

CPC-GRPO: Answer-Free Reinforcement Learning with Cross-Prompt Consensus Rewards

Anonymous ACL submission

Abstract

Reinforcement learning with verifiable rewards has improved reasoning in language models, but it typically relies on a ground-truth answer or an external verifier, which limits applicability and increases cost. We propose an answer-free training objective that derives rewards solely from the model’s own probabilities by exploiting prompt paraphrases as multiple semantic views of the same intent. For each paraphrase set, we generate candidate responses, rescore each response under the other paraphrased prompts via teacher forcing, and define a cross-prompt consensus reward that favors responses supported across views rather than those that fit a single phrasing. We optimize this reward using a policy update with an all-pairs objective and advantage broadcasting across prompt–response pairs. The framework naturally supports prefix-level training, enabling a controllable cost–signal trade-off. Experiments on RobustAlpacaEval and out-of-domain reasoning benchmarks (OpenBookQA, AQuA, HumanEval) show consistent gains over the pre-trained baselines on LLaMA3.2-3B and Qwen3-4B, alongside analyses demonstrating reward–performance alignment and the importance of design choices such as excluding self-view scores and ensembling-based candidates. All experiment code is available at our GitHub¹.

1 Introduction

Reinforcement learning (RL) is now a standard post-training tool for large language models (LLMs), and is used to bias generation toward more desirable outputs (Ouyang et al., 2022). Recent work further reports that RL can yield substantial gains on reasoning-centric tasks (Yu et al., 2025; Mroueh, 2025). In particular, RLVR (Reinforcement Learning with Verifiable Rewards)

shapes reasoning behavior using scalar rewards that are directly verifiable against a ground-truth answer (Lambert et al., 2025; Guo et al., 2025a). When combined with recent policy optimization methods such as GRPO (Zhihong Shao, 2024), DR.GRPO (Liu et al., 2025), and DAPO (Yu et al., 2025), RLVR has shown consistent improvements on mathematics and coding benchmarks.

Despite its clarity, RLVR has practical constraints. First, verification requires either ground-truth answers (e.g., exact match (EM)) or an external verifier (e.g., executor or grader), which makes it difficult to apply the same recipe to general instruction data where answer verification is not straightforward. Second, RLVR rewards are typically defined over a “reasoning → final answer” structure, so truncated responses cannot be evaluated. This pushes training toward long rollouts, while rollout generation and scoring can dominate RL cost in practice (Zhong et al., 2025; Hu et al., 2025; Yao et al., 2023).

We propose answer-free self-reward RL by leveraging paraphrase prompts of the same intent as multi-view observations. Given G paraphrase prompts, we generate responses and *rescore* each response under the other paraphrase conditions via teacher forcing to obtain cross-prompt scores. We then define a *consensus reward* as the multi-view average support, rather than the self-view score under a single prompt. To reduce prompt-specific overfitting, our default setting excludes the self-view term and applies group-wise normalization. Finally, we transform consensus rewards into group-relative advantages and perform an on-policy GRPO-style update, enabling answer-free policy optimization without ground-truth answers or external judges.

Because the reward is computed solely from model token probabilities, answer-free training requires neither labels nor verifiers. Moreover, the token-average form allows computing rewards on

¹https://anonymous.4open.science/r/Cross-Prompt_Consensus_GRPO-2660

082 response prefixes, enabling explicit performance–
083 cost trade-offs as a function of prefix length. We
084 evaluate answer-free training on general instruc-
085 tion data and out-of-domain reasoning benchmarks,
086 and analyze the performance–cost trade-off under
087 different prefix lengths.

088 2 Related Work

089 2.1 Reinforcement Learning in LLM

090 RL-based post-training is widely used to adapt
091 large language models by optimizing reward sig-
092 nals that reflect human preferences or task perfor-
093 mance (Ouyang et al., 2022). RLHF is a repre-
094 sentative approach, where a preference-based re-
095 ward model is trained and PPO-style updates are
096 applied (Schulman et al., 2017; Ziegler et al., 2019).
097 More recently, group-based optimization methods
098 such as GRPO update the policy using within-
099 group relative signals from multiple responses to
100 the same prompt, and have shown strong results
101 on reasoning-oriented benchmarks (Zhihong Shao,
102 2024). In the RLVR setting, GRPO can be com-
103 bined with verifiable rewards (e.g., exact match
104 (EM) or execution-based pass@k), but this typi-
105 cally requires a ground-truth answer or an external
106 verifier and often relies on sufficiently long rollouts
107 to compute rewards reliably (Guo et al., 2025a;
108 Yao et al., 2023). To reduce dependence on labeled
109 answers, prior work has explored reward construc-
110 tion without direct supervision. LLM-as-a-Judge
111 uses a strong external LLM to score or compare
112 outputs (Zheng et al., 2023), at the cost of addi-
113 tional inference and with open concerns about bias
114 and reproducibility. RLIF methods instead use in-
115 ternal model signals as intrinsic rewards. For ex-
116 ample, RENT uses entropy- (or confidence-) based
117 signals from the token distribution (Prabhudesai
118 et al., 2025), and Intuitor replaces external rewards
119 with self-certainty in a GRPO framework (Zhao
120 et al., 2025). However, self-certainty proxies can
121 be prompt- and task-dependent, and their alignment
122 with response quality is not guaranteed.

123 In contrast, we construct a reward from para-
124 phrase multi-view cross-prompt scoring, yielding a
125 consensus signal that does not require ground-truth
126 answers or external judges and does not rely solely
127 on self-certainty.

128 2.2 Prompt Paraphrase Ensemble

129 Ensembling is a classical strategy for reducing vari-
130 ance and improving stability. For LLMs, prompt

131 sensitivity implies that output quality can vary
132 substantially across semantically equivalent phras-
133 ings (Zhao et al., 2021). This has motivated multi-
134 prompt test-time methods that either aggregate to-
135 ken distributions during decoding or select a final
136 output from candidates generated by a prompt bank.
137 For example, M-Ped averages token distributions
138 across prompts via inner-batch ensembling (Guo
139 et al., 2025b), and multi-prompt MBR samples can-
140 didates from diverse prompts and selects outputs
141 using an MBR criterion (Heineman et al., 2024).
142 Bayesian Prompt Ensembles treat prompt varia-
143 tions as a source of uncertainty and aggregate out-
144 put probabilities to estimate calibrated uncertainty
145 for black-box LLMs (Tonolini et al., 2024).

146 Unlike prior work that primarily targets test-
147 time quality or uncertainty estimation, we use para-
148 phrase ensembles for *training-time* reward con-
149 struction. We define a multi-view consensus re-
150 ward via cross-prompt teacher-forcing rescaling,
151 enabling answer-free policy optimization without
152 ground-truth answers or external judges.

153 3 Methodology

154 In this section, we describe how to construct
155 an *answer-free* reward using a set of paraphrase
156 prompts and how to train with group-based pol-
157 icy optimization. The goal of our method is to (i)
158 define a reward signal that is computable without
159 ground-truth answer labels or external evaluators,
160 (ii) avoid over-reliance on self-certainty under a sin-
161 gle prompt, and (iii) concentrate probability mass
162 on responses that are *consistently supported* across
163 semantically equivalent multi-views.

164 To this end, we treat each paraphrase prompt as
165 a different observation (view) of the same intent,
166 and measure whether a particular response is not
167 only plausible under one prompt phrasing but also
168 explainable with high probability under other para-
169 phrase conditions via *cross-prompt scoring*. Each
170 response is scored by its average support (consen-
171 sus) across multi-views, and this score is converted
172 into a group-relative advantage. We further broad-
173 cast this advantage to all (prompt, response) pairs
174 used in the update, so that each conditional distribu-
175 tion $\pi_{\theta}(\cdot | x^{(i)})$ is encouraged to place more proba-
176 bility mass on regions that are commonly supported
177 across views (high-consensus regions) rather than
178 on prompt-specific idiosyncrasies. Since training is
179 on-policy, as updates proceed, sampled responses
180 increasingly reflect the preferences of the current

policy, and we may expect dynamics in which cross-prompt consensus is gradually strengthened. Below, we formalize the proposed method step by step.

3.1 Rollout with Paraphrase Ensemble

Each training sample consists of a set of G semantically equivalent paraphrase prompts for a single intent, $\mathcal{X} = \{x^{(i)}\}_{i=1}^G$. In on-policy rollouts, we generate a response for each prompt from the previous policy $\pi_{\theta_{\text{old}}}$:

$$y^{(j)} \sim \pi_{\theta_{\text{old}}}(\cdot | x^{(j)}), \quad j = 1, \dots, G. \quad (1)$$

In addition, to provide a response that is *more stably agreed upon* across multi-views, we generate a decoding-level ensemble (DLE) response $y^{(K)}$ ($K = G + 1$). At token step t , we define the ensemble next-token distribution as the average of the prompt-conditional distributions:

$$\bar{\pi}_{\theta_{\text{old}}}(\cdot | y_{<t}^{(K)}) = \frac{1}{G} \sum_{i=1}^G \pi_{\theta_{\text{old}}}(\cdot | x^{(i)}, y_{<t}^{(K)}). \quad (2)$$

Actual generation is performed by selecting the most probable token under the averaged distribution (e.g., greedy decoding):

$$y_t^{(K)} = \arg \max_v \bar{\pi}_{\theta_{\text{old}}}(v | y_{<t}^{(K)}). \quad (3)$$

Finally, the set of responses used for training is $\mathcal{Y} = \{y^{(j)}\}_{j=1}^K$ with $K = G + 1$.

3.2 Cross-prompt Consensus Reward

For each response $y^{(j)} = (y_1^{(j)}, \dots, y_{T_j}^{(j)})$, we measure support under other paraphrase conditions via cross-prompt scores. Concretely, we rescore $y^{(j)}$ under prompt $x^{(i)}$ via teacher forcing and define the token-mean probability score for pair (i, j) as:

$$s_{i,j} = \frac{1}{T_j} \sum_{t=1}^{T_j} \pi_{\theta_{\text{old}}}(y_t^{(j)} | x^{(i)}, y_{<t}^{(j)}). \quad (4)$$

This yields $G \times K$ cross-scores, which are computed solely from the model’s token probabilities without ground-truth labels or external evaluators. We discuss score variants that mitigate the influence of early tokens in Sec. 5.6.

We then define the consensus reward of response $y^{(j)}$ as its average support (consensus) across multi-views. In the default setting, for responses with

$j \leq G$, we exclude the self-view term ($i = j$) and take the average:

$$r_j = \frac{1}{G-1} \sum_{i=1, i \neq j}^G s_{i,j}, \quad j = 1, \dots, G. \quad (5)$$

This design is motivated by the fact that the self-view score is not, in a strict sense, a consensus signal; rather, it can act as a *prompt-specific fit* or self-certainty for a particular phrasing, and excluding it focuses the reward on cross-view consensus. Since the ensemble response $y^{(K)}$ has no specific self-view, we define it as the average over all views:

$$r_K = \frac{1}{G} \sum_{i=1}^G s_{i,K}, \quad K = G + 1. \quad (6)$$

Thus, a high r_j reflects *multi-view consensus* rather than self-certainty under a single prompt, and training encourages allocating more probability mass to responses with higher consensus.

3.3 All-pair Policy Optimization

We follow GRPO’s group-based policy optimization framework, but perform an all-pair GRPO-style optimization that utilizes *all* (prompt, response) combinations involved in computing the consensus reward for learning. The key is to (i) construct a group-relative signal from response-level consensus rewards and (ii) share this signal across (prompt, response) pairs that participate in cross-prompt rescoring, thereby strengthening high-consensus responses simultaneously under every view condition.

Group-relative advantage. Given K rewards $\{r_j\}_{j=1}^K$ for a sample, we define the group mean reward and advantages as:

$$\bar{r} = \frac{1}{K} \sum_{j=1}^K r_j, \quad A_j = r_j - \bar{r}. \quad (7)$$

For training stability, we normalize the advantages within each sample by their standard deviation:

$$\tilde{A}_j = \frac{A_j}{\sigma(A) + \varepsilon}, \quad (8)$$

where $\sigma(A)$ is the standard deviation of $\{A_j\}_{j=1}^K$ and ε is a small constant for numerical stability.

Algorithm 1 Paraphrase-Consensus GRPO

Require: Paraphrase prompts $\mathcal{X} = \{x^{(i)}\}_{i=1}^G$, policy π_θ , prefix cap L , clip ϵ (optional DLE)

- 1: **for** each policy update **do**
- 2: $\theta_{\text{old}} \leftarrow \theta$
- 3: **for** each sample \mathcal{X} in a minibatch **do**
- 4: **Rollout:** sample $y^{(j)} \sim \pi_{\theta_{\text{old}}}(\cdot | x^{(j)})$ for $j = 1..G$ (Eq. (1))
- 5: **if** DLE enabled **then**
- 6: generate $y^{(K)}$ by decoding-level ensemble (Eq. (2)–(3)); set $K \leftarrow G + 1$
- 7: **else**
- 8: set $K \leftarrow G$
- 9: **end if**
- 10: set $\tilde{T}_j \leftarrow \min(T_j, L)$ for all j (prefix-level training)
- 11: **Cross-score:** compute $s_{i,j}$ for all $i \leq G, j \leq K$ by teacher forcing up to \tilde{T}_j (Eq. (4))
- 12: **Reward:** compute consensus rewards $\{r_j\}_{j=1}^K$ (Eq. (5)–(6))
- 13: **Advantage:** compute normalized $\{\tilde{A}_j\}$ (Eq. (7)–(8))
- 14: **Broadcast:** set $\tilde{A}_{i,j} \leftarrow \tilde{A}_j$ for valid pairs \mathcal{P} (Eq. (11), (9))
- 15: accumulate clipped loss over $(i, j) \in \mathcal{P}, t \leq \tilde{T}_j$ (Eq. (12))
- 16: **end for**
- 17: update θ by minimizing the accumulated loss
- 18: **end for**

Advantage broadcast. Although rewards/advantages are defined at the response level, policy optimization is carried out over the (prompt, response) pairs used to compute cross-prompt scores. We therefore broadcast the normalized advantage for response $y^{(j)}$ to every valid (prompt, response) pair:

$$\tilde{A}_{i,j} := \tilde{A}_j. \quad (9)$$

This broadcast strengthens high-consensus responses under each view condition simultaneously, and consequently encourages $\pi_\theta(\cdot | x^{(i)})$ to place more probability mass on regions that are commonly supported across views (high-consensus regions) rather than on prompt-specific idiosyncrasies.

PPO clipped objective. For each (prompt, response) pair $(x^{(i)}, y^{(j)})$, we define the token-level probability ratio at step t as:

$$\rho_{i,j,t}(\theta) = \frac{\pi_\theta(y_t^{(j)} | x^{(i)}, y_{<t}^{(j)})}{\pi_{\theta_{\text{old}}}(y_t^{(j)} | x^{(i)}, y_{<t}^{(j)})}. \quad (10)$$

To reflect the exclusion of self-views, we define the set of valid (prompt, response) pairs as:

$$\mathcal{P} = \{(i, j) \mid 1 \leq i \leq G, 1 \leq j \leq K, i \neq j\}. \quad (11)$$

In the default setting ($K = G + 1$), G self-view pairs are excluded for $j \leq G$, whereas the ensemble response ($j = K$) has no excluded self-view, yielding

$$|\mathcal{P}| = G(G - 1) + G = G^2.$$

Finally, the policy update can be written as minimizing the objective obtained by averaging PPO’s clipped surrogate over all pairs:

$$\mathcal{L}(\theta) = -\frac{1}{|\mathcal{P}|} \sum_{(i,j) \in \mathcal{P}} \frac{1}{\tilde{T}_j} \sum_{t=1}^{\tilde{T}_j} \min \left(\rho_{i,j,t}(\theta) \tilde{A}_{i,j}, \text{clip}(\rho_{i,j,t}(\theta), 1 - \epsilon, 1 + \epsilon) \tilde{A}_{i,j} \right), \quad (12)$$

where ϵ is the PPO clipping coefficient.

3.4 Prefix-level Reward and Computational Trade-off

Because our reward is defined as a cumulative (average) function of token probabilities, it can be computed from response prefixes rather than requiring complete responses. Unlike RLVR, which must extract and verify the final answer, this means that learning signals can be formed without generating responses to completion.

Let L denote the maximum prefix length, and define the used length of each response as $\tilde{T}_j = \min(T_j, L)$. Accordingly, rollouts are restricted to generating only up to $y_{\leq \tilde{T}_j}^{(j)}$, and the inner sums in cross-prompt scoring (4) and the PPO loss (12) can likewise be performed only for $t = 1, \dots, \tilde{T}_j$. Since the self-attention computation of Transformer-based LLMs increases approximately as $O(L^2)$ with sequence length (Vaswani et al., 2023; Dao, 2023), controlling L allows us to directly manage the performance–cost trade-off. Algorithm 1 summarizes the overall training procedure of the proposed method.

Base Model	Method	Para.	Train	RAE (In-domain)		Reasoning (Out-of-domain)		
				Trained	Un-trained	OBQA	AQuA	HumanEval
LLaMA3.2-3B	pre-trained	X	–	<u>21.42</u>	<u>19.25</u>	58.66	52.36	60.06
	self-refine	X	–	9.49	8.73	34.80	43.37	49.39
	swarm-distillation (st)	✓	FT	18.08	17.23	53.11	47.44	60.06
	swarm-distillation (ot)	✓	FT	13.16	12.77	47.37	47.70	51.83
	INTUITOR	X	RL	19.68	18.21	<u>59.23</u>	50.13	59.45
	CPC-GRPO	✓	RL	22.89	21.04	59.94	53.35	61.59
Qwen3-4B	pre-trained	X	–	61.64	63.65	82.03	85.50	91.77
	self-refine	X	–	10.14	8.16	80.80	82.09	89.94
	swarm-distillation (st)	✓	FT	15.09	14.18	75.11	75.13	83.54
	swarm-distillation (ot)	✓	FT	30.56	33.31	77.14	80.45	90.24
	INTUITOR	X	RL	58.79	62.75	<u>81.57</u>	<u>85.63</u>	92.38
	CPC-GRPO	✓	RL	64.78	66.19	82.03	85.70	<u>92.07</u>

Table 1: Results on RobustAlpacaEval (RAE; in-domain) and out-of-domain reasoning benchmarks (OBQA/AQuA/HumanEval). RAE uses AlpacaEval weighted win rate; “Trained/Un-trained” follows the RAE split.

Base Model	Method	RAE Tra.		RAE Un-Tra.		OBQA		AQuA		HumanEval	
		Best	Worst	Best	Worst	Best	Worst	Best	Worst	Best	Worst
LLaMA3.2-3B	pre-trained	<u>40.06</u>	5.27	55.16	2.67	<u>89.80</u>	22.80	79.53	29.92	82.93	37.20
	self-refine	19.56	1.14	29.59	0.08	76.40	3.80	<u>80.31</u>	10.63	68.90	29.88
	swarm-distillation (st)	34.26	4.46	47.24	1.12	86.80	19.20	78.35	24.02	85.98	34.15
	swarm-distillation (ot)	25.71	3.66	35.36	0.01	78.40	17.60	79.13	22.05	77.44	26.22
	INTUITOR	37.34	<u>5.47</u>	50.62	0.99	88.60	25.20	77.17	<u>24.80</u>	83.54	<u>35.37</u>
	CPC-GRPO	42.73	8.16	<u>52.12</u>	<u>2.00</u>	90.20	<u>23.60</u>	84.65	<u>24.80</u>	89.02	34.15
Qwen3-4B	pre-trained	<u>87.37</u>	<u>33.56</u>	<u>94.06</u>	18.14	<u>94.80</u>	67.40	92.52	74.02	98.17	85.37
	self-refine	21.1	1.99	25.66	0.35	94.60	59.00	92.52	62.60	<u>97.56</u>	82.32
	swarm-distillation (st)	32.96	1.59	44.20	0.27	91.80	51.00	90.16	59.45	95.73	71.34
	swarm-distillation (ot)	54.61	7.28	71.93	2.53	92.80	58.60	89.37	68.50	96.95	83.54
	INTUITOR	83.75	30.15	92.26	17.66	94.60	66.40	<u>93.31</u>	75.98	96.95	87.80
	CPC-GRPO	89.10	39.67	94.78	22.49	95.20	<u>66.80</u>	94.09	<u>75.20</u>	98.17	<u>85.98</u>

Table 2: Prompt sensitivity across paraphrase variants, reported as best/worst over prompt variants per example.

4 Experiment

4.1 Datasets

We study an *answer-free* setting where training does not require ground-truth answers or external evaluators. For in-domain training, we use RobustAlpacaEval (RAE) (Cao et al., 2024), where each example provides one original prompt and ten semantically equivalent paraphrases (11 variants). We fix the paraphrase group size to $G=4$ by selecting the first four variants in the predefined index order for training. At evaluation time, we report AlpacaEval weighted win rate separately on (i) the *Trained* prompt subset and (ii) the remaining *Un-trained* variants, measuring generalization to unseen phrasings.

For out-of-domain evaluation, we use OpenBookQA, AQuA, and HumanEval (Mihaylov et al., 2018; Ling et al., 2017; Chen et al., 2021). We construct instruction-style prompt templates using

PromptSource and treat templates as prompt variants (views). Dataset statistics and construction details are provided in Appendix A.

4.2 Baselines

Since RAE does not provide ground-truth answers, we compare against methods that do not rely on answer verification (EM, pass@k) or external reward models. We evaluate: (1) **Pre-trained**, the base model without training; (2) **Self-refinement (test-time)** (Cao et al., 2024), which rewrites the prompt and answers the rewritten prompt; (3) **Swarm distillation** (Zhou et al., 2022), an unsupervised distillation baseline using paraphrase prompts as teachers (static/online); and (4) **INTUITOR** (Zhao et al., 2025), an RLIF-style GRPO baseline using self-certainty as intrinsic reward. All baselines are trained under the same setting as our method. Detailed protocols are in Appendix C.

4.3 Experimental Setting

We use LLaMA3.2-3B and Qwen3-4B as base models and conduct all FT/RL training with QLoRA (Detmers et al., 2023). All methods are trained for three epochs with a fixed paraphrase group size $G=4$. All experiments are evaluated under the same setting; due to limited resources, we report results for a single configuration, and provide multi-seed results in Appendix E.

For out-of-domain benchmarks, we provide all methods with the same minimal output-format rules and compute scores via rule-based parsing (Appendix D). For RAE, we use the official AlpacaEval weighted win rate with gpt4_turbo as the reference output and judge (Cao et al., 2024; Dubois et al., 2025). For OpenBookQA and AQuA, we report accuracy from parsed answers, and for HumanEval we report unit-test-based pass@1. Hyperparameters and additional details are in Appendix B.

4.4 Main Result

Table 1 summarizes results on RobustAlpacaEval (RAE; in-domain) and three out-of-domain reasoning benchmarks (OBQA/AQuA/HumanEval). Across both base models, CPC-GRPO (Cross-Prompt Consensus GRPO) achieves the strongest RAE weighted win rates on both the *Trained* and *Un-trained* prompt subsets, indicating that paraphrase-based consensus rewards can provide a usable post-training signal without ground-truth answers or external judges.

On LLaMA3.2-3B, CPC-GRPO improves RAE from 21.42/19.25 to 22.89/21.04 (*Trained/Un-trained*), and yields consistent gains on OBQA/AQuA/HumanEval (+1.28/+0.99/+1.53, respectively). On Qwen3-4B, CPC-GRPO improves RAE from 61.64/63.65 to 64.78/66.19, while changes on out-of-domain benchmarks are small (OBQA tied; AQuA and HumanEval change marginally). Relative to INTUITOR, CPC-GRPO is consistently stronger on RAE and remains competitive on the reasoning benchmarks.

In our configuration, self-refinement and swarm distillation do not improve overall performance and often reduce scores (Table 1), suggesting that prompt rewriting or unsupervised distillation alone is insufficient to match RL updates driven by a consensus reward.

Table 2 further reports oracle best and worst-case performance over prompt variants per exam-

Method	AVG	BEST	WORST
pre-trained	85.50	92.52	74.02
+ DLE	86.22	86.22	86.22
+ PCV	85.43	85.43	85.43
CPC-GRPO	85.70	94.09	75.20
+ DLE	88.19	88.19	88.19
+ PCV	87.01	87.01	87.01

Table 3: Test-time paraphrase ensembling on AQuA for Qwen3-4B.

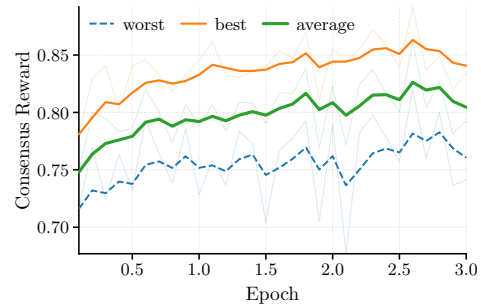


Figure 1: Training curves of the best/worst/mean consensus reward.

ple. CPC-GRPO generally increases the *Best* scores and improves *Worst* in many settings (e.g., Qwen3-4B RAE Trained Worst: 33.56→39.67), but the *Best–Worst* gap does not consistently shrink across tasks. This indicates that training improves robustness in aggregate while residual prompt sensitivity remains.

5 Analysis

5.1 Ensemble Performance

DLE averages the next-token distributions conditioned on the G paraphrase prompts and generates a single response by selecting tokens greedily from the averaged distribution (Eq. 2–3). **Paraphrase Consensus Voting (PCV)** computes the cross-prompt consensus reward (§3.2) for candidate responses generated from the G prompts, and then selects the response with the maximum reward. Since both methods combine multi-views to produce a single output, they remove performance dispersion (sensitivity) across prompt variants, yielding $AVG=BEST=WORST$.

Table 3 reports test-time ensemble performance on AQuA. For the pre-trained model, DLE provides a small gain (85.50→86.22), whereas PCV is comparable or slightly worse (85.50→85.43). In contrast, for the model trained with our method, both techniques improve performance, and DLE in particular yields a larger gain (85.70→88.19). This

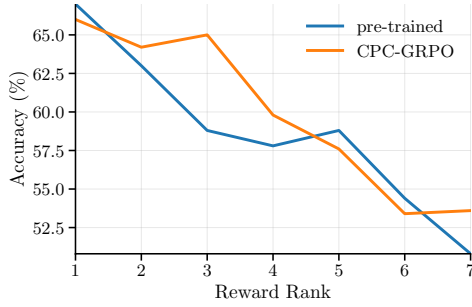


Figure 2: OpenBookQA accuracy of the consensus-reward rank for LLaMA.

suggests that as paraphrase consensus is strengthened during training, test-time methods that combine paraphrase information can operate more effectively.

5.2 Learning Curve

To analyze training dynamics, we track the *best/mean/worst* of the consensus reward (§3.2) for the response set \mathcal{Y} generated at each policy update. Concretely, given $\{r_j\}_{j=1}^K$ for a sample, we define best as $\max_j r_j$, worst as $\min_j r_j$, and mean as $\frac{1}{K} \sum_j r_j$, and record these statistics over the course of training.

In Figure 1, all three metrics tend to increase as training proceeds, supporting that the proposed optimization indeed operates in a direction that strengthens cross-prompt consensus. However, the gap between best and worst does not substantially shrink throughout training, indicating that a non-trivial level of prompt sensitivity remains even after training.

5.3 Alignment Between Consensus Reward and Performance

In this section, we analyze whether the consensus reward can serve as a *proxy for response quality* at test time. For each example, we generate candidate responses for each prompt variant and compute a cross-prompt-based consensus reward for each candidate, then sort candidates by their reward rank. We then measure task performance when selecting the response at rank r , to evaluate how well reward-based ranking aligns with actual performance.

Figure 2 shows accuracy on OpenBookQA by reward rank (with 1 denoting the highest reward). For the pre-trained model, accuracy is higher for top-ranked candidates and tends to decrease as rank worsens, suggesting that the consensus reward can act as a signal that distinguishes the relative quality of candidates even without ground-truth labels.

L	AQuA	HumanEval	Time (\times)
pre-trained	52.36	60.06	–
32	52.69	60.06	$\times 0.08$
64	52.69	61.59	$\times 0.14$
128	53.48	61.28	$\times 0.26$
256	53.22	63.11	$\times 0.45$
512	52.89	61.59	$\times 0.73$
1024	53.15	62.20	$\times 0.86$
2048	53.35	61.59	$\times 1.00$

Table 4: Effect of prefix length L and training time (normalized to CPC-GRPO as $\times 1.00$).

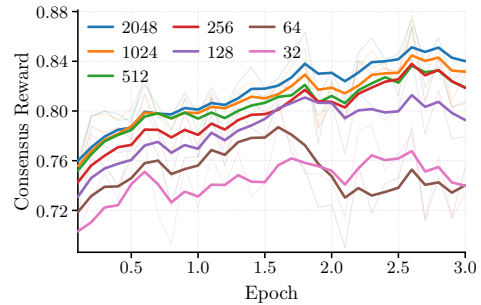


Figure 3: Training curves of the mean consensus reward (by prefix length L).

After training CPC-GRPO, this alignment trend largely persists, and we additionally observe that performance among the top ranks (rank 1–3) becomes more similar. This indicates that as training progresses, the quality of candidates receiving high rewards can improve together, which may reduce performance gaps among top-ranked candidates.

5.4 Prefix-level Training

Table 4 reports the performance–cost trade-off as a function of the response prefix length L used during training. With very short prefixes ($L=32, 64$), training cost is substantially reduced (to $\times 0.08$ and $\times 0.14$, respectively), but performance gains are limited or may be more variable. In contrast, from moderate lengths ($L \geq 128$), AQuA/HumanEval performance improves more stably while training cost remains significantly reduced; in particular, at $L=256$ we reduce wall-clock time to $\times 0.45$ while improving HumanEval to 63.11. Moreover, the longest prefix ($L=2048$) is not always optimal, showing that an appropriate L can be selected depending on the target task and resource constraints.

The prefix-conditioned learning curves in Figure 3 further confirm that the mean reward increases relatively stably for $L \geq 256$, whereas shorter prefixes yield noisier trajectories.

Base Model	Method	RAE (In-domain)		Reasoning (Out-of-domain)		
		Trained	Un-trained	OBQA	AQuA	HumanEval
LLaMA3.2-3B	CPC-GRPO	22.89	21.04	59.94	53.35	61.59
	CPC-GRPO (log-likelihood)	25.48	20.66	59.00	54.46	64.33
Qwen3-4B	CPC-GRPO	64.78	66.19	82.03	85.70	92.07
	CPC-GRPO (log-likelihood)	73.95	72.70	81.69	85.37	90.24

Table 5: Ablation of default CPC-GRPO vs. likelihood-based variant. **bold** denotes the best.

Method	AQuA	OBQA	Time (\times)
Intuitior	50.13	59.23	0.38
CPC-GRPO	53.35	59.94	1.00
w/o DLE	52.56	59.23	0.59
w/ Self-view	51.84	59.29	0.61

Table 6: Ablation Study on AQuA and OBQA with training time (normalized to CPC-GRPO as $\times 1.00$).

5.5 Ablation Study

Table 6 reports ablations of two key components in CPC-GRPO on LLaMA. Removing the decoding-level ensemble response (*w/o DLE*) reduces performance on both benchmarks (AQuA: 53.35 \rightarrow 52.56; OBQA: 59.94 \rightarrow 59.23), suggesting that the additional DLE candidate strengthens the consensus-oriented learning signal. In contrast, including self-view scores (*w/ Self-view*) decreases performance (AQuA: 53.35 \rightarrow 51.84; OBQA: 59.94 \rightarrow 59.29), consistent with self-view emphasizing prompt-specific fit over cross-view agreement. From an efficiency standpoint, *w/o DLE* reduces training time to $\times 0.59$ but sacrifices performance, indicating a performance–cost trade-off. Overall, using DLE while excluding self-view is the most effective setting in our experiments.

Why $O(G^2)$ rescoreing does not translate to quadratic wall-clock time. CPC-GRPO performs teacher-forced cross-prompt rescoreing over G prompt variants, which introduces $O(G^2)$ forward passes. However, system-level analyses of RLHF/RLVR pipelines report that autoregressive rollout generation typically dominates iteration time, often accounting for the majority of wall-clock runtime (Zheng et al., 2025; Gao et al., 2025). Since CPC-GRPO keeps the rollout budget comparable to baselines and adds computation mainly to rescoreing/optimization, the added $O(G^2)$ forwards are amortized by the rollout-dominated runtime, resulting in the non-quadratic overhead observed in Table 6.

5.6 Log-Likelihood Consensus Reward

Table 5 ablates the scoring function used in the consensus reward. While the default CPC-GRPO uses the mean of token probabilities $s_{i,j}$ (Eq. (4)), the variant replaces it with the log-likelihood:

$$s_{i,j}^{\log} = \frac{1}{T_j} \sum_{t=1}^{T_j} \log \pi_{\theta_{\text{old}}} \left(y_t^{(j)} \mid x^{(i)}, y_{<t}^{(j)} \right). \quad (13)$$

The consensus reward aggregation and the update remain unchanged; we only substitute $s_{i,j}^{\log}$ for $s_{i,j}$.

This log-scale scoring is motivated by the observation that cross-prompt teacher-forcing can assign extremely small probabilities to early tokens under some views, making probability-mean scores saturate near zero and reducing discriminability. Using log-probabilities expands the dynamic range in the low-probability regime and is aligned with standard maximum-likelihood training.

Empirically, the log-likelihood variant substantially improves in-domain RAE, but its gains on *Un-trained* prompts and out-of-domain correctness are mixed (Table 5). This suggests that optimizing cross-view likelihood can primarily strengthen linguistic plausibility across views, which may not consistently translate to correctness-focused evaluation under distribution shift.

6 Conclusion

We study answer-free RL for LLM post-training without ground-truth answers or external evaluators. Using paraphrase prompts as multi-view observations, we define a cross-prompt consensus reward via teacher-forcing rescoreing and optimize it with an all-pair GRPO-style update using group-relative advantages. Our method improves RobustAlpacaEval on both *Trained* and *Un-trained* prompts and shows non-degrading trends on out-of-domain reasoning benchmarks. Since the reward is computable on prefixes, we can explicitly trade off training cost and performance by adjusting the prefix length L .

559 Limitations

560 From a computational perspective, our method
561 requires cross-prompt rescoring (teacher-forcing
562 scoring) with respect to the paraphrase group size
563 G . Thus, per-sample scoring expands to approxi-
564 mately $G \times K$ (prompt, response) combinations,
565 and under the default setting ($K = G + 1$) it can
566 incur additional forward computation on the order
567 of $O(G^2)$. This cost can be substantial in particular
568 when accumulating token-level log-probabilities
569 over long responses. However, as shown in this pa-
570 per, limiting the prefix length L reduces the number
571 of processed tokens and can meaningfully reduce
572 training wall-clock time, and G and L serve as pri-
573 mary levers for controlling the performance–cost
574 trade-off. We also evaluate many prompt variants in
575 our experiments for analysis purposes (separating
576 *Trained/Un-trained* and measuring prompt sensi-
577 tivity), but in deployment one may operate with
578 single-prompt generation or with limited test-time
579 ensembling.

580 Selecting the prefix length L is non-trivial. L
581 simultaneously affects both the magnitude and vari-
582 ance of the learning signal, and its effect can de-
583 pend on a combination of factors, including (i) the
584 model’s paraphrase robustness, (ii) the sharpness
585 of token distributions, and (iii) the influence of
586 early tokens induced by task and prompt formats.
587 Therefore, it is difficult to guarantee that a fixed L
588 is optimal across all models and datasets, and in
589 practice additional tuning over L may be required.
590 Future work should explore automated strategies
591 such as scheduling that gradually increases L over
592 training, or adaptive selection of L based on reward
593 variance/stability.

594 References

595 Stephen H. Bach, Victor Sanh, Zheng-Xin Yong, Al-
596 bert Webson, Colin Raffel, Nihal V. Nayak, Abheesht
597 Sharma, Taewoon Kim, M Saiful Bari, Thibault
598 Fevry, Zaid Alyafeai, Manan Dey, Andrea Santilli,
599 Zhiqing Sun, Srulik Ben-David, Canwen Xu, Gun-
600 jan Chhablani, Han Wang, Jason Alan Fries, and 8
601 others. 2022. [Promptsources: An integrated develop-
602 ment environment and repository for natural language
603 prompts](#). *Preprint*, arXiv:2202.01279.

604 Bowen Cao, Deng Cai, Zhisong Zhang, Yuexian Zou,
605 and Wai Lam. 2024. [On the worst prompt perfor-
606 mance of large language models](#). In *The Thirty-
607 eighth Annual Conference on Neural Information
608 Processing Systems*.

609 Mark Chen, Jerry Tworek, Heewoo Jun, Qiming Yuan,

Henrique Ponde de Oliveira Pinto, Jared Kaplan, 610
Harri Edwards, Yuri Burda, Nicholas Joseph, Greg 611
Brockman, Alex Ray, Raul Puri, Gretchen Krueger, 612
Michael Petrov, Heidy Khlaaf, Girish Sastry, Pamela 613
Mishkin, Brooke Chan, Scott Gray, and 39 others. 614
2021. [Evaluating large language models trained on
code](#). 615
616

Tri Dao. 2023. [Flashattention-2: Faster attention with
better parallelism and work partitioning](#). *Preprint*,
arXiv:2307.08691. 617
618
619

Tim Dettmers, Artidoro Pagnoni, Ari Holtzman, and 620
Luke Zettlemoyer. 2023. [Qlora: Efficient finetuning
of quantized llms](#). *Preprint*, arXiv:2305.14314. 621
622

Yann Dubois, Balázs Galambosi, Percy Liang, and Tat- 623
sunori B. Hashimoto. 2025. [Length-controlled al-
pacaeval: A simple way to debias automatic evalua-
tors](#). *Preprint*, arXiv:2404.04475. 624
625
626

Wei Gao, Yuheng Zhao, Dakai An, Tianyuan Wu, Lunxi 627
Cao, Shaopan Xiong, Ju Huang, Weixun Wang, Siran 628
Yang, Wenbo Su, Jiamang Wang, Lin Qu, Bo Zheng, 629
and Wei Wang. 2025. [Rollpacker: Mitigating long-
tail rollouts for fast, synchronous rl post-training](#).
Preprint, arXiv:2509.21009. 630
631
632

Daya Guo, Dejian Yang, Haowei Zhang, Junxiao Song, 633
Peiyi Wang, Qihao Zhu, Runxin Xu, Ruoyu Zhang, 634
Shirong Ma, Xiao Bi, Xiaokang Zhang, Xingkai Yu, 635
Yu Wu, Z. F. Wu, Zhibin Gou, Zhihong Shao, Zhu- 636
oshu Li, Ziyi Gao, Aixin Liu, and 175 others. 2025a. 637
[Deepseek-r1 incentivizes reasoning in llms through
reinforcement learning](#). *Nature*, 645(8081):633–638. 638
639

Jiaxin Guo, Daimeng Wei, Yuanchang Luo, Hengchao 640
Shang, Zongyao Li, Jinlong Yang, Zhanglin Wu, 641
Zhiqiang Rao, Shimin Tao, and Hao Yang. 2025b. 642
[M-ped: Multi-prompt ensemble decoding for large
language models](#). In *Findings of the Association
for Computational Linguistics: EMNLP 2025*, pages 643
16693–16711, Suzhou, China. Association for Com- 644
putational Linguistics. 645
646
647

David Heineman, Yao Dou, and Wei Xu. 2024. [Improv-
ing minimum Bayes risk decoding with multi-prompt](#).
In *Proceedings of the 2024 Conference on Empirical
Methods in Natural Language Processing*, pages 648
22525–22545, Miami, Florida, USA. Association for 649
Computational Linguistics. 650
651
652
653

Jian Hu, Xibin Wu, Wei Shen, Jason Klein Liu, Weixun 654
Wang, Songlin Jiang, Haoran Wang, Hao Chen, Bin 655
Chen, Wenkai Fang, Xianyu, Yu Cao, Haotian Xu, 656
and Yiming Liu. 2025. [OpenRLHF: A ray-based
easy-to-use, scalable and high-performance RLHF
framework](#). In *Proceedings of the 2025 Conference
on Empirical Methods in Natural Language Process-
ing: System Demonstrations*, pages 656–666, Suzhou,
China. Association for Computational Linguistics. 661
662

Nathan Lambert, Jacob Morrison, Valentina Pyatkin, 663
Shengyi Huang, Hamish Ivison, Faeze Brahman, 664
Lester James Validad Miranda, Alisa Liu, Nouha 665
Dziri, Xinxin Lyu, Yuling Gu, Saumya Malik, Victoria 666

667	Graf, Jena D. Hwang, Jiangjiang Yang, Ronan Le Bras, Oyvind Tafjord, Christopher Wilhelm, Luca Soldaini, and 4 others. 2025. Tulu 3: Pushing frontiers in open language model post-training . In <i>Second Conference on Language Modeling</i> .	
668		
669		
670		
671		
672	Wang Ling, Dani Yogatama, Chris Dyer, and Phil Blunsom. 2017. Program induction by rationale generation: Learning to solve and explain algebraic word problems . In <i>Proceedings of the 55th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)</i> , pages 158–167, Vancouver, Canada. Association for Computational Linguistics.	
673		
674		
675		
676		
677		
678		
679	Zichen Liu, Changyu Chen, Wenjun Li, Penghui Qi, Tianyu Pang, Chao Du, Wee Sun Lee, and Min Lin. 2025. Understanding r1-zero-like training: A critical perspective . In <i>Second Conference on Language Modeling</i> .	
680		
681		
682		
683		
684	Todor Mihaylov, Peter Clark, Tushar Khot, and Ashish Sabharwal. 2018. Can a suit of armor conduct electricity? a new dataset for open book question answering . In <i>Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing</i> , pages 2381–2391, Brussels, Belgium. Association for Computational Linguistics.	
685		
686		
687		
688		
689		
690		
691	Youssef Mroueh. 2025. Reinforcement learning with verifiable rewards: Grpo’s effective loss, dynamics, and success amplification . <i>Preprint</i> , arXiv:2503.06639.	
692		
693		
694		
695	Long Ouyang, Jeff Wu, Xu Jiang, Diogo Almeida, Carroll L. Wainwright, Pamela Mishkin, Chong Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, John Schulman, Jacob Hilton, Fraser Kelton, Luke Miller, Maddie Simens, Amanda Askell, Peter Welinder, Paul Christiano, Jan Leike, and Ryan Lowe. 2022. Training language models to follow instructions with human feedback . In <i>Proceedings of the 36th International Conference on Neural Information Processing Systems, NIPS ’22</i> , Red Hook, NY, USA. Curran Associates Inc.	
696		
697		
698		
699		
700		
701		
702		
703		
704		
705		
706	Mihir Prabhudesai, Lili Chen, Alex Ippoliti, Katerina Fragkiadaki, Hao Liu, and Deepak Pathak. 2025. Maximizing confidence alone improves reasoning . <i>Preprint</i> , arXiv:2505.22660.	
707		
708		
709		
710	John Schulman, Filip Wolski, Prafulla Dhariwal, Alec Radford, and Oleg Klimov. 2017. Proximal policy optimization algorithms . <i>Preprint</i> , arXiv:1707.06347.	
711		
712		
713	Francesco Tonolini, Nikolaos Aletras, Jordan Massiah, and Gabriella Kazai. 2024. Bayesian prompt ensembles: Model uncertainty estimation for black-box large language models . In <i>Findings of the Association for Computational Linguistics: ACL 2024</i> , pages 12229–12272, Bangkok, Thailand. Association for Computational Linguistics.	
714		
715		
716		
717		
718		
719		
720	Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N. Gomez, Lukasz Kaiser, and Illia Polosukhin. 2023. Attention is all you need . <i>Preprint</i> , arXiv:1706.03762.	
721		
722		
723		
	Zhewei Yao, Reza Yazdani Aminabadi, Olatunji Ruwase, Samyam Rajbhandari, Xiaoxia Wu, Ammar Ahmad Awan, Jeff Rasley, Minjia Zhang, Conglong Li, Connor Holmes, Zhongzhu Zhou, Michael Wyatt, Molly Smith, Lev Kurilenko, Heyang Qin, Masahiro Tanaka, Shuai Che, Shuaiwen Leon Song, and Yuxiong He. 2023. Deepspeed-chat: Easy, fast and affordable rlhf training of chatgpt-like models at all scales . <i>Preprint</i> , arXiv:2308.01320.	724 725 726 727 728 729 730 731 732
	Qiyang Yu, Zheng Zhang, Ruofei Zhu, Yufeng Yuan, Xiaochen Zuo, Yu Yue, Weinan Dai, Tiantian Fan, Gaohong Liu, Lingjun Liu, Xin Liu, Haibin Lin, Zhiqi Lin, Bole Ma, Guangming Sheng, Yuxuan Tong, Chi Zhang, Mofan Zhang, Wang Zhang, and 16 others. 2025. Dapo: An open-source llm reinforcement learning system at scale . <i>Preprint</i> , arXiv:2503.14476.	733 734 735 736 737 738 739 740
	Tony Zhao, Eric Wallace, Shi Feng, Dan Klein, and Sameer Singh. 2021. Calibrate before use: Improving few-shot performance of language models . In <i>International Conference on Machine Learning</i> .	741 742 743 744
	Xuandong Zhao, Zhewei Kang, Aosong Feng, Sergey Levine, and Dawn Song. 2025. Learning to reason without external rewards . <i>arXiv preprint arXiv:2505.19590</i> .	745 746 747 748
	Haizhong Zheng, Yang Zhou, Brian R. Bartoldson, Bhavya Kailkhura, Fan Lai, Jiawei Zhao, and Beidi Chen. 2025. Act only when it pays: Efficient reinforcement learning for LLM reasoning via selective rollouts . In <i>The Thirty-ninth Annual Conference on Neural Information Processing Systems</i> .	749 750 751 752 753 754
	Lianmin Zheng, Wei-Lin Chiang, Ying Sheng, Siyuan Zhuang, Zhanghao Wu, Yonghao Zhuang, Zi Lin, Zhuohan Li, Dacheng Li, Eric Xing, Hao Zhang, Joseph E. Gonzalez, and Ion Stoica. 2023. Judging LLM-as-a-judge with MT-bench and chatbot arena . In <i>Thirty-seventh Conference on Neural Information Processing Systems Datasets and Benchmarks Track</i> .	755 756 757 758 759 760 761
	Qihao Zhu Runxin Xu Junxiao Song Mingchuan Zhang Y.K. Li Y. Wu Daya Guo Zhihong Shao, Peiyi Wang. 2024. Deepseekmath: Pushing the limits of mathematical reasoning in open language models .	762 763 764 765
	Yinmin Zhong, Zili Zhang, Bingyang Wu, Shengyu Liu, Yukun Chen, Changyi Wan, Hanpeng Hu, Lei Xia, Ranchen Ming, Yibo Zhu, and Xin Jin. 2025. Optimizing RLHF training for large language models with stage fusion . In <i>22nd USENIX Symposium on Networked Systems Design and Implementation (NSDI 25)</i> , pages 489–503, Philadelphia, PA. USENIX Association.	766 767 768 769 770 771 772 773
	Chunting Zhou, Junxian He, Xuezhe Ma, Taylor Berg-Kirkpatrick, and Graham Neubig. 2022. Prompt consistency for zero-shot task generalization . In <i>Findings of the Association for Computational Linguistics: EMNLP 2022</i> , pages 2613–2626, Abu Dhabi, United Arab Emirates. Association for Computational Linguistics.	774 775 776 777 778 779 780

781 Daniel M. Ziegler, Nisan Stiennon, Jeffrey Wu, Tom B.
782 Brown, Alec Radford, Dario Amodei, Paul Chris-
783 tiano, and Geoffrey Irving. 2019. *Fine-tuning lan-
784 guage models from human preferences*. *arXiv
785 preprint arXiv:1909.08593*.

A Detailed Datasets 786

RobustAlpacaEval (RAE). RAE is an 787
instruction-following benchmark based on 788
TinyAlpacaEval. Each example consists of (i) one 789
original prompt and (ii) ten semantically equivalent 790
paraphrase prompts (11 prompt variants in total). 791
The paraphrases are automatically generated using 792
GPT-4 and then manually reviewed/edited to 793
ensure meaning preservation and fluency (Cao 794
et al., 2024). In the main experiments, we use the 795
first $G=4$ variants under a fixed index order for 796
training (*Trained prompts*), and at evaluation time 797
we compute scores separately on (i) the *Trained* 798
subset and (ii) the remaining *Un-trained* subset 799
over the same set of examples. 800

Out-of-domain benchmarks. For evaluating 801
generalization, we use OpenBookQA, AQuA, and 802
HumanEval (Mihaylov et al., 2018; Ling et al., 803
2017; Chen et al., 2021). For each dataset, we 804
construct instruction-style prompts using templates 805
from PromptSource (Bach et al., 2022), and treat 806
different templates as prompt variants (= views), 807
forming a set of prompt variants per example. We 808
report the number of examples and prompt vari- 809
ants for each dataset in Table 7. We will release the 810
dataset construction and preprocessing scripts on 811
GitHub. 812

B Detailed Experimental Setting 813

Training. All fine-tuning (FT) and RL-based 814
methods are trained with QLoRA, and we fix the 815
paraphrase group size to $G=4$ throughout training. 816
For RL training, we perform rollouts and cross- 817
prompt scoring required for policy updates; to re- 818
duce rollout cost, we can use a multi-step update 819
that reuses a single rollout across multiple opti- 820
mization steps. The main hyperparameters of our 821
default setting are summarized in Table 8. Since 822
there is no verification-based reward, we do not use 823
a KL term in the default setting ($\beta=0$), and we fix 824
the learning rate to 1×10^{-5} . 825

Hardware. Training is performed on RTX 4090 826
 $\times 2$, and evaluation is performed on RTX 4090 \times 827
1. All experiments are run under the same machine 828
environment. 829

Evaluation metrics. We evaluate RAE using the 830
AlpacaEval weighted win rate, and set gpt4_turbo 831
as the reference output and judge (Dubois et al., 832
2025; Cao et al., 2024). We also separately re- 833

Dataset	#Examples	#Prompt Variants	Total Prompts
RobustAlpacaEval (RAE)	100	1+10	1100
AQuA	254	6	1524
OpenBookQA	500	7	3500
HumanEval	164	2	328

Table 7: Dataset statistics. We report the number of examples, the number of prompt variants per example, and the resulting total number of prompts.

Hyperparameter	Value
Prompt views (G)	4
Quantization	4-bit (QLoRA)
LoRA rank (r)	16
LoRA α	32
LoRA dropout	0.05
Optimizer	AdamW
Weight decay	0.0
Epochs	3
Batch size (per device)	1
Prefix length (L ; main)	2048
PPO clip ϵ	0.2
Learning rate	1×10^{-5}
KL coefficient β	0
LR scheduler	cosine
Warmup ratio	0.03

Table 8: Hyperparameter settings used in our experiments.

port results on the *Trained* prompt subset and the *Un-trained* subset over the same set of examples. For OpenBookQA and AQuA, we extract answer choices via rule-based parsing and compute accuracy. For HumanEval, we extract code blocks using rule-based procedures and compute unit-test-based pass@1. All methods are provided with the same output-format rules, and the rule templates and parsing procedures are described in Appendix D.

Decoding for rollouts and evaluation. Unless otherwise specified, we use the same decoding configuration across all methods for fair comparison. For rollouts, we generate one response per view (prompt variant) with a fixed maximum completion length and standard stopping criteria. For the decoding-level ensemble (DLE), we compute the per-step averaged token distribution across views and apply greedy decoding. Exact decoding parameters (e.g., temperature, top- p , max new tokens, and stop sequences) are fixed across runs and included in the released configuration files.

C Detailed Baselines

In this section, we summarize the key configurations and the training/inference protocols of the

baselines compared in the main text. All baselines use the same prompt format and decoding settings as in the main paper, and for the out-of-domain benchmarks (OpenBookQA/AQuA/HumanEval) we follow the output-format and parsing rules in Appendix D. (Complete hyperparameters and full template specifications are included in the released code and configuration files.)

Pre-trained. We evaluate the base model without any additional training. On RAE, we compute the AlpacaEval weighted win rate for each prompt variant, and the *Trained/Un-trained* subset split follows the main setting.

Self-refinement (test-time). We follow the self-refinement idea used in the RAE paper, but in this work we apply only a test-time template without any additional training (Cao et al., 2024). The model (i) rewrites the input prompt into a clearer form without changing its meaning, and then (ii) generates the final response conditioned on the rewritten prompt. The output consists of two sections, REWRITTEN PROMPT: and RESPONSE:, and we use only RESPONSE as the final output for evaluation (Figure 5). This baseline can improve performance via prompt clarification without training, but it has limitations in that the rewriting process may introduce subtle meaning shifts or amplify model biases.

Swarm distillation (unsupervised SFT with a pseudo-target pool). Swarm distillation is an unsupervised distillation method that treats multiple outputs obtained from a paraphrase prompt set for the same intent as a *pseudo-target pool*, and trains each prompt to imitate a pseudo-target sampled from the pool (Zhou et al., 2022; Cao et al., 2024). For each RAE example e , we form G prompts $\mathcal{X}_e = \{x^{(i)}\}_{i=1}^G$ and denote the corresponding set of candidate outputs as $\mathcal{Y}_e = \{y^{(i)}\}_{i=1}^G$. During training, for each prompt $x^{(i)}$, we uniformly sample a pseudo-target \tilde{y} from \mathcal{Y}_e and minimize the standard instruction SFT objective (cross-entropy

applied only on the answer-token span):

$$\mathcal{L}_{\text{swarm}}(\theta) = - \sum_e \sum_{i=1}^G E_{\tilde{y} \sim \text{Unif}(\mathcal{Y}_e)} \sum_{t=1}^{|\tilde{y}|} \log \pi_{\theta}(\tilde{y}_t | x^{(i)}, \tilde{y}_{<t}). \quad (14)$$

In this work, we evaluate a simple form of the baseline by using only *uniform sampling* from the pseudo-target pool, without separately tuning additional selection mechanisms or score-based aggregation. This approach can leverage expression diversity across paraphrases as supervision, but if pseudo-targets include teacher errors, it can suffer from confirmation bias where errors are amplified.

Static-teacher (st). Static-teacher keeps the pseudo-target pool \mathcal{Y}_e fixed throughout training. In our experiments, we use a pre-generated teacher output file (outputs for each prompt variant) as the fixed pool, and at each step we randomly sample a pseudo-target from \mathcal{Y}_e and train using Eq. (14).

Online-teacher (ot). Online-teacher is an *iterative self-training* variant that repeatedly refreshes the pseudo-target pool using responses generated by the current model during training. For each example e and prompt $x^{(i)}$, we generate a new response under the current policy and replace the pool element for that view:

$$y_e^{(i)} \leftarrow \text{Gen}(\pi_{\theta}; x^{(i)}). \quad (15)$$

We then uniformly sample pseudo-targets from the updated \mathcal{Y}_e and continue training with Eq. (14). In our implementation, we use deterministic decoding (e.g., greedy) to reduce variance in generation, and pool updates are performed only for examples/views that appear in the current minibatch. Online-teacher can quickly reflect the student’s latest distribution, but it poses risks of reduced diversity or mode collapse.

INTUITOR. INTUITOR is an RLIF-style method that optimizes with GRPO-style RL using intrinsic rewards derived from the model’s internal self-certainty signals, without external reward models or answer verification (Zhao et al., 2025). In this work, we use the public implementation as-is, while applying the same prompt format/decoding settings and evaluation rules as in the main paper for fair comparison. The detailed reward definition and training pipeline follow the original paper and

the released code (Zhao et al., 2025). This baseline provides an answer-free training signal, but the alignment between certainty signals and true task performance can vary with tasks and prompts.

D Formatting RULE

Test-time RULE To evaluate the out-of-domain datasets in this paper, we must parse the answer portion from each model response to compute the metrics. Therefore, for each dataset we prepend a minimal set of rules to the input so that the model outputs the answer in a prescribed format. The rules used in our experiments are shown in Figure 4.

Self-refinement baseline format. Self-refinement requires the model to first rewrite the input prompt and then produce a response to the rewritten prompt. At evaluation time, we use only the text after RESPONSE: as the final output (the rewritten prompt is excluded from evaluation). In RobustAlpacaEval, many cases skip the reasoning stage, so we provide an additional format specifically for RobustAlpacaEval. The formatting rule used in our experiments is shown in Figure 5.

E Generalization Verification

In the main text, due to resource constraints, we primarily analyze results from a single run per setting. To further examine whether the observed gains of the proposed method (CPC-GRPO) overly depend on a specific run, we repeat training with different random seeds and measure performance under the same evaluation protocol.

As shown in Table 9, for both base models, most metrics exhibit comparable improvements over the pre-trained baseline on in-domain (RAE) and out-of-domain reasoning benchmarks. In particular, for Qwen3-4B, all three seeds consistently maintain improved performance on both *Trained* and *Untrained* subsets, and the variation on out-of-domain benchmarks is limited. For LLaMA3.2-3B, seed-wise variability is relatively larger, but improvements are generally preserved on OBQA and HumanEval. Overall, these results indicate that the gains of CPC-GRPO are not confined to a particular case and exhibit broadly reproducible trends.

F Performance of RobustAlpacaEval

For follow-up work, we additionally include a performance table in the same format as (Cao et al.,

AQuA (multiple-choice)

RULE:

- You may reason briefly.
- Stay within 2048 tokens.
- Your response must end with exactly: The correct answer is <LETTER>
- <LETTER> must be one of A, B, C, D, or E.

OpenBookQA (multiple-choice/free-form)

RULE:

- You may reason briefly.
- Stay within 2048 tokens.
- Your response must end with exactly: The correct answer is <ANSWER>

HumanEval (code generation)

RULE:

- You may reason briefly.
- Stay within 2048 tokens.
- After briefly reasoning, output ONLY one code block: “python ... “
- The block must contain ONLY the complete function definition with the exact same signature.

Figure 4: Output constraints used for out-of-domain evaluation to enable deterministic parsing. The same rules are prepended for all methods.

989 2024). In this work, we report RobustAlpacaEval
990 results by separating paraphrases used for training
991 from the remaining paraphrases, whereas Cao et al.
992 (2024) reports performance over the full set of 11
993 prompt versions. Accordingly, we report results
994 aggregated over the full paraphrase set in Table 10.

RobustAlpacaEval

Rewrite the PROMPT to be clearer without changing meaning (do not add/remove requirements),
then respond to the rewritten prompt. Output exactly two sections.
ORIGINAL PROMPT: {...}

Out-of-Domain

Rewrite the PROMPT to be clearer without changing meaning (do not add/remove requirements),
then respond to the rewritten prompt. Output exactly two sections.

ANSWER FORMAT

REWRITTEN PROMPT: {Rewritten prompt}

RESPONSE: {Response of rewritten prompt}

ORIGINAL PROMPT: {...}

Figure 5: Formatting rule used for the self-refinement baseline. We only evaluate the content after RESPONSE: .

Base Model	Method	RAE (In-domain)		Reasoning (Out-of-domain)		
		Trained	Un-trained	OBQA	AQuA	HumanEval
LLaMA3.2-3B	pre-trained	21.42	19.25	58.66	52.36	60.06
	CPC-GRPO (seed1)	22.89	21.04	59.94	53.35	<u>61.59</u>
	CPC-GRPO (seed2)	20.12	<u>20.08</u>	59.60	<u>54.46</u>	<u>61.59</u>
	CPC-GRPO (seed3)	<u>21.56</u>	19.96	<u>59.71</u>	54.59	62.50
Qwen3-4B	pre-trained	61.64	63.65	82.03	85.50	91.77
	CPC-GRPO (seed1)	64.78	66.19	82.03	<u>85.70</u>	92.07
	CPC-GRPO (seed2)	<u>64.17</u>	<u>66.55</u>	<u>81.97</u>	85.50	92.07
	CPC-GRPO (seed3)	63.16	66.61	<u>81.57</u>	86.29	<u>91.16</u>

Table 9: Seed sensitivity of CPC-GRPO. Pre-trained scores are taken from Table 1. Within each base model and each column, **bold** denotes the best and underline denotes the second-best among CPC-GRPO seeds (ties are marked).

Base Model	Method	Para.	Train	RAE (All prompt versions)				
				Orig.	Best	Worst	Avg.	Std.
LLaMA3.2-3B	pre-trained	X	–	30.88	64.05	1.09	20.90	23.24
	self-refine	X	–	12.36	35.04	0.01	9.27	13.08
	swarm-distillation (st)	✓	FT	25.22	55.76	0.28	18.19	20.54
	swarm-distillation (ot)	✓	FT	17.15	43.59	0.00	13.28	16.29
	INTUITOR	X	RL	28.82	60.22	0.68	19.58	21.84
	CPC-GRPO	✓	RL	34.84	60.75	1.96	22.80	22.12
Qwen3-4B	pre-trained	X	–	68.77	95.64	11.40	63.58	31.89
	self-refine	X	–	13.40	35.49	0.35	9.17	12.62
	swarm-distillation (st)	✓	FT	16.60	53.11	0.04	14.72	19.48
	swarm-distillation (ot)	✓	FT	39.10	80.40	0.58	33.09	30.43
	INTUITOR	X	RL	70.05	95.45	9.46	62.31	32.29
	CPC-GRPO	✓	RL	68.62	96.04	15.73	66.03	30.30

Table 10: RobustAlpacaEval (RAE) statistics aggregated over all prompt versions (original + paraphrases).