we penalize -1.0 for any incomplete generations.

$$R_{\text{completion}}(o) = \begin{cases} -1.0, & \text{if } o \text{ is incomplete} \\ 0, & \text{otherwise} \end{cases} \quad (2.5)$$

This addition is crucial since rather than a reward gradient to address verbosity, a strict penalty term ensures explicit discouragement for truncated response.

2.3 Budget Forcing

In addition to fine-tuning-based methods, we adopt budget forcing as a test-time strategy to control and scale a model’s reasoning capability at inference time.

2.3.1 Overview of Budget Forcing

Formally, let T_{think} denote the number of tokens allocated to the reasoning phase. Budget forcing controls T_{think} by intervening during decoding to limit or extend the model’s generation. If the model generates more than a predefined threshold T_{max} , we terminate the thinking phase forcibly by appending a delimiter token, thereby shifting the model to answer generation. Conversely, if we aim to extend the thinking duration, we can instead append the "Wait" token to prompt the model to continue its reasoning.

2.3.2 Reinforcement Learning as a Corrective Mechanism

Numerous studies have already shown the necessity of employing RL on top of SFT for optimal performance (Setlur et al., 2025; Hou et al., 2025; Havrilla et al., 2024; X. Zhang et al., 2025; Chu et al., 2025). Employing budget forcing with SFT alone will predictably present some challenges. X. Zhang et al. (2025) have found that SFT actually degrades model performance on smaller LLMs



Figure 2.2: Average completion lengths during RL training (smoothed over 50 steps). During the first 2000 steps, the model consistently hits the token budget limit, but then it decreases and settles at around 350 tokens.

due to the long-context samples used in the dataset. We observed several distinct issues that comes with long-context SFT on smaller LLMs. Most notably is verbose responses, causing the model to backtrack and overthink its response (See Figure A.3). Moreover, even when prompted for the final answer, the model sometimes keeps self-correcting and ended up using all its token budget and cutting off its generation mid-sentence (See Figure A.4). Another issue found was language mixing and redundant responses which could be seen cluttering the reasoning traces (See Figure A.5 and Figure A.6).

RL have been shown to alleviate this problem (X. Zhang et al., 2025; Liu et al., 2025; J. Zhang & Zuo, 2025) by designing appropriate reward models to penalize those behavior. Moreover, Guo et al. (2025) have shown that employing RL resolves the issue of language mixing. We will introduce several metrics to assess these behaviors across our models.

2.3.3 Metrics

1. **Scaling:** Repeated applications of appending the "Wait" token can linearly increase the average number of reasoning tokens and thus push the model to deeper levels of analysis (Muenighoff et al., 2025; X. Zhang et al., 2025; Sui et al., 2025). Quantitatively evaluating this scaling effect, we adopt the metric proposed by

Muennighoff et al. (2025):

$$\text{Scaling} = \frac{1}{\binom{|\mathcal{A}|}{2}} \sum_{a,b \in \mathcal{A}, b > a} \frac{f(b) - f(a)}{b - a} \quad (2.6)$$

This captures how much accuracy improves as more reasoning tokens are allowed, where \mathcal{A} is the set of token budgets used during evaluation and $f(a)$ is the accuracy at budget level a . It's the average slope of the accuracy-vs-compute curve—higher means better utilization of extra compute.

2. **Incomplete Generations:** We will also measure the model's ability to comply with the token budget in giving a full-length, completed response for each thinking step and final answer step.

$$\text{Incomplete Gen.} = \frac{1}{|\mathcal{S}_a|} \sum_{s \in \mathcal{S}_a} f_t(s) \quad (2.7)$$

This measures the proportion of incomplete steps, where $\mathcal{S}_a = \{s_0 \dots s_L\}$ is the set of all steps at budget level a and f_t is a function which checks whether it is incomplete, which could either happen due to a truncated step, i.e. when the appended "Wait" token cut off the previous step before giving its boxed answer, or at the final answer segment when the model continued reasoning even after being prompted to finalize its answer and ran out of token budget. And it is worth noting that all cases of redundancy are also cases of incomplete generations.

3. **Average Token Length:** We measure the model's verbosity from the average complete output token length $|o|$ per sample, where $\mathcal{O}_a = \{o_0 \dots o_L\}$ is the set of all complete outputs, corresponding to 1319 samples from the GSM8K test set. A lower response length does not always mean better performance, but it is indicative of how well the model utilizes its reasoning budget. We will discuss this at a later section.

$$\text{Average Token Length} = \frac{1}{|\mathcal{O}_a|} \sum_{o \in \mathcal{O}_a} |o| \quad (2.8)$$

4. **Accuracy:** We use the zero-shot accuracy as our main metric to measure model performance.

3 Experiment

In this section we describe the configuration of our training procedure and the inference scheme we used during evaluation.

3.1 Models

We conducted the experiment following the three-stage pipeline using the Qwen2.5-1.5B-Instruct model (Qwen et al., 2025), using said model as our baseline. However, our main comparison will be between the SFT-only model versus the SFT+RL model.

3.2 Dataset

For the SFT stage we will be using the tokenized slk-1.1 dataset by Muennighoff et al. (2025), consisting of 1K math questions focused on quality, diversity and difficulty with a average response length of 10K tokens. Importantly, the reasoning traces contain keywords such as "Wait" or "Alternatively" followed by revisions to the previous step to encourage the desired self-correcting behavior during training. However, due to resource limitations we only use a subset of samples with a maximum token length of 10K, which amounts to 525 samples in total. And for the RL stage we will use a subset of 1K training samples from the GSM8K dataset by OpenAI (Cobbe et al., 2021b).

3.3 Training Setup

We conducted the SFT training using two NVIDIA A100 40GB GPUs and had FSDP enabled (Y. Zhao et al., 2023), and we used only a single GPU for the RL training. In both stages we trained the model using the TRL library (von Werra et al., 2020), and we implemented Flash Attention 2 (Dao, 2023) and loaded the model in BF16 to minimize VRAM usage. The hyperparameter configurations are shown in Table A.1.

3.4 Evaluation

We evaluated all our models under these two budget forcing inference scheme:

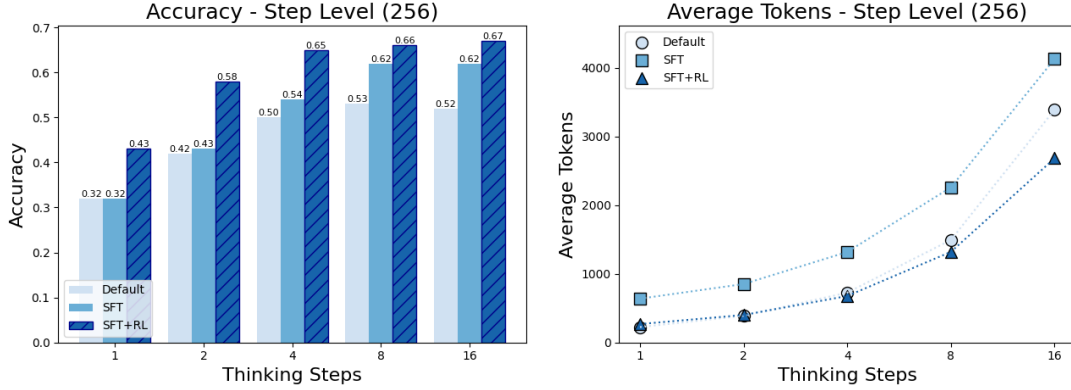


Figure 2.3: Zero-shot accuracy score (Left) and average tokens per sample (Right) of all three models under the Step-Level condition on the GSM8K test set. The SFT+RL model scores the highest in accuracy in all five thinking steps, while the SFT model performs similarly with the Default model in 1-2 thinking steps but gradually surpasses it. However, the SFT model outputs the most tokens while the SFT+RL model shows similar token consumption as the Default model.

1. **Baseline Full:** For each question we allow the model to generate an answer with the maximum token length for the Qwen2.5 family of models (32K) without apportioning its generation for thinking.
2. **Step-Level:** We evaluated the model multiple times with a thinking steps, where $a \in [1, 2, 4, 8, 16]$. Each step the model is allotted a thinking budget of 218 tokens and an additional 512 for the final answer, hence each corresponding to $\mathcal{A} = \{218, 512, 1024, 2048, 4096\}$ thinking token budgets. We chose 256 tokens as our per-step token budget because we found this value to have the highest accuracy score compared to 128 and 512 tokens. (See Table A.2).

For each condition and model, we measure the metrics outlined in Section 2.3.3 (except for Scaling in the Baseline Full condition). We used vLLM (Kwon et al., 2023) to speed up inference on a single A100 40GB GPU. For the sampling parameters we used greedy decoding and set the rest to the default values.

4 Results and Discussions

Figure 2.3 compares the accuracy on all three of the models. SFT+RL performs the best with a higher

accuracy on all five thinking steps. Under the Step-Level condition, it gained an average accuracy over the default model of **14.0%** across all conditions and **9.2%** over the SFT model. Moreover, we see a **14.1%** decrease in the average token length compared to the default model and a significant drop of **41.8%** compared to the SFT model. However, we also observe the SFT model showing the steepest performance gain as thinking steps increase. In Table 4.1 the SFT model has the highest scaling factor compared with the rest, with an increase of **3.73%** accuracy per extra reasoning step, suggesting that pure SFT alone benefits more from budget forcing, while we see the default model improving more moderately (scaling factor **2.77%**) rather than failing to benefit.

On the other hand, the SFT model was found to be incredibly inefficient and verbose. Given a single thinking step, **75.02%** of its responses reached the token limit and were truncated, and **85.28%** given 16 thinking steps. This is expected because of the long-context data on which it was trained, and more to the fact that smaller LLMs tend to perform worse after SFT (X. Zhang et al., 2025). We see in Figure 2.3 the SFT model scored the same in accuracy with the baseline model using a single thinking step, then gradually surpasses it with more steps. This shows that SFT increases its accuracy *given* more inference token budget. Hence why we see in the Baseline Full inference scheme

Model	Steps	Accuracy	Avg Tokens	Incomplete Gen. (%)	Scaling (%)
Default	BF	0.58	508	0.76	2.77
	1	0.32	224	20.09	
	2	0.42	389	18.15	
	4	0.50	726	18.13	
	8	<u>0.53</u>	1493	25.90	
	16	0.52	3389	51.18	
SFT	BF	0.67	5158	1.44	3.73
	1	0.32	<u>638</u>	<u>75.02</u>	
	2	0.43	851	79.73	
	4	0.54	1317	84.67	
	8	<u>0.62</u>	2258	86.97	
	16	<u>0.62</u>	4128	85.28	
SFT+RL	BF	0.65	669	0.88	3.32
	1	0.43	<u>268</u>	36.69	
	2	0.58	401	31.36	
	4	0.65	677	<u>26.96</u>	
	8	0.66	1319	27.35	
	16	0.67	2681	28.77	

Table 4.1: Performance metrics on the GSM8K dataset for 256 thinking tokens. For Incomplete Gen. and Average Tokens, lower is better. Numbers with boldface are the best across models, while the ones underlined are conditional for each model (excluding BF = Baseline Full).

the model performing at its best with the highest accuracy on par with the SFT+RL model and only a handful of incomplete generations.

However, we observe in both models a performance cap just beneath 70%. We suspect that this is mainly due to the model’s baseline capability. Despite the fact that it was adequately supplied with reasoning traces during SFT and further trained using conciseness rewards during RL, the SFT+RL model showed no further improvement beyond 68% accuracy. We hypothesized that it is due to the insufficient parameter size to fully capture the self-correcting behavior present in the SFT dataset. In other words, the model wasn’t large enough to generalize its reasoning and overfitted on the training data. We confirmed this by evaluating on the AIME and MATH500 benchmarks, and found that they didn’t score well enough to be worth including here.

5 Conclusion

In this paper, we have demonstrated that integrating reinforcement learning with supervised

Metric	Δ Default	Δ SFT
Accuracy	+0.14	+0.092
Avg. Tokens (%)	-14.1	-41.9
Incomplete Gen. (%)	+3.54	-52.11
Scaling (%)	+0.55	-0.41

Table 4.2: Mean improvements of the SFT+RL model over the Default and SFT model, excluding the Baseline Full condition (reductions in Avg. Tokens (%) and Incomplete Gen. (%) are better).

fine-tuning substantially enhances the effectiveness of budget forcing for mathematical reasoning, mitigating some setbacks introduced by long-context training data. On the GSM8K test set, our SFT+RL model achieved a zero-shot accuracy of **67%** at 16-step level, yielding an average gain of **14.0%** accuracy over the Default model and **9.2%** over the SFT model. It also reduced average token length to **1,069** tokens on average—**41.9%** lower than SFT and **14.1%** lower than Default. How-

ever, incomplete generations were not eliminated: under the 256-token budget they remained between **26.96%** and **36.69%**. We also see that the accuracy of all models scales positively with the number of thinking steps, with the SFT model benefiting most per additional step, whereas SFT+RL primarily lifts the baseline accuracy but offers a smaller marginal return. Finally, we see substantial accuracy improvements from the Baseline-Full condition, with an accuracy surpassing even the 16-step run for the Default and SFT model with a few incomplete generations.

Despite these gains, our study has several limitations. First, the 1.5B parameter backbone may not fully capture the complexity of self-correcting behavior in the SFT dataset, suggesting potential overfitting to the relatively small training sets. Second, evaluation was confined to GSM8K because our 1.5B model was not capable enough to score meaningful results on other benchmarks even with a high inference token budget. We hope that future work will explore scaling to larger models and extending evaluations to multiple reasoning domains to investigate the broader potential of budget forcing and other test-time scaling methods.

References

- Brown, B., Juravsky, J., Ehrlich, R., Clark, R., Le, Q. V., Ré, C., & Mirhoseini, A. (2024). Large language monkeys: Scaling inference compute with repeated sampling. *arXiv preprint arXiv:2407.21787*.
- Chen, J., Ren, J., Chen, X., Yang, C., Sun, R., & Arik, S. Ö. (2025). Sets: Leveraging self-verification and self-correction for improved test-time scaling. *arXiv preprint arXiv:2501.19306*.
- Chen, X., Xu, J., Liang, T., He, Z., Pang, J., Yu, D., ... Yu, D. (2025). *Do not think that much for $2+3=?$ on the overthinking of o1-like llms*. Retrieved from <https://arxiv.org/abs/2412.21187>
- Chu, T., Zhai, Y., Yang, J., Tong, S., Xie, S., Schuurmans, D., ... Ma, Y. (2025). Sft memorizes, rl generalizes: A comparative study of foundation model post-training. *arXiv preprint arXiv:2501.17161*.
- Cobbe, K., Kosaraju, V., Bavarian, M., Chen, M., Jun, H., Kaiser, L., ... others (2021a). Training verifiers to solve math word problems. *arXiv preprint arXiv:2110.14168*.
- Cobbe, K., Kosaraju, V., Bavarian, M., Chen, M., Jun, H., Kaiser, L., ... Schulman, J. (2021b). Training verifiers to solve math word problems. *arXiv preprint arXiv:2110.14168*.
- Dao, T. (2023). *Flashattention-2: Faster attention with better parallelism and work partitioning*. Retrieved from <https://arxiv.org/abs/2307.08691>
- DeepMind. (2025, Apr). Retrieved from <https://deepmind.google/technologies/gemini/>
- Gao, K., Cai, H., Shuai, Q., Gong, D., & Li, Z. (2024). Embedding self-correction as an inherent ability in large language models for enhanced mathematical reasoning. *arXiv preprint arXiv:2410.10735*.
- Guo, D., Yang, D., Zhang, H., Song, J., Zhang, R., Xu, R., ... others (2025). Deepseek-r1: Incentivizing reasoning capability in llms via reinforcement learning. *arXiv preprint arXiv:2501.12948*.
- Havrilla, A., Du, Y., Raparthy, S. C., Nalmpantis, C., Dwivedi-Yu, J., Zhuravinskyi, M., ... Raileanu, R. (2024). Teaching large language models to reason with reinforcement learning. *arXiv preprint arXiv:2403.04642*.
- Hoffmann, J., Borgeaud, S., Mensch, A., Buchatskaya, E., Cai, T., Rutherford, E., ... others (2022). Training compute-optimal large language models. *arXiv preprint arXiv:2203.15556*.
- Hou, Z., Lv, X., Lu, R., Zhang, J., Li, Y., Yao, Z., ... Dong, Y. (2025). Advancing language model reasoning through reinforcement learning and inference scaling. *arXiv preprint arXiv:2501.11651*.
- Huang, C., Huang, L., Leng, J., Liu, J., & Huang, J. (2025). Efficient test-time scaling via self-calibration. *arXiv preprint arXiv:2503.00031*.

- Huang, J., Chen, X., Mishra, S., Zheng, H. S., Yu, A. W., Song, X., & Zhou, D. (2023). Large language models cannot self-correct reasoning yet. *arXiv preprint arXiv:2310.01798*.
- Kamoi, R., Zhang, Y., Zhang, N., Han, J., & Zhang, R. (2024). When can llms actually correct their own mistakes? a critical survey of self-correction of llms. *Transactions of the Association for Computational Linguistics*, 12, 1417–1440.
- Kaplan, J., McCandlish, S., Henighan, T., Brown, T. B., Chess, B., Child, R., ... Amodei, D. (2020). Scaling laws for neural language models. *arXiv preprint arXiv:2001.08361*.
- Kumar, A., Zhuang, V., Agarwal, R., Su, Y., Co-Reyes, J. D., Singh, A., ... others (2024). Training language models to self-correct via reinforcement learning. *arXiv preprint arXiv:2409.12917*.
- Kwon, W., Li, Z., Zhuang, S., Sheng, Y., Zheng, L., Yu, C. H., ... Stoica, I. (2023). Efficient memory management for large language model serving with pagedattention. In *Proceedings of the acm sigops 29th symposium on operating systems principles*.
- Lightman, H., Kosaraju, V., Burda, Y., Edwards, H., Baker, B., Lee, T., ... Cobbe, K. (2023). Let’s verify step by step. In *The twelfth international conference on learning representations*.
- Liu, Z., Chen, C., Li, W., Qi, P., Pang, T., Du, C., ... Lin, M. (2025). Understanding rl-zero-like training: A critical perspective. *arXiv preprint arXiv:2503.20783*.
- Luo, L., Liu, Y., Liu, R., Phatale, S., Lara, H., Li, Y., ... others (2024). Improve mathematical reasoning in language models by automated process supervision. *arXiv preprint arXiv:2406.06592*, 2.
- Ma, R., Wang, P., Liu, C., Liu, X., Chen, J., Zhang, B., ... Li, J. (2025). S²r: Teaching llms to self-verify and self-correct via reinforcement learning. *arXiv preprint arXiv:2502.12853*.
- Madaan, A., Tandon, N., Gupta, P., Hallinan, S., Gao, L., Wiegrefe, S., ... others (2023). Self-refine: Iterative refinement with self-feedback. *Advances in Neural Information Processing Systems*, 36, 46534–46594.
- Muennighoff, N., Yang, Z., Shi, W., Li, X. L., Fei-Fei, L., Hajishirzi, H., ... Hashimoto, T. (2025). s1: Simple test-time scaling. *arXiv preprint arXiv:2501.19393*.
- OpenAI. (2024). Retrieved from <https://openai.com/index/learning-to-reason-with-llms/>
- OpenAI. (2025). Retrieved from <https://openai.com/index/openai-o3-mini/>
- Qwen, :, Yang, A., Yang, B., Zhang, B., Hui, B., ... Qiu, Z. (2025). *Qwen2.5 technical report*. Retrieved from <https://arxiv.org/abs/2412.15115>
- Schulman, J., Wolski, F., Dhariwal, P., Radford, A., & Klimov, O. (2017). *Proximal policy optimization algorithms*. Retrieved from <https://arxiv.org/abs/1707.06347>
- Setlur, A., Nagpal, C., Fisch, A., Geng, X., Eisenstein, J., Agarwal, R., ... Kumar, A. (2024). Rewarding progress: Scaling automated process verifiers for llm reasoning. *arXiv preprint arXiv:2410.08146*.
- Setlur, A., Rajaraman, N., Levine, S., & Kumar, A. (2025). Scaling test-time compute without verification or rl is suboptimal. *arXiv preprint arXiv:2502.12118*.
- Shao, Z., Wang, P., Zhu, Q., Xu, R., Song, J., Bi, X., ... others (2024). Deepseekmath: Pushing the limits of mathematical reasoning in open language models. *arXiv preprint arXiv:2402.03300*.
- Shen, S., Huang, F., Zhao, Z., Liu, C., Zheng, T., & Zhu, D. (2025). Long is more important than difficult for training reasoning models. *arXiv preprint arXiv:2503.18069*.
- Shinn, N., Cassano, F., Gopinath, A., Narasimhan, K., & Yao, S. (2023). Reflexion: Language agents with verbal reinforcement learning. *Advances in Neural Information Processing Systems*, 36, 8634–8652.

- Snell, C., Lee, J., Xu, K., & Kumar, A. (2024). Scaling llm test-time compute optimally can be more effective than scaling model parameters. *arXiv preprint arXiv:2408.03314*.
- Sui, Y., Chuang, Y.-N., Wang, G., Zhang, J., Zhang, T., Yuan, J., ... others (2025). Stop overthinking: A survey on efficient reasoning for large language models. *arXiv preprint arXiv:2503.16419*.
- Sun, Y., Zhou, G., Wang, H., Li, D., Dziri, N., & Song, D. (2025). Climbing the ladder of reasoning: What llms can-and still can't-solve after sft? *arXiv preprint arXiv:2504.11741*.
- Team, K., Du, A., Gao, B., Xing, B., Jiang, C., Chen, C., ... others (2025). Kimi k1. 5: Scaling reinforcement learning with llms. *arXiv preprint arXiv:2501.12599*.
- von Werra, L., Belkada, Y., Tunstall, L., Beeching, E., Thrush, T., Lambert, N., ... Gallouédec, Q. (2020). *Trl: Transformer reinforcement learning*. <https://github.com/huggingface/trl>. GitHub.
- Wang, P., Li, L., Shao, Z., Xu, R., Dai, D., Li, Y., ... Sui, Z. (2023). Math-shepherd: Verify and reinforce llms step-by-step without human annotations. *arXiv preprint arXiv:2312.08935*.
- Wang, X., Wei, J., Schuurmans, D., Le, Q., Chi, E., Narang, S., ... Zhou, D. (2022). Self-consistency improves chain of thought reasoning in language models. *arXiv preprint arXiv:2203.11171*.
- Wei, J., Wang, X., Schuurmans, D., Bosma, M., Xia, F., Chi, E., ... others (2022). Chain-of-thought prompting elicits reasoning in large language models. *Advances in neural information processing systems*, 35, 24824–24837.
- Wu, Y., Sun, Z., Li, S., Welleck, S., & Yang, Y. (2024). Inference scaling laws: An empirical analysis of compute-optimal inference for problem-solving with language models. *arXiv preprint arXiv:2408.00724*.
- Yang, A., Zhang, B., Hui, B., Gao, B., Yu, B., Li, C., ... Zhang, Z. (2024). Qwen2.5-math technical report: Toward mathematical expert model via self-improvement. *arXiv preprint arXiv:2409.12122*.
- Ye, Y., Huang, Z., Xiao, Y., Chern, E., Xia, S., & Liu, P. (2025). Limo: Less is more for reasoning. *arXiv preprint arXiv:2502.03387*.
- Yu, L., Jiang, W., Shi, H., Yu, J., Liu, Z., Zhang, Y., ... Liu, W. (2023). Metamath: Bootstrap your own mathematical questions for large language models. *arXiv preprint arXiv:2309.12284*.
- Zeng, Z., Cheng, Q., Yin, Z., Zhou, Y., & Qiu, X. (2025). Revisiting the test-time scaling of o1-like models: Do they truly possess test-time scaling capabilities? *arXiv preprint arXiv:2502.12215*.
- Zhang, J., & Zuo, C. (2025). Grpo-lead: A difficulty-aware reinforcement learning approach for concise mathematical reasoning in language models. *arXiv preprint arXiv:2504.09696*.
- Zhang, X., Huang, Z., Ni, C., Xiong, Z., Chen, J., & Oymak, S. (2025). Making small language models efficient reasoners: Intervention, supervision, reinforcement. *arXiv preprint arXiv:2505.07961*.
- Zhao, E., Awasthi, P., & Gollapudi, S. (2025). Sample, scrutinize and scale: Effective inference-time search by scaling verification. *arXiv preprint arXiv:2502.01839*.
- Zhao, Y., Gu, A., Varma, R., Luo, L., Huang, C.-C., Xu, M., ... Li, S. (2023). *Pytorch fsdp: Experiences on scaling fully sharded data parallel*. Retrieved from <https://arxiv.org/abs/2304.11277>

A Appendix

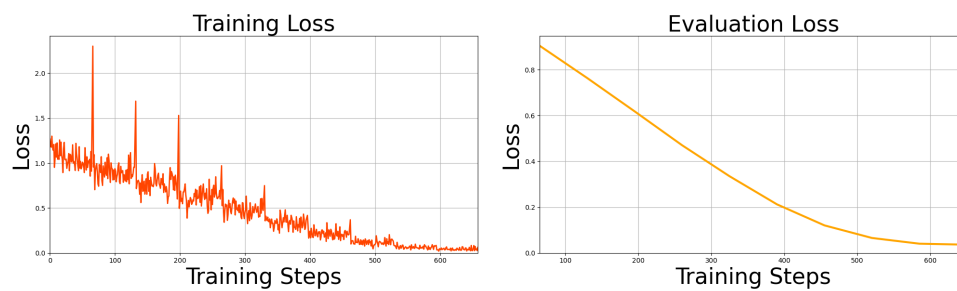


Figure A.1: Training and evaluation loss during SFT training. We achieved a final loss of 5.4% and 3.6% respectively.

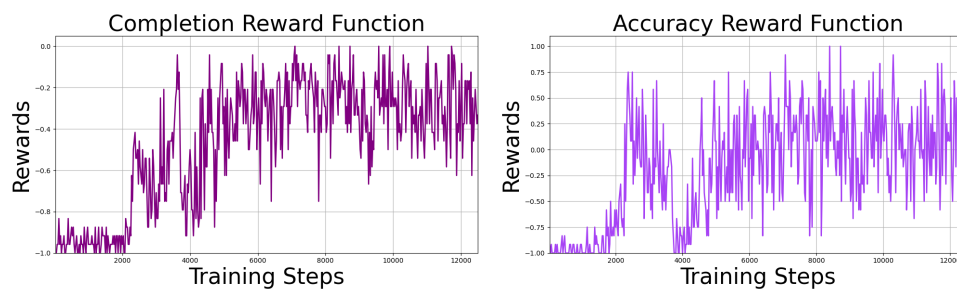


Figure A.2: Completion and Accuracy Rewards plots over 12500 training steps. We achieved a final reward of -0.33 and 0.06 respectively.

Question:

Every day, Wendi feeds each of her chickens three cups of mixed chicken feed, containing seeds, mealworms and vegetables to help keep them healthy. She gives the chickens their feed in three separate meals. In the morning, she gives her flock of chickens 15 cups of feed. In the afternoon, she gives her chickens another 25 cups of feed. How many cups of feed does she need to give her chickens in the final meal of the day if the size of Wendi's flock is 20 chickens?

Reasoning:

So each evening meal is **20** cups.

...

...

Wait, perhaps this is related to another problem piece. Let's consider that in English, "meal" can sometimes refer to a specific meal of the day. For example, "breakfast meal," "lunch meal," "dinner meal." But in this case, "meal" is given for the day, not specific to times. Alternatively, "the final meal" could refer to the last meal of the day, which would be evening if divided into three meals. But without mention of evening, the problem is incomplete.

...

...

Final Answer:

****Conclusion**:**

The remaining feed, 20 cups, is divided into three meals for the evening. Therefore, each evening meal amounts to:

$$\left[\frac{20 \text{ cups}}{3 \text{ meals}} \approx 6.667 \text{ cups/meal} \right]$$

However, the problem does not specify the exact number. Given the ambiguity and typical English usage where "final meal" refers to the last meal of the day (evening), the answer is the evening meal.

$$\left[\boxed{6.667 \text{ cups}} \text{ or } \boxed{20 \text{ cups}} \right]$$

But if required to use the final answer in $\boxed{\{\}}:$
 $\boxed{6.667}$

Figure A.3: Example of backtracking and overthinking. The model initially arrives at the correct answer, but begins to overthink and ventures into irrelevant reasoning trajectories which ended up giving an incorrect final answer.

Hyperparameter	SFT	RL
Learning rate	3×10^{-5}	5×10^{-5}
Scheduler	Cosine with 20% warmup	Cosine with 20% warmup
Per-device batch size	1	16
Epoch	10	2
Gradient accumulation	4	8
Maximum sequence length	10000	512
Optimizer	AdamW [0.90, 0.999]	
Weight decay	1×10^{-3}	
Maximum gradient normalization		0.1
G		6

Table A.1: Hyperparameters for the SFT and RL training.

<p>**Final Answer**</p> <p>\boxed{18}</p> <p>However, the problem states that she sells the remainder at the farmers' market daily. If she has 0 eggs on day 3, she can't sell 0 eggs (since she has none left), which would result in a revenue of 0. But the problem says "for \$2 per fresh duck egg." Therefore, if at some point she has no eggs left, her sales would have to stop. But perhaps in that case, she doesn't sell any eggs that day. So, her daily sales vary depending on the day:</p> <ul style="list-style-type: none"> - Day 1: 9 eggs sold for \$18 - Day 2: 9 eggs sold for \$18 - Day 3: 2 eggs sold for \$4 (but then she has 0 eggs left, skips day 4...) <p>But the problem doesn't specify what happens on days like day 3 where she has no eggs left. If we assume that on those days she doesn't sell anything, then her daily revenue is \$18. Otherwise, if she sells 9 eggs each of the first two days and 9 eggs on the third, her daily revenue is \$18 each for two days and \$4 on the third, which would be inconsistent. Therefore, perhaps on the third day, after selling 2 eggs, she doesn't make any more sales for the day. So, the revenue on day 3 is \$4, making her total daily revenue \$18 on days 1 and 2, and \$4 on day 3. But this is strange.</p> <p>Alternatively, perhaps each day she has the same number of eggs left, which is $16 - 7x$ eggs over x days. But if we take $x=2$, she has $16 - 7*2 = 16 - 14 = 2$ eggs left, which she sells for \$4, so per day that would be \$2. But if $x=3$, she has $16 - 7*3 = 16 - 21 = -5$ which is negative. That's not possible. Therefore, the number of days must be limited. So perhaps the problem is expecting a daily fixed revenue of \$18. However, given that she makes \$18 per day, with 9 eggs sold each of the first two days and 2 eggs sold on the third, but she can't have negative eggs, so maybe after selling 9 on day 1 and day 2,</p>

Figure A.4: Example of incomplete generation in the final answer segment. Despite being prompted for the final answer, the model continues to reason until it hits the maximum token limit.

Original total + points from fourth round = $81 + 2 = 83$. So, the total points after the extra points had been distributed is 83. So the answer should be 83.

But wait, am I missing something? Is the fourth round point addition for all players or just those who qualified? Wait, the problem says "in the fourth round, every player gets an extra point if they have over 20 points." So, it's every player regardless of whether they had their points calculated based on前三轮还是四轮。但是问题中说after three rounds,然后在第四轮给符合条件的加点。所以可能意思是前三轮已经计算了他们的点数，第四轮只根据前三轮的点数来判断是否加点，所以只有Mike和Tony在第四轮加点。

Figure A.5: Example of language mixing in the reasoning trace.

[illegible]

Figure A.6: Example of redundancy in the reasoning trace. This will always result in incomplete generations.

(a) 128 thinking tokens

Model	Steps	Accuracy	Avg Tokens	Incomplete Gen. (%)	Scaling (%)
Default	1	0.34	<u>206</u>	46.10	0.599
	2	0.33	300	<u>44.57</u>	
	4	0.38	521	64.22	
	8	0.39	988	68.98	
	16	<u>0.40</u>	1987	79.59	
SFT	1	0.23	<u>525</u>	<u>78.85</u>	3.68
	2	0.31	646	84.56	
	4	0.43	856	87.78	
	8	0.52	1330	91.83	
	16	<u>0.58</u>	2325	94.77	
SFT+RL	1	0.18	<u>180</u>	<u>51.59</u>	<u>5.84</u>
	2	0.36	270	60.00	
	4	0.56	469	55.36	
	8	0.62	889	53.34	
	16	<u>0.63</u>	1790	57.80	

(b) 512 thinking tokens

Model	Steps	Accuracy	Avg Tokens	Incomplete Gen. (%)	Scaling (%)
Default	1	0.27	<u>257</u>	3.98	<u>3.74</u>
	2	0.48	437	<u>3.13</u>	
	4	0.51	833	4.76	
	8	0.52	1919	15.13	
	16	<u>0.53</u>	5494	45.85	
SFT	1	0.48	<u>835</u>	66.91	1.99
	2	0.57	1306	73.62	
	4	<u>0.62</u>	2193	75.98	
	8	<u>0.62</u>	3868	72.79	
	16	0.61	6919	<u>62.99</u>	
SFT+RL	1	0.56	<u>339</u>	8.26	1.38
	2	0.63	491	<u>7.61</u>	
	4	0.65	849	8.95	
	8	<u>0.66</u>	1677	11.09	
	16	0.65	3466	13.22	

Table A.2: Performance metrics on the GSM8K dataset for 128 (a) and 512 (b) thinking tokens. As thinking tokens increase, we see a higher baseline accuracy and 2x-3x increase in token usage while substantially reducing incomplete generations for the Default and SFT+RL model. Furthermore, scaling factor drops for the SFT and SFT+RL model, but increases for the Default model.