

MAXIMIZING CONFIDENCE ALONE IMPROVES REASONING

Anonymous authors

Paper under double-blind review

ABSTRACT

Reinforcement learning (RL) has enabled machine learning models to achieve significant advances in many fields. Most recently, RL has empowered frontier language models to solve challenging math, science, and coding problems. However, central to any RL algorithm is the reward function, and reward engineering is a notoriously difficult problem in any domain. In this paper, we propose **RENT: Reinforcement Learning via Entropy Minimization** – a fully unsupervised RL method that requires no external reward or ground-truth answers, and instead uses the model’s entropy of its underlying distribution as an intrinsic reward. We find that by reinforcing the chains of thought that yield high model confidence on its generated answers, the model improves its reasoning ability. In our experiments, we showcase these improvements on an extensive suite of commonly-used reasoning benchmarks, including GSM8K, MATH500, AMC, AIME, and GPQA, and models of varying sizes from the Qwen, Mistral, and Llama families. The generality of our unsupervised learning method lends itself to applicability in a wide range of domains where external supervision is unavailable.

1 INTRODUCTION

Imagine you’re taking an exam. Once it begins, no new information is available and no external help can be sought. With only your own reasoning to rely on, how do you tackle a difficult problem? You might make an initial attempt, assess your confidence in the answer, and revise your reasoning until you feel sufficiently certain. Of course, confidence is not a guarantee of correctness – but in the absence of feedback, it is often the only intrinsic signal we have to guide further thought. In such settings, humans tend to optimize for confidence, or equivalently, to reduce uncertainty. In machine learning, uncertainty is commonly quantified via entropy – a measure of how peaked

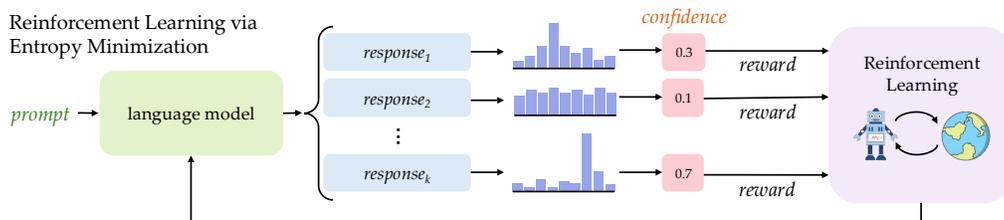


Figure 1: Overview of RENT: Reinforcement Learning via Entropy Minimization. For each response, we use the model’s underlying confidence (negative entropy) as a reward for reinforcement learning. This enables the model to learn without any external reward or ground-truth answers.

or diffuse a probability distribution is. Language models output distributions over tokens, and the entropy of these distributions reflects the model’s confidence: lower entropy implies more confident predictions. Yet despite the growing use of language models in reasoning tasks, current approaches to improvement still rely heavily on external supervision, rewarding correctness with respect to ground-truth labels Guo et al. (2025); Shao et al. (2024). This dependence is often impractical, particularly in real-world or open-ended scenarios where supervision is scarce or unavailable.

To address this, we propose **RENT: Reinforcement Learning via Entropy Minimization** – a fully unsupervised reinforcement learning method that improves reasoning performance by using the model’s own confidence as a reward. Specifically, we define the reward as the negative entropy of the model’s predicted token distributions. This signal is dense, general, and easy to compute, requiring no ground-truth answers. Importantly, not all parts of the response contribute equally to final performance. Through empirical analysis, we find that minimizing entropy over tokens near the end of the reasoning chain, especially those corresponding to the final answer, correlates most strongly with improved accuracy. In contrast, early tokens in the response show little correlation. This suggests that as the model approaches its final answer, it increasingly relies on its own confidence to guide reasoning, so encouraging confidence in these final steps is key.

Concretely, we frame the problem through the lens of test-time training (TTT), where a model continues to adapt even after deployment, using only information available at inference. Unlike conventional training, which relies on labeled data and offline optimization, TTT updates model parameters or behavior directly from the input distributions observed during testing. We demonstrate RENT’s effectiveness across diverse reasoning benchmarks, including GSM8K (Cobbe et al., 2021), MATH500 (Hendrycks et al., 2021; Lightman et al., 2023), AMC and AIME (Li et al., 2024), and GPQA (Rein et al., 2024). Our method scales across model families (Qwen, Mistral, and Llama) and sizes and consistently improves performance.

2 RELATED WORK

2.1 REINFORCEMENT LEARNING FOR REASONING

Initially, reinforcement learning (RL) for language models was mostly used for learning from human preferences (Christiano et al., 2017) and, traditionally, the RL optimization was done with algorithms such as PPO (Schulman et al., 2017). With the capabilities of language models continuing to improve, researchers have begun to explore the possibility of using RL to improve the performance of language models on reasoning tasks such as math (Cobbe et al., 2021; Hendrycks et al., 2021; Li et al., 2024), science (He et al., 2024; Rein et al., 2024), or coding (Li et al., 2022; Chaudhary, 2023) problems. In these settings, the model is prompted to generate a chain-of-thought (Wei et al., 2022) and final answer, and receives a reward based on how closely its final answer matches the ground-truth answer. These efforts present RL as an alternative to search-based approaches to chain-of-thought reasoning such as Tree of Thoughts (Yao et al., 2023) and Graph of Thoughts (Besta et al., 2024). Related lines of work include training a reward model to give feedback for every step in the chain of thought, and training RL models to encourage self-correcting behaviors in language models. Examples of RL methods in this space include Zelikman et al. (2022); Singh et al. (2023); Kumar et al. (2024); Qu et al. (2024); Uesato et al. (2022); Lightman et al. (2023); Wang et al. (2023). At scale, DeepSeek (Guo et al., 2025; Shao et al., 2024) proposed an open-source model that showed OpenAI o1 (Jaech et al., 2024)-level reasoning by performing RL in this manner, using a new algorithm GRPO (Shao et al., 2024).

2.2 CONFIDENCE AND CALIBRATION

Confidence measures quantify how certain a model is that its generated output is correct (Yoon et al., 2025; Spiess et al., 2024). In order to evaluate the confidence of machine learning models, it is necessary also to discuss *calibration* (Kalai and Vempala, 2024; Virk et al., 2024) - i.e., how aligned those confidences are with actual correctness. As language models are increasingly trusted to make important decisions, providing users with a reliable confidence measure would be useful in many situations (Geng et al., 2023; Manakul et al., 2023; Varshney et al., 2023; Hou et al., 2023; Jiang et al., 2023; Han et al., 2024; Spiess et al., 2024). As such, researchers have developed various confidence metrics for modern deep learning models and studied the extent to which they are calibrated. These include both methods that assume access to the model’s weights (Gupta et al., 2024; Kadavath et al., 2022; Xu et al., 2024; Spiess et al., 2024) and methods that estimate confidence via prompting alone (Xie et al., 2024; Geng et al., 2023; Tian et al., 2023; Xiong et al., 2023; Yang et al., 2024). In our paper, we use the model’s confidence to iteratively improve its own performance via reinforcement learning.

2.3 TEST-TIME ADAPTATION

Test-time adaptation is where a model is updated using data from the test distribution, without access to ground-truth labels. The goal is to improve performance in scenarios where there is a distribution shift between training and testing environments. Methods for adapting without labels include normalization techniques that recalibrate feature statistics at test time (Quiñero-Candela et al., 2008; Sun et al., 2017; Maria Carlucci et al., 2017; Schneider et al., 2020; Nado et al., 2020). The most relevant work to ours is Tent (Wang et al., 2020), which performs entropy minimization on model predictions during test time. This approach assumes that predictions on test data should be low in entropy if the model is well-adapted to the new distribution. Tent builds on earlier work that uses entropy minimization as a regularization strategy in semi-supervised learning (Grandvalet and Bengio, 2004; Lee et al., 2013; Berthelot et al., 2019) and domain adaptation contexts (Maria Carlucci et al., 2017; Shu et al., 2018; Saito et al., 2019), where encouraging confident predictions has proven effective for improving generalization. Recently, TTRL (Zuo et al., 2025) proposed test-time RL using majority voting as a reward. Compared to entropy, majority voting is a sparse reward and less general; for example, it cannot be applied to long-form free-response questions.

2.4 UNSUPERVISED REINFORCEMENT LEARNING

Unsupervised RL trains agents using intrinsic rewards like novelty, entropy, or mutual information, enabling skill acquisition without extrinsic feedback. Prior methods include ICM and RND for prediction error Pathak et al. (2017); Burda et al. (2018), APT Liu and Abbeel (2021) and ProtoRL Yarats et al. (2021) for entropy maximization, and DIAYN, APS, and SMM for skill discovery via mutual information Eysenbach et al. (2018); Liu and Abbeel (2021); Kim et al. (2023). An interesting observation is that while exploration methods primarily maximize entropy, we instead minimize it by reinforcing high-confidence outputs, and find that for language models, this leads to better reasoning performance without any external supervision.

3 METHOD

3.1 REINFORCEMENT LEARNING FOR LANGUAGE MODELS

The goal of reinforcement learning (RL) is to train a policy which generates actions that maximize the cumulative expected reward. In the context of modern language models, the policy π is a language model and the actions y_{pred} are sampled from the distribution over a discrete vocabulary. The task is formulated as a one-step RL problem in which the model generates $y_{\text{pred}} = \pi(x)$, where x is sampled from the dataset $\mathcal{D} = \{(x, y_{\text{target}})\}$, and receives some reward for the generation. Typically, the ground-truth answer y_{target} is used to give the model a reward $r = \mathcal{R}(y_{\text{target}}, y_{\text{pred}})$. One reward function which is currently used is simple string matching, where $\mathcal{R}(y_{\text{target}}, y_{\text{pred}}) = \mathbb{1}\{y_{\text{target}} = y_{\text{pred}}\}$. Our work focuses on instead doing *unsupervised* reinforcement learning, which does not require external supervision for the reward. Specifically, y_{target} is not used in the reward $r = \mathcal{R}(y_{\text{pred}})$ and we do not assume access to this at any point in training.

3.2 GROUP RELATIVE POLICY OPTIMIZATION (GRPO)

To optimize the policy, we adopt Group Relative Policy Optimization (GRPO; Shao et al., 2024), a critic-free variant of PPO that measures *relative* performance within a group of rollouts for the same input, rather than against a separate baseline policy.

For each training query x , we sample a group of G candidate responses $\{y^{(1)}, \dots, y^{(G)}\} \sim \pi_{\text{old}}(\cdot | x)$, and compute their rewards $r^{(i)} = \mathcal{R}(x, y^{(i)})$. GRPO normalizes these rewards within the group to form a group-relative advantage:

$$\mu_x = \frac{1}{G} \sum_{i=1}^G r^{(i)}, \quad \sigma_x = \sqrt{\frac{1}{G} \sum_{i=1}^G (r^{(i)} - \mu_x)^2}, \quad A^{(i)} = \frac{r^{(i)} - \mu_x}{\sigma_x + \varepsilon},$$

where μ_x and σ_x denote the mean and standard deviation of rewards in the group, and ε is a small constant for numerical stability. Intuitively, $A^{(i)}$ quantifies how much better or worse a response is relative to its peers for the same prompt.

Using these group-relative advantages, GRPO applies a PPO-style clipped objective to update the current policy π_θ away from π_{old} :

$$\mathcal{L}_{\text{GRPO}}(\theta) = \mathbb{E}_{x, \{y^{(i)}\} \sim \pi_{\text{old}}} \left[\frac{1}{G} \sum_{i=1}^G \min\left(\rho_\theta^{(i)} A^{(i)}, \text{clip}\left(\rho_\theta^{(i)}, 1 - \epsilon, 1 + \epsilon\right) A^{(i)}\right) \right],$$

where $\rho_\theta^{(i)} = \frac{\pi_\theta(y^{(i)}|x)}{\pi_{\text{old}}(y^{(i)}|x)}$ is the importance ratio and ϵ is the PPO clipping parameter. For more details, we refer the reader to Shao et al. (2024).

3.3 ENTROPY REWARD

For a given prompt x , the model generates a response $y_{\text{pred}} = y_{\text{pred},1}, \dots, y_{\text{pred},T} = \pi(x)$, where T is the number of tokens in the response. At each token $t \in \{1, \dots, T\}$, the model outputs a probability distribution p_t over the vocabulary \mathcal{V} , i.e., $p_t(v) = P(y_t = v | x, y_{<t})$. The entropy of this distribution measures the model’s uncertainty in predicting the next token and is given by:

$$H(p_t) = - \sum_{v \in \mathcal{V}} p_t(v) \log p_t(v)$$

To compute the total entropy of the response, we average the entropies across all tokens. The total entropy $H(\pi(x))$ provides a measure of the overall uncertainty in the model’s response. Higher entropy indicates greater uncertainty or more diverse token predictions, while lower entropy suggests more confident and peaked distributions at each token. We use the negative entropy of the predicted token distribution as a reward signal:

$$\mathcal{R}(y_{\text{pred}}) = -H(\pi(x)) = \frac{1}{T} \sum_{t=1}^T \sum_{v \in \mathcal{V}} p_t(v) \log p_t(v)$$

This reward encourages the model to produce more confident and peaked distributions over the vocabulary, effectively promoting lower uncertainty in its predictions. Within the RL framework, the learning objective becomes maximizing the expected reward over the data distribution:

$$\max_{\pi} \mathbb{E}_{x \sim \mathcal{D}} [\mathbb{E}_{y_{\text{pred}} \sim \pi(x)} [\mathcal{R}(y_{\text{pred}})]]$$

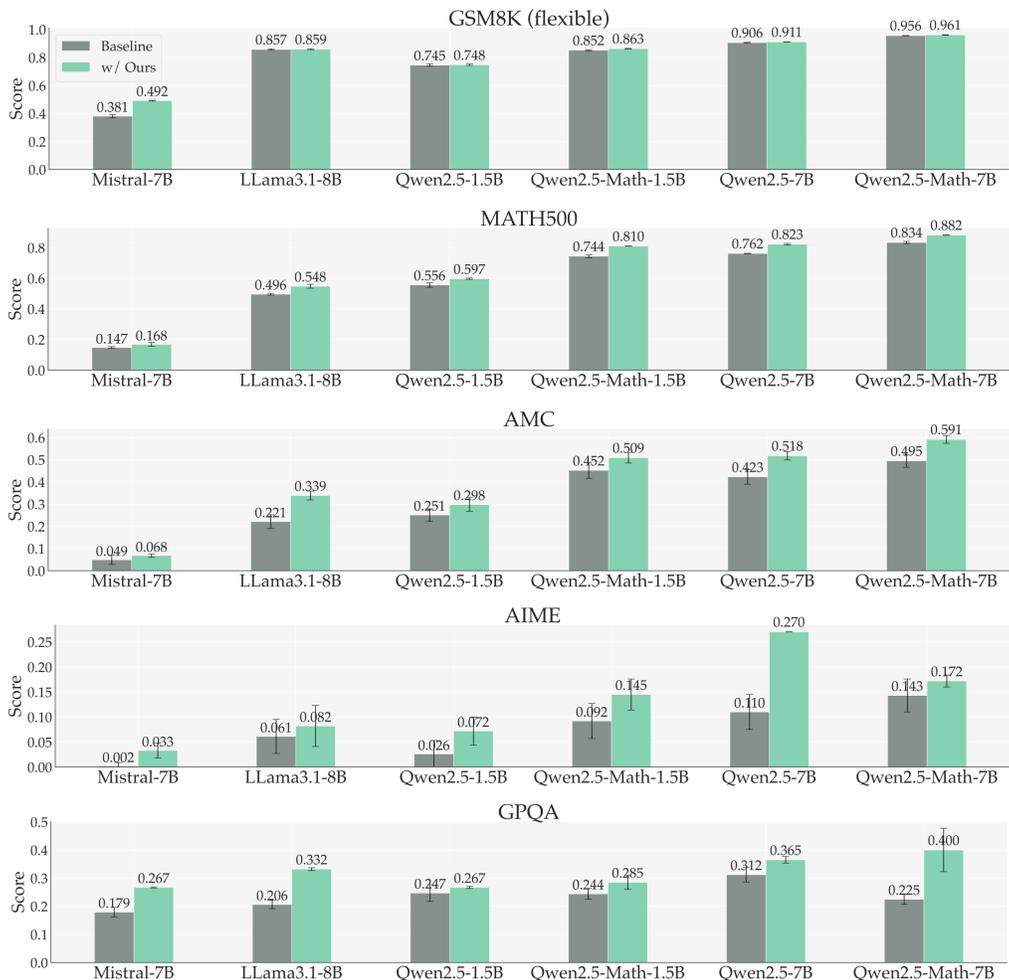
By optimizing this objective, the model learns to generate responses with lower entropy without relying on external supervision or labeled target responses.

4 EXPERIMENTS

4.1 EXPERIMENTAL SETUP

Benchmarks. We train a model with reinforcement learning on each dataset independently. We conduct our experiments on the following commonly-used benchmarks for evaluating the reasoning capabilities of large language models:

- **GSM8K** (Cobbe et al., 2021): GSM8K contains 8792 grade-school math word problems. The train set contains roughly 7473 problems and the test set contains roughly 1319 problems.
- **MATH500** (Hendrycks et al., 2021; Lightman et al., 2023): MATH (Hendrycks et al., 2021) is a dataset containing competition math problems spanning seven categories. It contains 12500 problems, of which 7500 are used for training and 5000 are used for testing. MATH500 (Lightman et al., 2023) is a subset of the MATH test set created by OpenAI by sampling uniformly at random from the test set.
- **AMC** (Li et al., 2024): The American Mathematics Competitions (AMCs) are competitions given to high school students. The specific dataset we use is comprised of 83 problems from the 2022 and 2023 AMC12 exams, which are given to 12th grade students. Although the original problems are in multiple-choice format, the dataset presents modified versions of the problem which expect an integer solution.



Note: Error bars may be large as they indicate the standard deviation over the samples, *not* training seeds.

Figure 2: Performance on GSM8K, MATH500, AMC, AIME, and GPQA. The standard deviations reported are over 5, 5, 32, 64, and 10 samples, respectively. Across benchmarks and models, we find that entropy minimization alone is an effective reward for improving the reasoning ability of language models. All models are Instruct models; the "Instruct" is omitted for brevity.

- AIME24 (Li et al., 2024): The American Invitational Mathematics Examination (AIME) is a prestigious high school mathematics competition. It consists of 15 questions meant to be completed in 3 hours and is given to top-scoring students on the AMC exam. Each year, there are two versions of the exam which consist of distinct questions. We train on the 30 problems from both versions of the 2024 exam.
- GPQA (Rein et al., 2024): GPQA is a dataset of 448 multiple-choice problems in biology, physics, and chemistry at the PhD level. They are intended to be "Google-proof" in the sense that they require advanced reasoning skills.

Since we are interested in test-time adaptation, and we do not assume access to the ground-truth answer, we use the same dataset for training and evaluation. Additionally, some benchmarks do not have standardized train sets. The exception is GSM8K, where we use the standard train and test sets; this shows that generalization does occur and RENT is not merely overfitting to the test set.

270
271
272
273
274
275
276
277
278
279
280
281
282
283
284
285
286
287
288
289
290
291
292
293
294
295
296
297
298
299
300
301
302
303
304
305
306
307
308
309
310
311
312
313
314
315
316
317
318
319
320
321
322
323

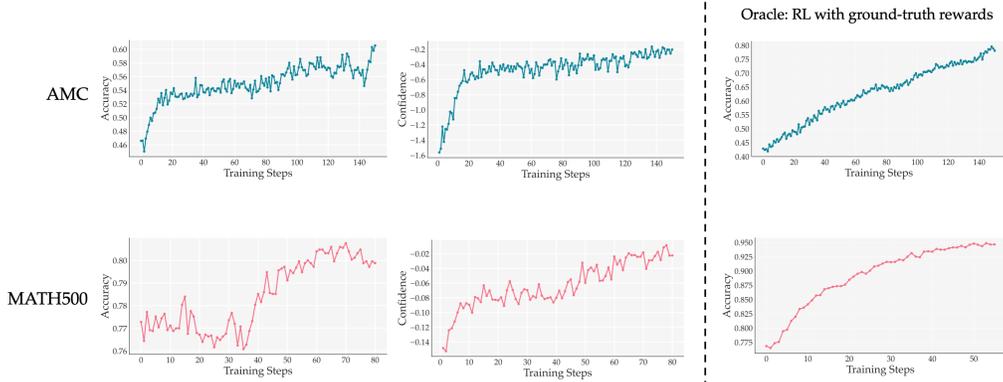


Figure 3: Accuracy and confidence over the course of training. The trends indicate that accuracy and confidence are indeed highly correlated and therefore it is natural to use confidence as a reward.

Models. We conduct experiments on a wide range of models from different model families and with varying parameter counts: Mistral-7B-Instruct-v0.3, Llama3.1-8B-Instruct, Qwen2.5-1.5B-Instruct, Qwen2.5-Math-1.5B-Instruct, Qwen2.5-7B-Instruct, and Qwen2.5-Math-7B-Instruct.

Implementation details. For the RL optimization we use GRPO (Shao et al., 2024) with a learning rate of 1×10^{-6} and the Adam optimizer. The batch sizes and sampling hyperparameters may vary among models and datasets. We provide a full list of hyperparameters in the Appendix.

4.2 MAIN RESULTS

Figure 2 shows the performance of models before and after entropy minimization on GSM8K, MATH500, AMC, AIME24, and GPQA. We report standard deviations reported are over 5, 5, 32, 64, and 10 samples, respectively. Note that all models are Instruct models (e.g., Qwen2.5-1.5B refers to Qwen2.5-1.5B-Instruct). Across model families, model sizes, and benchmarks, entropy minimization allows large language models to improve their reasoning skills, without any external supervision. On the Math models such as Qwen2.5-Math-1.5B and Qwen2.5-Math-7B, the base model often struggles at following instructions and therefore the initial score is zero or near zero, and therefore the boost from entropy minimization is quite large. On models that are already proficient at instruction following, we can still see strong performance improvements from entropy minimization. Given the potential pitfall of overconfidence in language models, we performed extensive experimentation to ensure empirically that entropy minimization is a robust and generalizable reward function across datasets and models.

4.3 IS IT JUST FORMATTING?

It is a well-known issue with reasoning benchmarks that language models can lose points simply because they do not know how to put their answers in the right format. For example, MATH500 expects final answers to be placed in "boxed". A nontrivial amount of engineering effort has gone into both designing prompts that encourage correct formatting and implementing parsers that effectively extract answers from language model responses, in attempts to mitigate this issue. Therefore, one might wonder if, instead of learning to perform complex reasoning, RENT merely encourages the model to put its answers in the right format. Table 1 shows that this is not the case. Models trained with the RENT reward outperform only using a format reward, which simply assigns a binary reward based on whether the correct format is followed in the response. In some cases, the performance of our method is similar to (or even slightly worse than) just using format reward, but of course it is expected that unsupervised RL methods might not always lead to significant improvements. For example, if the benchmark is extremely easy and the model only needs to learn the right format to achieve near-perfect scores, RENT would not outperform format reward. Or, if the benchmark is so hard that it is beyond the model’s capabilities altogether, neither method would perform well.

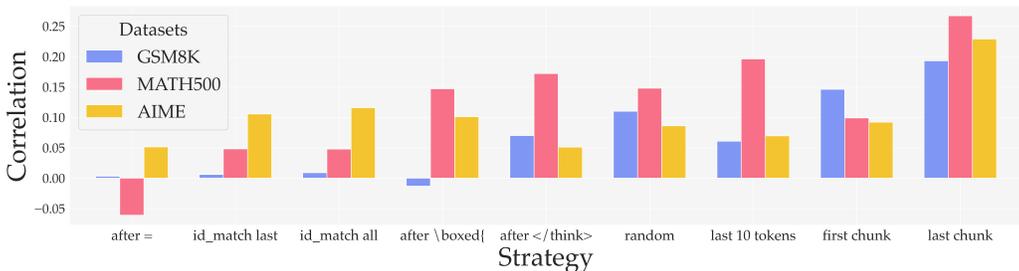


Figure 4: Evaluation (by computing correlation between accuracy and confidence) of various strategies for selecting which tokens to minimize the entropy over. We find the highest correlation between accuracy and confidence in the last few tokens of the response.

Table 1: Comparison to RL with a format reward. The best result on each benchmark is indicated in bold. RENT generally outperforms only using a format reward.

	GSM8K	MATH500	AMC	AIME	GPQA
<i>Mistral-7B-Instruct-v0.3</i>	0.381	0.147	0.049	0.002	0.179
w/ Format reward only	0.393	0.150	0.051	0.015	0.240
w/ RENT (Ours)	0.492	0.168	0.068	0.033	0.267
<i>LLama3.1-8B-Instruct</i>	0.857	0.496	0.221	0.061	0.206
w/ Format reward only	0.866	0.533	0.265	0.086	0.282
w/ RENT (Ours)	0.859	0.548	0.339	0.082	0.332
<i>Qwen2.5-1.5B-Instruct</i>	0.745	0.548	0.251	0.026	0.247
w/ Format reward only	0.754	0.558	0.259	0.054	0.271
w/ RENT (Ours)	0.748	0.597	0.298	0.072	0.267
<i>Qwen2.5-Math-1.5B-Instruct</i>	0.852	0.744	0.452	0.092	0.244
w/ Format reward only	0.857	0.756	0.490	0.117	0.276
w/ RENT (Ours)	0.863	0.810	0.504	0.145	0.285
<i>Qwen2.5-7B-Instruct</i>	0.906	0.762	0.423	0.110	0.311
w/ Format reward only	0.913	0.774	0.458	0.156	0.338
w/ RENT (Ours)	0.911	0.823	0.518	0.270	0.365
<i>Qwen2.5-Math-7B-Instruct</i>	0.956	0.834	0.495	0.143	0.225
w/ Format reward only	0.957	0.873	0.560	0.154	0.340
w/ RENT (Ours)	0.967	0.882	0.591	0.167	0.400

However, across datasets and model sizes, we find a consistent improvement over using the format reward and this assures us that the model is actually learning to think through difficult problems and improve its ability to reason.

4.4 CORRELATION BETWEEN ENTROPY AND ACCURACY

Figure 3 shows the accuracy and confidence throughout training Qwen2.5-Math-7B and Qwen2.5-7B-Instruct on the AMC and MATH500 datasets respectively. Critically, as the model improves its confidence via RENT, the accuracy of the model improves as well. This demonstrates the significant correlation between answer confidence and answer accuracy, supporting our initial hypothesis.

4.5 QUALITATIVE SAMPLES

Table 3 shows a qualitative sample from GSM8K. See the Appendix for a qualitative sample from AIME. The qualitative samples verify that the model indeed learns meaningful reasoning skills via entropy minimization. It is not merely learning to format its answer correctly, or otherwise collapsing to some other reward-hacking behavior.

4.6 COMPARISON TO CONCURRENT WORK

In this section, we compare RENT to concurrent papers which use intrinsic rewards. We evaluate on GSM8K, MATH500, AMC, AIME, and GPQA and run all experiments with Qwen2.5-7B-Instruct as the baseline model. Table 2 shows comparisons to the following methods:

- Test-Time Reinforcement Learning (TTRL) (Zuo et al., 2025) assigns a reward of 1 to the majority answer and 0 to all other answers. In our experiments, we reimplemented this majority voting reward in our codebase.
- Intuitor (Zhao et al., 2025) uses the forward KL divergence between a uniform distribution and the model’s distribution as the reward. In contrast, we use entropy, which is the reverse KL divergence from the uniform distribution. Intuitor is mode-seeking while RENT is mode-covering. We ran the publicly available Intuitor code (which is also implemented on top of verl framework Sheng et al. (2024)) with the same batch size, epochs and evaluation strategy as RENT for fair comparison.
- Shao et al. (2025) suggested that even random or “spurious” rewards could be used to improve reasoning. To compare against spurious rewards, we modify our code to set the reward for every generation randomly to 0 or 1 with equal probability. Intuitively, we believe spurious rewards might work because gradients from correct examples contribute to learning, while gradients from incorrect examples might cancel each other out. Our hypothesis is supported by work such as Rolnick et al. (2017), which shows that learning can still happen even when diluting datasets with incorrect labels.

Empirically, we find that RENT is the best of the four methods on average. Compared to TTRL and Intuitor, performance is similar on most benchmarks except AIME, where RENT outperforms both by a large margin. This is especially interesting since AIME is the hardest benchmark in our evaluations (i.e., the initial accuracy of the model is the lowest). Spurious rewards are not competitive with the other three methods; we conclude that random reward values are not enough and it is indeed beneficial to use meaningful unsupervised rewards that explicitly encourage some measure of confidence. Additionally, we discuss more concurrent works in the Appendix.

Table 2: Comparison of RENT with three concurrent papers: TTRL (Zuo et al., 2025), Intuitor (Zhao et al., 2025), and Spurious Rewards (Shao et al., 2025). The best result on each benchmark is indicated in bold. RENT is the best-performing method on MATH500, AIME, GPQA, and is the best on average.

	GSM8K	MATH500	AMC	AIME	GPQA	Average
<i>Qwen2.5-7B-Instruct</i>						
w/ TTRL	0.933	0.822	0.521	0.172	0.346	0.559
w/ Intuitor (forward KL)	0.929	0.783	0.525	0.200	0.337	0.555
w/ Spurious Rewards	0.910	0.774	0.459	0.156	0.342	0.528
w/ RENT (Ours)	0.911	0.823	0.518	0.270	0.365	0.577

4.7 WHICH TOKENS TO MINIMIZE ENTROPY OVER

We investigated which response tokens are the best to minimize entropy over. To do this, we evaluated the following response token selection strategies:

- **last chunk** - split the response into k chunks, and select all tokens in the last chunk
- **first chunk** - split the response into k chunks, and select all tokens in the first chunk
- **last 10 tokens** - select the last k response tokens
- **random** - select k random tokens from the response
- **after ;/think_i** - select all tokens after the thought process (after the “;/think_i” tag)
- **after \boxed{ }** - select all tokens after the boxed answer
- **id_match_last** - select the tokens in the model’s final answer

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

Table 3: The model splits the problem into its key steps and correctly calculates the target age

Prompt	system Please reason step by step, and put your final answer within <code>\boxed{}</code> . user Jaymee is 2 years older than twice the age of Shara. If Shara is 10 years old, how old is Jaymee? Let’s think step by step and output the final answer within <code>"\boxed{}</code> ". assistant
Response	To determine Jaymee’s age, we need to follow these steps: 1. Identify Shara’s age. 2. Calculate twice Shara’s age. 3. Add 2 years to twice Shara’s age. Step 1: Shara’s age is given as 10 years old. Step 2: Twice Shara’s age is calculated as: $2 \times 10 = 20$ Step 3: Jaymee is 2 years older than twice Shara’s age, so we add 2 years to 20: $20 + 2 = 22$ Therefore, Jaymee is <code>\boxed{22}</code> years old.

- **id_match_all** - select all tokens in the model’s final answer and any previous occurrences
- **after =** - minimize all tokens after the equality symbol (“=”)

Figure 4 shows the initial correlation between negative entropy and accuracy for each of these strategies on three datasets. While most token selection strategies do result in a positive correlation between confidence and accuracy, we note that the “last chunk” strategy has a significantly higher correlation compared to the “first chunk” strategy. This suggests that the most important tokens to minimize entropy over are tokens that occur later in the response. Furthermore, based on the low correlation results from the “id_match_last” and “id_match_all” strategies, we find that it is not sufficient to simply minimize the entropy of the final answer tokens; this suggests that, counterintuitively, the token-level confidence of the final answer tokens is not well-calibrated to its true response confidence/accuracy.

5 LIMITATIONS

Fundamentally, unsupervised learning alone is relatively limited compared to methods which are able to use external supervision for learning. Therefore, it is not surprising that our method is not able to match the performance of methods that have access to the ground-truth answers. It is, of course, a possibility for the model to be confidently wrong. Overconfidence is a well-known issue with language models and these calibration errors can cause RENT to fail catastrophically. It could be dangerous to deploy such an unsupervised learning method in the real world without any safeguards. However, we generally find empirically that confidence does correlate with accuracy and the performance does improve by using confidence alone. This indicates that even if the model is overconfident on some answers, it is well-calibrated overall.

6 CONCLUSION

We presented RENT, an unsupervised reinforcement learning method which uses entropy as a reward. In our experiments, we showed that by simply minimizing entropy, we can improve the reasoning performance of language models on GSM8K, MATH500, AMC, AIME, and GPQA. Our reward function is general and can be applied on a wide range of domains. We are excited about the possibility of using entropy minimization and, more broadly, unsupervised RL to improve the capabilities of machine learning models in regimes where external supervision is unavailable.

7 REPRODUCIBILITY STATEMENT

We described the experimental setup in Section 4.1 and provided a list of hyperparameters in the Appendix. Furthermore, we have submitted the source code as supplementary material.

REFERENCES

- Shivam Agarwal, Zimin Zhang, Lifan Yuan, Jiawei Han, and Hao Peng. The unreasonable effectiveness of entropy minimization in llm reasoning. *arXiv preprint arXiv:2505.15134*, 2025.
- David Berthelot, Nicholas Carlini, Ian Goodfellow, Nicolas Papernot, Avital Oliver, and Colin Raffel. Mixmatch: A holistic approach to semi-supervised learning. *Advances in neural information processing systems*, 32, 2019.
- Maciej Besta, Nils Blach, Ales Kubicek, Robert Gerstenberger, Michal Podstawski, Lukas Gianinazzi, Joanna Gajda, Tomasz Lehmann, Hubert Niewiadomski, Piotr Nyczyk, et al. Graph of thoughts: Solving elaborate problems with large language models. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 38, pages 17682–17690, 2024.
- Yuri Burda, Harrison Edwards, Amos Storkey, and Oleg Klimov. Exploration by random network distillation. *arXiv preprint arXiv:1810.12894*, 2018.
- Sahil Chaudhary. Code alpaca: An instruction-following llama model for code generation. <https://github.com/sahil280114/codealpaca>, 2023.
- Daixuan Cheng, Shaohan Huang, Xuekai Zhu, Bo Dai, Wayne Xin Zhao, Zhenliang Zhang, and Furu Wei. Reasoning with exploration: An entropy perspective. *arXiv preprint arXiv:2506.14758*, 2025.
- Paul F Christiano, Jan Leike, Tom Brown, Miljan Martic, Shane Legg, and Dario Amodei. Deep reinforcement learning from human preferences. *Advances in neural information processing systems*, 30, 2017.
- 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.
- Ganqu Cui, Yuchen Zhang, Jiacheng Chen, Lifan Yuan, Zhi Wang, Yuxin Zuo, Haozhan Li, Yuchen Fan, Huayu Chen, Weize Chen, et al. The entropy mechanism of reinforcement learning for reasoning language models. *arXiv preprint arXiv:2505.22617*, 2025.
- Benjamin Eysenbach, Abhishek Gupta, Julian Ibarz, and Sergey Levine. Diversity is all you need: Learning skills without a reward function. *arXiv preprint arXiv:1802.06070*, 2018.
- Jiahui Geng, Fengyu Cai, Yuxia Wang, Heinz Koepl, Preslav Nakov, and Iryna Gurevych. A survey of confidence estimation and calibration in large language models. *arXiv preprint arXiv:2311.08298*, 2023.
- Yves Grandvalet and Yoshua Bengio. Semi-supervised learning by entropy minimization. *Advances in neural information processing systems*, 17, 2004.
- 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 llms via reinforcement learning. *arXiv preprint arXiv:2501.12948*, 2025.
- Neha Gupta, Harikrishna Narasimhan, Wittawat Jitkittum, Ankit Singh Rawat, Aditya Krishna Menon, and Sanjiv Kumar. Language model cascades: Token-level uncertainty and beyond. *arXiv preprint arXiv:2404.10136*, 2024.
- Jiuzhou Han, Wray Buntine, and Ehsan Shareghi. Towards uncertainty-aware language agent. *arXiv preprint arXiv:2401.14016*, 2024.

- 540 Chaoqun He, Renjie Luo, Yuzhuo Bai, Shengding Hu, Zhen Leng Thai, Junhao Shen, Jinyi Hu,
541 Xu Han, Yujie Huang, Yuxiang Zhang, et al. Olympiadbench: A challenging benchmark for
542 promoting agi with olympiad-level bilingual multimodal scientific problems. *arXiv preprint*
543 *arXiv:2402.14008*, 2024.
- 544 Dan Hendrycks, Collin Burns, Saurav Kadavath, Akul Arora, Steven Basart, Eric Tang, Dawn Song,
545 and Jacob Steinhardt. Measuring mathematical problem solving with the math dataset. *arXiv*
546 *preprint arXiv:2103.03874*, 2021.
- 547 Bairu Hou, Yujian Liu, Kaizhi Qian, Jacob Andreas, Shiyu Chang, and Yang Zhang. Decomposing
548 uncertainty for large language models through input clarification ensembling. *arXiv preprint*
549 *arXiv:2311.08718*, 2023.
- 551 Aaron Jaech, Adam Kalai, Adam Lerer, Adam Richardson, Ahmed El-Kishky, Aiden Low, Alec
552 Helyar, Aleksander Madry, Alex Beutel, Alex Carney, et al. Openai o1 system card. *arXiv*
553 *preprint arXiv:2412.16720*, 2024.
- 554 Zhengbao Jiang, Frank F Xu, Luyu Gao, Zhiqing Sun, Qian Liu, Jane Dwivedi-Yu, Yiming Yang,
555 Jamie Callan, and Graham Neubig. Active retrieval augmented generation. In *Proceedings of*
556 *the 2023 Conference on Empirical Methods in Natural Language Processing*, pages 7969–7992,
557 2023.
- 558 Fangkai Jiao, Geyang Guo, Xingxing Zhang, Nancy F Chen, Shafiq Joty, and Furu Wei. Preference
559 optimization for reasoning with pseudo feedback. *arXiv preprint arXiv:2411.16345*, 2024.
- 561 Saurav Kadavath, Tom Conerly, Amanda Askell, Tom Henighan, Dawn Drain, Ethan Perez,
562 Nicholas Schiefer, Zac Hatfield-Dodds, Nova DasSarma, Eli Tran-Johnson, et al. Language mod-
563 els (mostly) know what they know. *arXiv preprint arXiv:2207.05221*, 2022.
- 564 Adam Tauman Kalai and Santosh S Vempala. Calibrated language models must hallucinate. In
565 *Proceedings of the 56th Annual ACM Symposium on Theory of Computing*, pages 160–171, 2024.
- 566 Seongun Kim, Kywoon Lee, and Jaesik Choi. Variational curriculum reinforcement learning for
567 unsupervised discovery of skills. *arXiv preprint arXiv:2310.19424*, 2023.
- 568 Aviral Kumar, Vincent Zhuang, Rishabh Agarwal, Yi Su, John D Co-Reyes, Avi Singh, Kate Baumli,
569 Shariq Iqbal, Colton Bishop, Rebecca Roelofs, et al. Training language models to self-correct via
570 reinforcement learning. *arXiv preprint arXiv:2409.12917*, 2024.
- 571 Dong-Hyun Lee et al. Pseudo-label: The simple and efficient semi-supervised learning method for
572 deep neural networks. In *Workshop on challenges in representation learning, ICML*, volume 3,
573 page 896. Atlanta, 2013.
- 574 Jia Li, Edward Beeching, Lewis Tunstall, Ben Lipkin, Roman Soletskyi, Shengyi Huang, Kashif
575 Rasul, Longhui Yu, Albert Q Jiang, Ziju Shen, et al. Numinamath: The largest public dataset in
576 ai4maths with 860k pairs of competition math problems and solutions. *Hugging Face repository*,
577 13:9, 2024.
- 578 Yujia Li, David Choi, Junyoung Chung, Nate Kushman, Julian Schrittwieser, Rémi Leblond, Tom
579 Eccles, James Keeling, Felix Gimeno, Agustin Dal Lago, Thomas Hubert, Peter Choy, Cyprien
580 de Masson d’Autume, Igor Babuschkin, Xinyun Chen, Po-Sen Huang, Johannes Welbl, Sven
581 Goyal, Alexey Cherepanov, James Molloy, Daniel Mankowitz, Esme Sutherland Robson, Push-
582 meet Kohli, Nando de Freitas, Koray Kavukcuoglu, and Oriol Vinyals. Competition-level code
583 generation with alphacode. *arXiv preprint arXiv:2203.07814*, 2022.
- 584 Hunter Lightman, Vineet Kosaraju, Yuri Burda, Harrison Edwards, Bowen Baker, Teddy Lee, Jan
585 Leike, John Schulman, Ilya Sutskever, and Karl Cobbe. Let’s verify step by step. In *The Twelfth*
586 *International Conference on Learning Representations*, 2023.
- 587 Hao Liu and Pieter Abbeel. Behavior from the void: Unsupervised active pre-training. *Advances in*
588 *Neural Information Processing Systems*, 34:18459–18473, 2021.
- 589 Potsawee Manakul, Adian Liusie, and Mark JF Gales. Selfcheckgpt: Zero-resource black-box hallu-
590 cination detection for generative large language models. *arXiv preprint arXiv:2303.08896*, 2023.

- 594 Fabio Maria Carlucci, Lorenzo Porzi, Barbara Caputo, Elisa Ricci, and Samuel Rota Buló. Autodial:
595 Automatic domain alignment layers. In *Proceedings of the IEEE international conference on*
596 *computer vision*, pages 5067–5075, 2017.
- 597 Zachary Nado, Shreyas Padhy, D Sculley, Alexander D’Amour, Balaji Lakshminarayanan, and
598 Jasper Snoek. Evaluating prediction-time batch normalization for robustness under covariate
599 shift. *arXiv preprint arXiv:2006.10963*, 2020.
- 600 Deepak Pathak, Pulkit Agrawal, Alexei A Efros, and Trevor Darrell. Curiosity-driven exploration by
601 self-supervised prediction. In *International conference on machine learning*, pages 2778–2787.
602 PMLR, 2017.
- 603 Yuxiao Qu, Tianjun Zhang, Naman Garg, and Aviral Kumar. Recursive introspection: Teaching
604 language model agents how to self-improve. *Advances in Neural Information Processing Systems*,
605 37:55249–55285, 2024.
- 606 Joaquin Quiñero-Candela, Masashi Sugiyama, Anton Schwaighofer, and N Lawrence. Covariate
607 shift and local learning by distribution matching. *Dataset Shift in Machine Learning*, pages 131–
608 160, 2008.
- 609 David Rein, Betty Li Hou, Asa Cooper Stickland, Jackson Petty, Richard Yuanzhe Pang, Julien Di-
610 rani, Julian Michael, and Samuel R Bowman. Gpqa: A graduate-level google-proof q&a bench-
611 mark. In *First Conference on Language Modeling*, 2024.
- 612 David Rolnick, Andreas Veit, Serge Belongie, and Nir Shavit. Deep learning is robust to massive
613 label noise. *arXiv preprint arXiv:1705.10694*, 2017.
- 614 Kuniaki Saito, Donghyun Kim, Stan Sclaroff, Trevor Darrell, and Kate Saenko. Semi-supervised do-
615 main adaptation via minimax entropy. In *Proceedings of the IEEE/CVF international conference*
616 *on computer vision*, pages 8050–8058, 2019.
- 617 Steffen Schneider, Evgenia Rusak, Luisa Eck, Oliver Bringmann, Wieland Brendel, and Matthias
618 Bethge. Improving robustness against common corruptions by covariate shift adaptation. *Ad-
619 vances in neural information processing systems*, 33:11539–11551, 2020.
- 620 John Schulman, Filip Wolski, Prafulla Dhariwal, Alec Radford, and Oleg Klimov. Proximal policy
621 optimization algorithms. *arXiv preprint arXiv:1707.06347*, 2017.
- 622 Rulin Shao, Shuyue Stella Li, Rui Xin, Scott Geng, Yiping Wang, Sewoong Oh, Simon Shaolei
623 Du, Nathan Lambert, Sewon Min, Ranjay Krishna, et al. Spurious rewards: Rethinking training
624 signals in rlvr. *arXiv preprint arXiv:2506.10947*, 2025.
- 625 Zhihong Shao, Peiyi Wang, Qihao Zhu, Runxin Xu, Junxiao Song, Xiao Bi, Haowei Zhang,
626 Mingchuan Zhang, YK Li, Y Wu, et al. Deepseekmath: Pushing the limits of mathematical
627 reasoning in open language models. *arXiv preprint arXiv:2402.03300*, 2024.
- 628 Guangming Sheng, Chi Zhang, Zilingfeng Ye, Xibin Wu, Wang Zhang, Ru Zhang, Yanghua Peng,
629 Haibin Lin, and Chuan Wu. Hybridflow: A flexible and efficient rlhf framework. *arXiv preprint*
630 *arXiv: 2409.19256*, 2024.
- 631 Rui Shu, Hung H Bui, Hirokazu Narui, and Stefano Ermon. A dirt-t approach to unsupervised
632 domain adaptation. *arXiv preprint arXiv:1802.08735*, 2018.
- 633 Avi Singh, John D Co-Reyes, Rishabh Agarwal, Ankesh Anand, Piyush Patil, Xavier Garcia, Peter J
634 Liu, James Harrison, Jaehoon Lee, Kelvin Xu, et al. Beyond human data: Scaling self-training
635 for problem-solving with language models. *arXiv preprint arXiv:2312.06585*, 2023.
- 636 Claudio Spiess, David Gros, Kunal Suresh Pai, Michael Pradel, Md Rafiqul Islam Rabin, Amin
637 Alipour, Susmit Jha, Prem Devanbu, and Toufique Ahmed. Calibration and correctness of lan-
638 guage models for code. *arXiv preprint arXiv:2402.02047*, 2024.
- 639 Baochen Sun, Jiashi Feng, and Kate Saenko. Correlation alignment for unsupervised domain adap-
640 tation. *Domain adaptation in computer vision applications*, pages 153–171, 2017.

- 648 Katherine Tian, Eric Mitchell, Allan Zhou, Archit Sharma, Rafael Rafailov, Huaxiu Yao, Chelsea
649 Finn, and Christopher D Manning. Just ask for calibration: Strategies for eliciting calibrated
650 confidence scores from language models fine-tuned with human feedback. *arXiv preprint*
651 *arXiv:2305.14975*, 2023.
- 652 Jonathan Uesato, Nate Kushman, Ramana Kumar, Francis Song, Noah Siegel, Lisa Wang, Antonia
653 Creswell, Geoffrey Irving, and Irina Higgins. Solving math word problems with process-and
654 outcome-based feedback. *arXiv preprint arXiv:2211.14275*, 2022.
- 656 Neeraj Varshney, Wenlin Yao, Hongming Zhang, Jianshu Chen, and Dong Yu. A stitch in time saves
657 nine: Detecting and mitigating hallucinations of llms by validating low-confidence generation.
658 *arXiv preprint arXiv:2307.03987*, 2023.
- 659 Yuvraj Virk, Premkumar Devanbu, and Toufique Ahmed. Enhancing trust in llm-generated code
660 summaries with calibrated confidence scores. *arXiv preprint arXiv:2404.19318*, 2024.
- 662 Dequan Wang, Evan Shelhamer, Shaoteng Liu, Bruno Olshausen, and Trevor Darrell. Tent: Fully
663 test-time adaptation by entropy minimization. *arXiv preprint arXiv:2006.10726*, 2020.
- 665 Peiyi Wang, Lei Li, Zhihong Shao, RX Xu, Damai Dai, Yifei Li, Deli Chen, Y Wu, and Zhifang
666 Sui. Math-shepherd: A label-free step-by-step verifier for llms in mathematical reasoning. *arXiv*
667 *preprint arXiv:2312.08935*, 2023.
- 668 Shenzhi Wang, Le Yu, Chang Gao, Chujie Zheng, Shixuan Liu, Rui Lu, Kai Dang, Xionghui Chen,
669 Jianxin Yang, Zhenru Zhang, et al. Beyond the 80/20 rule: High-entropy minority tokens drive
670 effective reinforcement learning for llm reasoning. *arXiv preprint arXiv:2506.01939*, 2025.
- 672 Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Fei Xia, Ed Chi, Quoc V Le, Denny
673 Zhou, et al. Chain-of-thought prompting elicits reasoning in large language models. *Advances in*
674 *neural information processing systems*, 35:24824–24837, 2022.
- 675 Liangru Xie, Hui Liu, Jingying Zeng, Xianfeng Tang, Yan Han, Chen Luo, Jing Huang, Zhen Li,
676 Suhang Wang, and Qi He. A survey of calibration process for black-box llms. *arXiv preprint*
677 *arXiv:2412.12767*, 2024.
- 679 Miao Xiong, Zhiyuan Hu, Xinyang Lu, Yifei Li, Jie Fu, Junxian He, and Bryan Hooi. Can llms
680 express their uncertainty? an empirical evaluation of confidence elicitation in llms. *arXiv preprint*
681 *arXiv:2306.13063*, 2023.
- 682 Tianyang Xu, Shujin Wu, Shizhe Diao, Xiaoze Liu, Xingyao Wang, Yangyi Chen, and Jing Gao.
683 Sayself: Teaching llms to express confidence with self-reflective rationales. In *Proceedings of*
684 *the 2024 Conference on Empirical Methods in Natural Language Processing*, pages 5985–5998,
685 2024.
- 687 Daniel Yang, Yao-Hung Hubert Tsai, and Makoto Yamada. On verbalized confidence scores for
688 llms. *arXiv preprint arXiv:2412.14737*, 2024.
- 689 Shunyu Yao, Dian Yu, Jeffrey Zhao, Izhak Shafran, Tom Griffiths, Yuan Cao, and Karthik
690 Narasimhan. Tree of thoughts: Deliberate problem solving with large language models. *Ad-*
691 *vances in neural information processing systems*, 36:11809–11822, 2023.
- 693 Denis Yarats, Rob Fergus, Alessandro Lazaric, and Lerrel Pinto. Reinforcement learning with proto-
694 typical representations. In *International Conference on Machine Learning*, pages 11920–11931.
695 PMLR, 2021.
- 697 Dongkeun Yoon, Seungone Kim, Sohee Yang, Sunkyoung Kim, Soyeon Kim, Yongil Kim, Eunbi
698 Choi, Yireun Kim, and Minjoon Seo. Reasoning models better express their confidence. *arXiv*
699 *preprint arXiv:2505.14489*, 2025.
- 700 Eric Zelikman, Yuhuai Wu, Jesse Mu, and Noah Goodman. Star: Bootstrapping reasoning with
701 reasoning. *Advances in Neural Information Processing Systems*, 35:15476–15488, 2022.

702 Qingyang Zhang, Haitao Wu, Changqing Zhang, Peilin Zhao, and Yatao Bian. Right question
703 is already half the answer: Fully unsupervised llm reasoning incentivization. *arXiv preprint*
704 *arXiv:2504.05812*, 2025.

705
706 Xuandong Zhao, Zhewei Kang, Aosong Feng, Sergey Levine, and Dawn Song. Learning to reason
707 without external rewards. *arXiv preprint arXiv:2505.19590*, 2025.

708
709 Yuxin Zuo, Kaiyan Zhang, Shang Qu, Li Sheng, Xuekai Zhu, Biqing Qi, Youbang Sun, Ganqu
710 Cui, Ning Ding, and Bowen Zhou. Ttrl: Test-time reinforcement learning. *arXiv preprint*
711 *arXiv:2504.16084*, 2025.

712
713
714
715
716
717
718
719
720
721
722
723
724
725
726
727
728
729
730
731
732
733
734
735
736
737
738
739
740
741
742
743
744
745
746
747
748
749
750
751
752
753
754
755

A APPENDIX

A.1 HYPERPARAMETERS

A full list of hyperparameters can be found in Table 4.

A.2 MORE QUALITATIVE SAMPLES

Table 5 shows a qualitative sample from AIME. The model indeed learns meaningful reasoning skills via entropy minimization. It is not merely learning to format its answer correctly, or otherwise collapsing to some other reward-hacking behavior.

Table 4: Hyperparameters.

Hyperparameter	Value
Max prompt length	1024
Max response length	3072
Batch size	64 GSM8K 500 MATH500
	80 AMC
	30 AIME
	64 Countdown
	196 GPQA
Policy mini batch size	32 GSM8K 32 MATH500
	80 AMC
	30 AIME
	32 Countdown
	32 GPQA
Policy micro batch size per GPU	8
Learning rate	1×10^{-6}
Weight decay	0.01
Learning rate warmup	Constant
Optimizer	Adam
Temperature	1.0 for train 0.8 for validation
Top k	-1
Top p	1
Number of samples per example n	5
Remove padding	True
Use KL loss	True
KL loss coefficient	0.001
Clip ratio	0.2
Grad clip	1.0

A.3 PASS@K CURVES

Figure 5, Figure 6, and Figure 7 show the pass@k performance curves for AIME, GSM8K, and MATH500 respectively.

A.4 RESPONSE LENGTH

Figure 8 shows the mean response length of Qwen2.5-7B-Instruct over the course of training on AIME. Overall, the response length increases during training.

A.5 MORE CONCURRENT WORKS

Cheng et al. (2025) introduce an entropy-augmented advantage term in RL to encourage longer, deeper chain-of-thought exploration (addressing exploitation-induced plateaus). This is different from our test-time scenario where we show performance gains with no external rewards.

Table 5: The model identifies the geometric series and succinctly solves the problem afterwards

<p>810</p> <p>811</p> <p>812</p> <p>813</p> <p>814</p> <p>815</p> <p>816</p> <p>817</p> <p>818</p> <p>819</p> <p>820</p> <p>821</p>	<p>system</p> <p>Please reason step by step, and put your final answer within $\boxed{\}$.</p> <p>user</p> <p>The infinite product</p> $\sqrt[3]{10} \cdot \sqrt[3]{\sqrt[3]{10}} \cdot \sqrt[3]{\sqrt[3]{\sqrt[3]{10}}} \dots$ <p>evaluates to a real number \sqrt{m}. What is the value of m? Let's think step by step and output the final answer within $\boxed{\}$.</p> <p>assistant</p>
<p>822</p> <p>823</p> <p>824</p> <p>825</p> <p>826</p> <p>827</p> <p>828</p> <p>829</p> <p>830</p> <p>831</p> <p>832</p> <p>833</p> <p>834</p> <p>835</p> <p>836</p> <p>837</p> <p>838</p> <p>839</p> <p>840</p> <p>841</p> <p>842</p> <p>843</p> <p>844</p> <p>845</p>	<p>Response</p> <p>To find the value of the infinite product $\sqrt[3]{10} \cdot \sqrt[3]{\sqrt[3]{10}} \cdot \sqrt[3]{\sqrt[3]{\sqrt[3]{10}}} \dots$, we can express each term in the product as a power of 10.</p> <p>The first term is $\sqrt[3]{10} = 10^{1/3}$. The second term is $\sqrt[3]{\sqrt[3]{10}} = (10^{1/3})^{1/3} = 10^{1/9}$. The third term is $\sqrt[3]{\sqrt[3]{\sqrt[3]{10}}} = (10^{1/9})^{1/3} = 10^{1/27}$. And so on.</p> <p>The infinite product can be written as:</p> $10^{1/3} \cdot 10^{1/9} \cdot 10^{1/27} \dots = 10^{1/3+1/9+1/27+\dots}$ <p>The exponent is an infinite geometric series with the first term $a = \frac{1}{3}$ and common ratio $r = \frac{1}{3}$. The sum S of an infinite geometric series is given by:</p> $S = \frac{a}{1-r} = \frac{\frac{1}{3}}{1-\frac{1}{3}} = \frac{\frac{1}{3}}{\frac{2}{3}} = \frac{1}{2}.$ <p>So the infinite product evaluates to:</p> $10^{1/2} = \sqrt{10}.$ <p>Thus, the value of m is 10. The final answer is:</p> $\boxed{10}$

846 Cui et al. (2025) study the dynamics of RL training for reasoning and identify the entropy collapse problem where policy entropy drops early, stalling performance. While this may be bad for a longer-term training-time scenario, our goal during test-time training is to maximize the immediate performance gains of the model, so we do not believe this should be a downside for our method.

850 Wang et al. (2025) analyze which tokens matter in reasoning by examining token-level entropy during RL with verifiable rewards. They find that only a small minority of high-entropy tokens drive most reasoning improvements, and show that restricting policy updates to those tokens yields large gains (a beyond-80/20 effect). We believe this is complementary work; it shows that certain tokens are more impactful during training, but does not use their entropy directly in any kind of reward. Additionally, the focus of this paper is training-time improvements using ground-truth labels, while we focus on unlabeled test-time improvements.

857 Zhang et al. (2025) propose EMPO, a fully unsupervised method that (like ours) minimizes a model's predictive entropy to incentivize better reasoning. EMPO uses a semantic confidence reward where they first sample many responses, group them into similar meaning, and penalize the model for the number of meaning clusters it has. This is distinct from our token-level entropy minimization, which does not require multiple samples per update, and instead uses the mean token entropy as a lower-level confidence reward.

863 Agarwal et al. (2025) similarly incorporates token level entropy minimization as an unsupervised reward for RL. However, we note that this paper is concurrent work to ours. Additionally, we

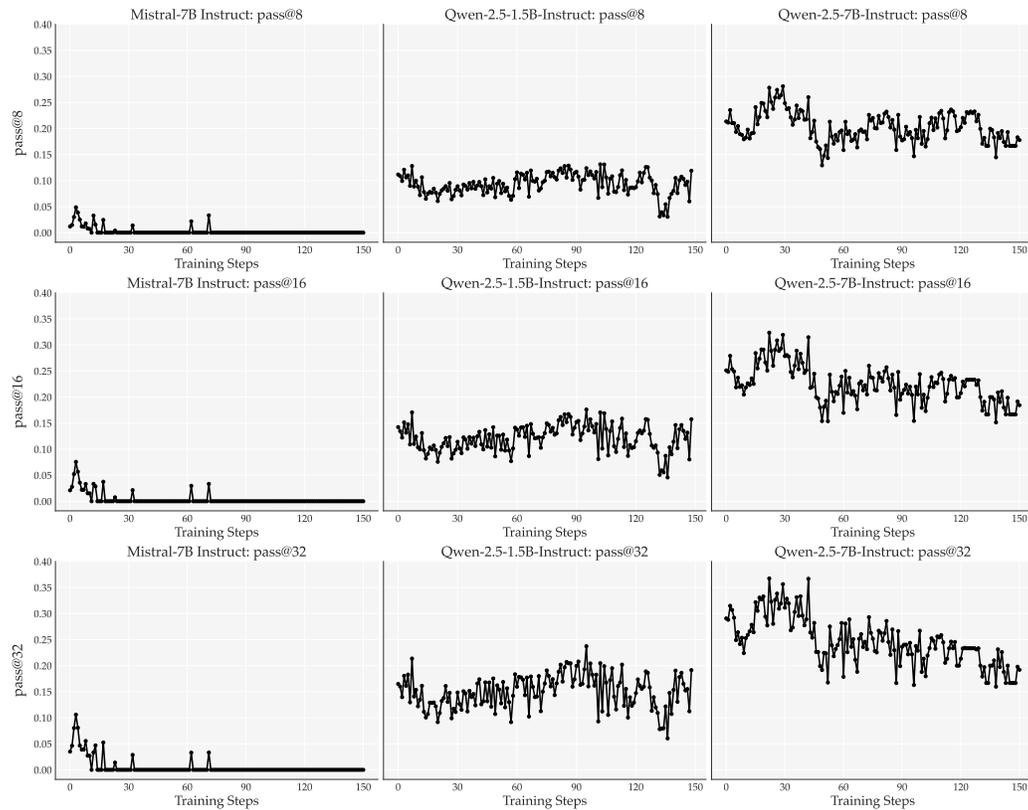


Figure 5: Pass@k curves on the AIME dataset.

believe our experiments on token minimization patterns offer unique, meaningful insights into the token minimization technique that distinguish our work from this paper.

Jiao et al. (2024) take a different approach by optimizing reasoning via pseudo feedback from frontier LLMs in a preference-based RL framework. In contrast, our method relies solely on the model’s own entropy as the reward to improve reasoning.

918
919
920
921
922
923
924
925
926
927
928
929
930
931
932
933
934
935
936
937
938
939
940
941
942
943
944
945
946
947
948
949
950
951
952
953
954
955
956
957
958
959
960
961
962
963
964
965
966
967
968
969
970
971

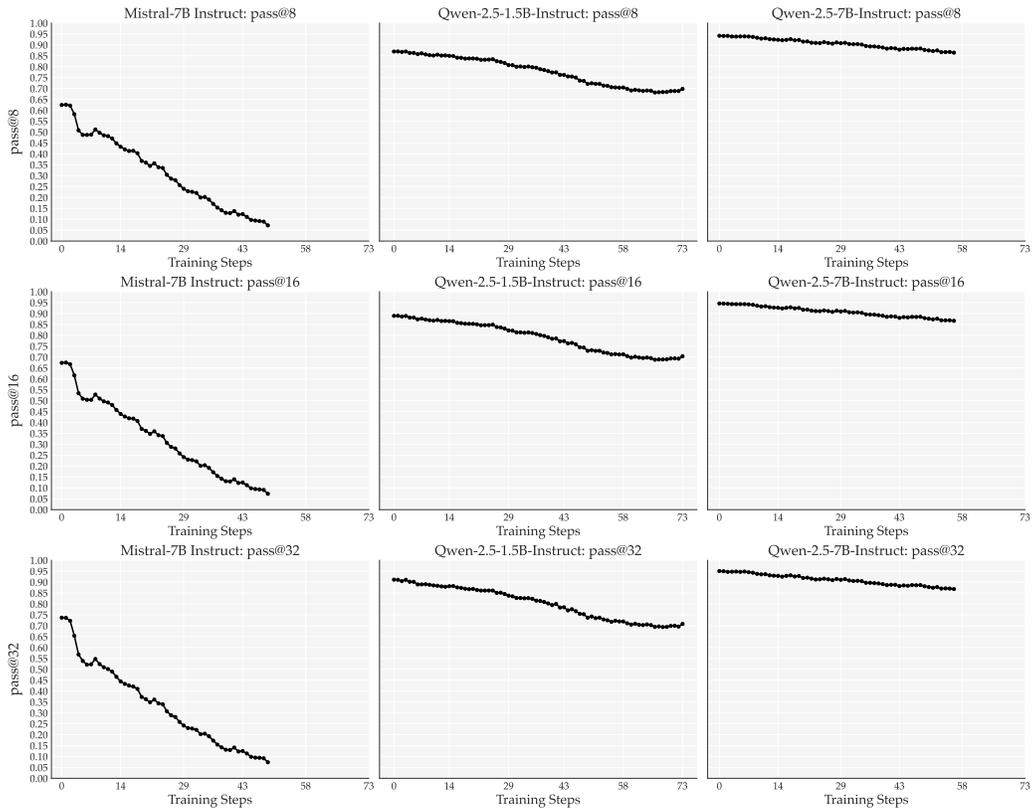


Figure 6: Pass@k curves on the GSM8K dataset.

972
 973
 974
 975
 976
 977
 978
 979
 980
 981
 982
 983
 984
 985
 986
 987
 988
 989
 990
 991
 992
 993
 994
 995
 996
 997
 998
 999
 1000
 1001
 1002
 1003
 1004
 1005
 1006
 1007
 1008
 1009
 1010
 1011
 1012
 1013
 1014
 1015
 1016
 1017
 1018
 1019
 1020
 1021
 1022
 1023
 1024
 1025

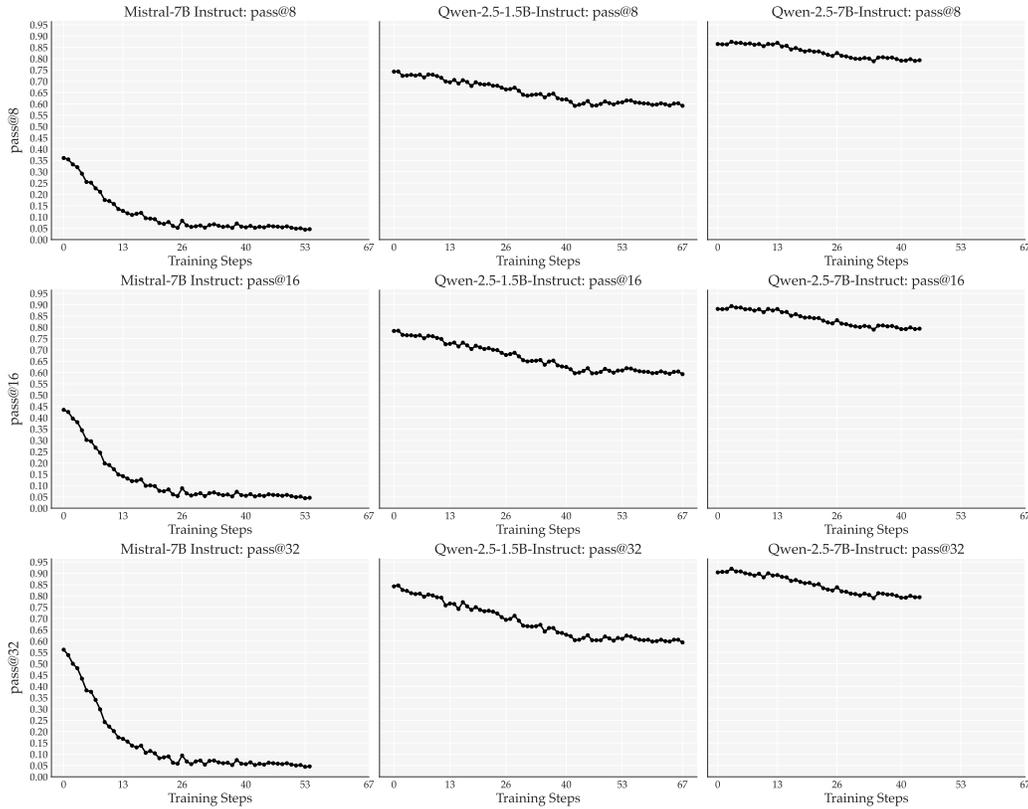


Figure 7: Pass@k curves on the MATH500 dataset.

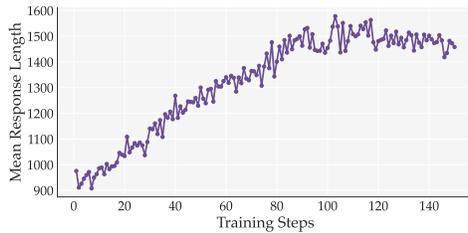


Figure 8: Mean response length of Qwen2.5-7B-Instruct over the course of training on AIME.