

# Think Earlier, Not Longer: Prompt Optimization via Reducing Unhealthy Exploration

Anonymous ACL submission

## Abstract

While large language models exhibit strong reasoning capabilities, prior work shows that their performance can be further enhanced by encouraging greater exploration. However, existing approaches overlook the presence of unhealthy exploration that increases exploration-related token usage without contributing to effective problem-solving. In this work, we show that prompt ambiguity can artificially prolong early-stage exploration, manifested as an elevated and delayed early-stage entropy peak. Although this uncertainty may be gradually resolved as reasoning progresses, reflected in the eventual convergence of the late-stage entropy peak, it does not meaningfully improve accuracy or self-consistency and instead substantially reduces reasoning efficiency. Motivated by these observations, we propose an entropy-dynamics-aware prompt optimization framework that trains a lightweight optimizer to generate concise clarifications. These clarifications aim to reduce ambiguity-induced early-stage uncertainty while preserving the model’s reasoning capabilities. Extensive experiments across multiple models, reasoning budgets, and benchmarks demonstrate that our approach consistently improves reasoning efficiency by up to 52%, by reducing unhealthy exploration without sacrificing accuracy.

## 1 Introduction

Recent work on large language model (LLM) reasoning has proposed a variety of techniques to improve reasoning performance by encouraging greater exploration during generation, such as test-time scaling (Muennighoff et al., 2025; Snell et al., 2025) and parallel reasoning (Wang et al., 2022; Yao et al., 2023; Besta et al., 2024). However, these approaches implicitly assume that increased exploration is uniformly beneficial, overlooking the presence of unhealthy exploration that substantially increases token usage without contributing to effective downstream problem-solving.

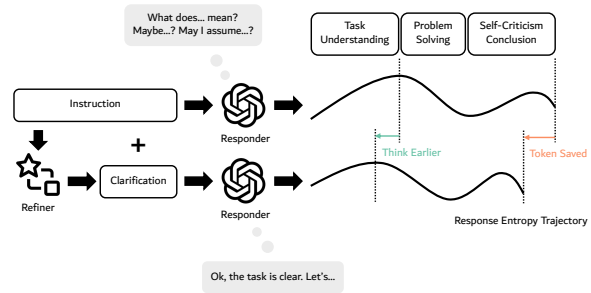


Figure 1: Early-stage exploration induced by prompt ambiguity can be effectively mitigated through concise clarification, thereby improving reasoning efficiency.

To expose this phenomenon, we conduct controlled prompt ambiguity experiments that isolate the effect of task underspecification while keeping the underlying reasoning task unchanged. Specifically, we systematically vary prompt ambiguity by masking technical terms with their abbreviations. Across multiple models and reasoning budgets, increased prompt ambiguity leads to a substantial inflation in token usage, while no improvement in either accuracy or self-consistency.

By decomposing token-level entropy trajectories into three stages: task understanding (early), problem solving (middle), and self-criticism (late), we observe that the early-stage entropy peak becomes both elevated and delayed. In contrast, the late-stage entropy peak remains highly convergent. This observation indicates that exploration induced by prompt ambiguity is progressively resolved rather than directly advancing problem-solving.

Motivated by these observations, we propose a prompt optimization framework that trains a lightweight optimizer via multi-turn reinforcement learning to generate concise clarifications for the original prompt. We design an entropy-peak-based reward that encourages the generated clarifications to both reduce and advance the early-stage entropy peak, while preserving the late-stage entropy peak. In other words, our approach enables the model to

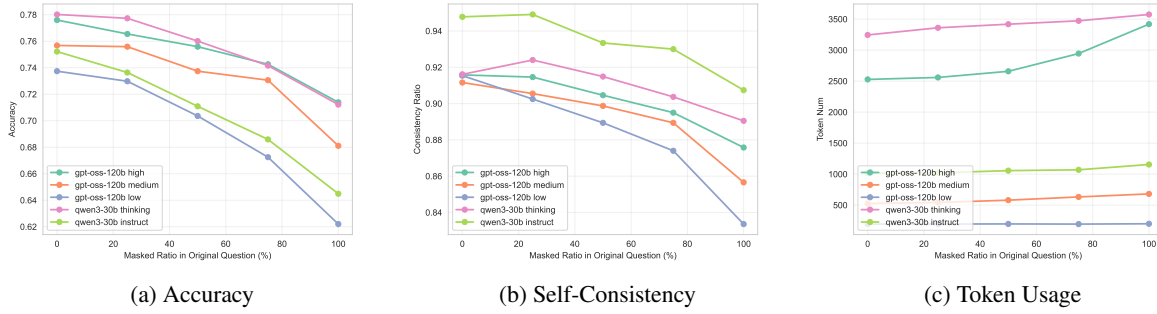


Figure 2: Accuracy, consistency, and token usage of GPT-oss-120B under three levels of reasoning effort, Qwen3-30B-Instruct and Qwen3-30B-Thinking as the mask ratio increases.

understand the task more efficiently, without harming its intrinsic multi-step reasoning capability.

Our experimental results show that the proposed prompt optimization method substantially improves reasoning efficiency, achieving up to 52.9% gains on in-domain benchmarks. Moreover, our method demonstrates strong generalization ability, improving reasoning efficiency by up to 38.54% on out-of-domain benchmarks.

In summary, the main contributions of our work are as follows:

- We conduct controlled prompt ambiguity experiments, revealing that ambiguity-driven unhealthy exploration artificially extends the early reasoning stage without improving problem-solving performance.
- We propose an entropy-dynamics-aware prompt optimization framework that trains a lightweight optimizer to generate concise clarifications, selectively reducing unhealthy early-stage uncertainty while preserving healthy reasoning.
- Extensive experiments across multiple models, reasoning budgets, and both in-domain and out-of-domain benchmarks demonstrate that our method consistently improves reasoning efficiency without sacrificing accuracy.

## 2 Entropy Dynamics Observation

### 2.1 Controlled Prompt Ambiguity

To elicit unhealthy exploration, we manipulate prompt ambiguity rather than reasoning depth. When key terms or constraints in a prompt are underspecified, the model expends additional tokens on task interpretation and linguistic clarification before committing to a concrete reasoning trajectory.

Specifically, we sample 1,000 instances from MMLU-Pro (Wang et al., 2024) and use GPT-5<sup>1</sup> to identify technical terms, which are then replaced by their initialisms (e.g., page fault → PF). By varying the proportion of masked terms, we systematically control the degree of prompt ambiguity while keeping the reasoning task itself unchanged.

We evaluate GPT-oss-120B (Agarwal et al., 2025a) under three levels of reasoning effort (high, medium, and low) across five mask ratios (0%, 25%, 50%, 75%, and 100%). The temperature is set to 0.7, and for each question we perform 16 rollouts. As shown in Figure 2, increasing prompt ambiguity consistently degrades accuracy and self-consistency, while increasing token usage.

### 2.2 Token-Level Entropy Trajectories

To better understand how prompt ambiguity affects the generation process, we analyze token-level entropy trajectories under medium and high reasoning effort. Figures 3 and 4 report average token-level entropy curves over both relative positions and absolute positions.

Across settings, the entropy trajectory consistently exhibits a three-stage structure:

- **Early Stage:** Entropy increases sharply, reflecting the reasoning model’s prompt interpretation and identification;
- **Middle Stage:** Entropy gradually decreases as the model commits to a specific reasoning trajectory and focuses on problem solving;
- **Final Stage:** Entropy shows a mild increase followed by a sharp drop, presenting reasoning verification, and answer finalization.

<sup>1</sup><https://cdn.openai.com/gpt-5-system-card.pdf>

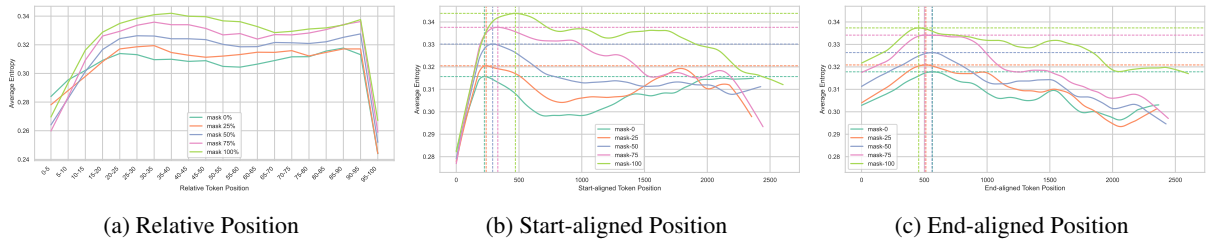


Figure 3: Token-level entropy of GPT-oss-120B (medium reasoning effort) across five mask ratios under three alignment schemes: relative position, which normalizes token indices by response length using 5% bins; start-aligned position, which uses absolute token indices aligned at the first generated token; and end-aligned position, which uses absolute token indices aligned at the final generated token (distance to termination).

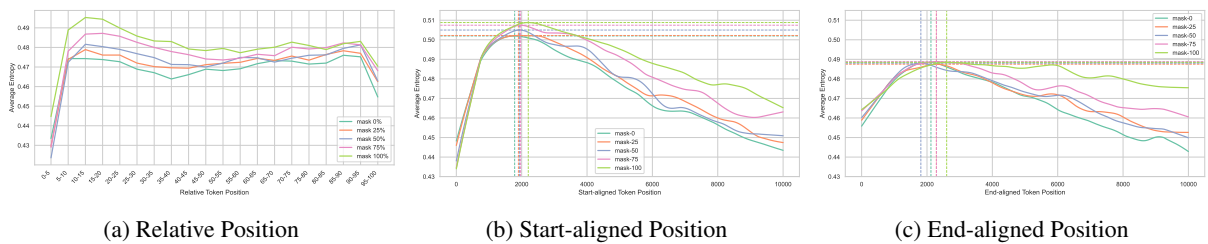


Figure 4: Token-level entropy of GPT-oss-120B (high reasoning effort).

### 2.3 Early-stage Uncertainty

A salient feature of the entropy trajectory is the early-stage entropy peak, which marks the transition from prompt interpretation to committed problem solving. In this phase, the model is actively forming a task understanding and may explore multiple plausible interpretations and solution directions. When the prompt is underspecified, this early-stage exploration is augmented by interpretive uncertainty induced by prompt ambiguity, manifesting as an elevated and delayed entropy peak.

As shown in Figures 3b and 4b, increasing the mask ratio consistently elevates the early-stage entropy peak and shifts it to later token positions. Increasing reasoning effort partially compensates for this increased early-stage uncertainty, not by resolving it efficiently, but by doing so at the cost of substantial token expansion. Under medium reasoning effort, the difference in early-stage peak height between mask ratio 0 and 100 is 0.027, accompanied by an increase of 335 generated tokens. Under high reasoning effort, the peak height difference is reduced to 0.006, yet the corresponding token increase rises to 437.

### 2.4 Late-Stage Convergence

As shown in Figures 3c and 4c, despite substantial differences in early-stage uncertainty, entropy trajectories exhibit a contrasting pattern near generation termination. When responses are aligned

by their termination points, entropy dynamics in the late stage become increasingly similar across different mask ratios.

Notably, higher reasoning effort leads to stronger late-stage convergence. Under medium reasoning effort, the difference in late-stage peak height across mask ratios is 0.019, representing a compression to approximately 70% of the corresponding early-stage difference. Under high reasoning effort, this difference further collapses to  $5.1 \times 10^{-5}$ , amounting to only 0.8% of the early-stage peak difference. This trend mirrors the observation that higher reasoning effort also exhibits greater robustness to the accuracy and self-consistency degradation introduced by prompt ambiguity.

Importantly, this convergence does not imply the disappearance of all early-stage exploration. As the model commits to a task interpretation, the ambiguity-induced unhealthy components component is progressively resolved, while structurally necessary healthy components persists and is carried forward into the late stage, where it supports answer verification and finalization.

## 3 Prompt Optimization

Building on the observations in Section 2, we design a framework that trains a lightweight prompt optimizer using entropy-dynamics-guided multi-turn reinforcement learning. The optimizer learns to generate a clarification that helps the reasoning

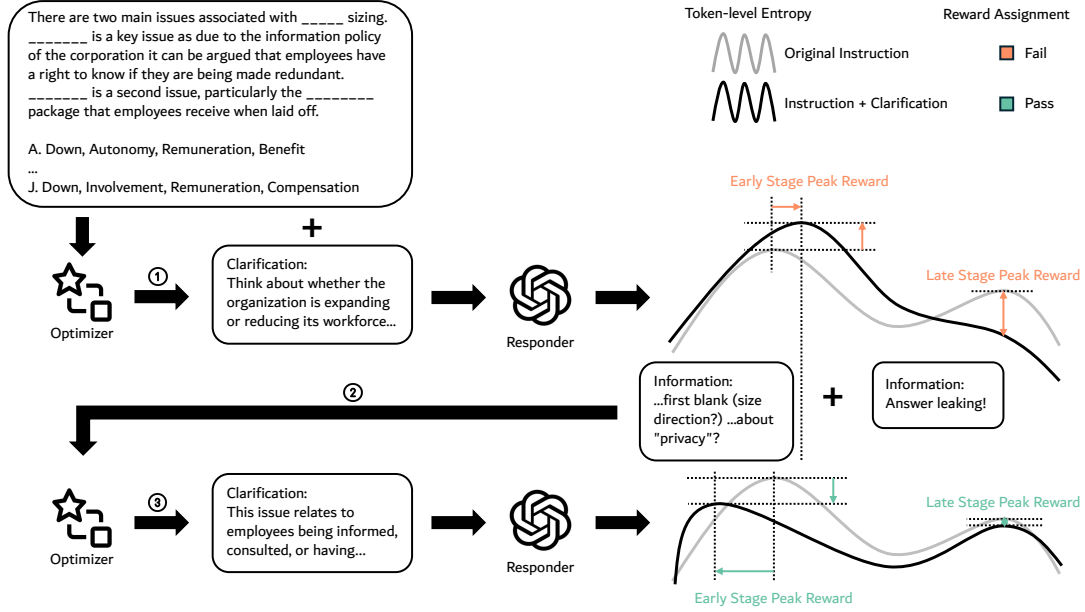


Figure 5: Overview of the proposed prompt optimizer training framework. The optimizer is first given the original prompt and generates a initial clarification (Step 1). This clarification is then appended to the original prompt and provided to the Responder. If the resulting response passes the entropy peak check, the clarification is considered qualified. Otherwise, feedback information is presented to the optimizer (Step 2), which then generates a refined clarification in the next turn (Step 3). This process repeats until a qualified clarification is produced or the maximum number of turns is reached. Finally, the last-turn clarification reward is used to update the optimizer.

model more accurately interpret the prompt. By reducing the prevalence of ambiguity-induced unhealthy exploration, this approach effectively mitigates unnecessary token expenditure and improves reasoning efficiency. An overview of the training process is shown in Figure 5.

### 3.1 Problem Formulation

To clarify our objective, we distinguish between two qualitatively different forms of exploration. Unhealthy exploration refers to uncertainty-driven token expenditure in the early stage of generation, where the model has not yet formed a clear understanding of the task. This form of exploration manifests as elevated and delayed early-stage entropy. Importantly, although this ambiguity-induced exploration is largely self-resolving as the model commits to a task interpretation, it is inefficient, contributing directly to response length inflation without improving reasoning quality.

In contrast, healthy exploration occurs after the task has been sufficiently specified, when the model actively explores multiple plausible reasoning paths. This form of exploration is a desirable component of complex reasoning and is not eliminated as task interpretation stabilizes. Instead, it also contributes to answer verification, and ultimately

leading to more reliable conclusions.

Our goal is to improve reasoning efficiency by selectively reducing unhealthy exploration while preserving healthy exploration. To this end, we introduce a lightweight prompt optimizer  $P_\theta$  that augments the original prompt  $x$  with a concise clarification:

$$x' = x \oplus P_\theta(x), \quad (1)$$

The clarification is restricted to clarifying ambiguous aspects of the task formulation and must not reveal solution steps, or answer choices.

### 3.2 Entropy Peak Reward

The prompt optimizer is guided solely by entropy dynamics, which reflect the internal reasoning behavior of the model, without any correctness-based supervision. This design ensures that the optimizer focuses on improving how the reasoning model commits to a solution, rather than what answer it produces.

Given an augmented prompt  $x'$ , the frozen reasoning model generates a response of length  $T$ . At each generation step  $t$ , the model produces a predictive distribution  $p_t(\cdot | x', y_{<t})$ , whose token-level entropy is defined as

$$e_t = - \sum_v p_t(v) \log p_t(v). \quad (2)$$

**Early Stage Peak Reward** The primary objective of prompt optimization is to reduce early stage peak reward and encourage earlier commitment to a reasoning trajectory.

We focus on the early portion of the response, corresponding to the prompt interpretation phase. Let  $T^{\text{early}} = \lfloor \alpha T \rfloor$ , where  $\alpha \in (0, 1)$  is a fixed ratio. We apply a smoothing operator  $\mathcal{S}(\cdot)$  to the entropy sequence and define the early-stage entropy peak as

$$\begin{aligned} p^{\text{early}} &= \arg \max_{1 \leq t \leq T^{\text{early}}} \mathcal{S}(e_t), \\ h^{\text{early}} &= \max_{1 \leq t \leq T^{\text{early}}} \mathcal{S}(e_t), \end{aligned} \quad (3)$$

where  $p^{\text{early}}$  denotes the peak position and  $h^{\text{early}}$  the peak height. Intuitively,  $p^{\text{early}}$  measures how long the model remains uncertain before committing to a reasoning trajectory, while  $h^{\text{early}}$  quantifies the severity of uncertainty.

To obtain robust estimates, we perform  $K$  independent rollouts of the reasoning model. For each rollout  $k$ , we extract a peak pair  $(p^{\text{early},(k)}, h^{\text{early},(k)})$  and aggregate them as

$$\begin{aligned} \bar{p}^{\text{early}} &= \frac{1}{K} \sum_{k=1}^K p^{\text{early},(k)}, \\ \bar{h}^{\text{early}} &= \frac{1}{K} \sum_{k=1}^K h^{\text{early},(k)}. \end{aligned} \quad (4)$$

For each prompt, baseline statistics  $\bar{p}_0^{\text{early}}$  and  $\bar{h}_0^{\text{early}}$  are precomputed using the original prompt. The effect of optimization is then measured by

$$\begin{aligned} \Delta p^{\text{early}} &= \bar{p}_0^{\text{early}} - \bar{p}^{\text{early}}, \\ \Delta h^{\text{early}} &= \bar{h}_0^{\text{early}} - \bar{h}^{\text{early}}. \end{aligned} \quad (5)$$

Finally, we define the early-stage peak reward as

$$R^{\text{early}} = \mathbb{I}(\Delta p^{\text{early}} > 0 \wedge \Delta h^{\text{early}} > 0), \quad (6)$$

which assigns a positive reward only when prompt optimization successfully reduces both the timing and magnitude of the early-stage entropy peak.

**Late Stage Peak Reward** The late-stage peak reward serves as a regularization signal rather than a primary optimization objective, ensuring that healthy exploration in later stages is preserved.

To enforce this constraint, we additionally monitor entropy behavior in the termination neighborhood. Let  $T^{\text{late}} = \lceil (1 - \gamma)T \rceil$ , where  $\gamma \in (0, 1)$

specifies the fraction of tokens considered as the late stage. Using the same smoothing operator  $\mathcal{S}(\cdot)$ , we define the late-stage entropy peak height for a single rollout as

$$h^{\text{late}} = \max_{T^{\text{late}} \leq t \leq T} \mathcal{S}(e_t). \quad (7)$$

To obtain robust estimates, we again perform  $K$  independent rollouts. For each rollout  $k$ , we extract a late-stage peak height  $h^{\text{late},(k)}$  and aggregate them as

$$\bar{h}^{\text{late}} = \frac{1}{K} \sum_{k=1}^K h^{\text{late},(k)}. \quad (8)$$

As with the early-stage statistics, we precompute a baseline late-stage peak height  $\bar{h}_0^{\text{late}}$  using the original prompt. The effect of optimization is then measured by

$$\Delta h^{\text{late}} = \bar{h}_0^{\text{late}} - \bar{h}^{\text{late}}. \quad (9)$$

Finally, we define the late-stage peak reward as

$$R^{\text{late}} = \mathbb{I}(\Delta h^{\text{late}} < \Delta h^{\text{early}}), \quad (10)$$

which assigns a positive reward only when that the late-stage entropy peak remains closer to that of the original prompt than the early-stage peak.

### 3.3 Multi-Turn Reinforcement Learning

Directly optimizing early-stage uncertainty using single-step reinforcement learning is challenging, as a scalar reward provides limited information. To address this limitation, we introduce a multi-turn reinforcement-learning framework in which the optimizer iteratively refines its clarifications based on the feedback derived from the reasoning model until a qualified clarification is produced or the maximum number of turns is reached.

In each training episode, the optimizer first generates an initial clarification, which is appended to the original prompt and passed to the reasoning model. The resulting response is then analyzed to determine whether it satisfies the entropy peak reward criteria. If the early-stage entropy peak conditions are not met, the optimizer receives the reasoning model’s generation tokens preceding the early-stage entropy peak. These tokens indicate where the model begins to exhibit uncertainty. If the late-stage entropy peak conditions are not satisfied, the optimizer instead receives feedback signals that do not leak information that could directly influence the reasoning process of the model. The

Model	Effort	Method	In-domain Benchmarks						Out-of-domain Benchmarks					
			MMLU-Pro			SuperGPQA			BBH			MedQA		
			Acc.↑	Tok.↓	Eff.↑	Acc.↑	Tok.↓	Eff.↑	Acc.↑	Tok.↓	Eff.↑	Acc.↑	Tok.↓	Eff.↑
GPT-oss-120B	High	Pure	<b>80.60</b>	2383.35	3.38	49.40	4500.31	1.09	84.40	1496.53	5.63	<b>91.10</b>	1372.63	6.63
		Early-Stop	58.70	<b>1285.58</b>	4.56	28.30	<b>2069.67</b>	1.36	57.70	<b>865.16</b>	6.66	85.70	1164.39	7.36
		EvoPrompt	80.10	2943.89	2.72	<b>50.40</b>	5635.07	0.89	-	-	-	-	-	-
	Medium	Pure	79.30	566.69	13.99	47.20	1055.71	4.47	<b>84.40</b>	421.18	20.03	88.00	292.65	30.07
		Early-Stop	74.50	508.60	14.64	39.30	772.76	5.08	65.70	<b>306.84</b>	21.41	88.00	291.64	30.17
		EvoPrompt	79.90	654.63	12.20	49.50	1020.56	4.85	-	-	-	-	-	-
		Ours	<b>80.60</b>	<b>472.89</b>	<b>17.04</b>	<b>48.60</b>	<b>725.61</b>	<b>6.69</b>	83.90	371.77	<b>22.56</b>	<b>88.30</b>	<b>261.31</b>	<b>33.79</b>
	Low	Pure	75.70	216.22	35.01	42.70	294.56	14.49	83.30	206.78	40.28	<b>85.40</b>	100.86	84.67
		Early-Stop	74.80	213.46	35.04	41.60	281.57	14.77	72.20	<b>181.01</b>	39.88	85.30	100.79	84.63
		EvoPrompt	<b>76.60</b>	245.61	31.18	44.90	315.00	14.25	-	-	-	-	-	-
		Ours	76.20	<b>186.15</b>	<b>40.93</b>	<b>45.40</b>	<b>252.01</b>	<b>18.01</b>	<b>83.50</b>	188.35	<b>44.33</b>	84.90	<b>100.04</b>	<b>84.86</b>
	Qwen3-30B	Thinking	Pure	81.40	3316.56	2.45	<b>53.30</b>	4765.22	1.11	<b>85.90</b>	1798.92	4.77	<b>84.90</b>	1982.31
Early-Stop			54.70	<b>2138.95</b>	2.55	33.30	<b>2669.79</b>	1.24	40.20	<b>1012.05</b>	3.97	83.40	1945.48	4.28
EvoPrompt			79.60	3670.03	2.16	51.80	4111.81	1.25	-	-	-	-	-	-
Ours			<b>81.70</b>	2491.31	<b>3.27</b>	53.10	3795.22	<b>1.39</b>	85.30	1406.23	<b>6.06</b>	84.80	<b>1795.78</b>	<b>4.72</b>
Instruct		Pure	76.50	1147.30	6.66	52.50	1986.15	2.64	83.30	564.12	14.76	<b>73.40</b>	55.71	131.75
Early-Stop	52.10	<b>464.20</b>	<b>11.22</b>	32.00	<b>858.16</b>	<b>3.72</b>	54.00	<b>339.47</b>	15.90	72.70	<b>51.44</b>	141.32		
EvoPrompt	78.60	1704.12	4.61	49.50	2474.67	2.00	-	-	-	-	-	-		
Ours	<b>80.10</b>	875.32	9.15	<b>52.60</b>	1497.04	3.51	<b>83.40</b>	463.99	<b>17.97</b>	73.20	51.68	<b>141.64</b>		

Table 1: Accuracy (Acc.), token usage (Tok.) and reasoning efficiency (Eff.) of reasoning models evaluated on both in-domain and out-of-domain benchmarks. Bolded values indicate the best performance for each reasoning model.

optimizer is encouraged to generate a more targeted clarification in the subsequent turn.

The multi-turn interaction is introduced solely as a training scaffold to stabilize optimization and improve credit assignment. At inference time, the prompt optimizer operates in a single-shot manner, generating at most one clarification per prompt and introducing no additional interaction overhead.

## 4 Experiments

### 4.1 Setup

**Training Setting** We adopt Qwen3-4B-Instruct as a lightweight prompt optimizer. Based on observations from our entropy dynamics experiments, we set the early-stage ratio to 0.5 and the late-stage ratio to 0.3. We reuse 1,000 samples from MMLU-Pro and additionally sample 1,000 instances from SuperGPQA (Du et al., 2025) to construct the training set. These datasets span a wide range of professional domains, enabling the prompt optimizer to learn robust clarification strategies.

The target reasoning models used during training include GPT-oss-120B under low, medium, and high reasoning effort settings, as well as both the thinking and instruct variants of Qwen3-30B (Yang et al., 2025). For entropy peak reward computation, we set the reasoning model temperature to 0.7 and perform 16 rollouts to estimate the average peak position and peak height. We approximate entropy using only the top-5 token log probabilities.

**Evaluation Setting** We set the reasoning temperature to 0.0 and perform a single rollout to ensure reproducibility. The prompt template used for evaluation is provided in Appendix C. To quantify reasoning efficiency, we compute:

$$\text{Reasoning Efficiency} = \frac{\text{Accuracy}}{\text{Token Usage}} \quad (11)$$

### 4.2 Empire Results

**In-domain Evaluation** The main baselines we compare against are Pure, Early-Stop (Sharma and Chopra, 2025) and EvoPrompt (Tong et al., 2025). For the in-domain benchmarks, we follow the original implementation of the Early-Stop method and EvoPrompt method. Specifically, for the former, the entropy threshold is computed as the mean entropy of each reasoning model on the training sets of MMLU-Pro and SuperGPQA, and the patience parameter is set to 50. For the latter, we obtain optimized prompts by applying the evolutionary algorithm on the same training sets from MMLU-Pro and SuperGPQA.

As shown in Table 1, our method achieves the highest reasoning efficiency on most reasoning models, without harming, and in some cases even slightly improving, accuracy. In contrast, the Early-Stop method also demonstrates effectiveness in improving reasoning efficiency, but it carries the risk of reducing accuracy, as it directly truncates subsequent tokens once the entropy remains below a

TU Token	Enrichment	SC Token	Enrichment
interpret	4.01	double-check	2.54
what does	1.85	verify	1.58
means	1.53	mismatch	1.47
assume	1.38	recompute	1.24

Table 2: Task-understanding (TU) tokens enriched around early-stage entropy peaks and self-criticism (SC) tokens enriched around late-stage entropy peaks.

threshold for a predefined number of consecutive steps. Meanwhile, EvoPrompt primarily focuses on improving accuracy. Although it successfully enhances accuracy for most reasoning models, it often increases token usage, which in turn leads to lower reasoning efficiency in most cases.

**Out-of-domain Evaluation** To assess generalization, we further evaluate our method on out-of-domain benchmarks, including BBH (Suzgun et al., 2023) and MedQA (Jin et al., 2021). For the Early-Stop method, we estimate the entropy threshold by averaging the thresholds obtained for each reasoning model on MMLU-Pro and SuperGPQA, while keeping the patience parameter fixed at 50. In contrast, EvoPrompt requires a development set for prompt optimization and is therefore not applicable in this out-of-domain setting.

As shown in Table 1, our method achieves a significant improvement in reasoning efficiency across different reasoning models, indicating that the learned clarification strategies are robust to substantial domain shifts. We further report the total token consumption and the corresponding reasoning efficiency in Appendix D.

## 5 Analysis

### 5.1 Tokens Associated with Entropy Peaks

To quantify the association between entropy peaks and specific token types, we measure token enrichment using a normalized occurrence ratio. Specifically, for each token  $w$ , we define

$$\text{Enrich}(w) = \frac{C_{\text{seg}}(w)/N_{\text{seg}}}{C_{\text{other}}(w)/N_{\text{other}}}, \quad (12)$$

where  $C_{\text{seg}}(w)$  and  $N_{\text{seg}}$  denote the number of occurrences of  $w$  and the total number of tokens within entropy-peak-adjacent segments, and  $C_{\text{other}}(w)$  and  $N_{\text{other}}$  are defined analogously for all remaining positions in the response.

As shown in Table 2, early-stage entropy peaks are selectively associated with tokens indicating

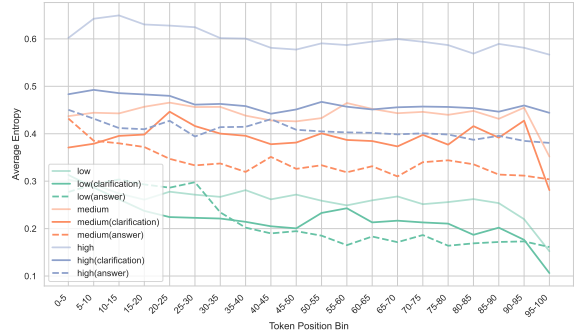


Figure 6: Relative-position entropy trajectories of GPT-oss-120B under high, medium and low reasoning effort, with clarification and answer conditioning.

task understanding and semantic clarification. In contrast, late-stage entropy peaks are characterized by the enrichment of self-criticism and verification-related tokens. This apparent enrichment of characteristic functional tokens provide empirical evidence that entropy peaks are not merely statistical artifacts, but are systematically aligned with distinct cognitive functions, further supporting our stage decomposition of the reasoning process.

### 5.2 Entropy Dynamics with Answer Conditioning

To investigate how answer leakage influences a model’s reasoning capability, we compare the token-level entropy dynamics of GPT-oss-120B with clarification-given prompting and answer-conditioned prompting, in which the answer is explicitly provided. As shown in Figure 6, under medium and high reasoning effort, answer-conditioned prompting leads to a sharp drop in entropy at the early stage, indicating a loss of exploratory behavior. Moreover, the late stage entropy peak almost entirely disappears. In contrast, when clarifications are provided, the model’s overall entropy trajectory remains much closer to that observed with the original prompting. This suggests that our prompt optimization method preserves the model’s intrinsic reasoning capability while avoiding answer leakage. Human evaluation in Appendix E further support this conclusion.

Under low reasoning effort, the behavior differs slightly. For both the original and clarification-given settings, entropy decreases in the early stage, reflecting the model’s limited exploratory capacity under constrained reasoning budgets. However, with answer-conditioned prompting, entropy initially increases before sharply dropping. We hy-

Model	Effort	Method	Early	Middle	Late
GPT-oss-120B	High	Pure	482.36	1564.49	336.48
		Ours	270.51	1435.80	245.97
	Medium	Pure	143.63	334.04	89.02
		Ours	67.69	304.69	70.49
Qwen3-30B	Thinking	Pure	881.13	1815.56	618.86
		Ours	445.55	1661.84	462.64
	Instruct	Pure	298.67	670.41	178.21
		Ours	222.51	651.29	161.50

Table 3: Average token usage of GPT-oss-120B and Qwen3-30B across early, middle and late stages.

461 pothesize that this entropy increase reflects a tran- 499  
462 sient conflict between the model’s internal reason- 500  
463 ing trajectory and the externally imposed answer- 501  
464 conditioned conclusion. This further highlights 502  
465 that answer conditioning fundamentally alters the 503  
466 model’s reasoning dynamics. 504

### 467 5.3 Stage-wise Token Usage Ablation 505

468 To examine whether the improvement in reasoning 506  
469 efficiency arises from reducing unhealthy explora- 507  
470 tion in the early stage, we conduct a token-usage 508  
471 stage ablation analysis by measuring the average 509  
472 number of tokens generated during the early, mid- 510  
473 dle, and late stages of generation. 511

474 As shown in Table 3, our method consistently re- 512  
475 duces early-stage token usage across all reasoning 513  
476 models. This suggests that the proposed strategy 514  
477 enables the model to converge to a stable semantic 515  
478 interpretation more efficiently, avoiding redundant 516  
479 reformulations and excessive exploratory phrasing 517  
480 at the beginning of generation. 518

481 In contrast, only modest reductions are observed 519  
482 in the middle and late stages. This indicates that 520  
483 sufficient exploration capacity for multi-step rea- 521  
484 soning is preserved. Such behavior aligns well with 522  
485 our design goal of avoiding over-regularization dur- 523  
486 ing the main reasoning phase. 524

## 487 6 Related Work 525

### 488 6.1 Entropy Dynamics 526

489 Entropy is widely used to probe the internal behav- 527  
490 ior of large language models. In instruction-tuned 528  
491 models, prior work (Kuhn et al., 2023; Nikitin et al., 529  
492 2024) primarily interprets entropy and related un- 530  
493 certainty measures as indicators of response reli- 531  
494 ability and model confidence. Farquhar et al. (Far- 532  
495 quhar et al., 2024) further introduce semantic en- 533  
496 tropy to measure disagreement in the semantic 534  
497 space of generated outputs and apply it to hallu- 535  
498 cination detection. Subsequent work (Han et al., 536

2024; Nguyen et al., 2025) extends this line by 499  
using semantic entropy and related measures to 500  
identify hallucinations in LLM outputs. 501

502 With the emergence of large reasoning mod- 502  
els, recent studies have begun to interpret en- 503  
tropy as a signal of reasoning convergence and 504  
to use it reactively to dynamically adjust reason- 505  
ing depth (Zhang et al., 2025) or as a confidence signal 506  
for early stopping (Sharma and Chopra, 2025). In 507  
addition, several works (Wang et al., 2025; Agar- 508  
wal et al., 2025b) incorporate entropy directly into 509  
training objectives, using it as a reward or regu- 510  
larization signal to prevent entropy collapse and 511  
thereby encourage exploration. 512

### 513 6.2 Prompt Optimization 513

514 Prompt optimization aims to improve task perfor- 514  
mance by rewriting, expanding, or searching over 515  
prompts, while keeping model parameters fixed. 516  
Early work in this direction treats prompts as dis- 517  
crete textual objects and relies on heuristic strate- 518  
gies (Schick and Schütze, 2020; Shin et al., 2020) 519  
or black-box search (Wallace et al., 2019) to ex- 520  
plore prompt variants. 521

522 More recent approaches leverage LLMs them- 522  
selves as prompt optimizers. In this paradigm, an 523  
LLM acts as a meta-optimizer (Yang et al., 2023), 524  
generating candidate prompts (Zhou et al., 2022), 525  
mutating them via evolutionary operators (Tong 526  
et al., 2025) with improved variants selected in a 527  
closed optimization loop. 528

529 Beyond using frozen LLMs as prompt optimiz- 529  
ers, recent work has explored explicitly training 530  
or adapting LLMs to better function as optimizers. 531  
This line of research treats optimization itself as 532  
a learnable behavior, where models are trained to 533  
propose improved solutions based on feedback sig- 534  
nals such as task rewards (Deng et al., 2022) or 535  
preference comparisons (Lin et al., 2024). 536

## 537 7 Conclusion 537

538 In this work, we show that a substantial portion of 538  
exploration in LLM reasoning is unhealthy: it is 539  
induced by prompt ambiguity, and inflates token 540  
usage without improving reasoning quality. We 541  
propose an entropy-dynamics-aware prompt opti- 542  
mization framework that generates concise clari- 543  
fications to reduce early-stage uncertainty while 544  
preserving healthy exploration. Extensive experi- 545  
ments show that our method consistently improves 546  
reasoning efficiency without sacrificing accuracy. 547

## 548 Limitations

549 Despite the effectiveness of our proposed prompt  
550 optimization framework, our work has several limi-  
551 tations that merit discussion:

- 552 • Although we evaluate across multiple mod-  
553 els, reasoning budgets, and in- and out-of-  
554 domain benchmarks, our experiments remain  
555 limited to a finite set of architectures and  
556 datasets. It remains unclear how well the pro-  
557 posed approach generalizes to other genera-  
558 tion paradigms, such as interactive dialogue,  
559 or open-ended creative generation.
- 560 • Although entropy peaks empirically align with  
561 distinct reasoning stages, entropy remains an  
562 indirect proxy for internal uncertainty and ex-  
563 ploration. Future work may incorporate com-  
564plementary signals, such as semantic entropy  
565 or representation-level measures, to better dis-  
566tinguish unhealthy from healthy exploration.
- 567 • Our explicit restriction that clarifications re-  
568 main concise and non-informative with re-  
569 spect to solution steps may limit the opti-  
570 mizer’s ability to handle prompts that are  
571 deeply underspecified or structurally flawed.  
572 In such cases, more substantial prompt refor-  
573 mulation may be required, which falls beyond  
574 the scope of the current framework.

## 575 References

576 Sandhini Agarwal, Lama Ahmad, Jason Ai, Sam Alt-  
577 man, Andy Applebaum, Edwin Arbus, Rahul K  
578 Arora, Yu Bai, Bowen Baker, Haiming Bao, and 1  
579 others. 2025a. gpt-oss-120b & gpt-oss-20b model  
580 card. *arXiv preprint arXiv:2508.10925*.

581 Shivam Agarwal, Zimin Zhang, Lifan Yuan, Jiawei Han,  
582 and Hao Peng. 2025b. The unreasonable effective-  
583 ness of entropy minimization in llm reasoning. *arXiv*  
584 *preprint arXiv:2505.15134*.

585 Maciej Besta, Nils Blach, Ales Kubicek, Robert Gersten-  
586 berger, Michal Podstawski, Lukas Gianinazzi, Joanna  
587 Gajda, Tomasz Lehmann, Hubert Niewiadomski, Pi-  
588 otr Nyczyk, and 1 others. 2024. Graph of thoughts:  
589 Solving elaborate problems with large language mod-  
590 els. In *Proceedings of the AAAI conference on artifi-*  
591 *cial intelligence*, volume 38, pages 17682–17690.

592 Mingkai Deng, Jianyu Wang, Cheng-Ping Hsieh, Yi-  
593 han Wang, Han Guo, Tianmin Shu, Meng Song, Eric  
594 Xing, and Zhiting Hu. 2022. Rlprompt: Optimizing  
595 discrete text prompts with reinforcement learning.

*In Proceedings of the 2022 Conference on Empiri-*  
*cal Methods in Natural Language Processing*, pages  
3369–3391. 596  
597  
598

Xinrun Du, Yifan Yao, Kaijing Ma, Bingli Wang,  
Tianyu Zheng, King Zhu, Minghao Liu, Yiming  
Liang, Xiaolong Jin, Zhenlin Wei, and 1 others. 2025.  
Supergpqa: Scaling llm evaluation across 285 gradu-  
ate disciplines. *arXiv preprint arXiv:2502.14739*. 599  
600  
601  
602  
603

Sebastian Farquhar, Jannik Kossen, Lorenz Kuhn, and  
Yarin Gal. 2024. Detecting hallucinations in large  
language models using semantic entropy. *Nature*,  
630(8017):625–630. 604  
605  
606  
607

Jiatong Han, Jannik Kossen, Muhammed Razzak, Lisa  
Schut, Shreshth A Malik, and Yarin Gal. 2024. Se-  
mantic entropy probes: Robust and cheap hallucina-  
tion detection in llms. In *ICML 2024 Workshop on*  
*Foundation Models in the Wild*. 608  
609  
610  
611  
612

Di Jin, Eileen Pan, Nassim Oufattole, Wei-Hung Weng,  
Hanyi Fang, and Peter Szolovits. 2021. What disease  
does this patient have? a large-scale open domain  
question answering dataset from medical exams. *Ap-*  
*plied Sciences*, 11(14):6421. 613  
614  
615  
616  
617

Lorenz Kuhn, Yarin Gal, and Sebastian Farquhar. 2023.  
Semantic uncertainty: Linguistic invariances for un-  
certainty estimation in natural language generation.  
*arXiv preprint arXiv:2302.09664*. 618  
619  
620  
621

Xiaoqiang Lin, Zhongxiang Dai, Arun Verma, See-  
Kiong Ng, Patrick Jaillet, and Bryan Kian Hsiang  
Low. 2024. Prompt optimization with human feed-  
back. *arXiv preprint arXiv:2405.17346*. 622  
623  
624  
625

Niklas Muennighoff, Zitong Yang, Weijia Shi, Xi-  
ang Lisa Li, Li Fei-Fei, Hannaneh Hajishirzi, Luke  
Zettlemoyer, Percy Liang, Emmanuel Candès, and  
Tatsunori B Hashimoto. 2025. s1: Simple test-time  
scaling. In *Proceedings of the 2025 Conference on*  
*Empirical Methods in Natural Language Processing*,  
pages 20286–20332. 626  
627  
628  
629  
630  
631  
632

Dang Nguyen, Ali Payani, and Baharan Mirzasoleiman.  
2025. Beyond semantic entropy: Boosting llm uncer-  
tainty quantification with pairwise semantic similar-  
ity. *arXiv preprint arXiv:2506.00245*. 633  
634  
635  
636

Alexander Nikitin, Jannik Kossen, Yarin Gal, and Pekka  
Marttinen. 2024. Kernel language entropy: Fine-  
grained uncertainty quantification for llms from se-  
mantic similarities. *Advances in Neural Information*  
*Processing Systems*, 37:8901–8929. 637  
638  
639  
640  
641

Timo Schick and Hinrich Schütze. 2020. Few-shot text  
generation with pattern-exploiting training. *arXiv*  
*preprint arXiv:2012.11926*. 642  
643  
644

Aman Sharma and Paras Chopra. 2025. Think  
just enough: Sequence-level entropy as a confi-  
dence signal for llm reasoning. *arXiv preprint*  
*arXiv:2510.08146*. 645  
646  
647  
648

649	Taylor Shin, Yasaman Razeghi, Robert L Logan IV,	with large language models. <i>Advances in neural</i>	705
650	Eric Wallace, and Sameer Singh. 2020. Autoprompt:	<i>information processing systems</i> , 36:11809–11822.	706
651	Eliciting knowledge from language models with		
652	automatically generated prompts. <i>arXiv preprint</i>	Jinghan Zhang, Xiting Wang, Fengran Mo, Yeyang	707
653	<i>arXiv:2010.15980</i> .	Zhou, Wanfu Gao, and Kunpeng Liu. 2025. Entropy-	708
		based exploration conduction for multi-step reason-	709
654	Charlie Victor Snell, Jaehoon Lee, Kelvin Xu, and Avi-	ing. <i>arXiv preprint arXiv:2503.15848</i> .	710
655	ral Kumar. 2025. Scaling llm test-time compute opti-		
656	mally can be more effective than scaling parameters	Yongchao Zhou, Andrei Ioan Muresanu, Ziwen Han,	711
657	for reasoning. In <i>The Thirteenth International Con-</i>	Keiran Paster, Silviu Pitis, Harris Chan, and Jimmy	712
658	<i>ference on Learning Representations</i> .	Ba. 2022. Large language models are human-level	713
		prompt engineers. In <i>The eleventh international con-</i>	714
659	Mirac Suzgun, Nathan Scales, Nathanael Schärli, Se-	<i>ference on learning representations</i> .	715
660	bastian Gehrmann, Yi Tay, Hyung Won Chung,		
661	Aakanksha Chowdhery, Quoc Le, Ed Chi, Denny	<b>A Baseline</b>	716
662	Zhou, and 1 others. 2023. Challenging big-bench		
663	tasks and whether chain-of-thought can solve them.	<b>Early Stop</b> (Sharma and Chopra, 2025) is an	717
664	In <i>Findings of the Association for Computational</i>	inference-time efficiency method that monitors	718
665	<i>Linguistics: ACL 2023</i> , pages 13003–13051.	sequence-level entropy during generation and ter-	719
		minates reasoning once the entropy remains below	720
666	Zeliang Tong, Zhuojun Ding, and Wei Wei. 2025. Evo-	a predefined threshold for a fixed number of steps.	721
667	prompt: Evolving prompts for enhanced zero-shot	By truncating later-stage generation, it aims to re-	722
668	named entity recognition with large language models.	duce unnecessary token usage while maintaining	723
669	In <i>Proceedings of the 31st International Conference</i>	acceptable accuracy.	724
670	<i>on Computational Linguistics</i> , pages 5136–5153.		
		<b>Evoprompt</b> (Tong et al., 2025) is a prompt opti-	725
671	Eric Wallace, Shi Feng, Nikhil Kandpal, Matt Gardner,	mization framework that applies evolutionary algo-	726
672	and Sameer Singh. 2019. Universal adversarial trig-	rithms to iteratively refine prompts based on task	727
673	gers for attacking and analyzing nlp. <i>arXiv preprint</i>	performance. It focuses on improving answer ac-	728
674	<i>arXiv:1908.07125</i> .	curacy by exploring prompt variants through mu-	729
		tation and selection, and requires a development	730
675	Chen Wang, Zhaochun Li, Jionghao Bai, Yuzhi Zhang,	set to evaluate candidate prompts and guide the	731
676	Shisheng Cui, Zhou Zhao, and Yue Wang. 2025. Ar-	evolutionary search process.	732
677	bitrary entropy policy optimization: Entropy is con-		
678	trollable in reinforcement fine-tuning. <i>arXiv preprint</i>	<b>B Benchmarks</b>	733
679	<i>arXiv:2510.08141</i> .		
		<b>MMLU-Pro</b> (Wang et al., 2024) is a large-	734
680	Xuezhi Wang, Jason Wei, Dale Schuurmans, Quoc Le,	scale multi-task benchmark designed to evaluate	735
681	Ed Chi, Sharan Narang, Aakanksha Chowdhery, and	professional-level language understanding across	736
682	Denny Zhou. 2022. Self-consistency improves chain	diverse academic and technical domains. Com-	737
683	of thought reasoning in language models. <i>arXiv</i>	pared to the original MMLU benchmark, it features	738
684	<i>preprint arXiv:2203.11171</i> .	increased difficulty and reduced annotation noise,	739
		making it well-suited for evaluating advanced rea-	740
685	Yubo Wang, Xueguang Ma, Ge Zhang, Yuansheng Ni,	soning capabilities of large language models.	741
686	Abhranil Chandra, Shiguang Guo, Weiming Ren,		
687	Aaran Arulraj, Xuan He, Ziyang Jiang, and 1 others.	<b>SuperGPQA</b> (Du et al., 2025) is a comprehensive	742
688	2024. Mmlu-pro: A more robust and challenging	benchmark covering 285 graduate-level disciplines,	743
689	multi-task language understanding benchmark. <i>Ad-</i>	with questions curated to require domain-specific	744
690	<i>vances in Neural Information Processing Systems</i> ,	knowledge and multi-step reasoning. Its breadth	745
691	37:95266–95290.	and difficulty make it a challenging testbed for	746
		assessing both reasoning robustness and general-	747
692	An Yang, Anfeng Li, Baosong Yang, Beichen Zhang,	ization across professional domains.	748
693	Binyuan Hui, Bo Zheng, Bowen Yu, Chang		
694	Gao, Chengen Huang, Chenxu Lv, and 1 others.	<b>Big-Bench Hard (BBH)</b> (Suzgun et al., 2023) a	749
695	2025. Qwen3 technical report. <i>arXiv preprint</i>	subset of the BIG-Bench benchmark consisting of	750
696	<i>arXiv:2505.09388</i> .	tasks that are empirically difficult for large lan-	751
		guage models. The benchmark emphasizes compo-	752
697	Chengrun Yang, Xuezhi Wang, Yifeng Lu, Hanxiao	sitional and multi-step reasoning, and is commonly	753
698	Liu, Quoc V Le, Denny Zhou, and Xinyun Chen.		
699	2023. Large language models as optimizers. In		
700	<i>The Twelfth International Conference on Learning</i>		
701	<i>Representations</i> .		
702	Shunyu Yao, Dian Yu, Jeffrey Zhao, Izhak Shafran,		
703	Tom Griffiths, Yuan Cao, and Karthik Narasimhan.		
704	2023. Tree of thoughts: Deliberate problem solving		

Model	Effort	Method	In-domain Benchmarks						Out-of-domain Benchmarks					
			MMLU-Pro			SuperGPQA			BBH			MedQA		
			Acc.↑	Tok.↓	Eff.↑	Acc.↑	Tok.↓	Eff.↑	Acc.↑	Tok.↓	Eff.↑	Acc.↑	Tok.↓	Eff.↑
GPT-oss-120B	High	Pure	<b>80.60</b>	2383.35	3.38	49.40	4500.31	1.09	84.40	1496.53	5.63	<b>91.10</b>	1372.63	6.63
		Ours	80.30	<b>1883.77</b>	<b>4.26</b>	<b>49.80</b>	<b>3469.76</b>	<b>1.43</b>	<b>85.40</b>	<b>1107.80</b>	<b>7.70</b>	90.20	<b>916.28</b>	<b>9.84</b>
	Medium	Pure	79.30	566.69	13.99	47.20	1055.71	4.47	<b>84.40</b>	421.18	20.03	88.00	292.65	30.07
		Ours	<b>80.60</b>	<b>522.75</b>	<b>15.41</b>	<b>48.60</b>	<b>776.53</b>	<b>6.25</b>	83.90	<b>387.09</b>	<b>21.67</b>	<b>88.30</b>	<b>280.23</b>	<b>31.50</b>
	Low	Pure	75.70	216.22	35.01	42.70	294.56	14.49	83.30	206.78	40.28	<b>85.40</b>	<b>100.86</b>	<b>84.67</b>
		Ours	<b>76.20</b>	<b>209.88</b>	<b>36.30</b>	<b>45.40</b>	<b>263.75</b>	<b>17.25</b>	<b>83.50</b>	<b>200.52</b>	<b>41.64</b>	84.90	116.78	72.70
Qwen3-30B	Thinking	Pure	81.40	3316.56	2.45	<b>53.30</b>	4765.22	1.11	<b>85.90</b>	1798.92	4.77	<b>84.90</b>	1982.31	4.28
		Ours	<b>81.70</b>	<b>2553.35</b>	<b>3.19</b>	53.10	<b>3855.48</b>	<b>1.37</b>	85.30	<b>1452.83</b>	<b>5.87</b>	84.80	<b>1328.35</b>	<b>6.38</b>
	Instruct	Pure	76.50	1147.30	6.66	52.50	1986.15	2.64	83.30	564.12	14.76	<b>73.40</b>	<b>55.71</b>	<b>131.75</b>
		Ours	<b>80.10</b>	<b>1025.40</b>	<b>7.81</b>	<b>52.60</b>	<b>1951.21</b>	<b>2.69</b>	<b>83.40</b>	<b>484.39</b>	<b>17.21</b>	73.20	67.80	107.96

Table 4: Total token usage of reasoning models.

used to evaluate the effectiveness of advanced reasoning strategies.

**MedQA** (Jin et al., 2021) is a medical question answering benchmark derived from professional medical licensing examinations. The dataset requires complex clinical reasoning and domain knowledge, making it a standard benchmark for evaluating reasoning performance in the medical domain.

## C Prompt

### Clarification Generation

*You are a helpful assistant that generates a hint to help a reasoning model understand the question and avoid ambiguity.*

*Target Reasoning Model: {{Target Model}}*

*Rules:*

- 1) Do NOT change the question text or options.
- 2) Conclude your response with label 'Clarification:', followed by one single your generated clarification.
- 3) No solution steps, no answer letter.

*Instruction: {{Instruction}}*

### Clarification Refinement (Early Stage Reward Failure)

*Your previous hint did not sufficiently reduce early-stage uncertainty. Below are the tokens generated by the reasoning model \*before\* its main uncertainty peak, indicating where the model started to become confused or ambiguous.*

*Early-stage Tokens: {{Tokens}}*

*Based on these tokens, identify what the reasoning model misunderstood or was uncertain about, and generate a NEW, more targeted hint to clarify the instruction, ensuring:*

- 1) Do NOT change the question text or options.

- 2) Conclude your response with label 'Clarification:', followed by one single your generated clarification.

- 3) No solution steps, no answer letter.

### Clarification Refinement (Late Stage Reward Failure)

*Your previous hint destabilized the reasoning model's late-stage processing. This means the hint may have been too specific, leaked information, or disrupted the model's natural reasoning flow.*

*Please generate a NEW, more targeted hint to clarify the instruction, ensuring:*

- 1) Do NOT change the question text or options.
- 2) Conclude your response with label 'Clarification:', followed by one single your generated clarification.
- 3) No solution steps, no answer letter.

## D Total Token Usage

To demonstrate that our prompt optimizer does not introduce excessive token overhead that would offset the reasoning efficiency gains of the reasoning model, we report both the total token usage and the corresponding reasoning efficiency.

We denote by  $T_o$  the number of tokens used by the prompt optimizer,  $P_o$  the parameter size of the optimizer,  $P_r$  the number of activated parameters of the reasoning model during inference, and  $T_r$  the number of tokens used by the reasoning model. The total token usage is computed as

$$T_{\text{total}} = T_o \cdot \frac{P_o}{P_r} + T_r. \quad (13)$$

Specifically,  $T_o = 4$  for our prompt optimizer trained on Qwen3-4B-Instruct. For GPT-OSS-120B, the number of activated parameters during

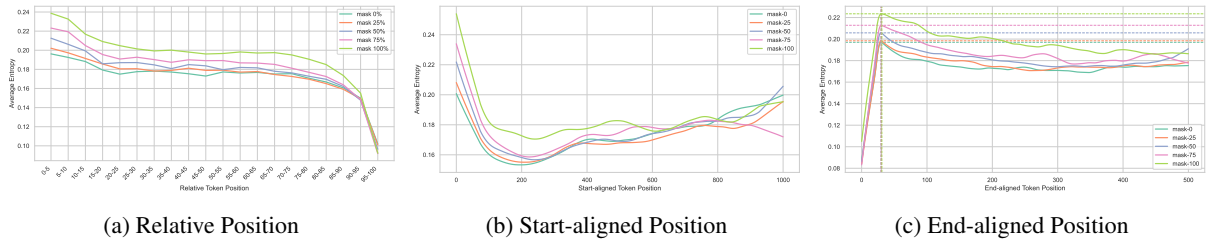


Figure 7: Token-level entropy of GPT-oss-120B under low reasoning effort.

inference is  $P_r = 5.1$ , as reported in our analysis, while for the Qwen3-30B series,  $P_r = 3$ . As shown in Table 4, even after accounting for the token usage introduced by the prompt optimizer, our proposed prompt optimization framework still achieves a significant improvement in reasoning efficiency.

## E Human Evaluation

Although we have shown that our prompt optimization framework preserves the reasoning capability of the underlying model, we further conduct a manual evaluation to assess whether the generated clarifications exhibit answer leakage.

Specifically, we randomly sample 100 generated clarifications from the evaluation set and invite three independent human annotators with backgrounds in computer science and machine learning to assess them. Each clarification is categorized into one of the following three levels:

- **No leakage:** The clarification purely disambiguates the task or constraints without revealing the answer, solution steps, or providing strong hints toward a specific option.
- **Partial leakage:** The clarification contains mild cues or narrowing signals that may reduce the reasoning space, but does not directly reveal the correct answer or solution.
- **Full leakage:** The clarification explicitly reveals the correct answer or provides decisive information that eliminates the need for reasoning.

All clarifications are annotated independently by the three annotators following detailed guidelines. The final label for each sample is determined by majority vote.

In addition to human evaluation, we perform an automatic assessment using an LLM-as-a-judge.

	No Leakage	Partial Leakage	Full Leakage
<b>Human</b>	83	17	0
<b>GPT-5</b>	89	11	0

Table 5: Human evaluation results on answer leakage in generated clarifications.

Specifically, we employ GPT-5 as an impartial evaluator and prompt it to classify each clarification into the same three leakage categories using an identical taxonomy and conservative decision rules. To reduce stochastic effects, the model is queried with temperature set to zero, and each clarification is evaluated independently.

As shown in Table 5, the majority of generated clarifications exhibit no answer leakage. A small portion falls into the partial leakage category, while no cases directly reveal the correct answer. The LLM-as-a-judge results are broadly consistent with human annotations, further supporting the reliability of our evaluation.

According to our analysis, direct answer leakage substantially alters the model’s reasoning dynamics, causing the late-stage entropy peak to disappear, which indicates a collapse of the verification process. In contrast, partial leakage, in the form of leaking reasoning cues, may conflict with the model’s original reasoning trajectory and instead lead to an elevated and delayed early-stage entropy peak. Both behaviors deviate from the desired entropy dynamics and are therefore penalized by our entropy-peak-based reward.

## F Additional Entropy Dynamics

We also present the entropy dynamics of GPT-oss-120B under low reasoning effort, together with results from the Qwen3-30B series.

As shown in Figure 7, unlike the medium- and high-reasoning-effort settings, GPT-oss-120B under low reasoning effort does not exhibit a pronounced early-stage entropy increase, even when prompt ambiguity is introduced. Nevertheless, ele-

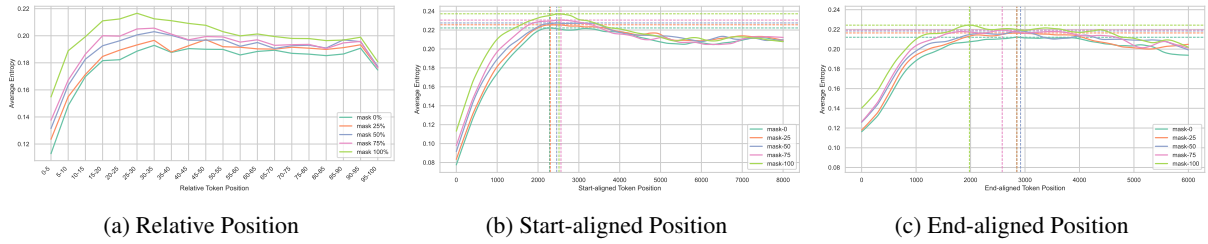


Figure 8: Token-level entropy of Qwen3-30B-Instruct.

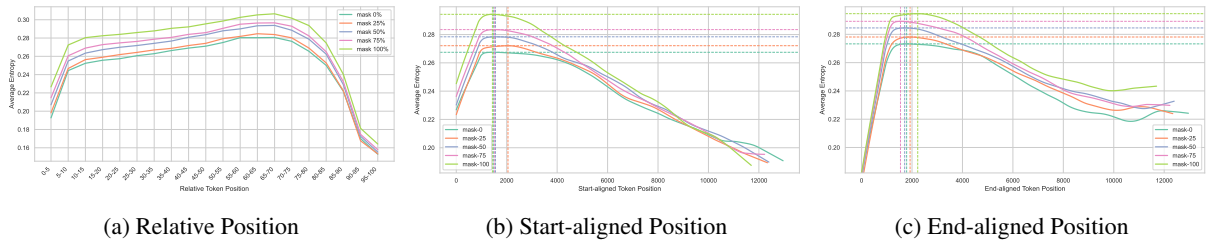


Figure 9: Token-level entropy of Qwen3-30B-Thinking.

892 vated entropy remains observable in the late stage  
 893 of generation, indicating that uncertainty is not  
 894 eliminated but manifests at a later phase. These ob-  
 895 servations reveal an asymmetric capacity allocation  
 896 under constrained reasoning budgets: early-stage  
 897 exploratory behaviors are sacrificed first, while late-  
 898 stage verification and answer finalization are com-  
 899 paratively preserved. With more ample reasoning  
 900 budgets, models are able to invest additional tokens  
 901 in early-stage exploration, which explains why to-  
 902 ken usage increases more rapidly with higher mask  
 903 ratios under higher reasoning effort.

904 As shown in Figure 8 and Figure 9, the Qwen3-  
 905 30B series exhibits trends similar to those observed  
 906 for GPT-oss-120B. Specifically, the entropy trajec-  
 907 tory follows a consistent pattern: it initially in-  
 908 creases, then decreases, rises again, and finally  
 909 drops sharply toward termination. Moreover, as  
 910 prompt ambiguity increases, the early-stage en-  
 911 tropy peak becomes both elevated and delayed,  
 912 while the late-stage entropy peak remains highly  
 913 convergent. These results further demonstrate the  
 914 generality of the observed dual-peak entropy dy-  
 915 namics across different model families, supporting  
 916 the robustness of our analysis.