

VRPO: Rethinking Value Modeling for Robust RL under Noisy Supervision in LLM Post-Training

Anonymous ACL submission

Abstract

Reinforcement Learning (RL) in real-world environments often suffers from ambiguous or incomplete reward supervision, which undermines policy stability and generalization. Such noise may cause models to ignore key information or even collapse in advantage estimation. We find that a strong value model is essential for absorbing unstable signals and producing reliable advantages, offering denser and more robust supervision than the reward model. To better optimize noisy supervision, we propose VRPO, a framework that enhances value modeling for robust RL in LLM post-training. VRPO integrates (1) auxiliary losses guided by entropy and perplexity from a frozen language model, and (2) a variational information bottleneck, enabling the value model to filter noise and capture key words. This design allows the value model to correct noise rewards and generate more reliable advantage estimates, transforming it from a passive predictor into an active noise regulator. Experiments on multi-turn dialogue, math reasoning, and science QA with both rule-based and model-based rewards show that VRPO consistently outperforms baselines such as PPO and GRPO. Our work highlights the central role of the value model in Robust RL and provides a principled and practical approach to policy optimization under noisy supervision.

1 Introduction

Reinforcement Learning has achieved remarkable success across a wide range of applications (Perolat et al., 2022; Chen et al., 2022; Bellemare et al., 2020; Zhou et al., 2020; Yue et al., 2024; Caregnato-Neto et al., 2024). However, deploying RL in real-world scenarios often involves noisy or imperfect supervision, particularly when optimization depends on human feedback or learned reward models. This challenge is especially evident in LLM Post-Training RL and preference-based Reinforcement Learning from Human Feedback

(RLHF) (Casper et al., 2023; Zhang et al., 2025), where reward signals are approximate and not directly derived from ground-truth annotations, such as those generated by Generalized Advantage Estimation (GAE) (Schulman et al., 2018).

To address this issue, recent works have proposed robust RL training methods by either denoising reward models (Cheng et al., 2025; Miao et al., 2024) or filtering corrupted data before policy updates (Cheng et al., 2024; Wang et al., 2024b). However, these methods implicitly assume that reward errors can be corrected during training, which often does not hold in practice: training rewards are frequently ambiguous, sparse, or fundamentally unreliable (Poiani et al., 2024). As a result, inaccurate reward signals can propagate throughout RL training, leading to instability in textual perception, loss of critical information during advantage estimation, and ultimately degraded policy optimization performance and convergence stability.

While noisy reward modeling has received considerable attention, the value model, which provides denser supervision and plays a central role in training, has been largely overlooked as a potential denoising mechanism. In this paper, we adopt an information-theoretic perspective (Hung et al., 2023) and propose an alternative approach: rather than relying solely on reward estimation or data filtering, we enhance RL robustness by directly optimizing the value model to absorb uncertainty and stabilize training. This perspective highlights the value model as a critical component for guiding language perception under noisy supervision.

We introduce **VRPO** (Value Model Boosting for Robust Policy Optimization), a novel framework that integrates a noise-resilient and semantically aware value model into RL training. VRPO combines two key innovations: (1) an auxiliary loss guided by entropy and perplexity from a frozen language model, and (2) a variational information bottleneck (IB) architecture for robust value model-

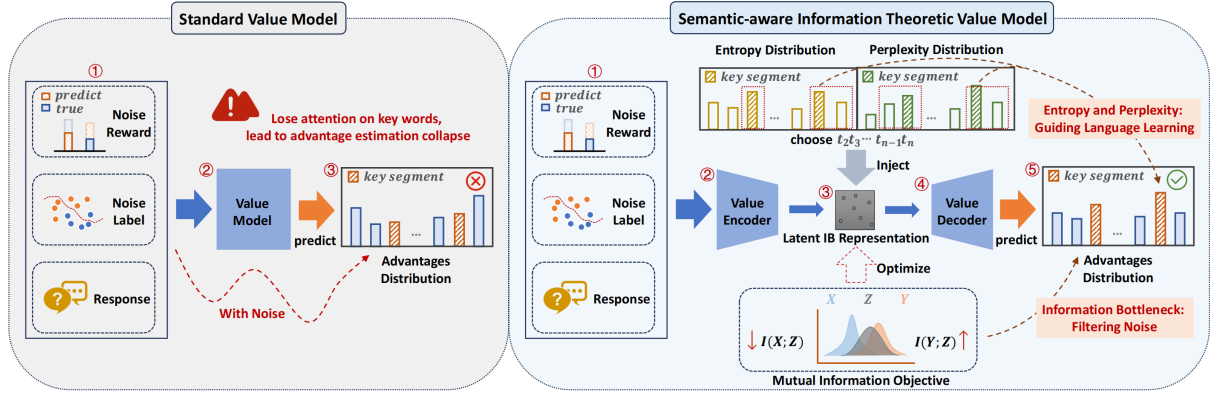


Figure 1: Comparison between the Standard Value Model and our Semantic-aware Information Theoretic Value Model. The standard model fits noisy rewards and labels, leading to unstable advantage estimates. VRPO enhances robustness by guiding learning with entropy and perplexity signals, and filtering noise via a variational information bottleneck, yielding more reliable value predictions and improved PPO stability.

ing. The IB mechanism constrains irrelevant information flow, enhancing tolerance to noise, while language signals from the frozen LM guide model to capture key words under noisy training. These auxiliary objectives further align the model’s internal feature space with the semantic space of language, effectively suppressing spurious noise while retaining task-relevant information. Moreover, VRPO alleviates overfitting (Gao et al., 2022) and mitigates contextual noise (Wang et al., 2022), thereby improving training robustness.

We evaluate VRPO across multiple tasks, including multi-turn dialogue, mathematical reasoning, and scientific question answering, under both rule-based and model-based noisy reward settings. Experimental results show that VRPO consistently outperforms baselines such as PPO and GRPO in noisy supervision scenarios. By incorporating information-theoretic regularization and semantic guidance into RL training, VRPO transforms the value model from a passive estimator into an active, noise-aware regulator. Through context-appropriate advantage estimation and robust key-signal extraction, our method significantly enhances stability and generalization in complex environments.

Our contributions are as follows:

- We propose entropy and perplexity-based auxiliary losses that mitigate advantage collapse from noisy rewards improving the semantic robustness of the value model.
- We introduce a variational information bottleneck structure that constrains irrelevant information flow.

- We empirically demonstrate that VRPO improves RL robustness and generalization in noisy RLHF and RL scenarios, offering a principled and scalable solution for real-world RL applications.

2 Related Work

Robust RL methods primarily focus on improving reward models or filtering noisy data. Reward-oriented approaches build noise-resistant discriminators via conservative gradients (Liang et al., 2024), ensemble models (Wang et al., 2024a), or uncertainty-aware losses (Wu et al., 2024). For example, data-centric methods like RIME (Cheng et al., 2024) use pre-trained denoisers with KL bounds or apply reward consensus to remove inconsistent preferences (Wang et al., 2024a). Our work shifts robustness efforts from reward correction to value model compensation. Through value model optimization, we partially absorb and filter out the noise propagation throughout the training process.

The Information Bottleneck (IB) framework offers a principled way to learn compact (Dai et al., 2018; Goyal et al., 2019; Hafner et al., 2019; Ibarz et al., 2018), robust representations by balancing compression and prediction (Wang et al., 2022) through variational approximations (Martini et al., 2024). In RL, IB has enhanced policy or reward robustness under noisy supervision (Miao et al., 2024) but has rarely been applied to value models. We fill this gap by integrating IB regularization into the value function, enabling it to absorb noise during RL training.

3 Method

3.1 Motivation

In standard Proximal Policy Optimization (PPO), policy updates depend on the joint contribution of reward signals and value estimates via advantage computation. The core optimization objective is:

$$\mathcal{L}_{\text{PPO}} = \mathbf{E}_t \left[\min \left(r_t(\theta) \hat{A}_t, \text{clip}(r_t(\theta), 1 - \epsilon, 1 + \epsilon) \hat{A}_t \right) \right], \quad (1)$$

where the advantage \hat{A}_t is typically computed through Generalized Advantage Estimation (GAE):

$$\hat{A}_t = \delta_t + (\gamma\lambda)\delta * t + 1 + (\gamma\lambda)^2\delta_{t+2} + \dots, \quad (2)$$

$$\delta_t = r_t + \gamma V(s_{t+1}) - V(s_t). \quad (3)$$

This definition highlights the strong coupling between reward signals r_t and the value function $V(s)$. When r_t is noisy, its error propagates through \hat{A}_t , leading to unstable training and degraded final policy performance.

We observe that under noisy supervision, the model often **misallocates advantage values** and fails to emphasize semantically important tokens. As illustrated in Figure 2 (right), for correct answers, the model tends to concentrate higher advantages toward the sequence end, resulting in **length hacking** (Chaudhari et al., 2024; Laidlaw et al., 2025), where longer responses are rewarded regardless of quality. For incorrect answers, advantages are incorrectly assigned to irrelevant tokens. The root cause of such errors lies in biased reward signals: noisy r_t induces distorted advantages under GAE.

To address this, we propose a complementary perspective: **enhancing the value model itself to compensate for biased advantage estimation under noisy rewards**. Specifically, we design value models guided by semantic signals from a frozen language model, ensuring robustness to linguistic noise and preserving attention to key tokens. Additionally, we incorporate an information bottleneck mechanism to constrain irrelevant information flow, thereby improving the value model’s noise tolerance and stabilizing training. As shown in Figures 2, this design not only preserves critical signals in advantage estimation but also suppresses fluctuations under irrelevant perturbations.

Compared with reward-only denoising approaches, our framework leverages denser value supervision to more effectively correct distorted

advantages, leading to superior training stability and final policy performance.

3.2 Information-Theoretic Value Modeling

To mitigate reward noise propagation in PPO, the value model must extract return-relevant information while filtering task-irrelevant details. This capability is essential for robust training under ambiguous, sparse, or inconsistent reward signals.

We address this by reformulating value modeling from an information-theoretic perspective. Let the input be a random variable X , the latent representation Z , and the return Y . Assuming Z follows a Gaussian distribution, we define:

$$I_{\text{bottleneck}} = I(X; Z), \quad I_{\text{value}} = I(Z; Y) \quad (4)$$

Here, $I_{\text{bottleneck}}$ measures information retained from the input, while I_{value} quantifies predictive information for returns. The objective is to learn a latent representation that maximizes return relevance while minimizing input redundancy:

$$\max J(\theta) = I(Z; Y) - \beta I(X; Z) \quad (5)$$

where β controls the trade-off between compression and relevance, and θ denotes model parameters.

Since mutual information is intractable in high-dimensional settings, we optimize a variational lower bound using training data $D = \{(x_i, y_i)\}_{i=1}^N$:

$$J(\phi, \psi) \geq J_{\text{VLB}}(\phi, \psi) = \mathbf{E}_{(x,y) \sim D} [J_{\text{value}} - \beta J_{\text{bottleneck}}] \quad (6)$$

$$J_{\text{value}} = \mathbf{E}_{z \sim p_\psi(z|x)} [\log q_\psi(y|z)] \quad (7)$$

$$J_{\text{bottleneck}} = \text{KL}[p_\phi(z|x) \| r(z)] \quad (8)$$

where $p_\phi(z|x)$ is the variational encoder, $q_\psi(y|z)$ the return predictor, and $r(z)$ the standard Gaussian prior. Parameters ϕ and ψ correspond to the encoder and decoder.

In practice, we model $p_\phi(z|x)$ as a diagonal Gaussian:

$$z = f_\phi^\mu(x) + f_\phi^\sigma(x) \cdot \epsilon, \quad \epsilon \sim \mathcal{N}(0, I) \quad (9)$$

and use a lightweight MLP q_ψ to predict returns from z . The final training objective becomes:

$$\mathcal{L}_{\text{value}} = \mathbf{E}_{(x,y) \sim D} [-\log q_\psi(z) + \beta \cdot \text{KL}(p_\phi(z|x) \| r(z))] \quad (10)$$

This structure encourages the value model to learn compact, reward-relevant latent representations while suppressing irrelevant or noisy input signals.

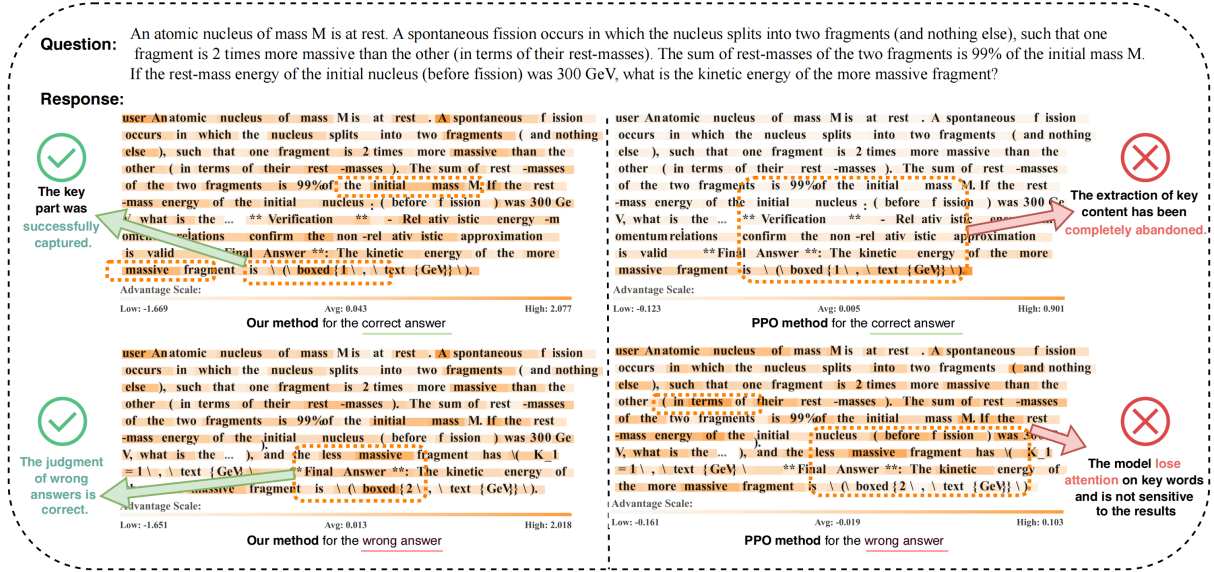


Figure 2: Token-level advantage estimation for the same response across different methods. Our method exhibits sharper focus on critical reasoning steps.

3.3 Enhancing Semantic Awareness for Noise-Resistant Value Modeling

While the information bottleneck in Section 3.2 helps filter out irrelevant features, it does not directly address semantic misalignment. Under noisy reward supervision, the value model may still lose attention on key words, resulting in prediction errors. To mitigate this, we introduce a semantic-level regularization mechanism that enhances the alignment between the model’s feature space and the language space.

Concretely, we incorporate a frozen LM head into the value model to express its internal token-level understanding. We then guide this understanding using auxiliary losses based on entropy and perplexity, encouraging semantic consistency with the actor model despite noisy rewards. Let $P_V(y_t|x)$ denote the token-level prediction distribution from the LM head. We define the auxiliary losses as:

$$L_{\text{ent}} = \sum_{t \in T_{\text{ent}}} H[P_V(y_t|x)], \quad (11)$$

$$L_{\text{ppl}} = \sum_{t \in T_{\text{ppl}}} -\log P_V(y_t = y_t^*|x) \quad (12)$$

where $H[\cdot]$ is entropy, and $T_{\text{ent}}, T_{\text{ppl}}$ are dynamically selected token subsets with high entropy and perplexity respectively:

$$T_{\text{ent}} = \{t : H[P_V(y_t|x)] > \hat{T}_{\text{entropy}}\}, \quad (13)$$

$$T_{\text{ppl}} = \{t : -\log P_V(y_t = y_t^*|x) > \hat{T}_{\text{perplexity}}\} \quad (14)$$

where \hat{T}_{entropy} and $\hat{T}_{\text{perplexity}}$ represent the threshold of entropy and perplexity. The final semantic regularization loss is:

$$L_{\text{semantic}} = \lambda_{\text{ent}} L_{\text{ent}} + \lambda_{\text{ppl}} L_{\text{ppl}} \quad (15)$$

where λ_{ent} and λ_{ppl} represent the weight of the loss.

This method provides the value model with stable, semantically meaningful supervision, helping it stay anchored to the input’s linguistic structure and resist reward noise. Combined with the information bottleneck from Section 3.2, it reinforces both semantic alignment and robustness.

4 Experiments

4.1 Setup

We evaluate VRPO in two distinct but complementary settings: **rule-based rewards** and **model-based rewards**, both designed to test robustness under noisy supervision.

Rule-based reward setting This scenario mainly considers the noise in unsupervised learning. We simulate supervision by performing majority voting over multiple outputs from pretrained models. Two variants are considered: (i) *training-time augmentation* where noisy pseudo-labels from the training set supervise RL training and (ii) *test-time augmentation*, where noisy pseudo-labels from the test set supervise RL training. Models are cold-started via supervised fine-tuning and constrained to specific answer formats. The results of test-time augmentation can be found in the appendix.

Model-based reward setting We conduct experiments on multi-turn dialogue, using a separate reward model trained on annotated dialogue data. A policy model is optimized via RL under this noisy reward.

Baselines In both settings, we compare VRPO with PPO and GRPO under identical initialization. In addition, we include comparisons with several strong robust reinforcement learning baselines, including Dr.GRPO, Reinforce++, KTAE(Sun et al., 2025), and λ -GRPO(Wang et al., 2025). All models are based on Qwen2.5 (Qwen et al., 2025), Qwen3 (Yang et al., 2025) or Llama3.1 variants, fine-tuned appropriately.

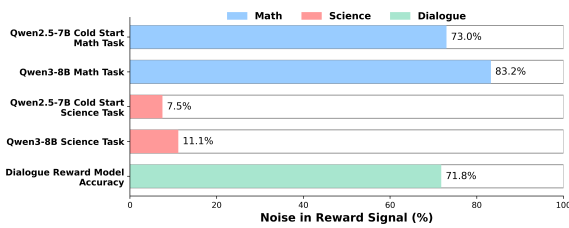


Figure 3: Noise statistics in the various tasks. A significant portion of rewards contains inaccuracies.

Training Noise Noise is pervasive in RL training, especially under model-based feedback. Figure 3 illustrates the error rate in various reward signals. The noise in the rule-based reward from the majority vote generated data.

Implementation Details In the majority voting during test-time optimization, 32 samples are used, while 5 samples are used for voting during training-time optimization. The weights for the auxiliary losses (entropy and perplexity) are set to 0.5, making 80% of the labels effective. The experiments are conducted on 8 NVIDIA A100 80GB GPUs. For further dataset construction, training details, and model initialization, please refer to Appendix.

4.2 Evaluation Metrics

To evaluate the effectiveness of our value model, we consider both task-level performance and value estimation capabilities across different task domains.

4.2.1 Task-Level Metrics

For mathematical reasoning and scientific knowledge tasks, we adopt accuracy as the primary metric to measure whether the model produces correct answers. For multi-turn dialogue tasks, we employ a tripartite evaluation framework that captures

both task execution and communicative quality: (1) Task Completion Rate (TCR), assessing whether the model successfully fulfills the user-intended objective; (2) Ask Completion Rate (ACR), measuring whether detailed aspects of the user query are adequately addressed; and (3) Goal Completion Rate (GCR), evaluating the fluency, coherence, and appropriateness of the model’s responses. This comprehensive evaluation allows us to examine not only whether the model achieves the intended outcomes, but also how effectively it aligns with human communication standards.

4.2.2 Value Model Performance Metrics

To directly evaluate the accuracy and robustness of the value model, we adopt the following criteria:

Explained Variance We use the explained variance to measure how much of the empirical return is captured by the predicted value:

$$\text{Explained Variance} = 1 - \frac{\text{Var}[\hat{R} - V]}{\text{Var}[\hat{R}]} \quad (16)$$

where \hat{R} denotes the empirical return and V is the predicted value. A value near 1 indicates high explanatory power, while values near or below zero indicate poor or misleading predictions.

Prediction Error Evaluation To assess the value model’s prediction accuracy, we adopt Mean Squared Error (MSE) and Mean Absolute Error (MAE) between the predicted value $V_{\theta}(s_t)$ and the reference target $V_{\text{target}}(s_t)$:

$$\mathcal{L}_{\text{MSE}} = \frac{1}{T} \sum_{t=0}^{T-1} (V_{\theta}(s_t) - V_{\text{target}}(s_t))^2, \quad (17)$$

$$\mathcal{L}_{\text{MAE}} = \frac{1}{T} \sum_{t=0}^{T-1} |V_{\theta}(s_t) - V_{\text{target}}(s_t)| \quad (18)$$

Here, T denotes the trajectory length, and $V_{\text{target}}(s_t)$ is computed via n -step Temporal Difference (TD) with Generalized Advantage Estimation (GAE):

$$V_{\text{target}}(s_t) = \sum_{k=0}^{n-1} \gamma^k r_{t+k} + \gamma^n V(s_{t+n}) \quad (19)$$

If $t + n \geq T$, the bootstrap term $\gamma^n V(s_{t+n})$ is set to zero to account for episode termination. MSE emphasizes large errors and captures prediction stability, while MAE reflects overall distributional deviation. As the rewards are obtained through simulation, **these metrics primarily serve as relative indicators of training quality rather than absolute performance benchmarks.**

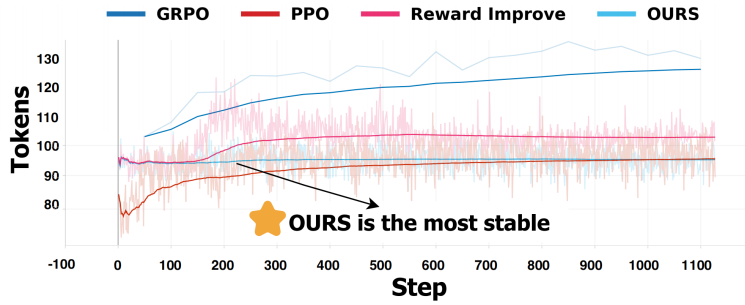


Figure 4: Average response length across training. Our method maintains stable lengths (94 to 95), while PPO and GRPO increase significantly due to length reward hacking.

4.2.3 Advantage Visualization

To qualitatively assess the robustness introduced by value model training, we visualize token-level advantage values across generated sequences. In reinforcement learning-based language generation, the advantage function quantifies the relative benefit of taking action a_t in state s_t , and is computed as:

$$A_t = Q(s_t, a_t) - V(s_t) \quad (20)$$

where $Q(s_t, a_t)$ represents the estimated return of taking action a_t in state s_t , and $V(s_t)$ denotes the baseline value of the current state.

Higher advantage values indicate that the corresponding token exhibits significantly better performance compared to the average policy behavior under the current context. These tokens function as positive learning signals and reflect directions along which the policy should be reinforced. Therefore, visualizing token-level advantages enables us to analyze which parts of the generation are being prioritized during optimization and how the learned policy aligns with high-reward behaviors.

4.3 Dialogue Task under Model-Based Reward RL Training

To assess the effectiveness of our approach in realistic settings with model-based reward supervision, we conduct multi-turn dialogue experiments on the **Honor-Dialogue Dataset**. This dataset, constructed by us, contains multi-domain, task-oriented dialogues collected from real-world scenarios. The task requires the model to act as a dialogue assistant, producing natural and effective responses. Under noisy dialogue reward model supervision, the model is trained with reinforcement learning on real multi-turn dialogues. Performance is evaluated through interactions with GPT-4o using the dialogue standards in Section 4.2.1. Dataset statistics are provided in the Appendix, and Fig-

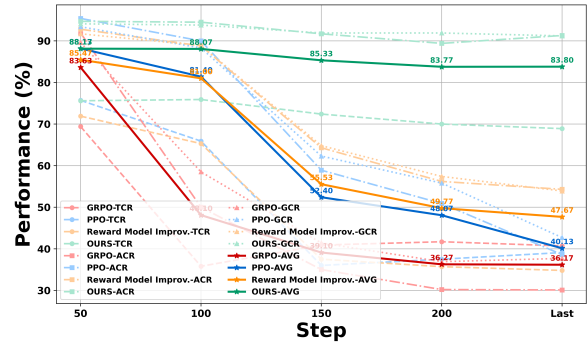


Figure 5: Training performance on dialogue tasks. VRPO (the top lines) improves task completion from 72.1% to 75.9%, avoids collapse under noise, and achieves a final average performance of 83.80%, outperforms GRPO (36.17%) and PPO (40.13%) in stability.

ure 6 shows the training dynamics. All real information in the data has undergone anonymization, and the use of the data has been authorized.

VRPO enhances training stability in real-world general scenarios Despite noisy supervision, our method effectively improves dialogue agent performance, surpassing the cold-start baseline from 85.87% to 88.13%. In contrast, PPO and GRPO exhibit severe performance collapse under noisy rewards, with GRPO performing the worst. This demonstrates our method’s stability and resilience in noisy settings. For further experimental data results, please refer to Appendix.

Strong Value Model Helps Reduce Reward Hacking A stronger value model helps mitigate reward model hacking by resisting length bias during training. As shown in Figure 4, the performance collapse largely stems from length-based reward exploitation. Our method stabilizes advantage estimation, preventing it from scaling with response length. In contrast, PPO and GRPO exhibit sharp length inflation, revealing greater vulnerability to

445 biased rewards.

446 **Value model better corrects model training bias**
447 **compared to Reward model** We applied our
448 method to the reward model for the experiment.
449 The results show that the value model provides
450 denser supervision, leading to superior correction
451 of the final advantage estimates. Compared to re-
452 ward model optimization methods, VRPO better
453 preserves training performance in dialogue tasks,
454 achieving an optimal average performance im-
455 provement of 75.90%, with a final average perfor-
456 mance exceeding 36.13%. During training, length
457 fluctuations are more stable, effectively addressing
458 reward biases caused by length during training.

459 **4.4 Math and Science Task under Rule-Based** 460 **Rewards RL Training**

461 To further assess the effectiveness of our approach
462 in realistic settings where supervision comes from
463 rule-based rewards, we conduct a series of exper-
464 iments on mathematical and scientific reasoning
465 tasks during training-time optimization. Specif-
466 ically, we use the Light-R1 dataset (Wen et al.,
467 2025) for math tasks and the SuperGPQA(Team
468 et al., 2025) dataset for scientific tasks. For each
469 dataset, we perform multiple inferences with the
470 original model, followed by voting to filter the
471 results before training. The evaluation is then
472 conducted on math and scientific QA datasets.
473 The test datasets include four math-focused tasks:
474 MATH500 (Hendrycks et al., 2021), AIME24,
475 Minerva-Math (Lewkowycz et al., 2022), AMC23,
476 as well as three scientific knowledge tasks: Sam-
477 pleQA (Wei et al., 2024), GPQA , and HLE (Hu-
478 manity’s Last Exam)(Phan et al., 2025). Table 8
479 reports the accuracy of models before and after RL
480 training, comparing weak and strong models under
481 noisy conditions. Due to the limitation of memory,
482 multiple rounds of generation are carried out here.
483 In each round, the generated length is at most 4096
484 tokens, and the result of the last round is used for
485 evaluation.

486 **VRPO achieves optimization under noisy su-**
487 **pervision across multiple domains** By utiliz-
488 ing a value model with stronger language percep-
489 tion, VRPO consistently outperforms other base-
490 line methods in multiple domains. This improve-
491 ment is particularly evident in math reasoning and
492 scientific knowledge tasks. For example, after
493 training with the Qwen3-8B model, the average
494 accuracy on math tasks increased from 58.40%

495 to 60.55%, and on scientific tasks from 2.96% to
496 3.82%. These results highlight the model’s ability
497 to extract relevant information from noisy or fuzzy
498 feedback and showcase its strong generalization
499 across various domains.

500 **VRPO shows significant effectiveness in scenar-**
501 **ios with severe training collapse** In math tasks
502 with a weaker model (Qwen2.5-7B cold-start),
503 model training is heavily impacted by noise, lead-
504 ing to a notable performance drop in GRPO, with
505 accuracy on MATH500 dropping to 50.40%. With
506 a stronger model (Qwen3-8B), noise impact de-
507 creases, and all methods improve. Our method still
508 provides stable gains, achieving 38.14% average ac-
509 curacy on all datasets, outperforming Reinforce++
510 by 1.54%. In the weaker model setup, the gain is
511 even higher (6.49%), demonstrating the advantage
512 of our method under more challenging conditions.

513 **VRPO helps guide the model’s advantage esti-**
514 **mates to focus on key words** Figure 2 visualizes
515 the advantage estimation for the same response
516 across different methods. While PPO disperses
517 attention across tokens, our model focuses on crit-
518 ical text information. It not only captures the key
519 elements of the response but also shows a more
520 appropriate perception and judgment of its answer.

521 **4.5 Ablation Experiment and Discussion**

522 **Enhanced Semantic Awareness rather than Se-**
523 **matic Understanding** As shown in Table 2, we
524 conduct an ablation study to evaluate how different
525 forms of language-awareness affect model robust-
526 ness under noisy reward supervision. Our results
527 suggest that directly solving tasks via cross-entropy
528 is suboptimal for the value model. Instead, regu-
529 larizing with entropy and perplexity, focusing on
530 semantic perception, leads to stronger generaliza-
531 tion and robustness. Specifically, minimizing these
532 losses improves the value model’s ability to identify
533 meaningful signals amid noisy feedback, boosting
534 average accuracy to 42.41%, compared to 39.51%
535 under standard cross-entropy training. This im-
536 provement stems from the model’s enhanced align-
537 ment with linguistic structure, enabling it to act not
538 as a task solver, but as a semantic signal regulator,
539 assessing responses with higher fidelity and filter-
540 ing unstable patterns. While actor improvements
541 help with comprehension, our results show that the
542 value model is more effective in absorbing reward
543 uncertainty, yielding a 4% gain over actor-centric
544 training on average.

Domain	Math				Factuality	Science	Knowledge	ALL
Method	MATH500	AIME24	Minerva-Math	AMC23	SampleQA	GPQA	HLE	AVG
<i>Qwen2.5-7B-Cold Start(Weak Model)</i>								
Base	71.60%	6.67%	18.75%	52.50%	2.36%	1.45%	3.29%	22.37%
GRPO	50.40%	6.67%	13.60%	22.50%	2.36%	2.36%	2.97%	14.41%
PPO	67.40%	13.33%	19.12%	35.83%	2.54%	2.54%	3.43%	20.60%
Reinforce++	63.80%	6.67%	10.66%	40.83%	2.64%	2.17%	3.34%	18.59%
Dr.GRPO	73.00%	13.33%	18.75%	53.33%	2.54%	1.27%	3.34%	23.65%
Ours	72.20%	23.33%	19.85%	50.00%	2.82%	3.44%	3.94%	25.08%
<i>Qwen3-8B(Strong Model)</i>								
Base	87.40%	41.67%	28.68%	75.83%	2.89%	3.10%	2.89%	34.64%
GRPO	89.20%	35.00%	28.68%	84.17%	3.03%	2.98%	3.24%	35.19%
PPO	86.40%	36.67%	30.51%	80.00%	2.82%	2.17%	3.29%	34.55%
Reinforce++	89.00%	45.00%	27.21%	84.17%	3.19%	4.35%	3.29%	36.60%
Dr.GRPO	87.80%	45.00%	27.94%	85.83%	2.50%	3.99%	3.10%	36.59%
Ours	90.20%	46.67%	31.99%	86.67%	3.21%	4.35%	3.89%	38.14%

Table 1: Accuracy (%) on train-time optimization with rule-based rewards across multiple reasoning benchmarks. Our method achieves consistent improvements under both weaker and stronger base models.

Method	MATH500	AIME24	Minerva-Math	AMC23	AVG
<i>Actor Model</i>					
Cross-Entropy Loss	73.20%	10.00%	19.12%	50.83%	38.29%
CE with Entropy/Perplexity Filtering	70.40%	13.33%	19.12%	50.83%	38.42%
Entropy + Perplexity Minimization	69.40%	10.00%	18.01%	54.17%	37.90%
<i>Value Model</i>					
Cross-Entropy Loss	73.00%	10.00%	18.38%	56.67%	39.51%
CE with Entropy/Perplexity Filtering	72.20%	10.00%	18.75%	57.50%	39.61%
Entropy + Perplexity Minimization(Ours)	74.40%	13.33%	20.22%	61.67%	42.41%

Table 2: Comparison of different training objectives on math reasoning tasks under test-time reward perturbation. Entropy- and perplexity-based regularization improves value model robustness (+2.9%) over standard cross-entropy, confirming its role as a semantic signal regulator.

Activation	MATH 500	AIME 24	Minerva Math	AMC 23	AVG
0%(IB-only)	72.20%	10.00%	20.22%	51.67%	38.52%
20%	75.00%	6.67%	18.38%	54.17%	38.56%
50%	73.40%	13.33%	19.49%	55.83%	40.51%
80%(Ours)	74.40%	13.33%	20.22%	61.67%	42.41%
100%	73.80%	16.67%	19.12%	55.00%	41.15%

Table 3: Impact of partial token activation on performance via entropy/perplexity-based loss, activating 80% of high-uncertainty tokens yields the best accuracy.

Partial Token Activation may Strengthen Semantic Learning We investigate how partially activating high-uncertainty tokens affects robustness in noisy reward training by selectively applying language-aware losses. As shown in Table 13, activating 80% of high-entropy/perplexity tokens yields the best performance (42.41%), demonstrating that focusing supervision on uncertain regions strengthens the value model’s alignment with language semantics. Interestingly, fully activating all tokens (100%) slightly reduces performance (to 41.15%), likely due to the inclusion of noisy or uninformative tokens, which destabilize value es-

timation. This reveals a trade-off: while broader activation promotes semantic sensitivity, indiscriminate supervision can introduce harmful variance. Partial token activation offers a balanced solution that enhances robustness without compromising core value function stability.

5 Conclusion

In this study, we rethink the often-overlooked role of the value model in RL frameworks, particularly under noisy reward supervision in LLM post-training. We propose VRPO, a novel training framework that improves the robustness of RL by incorporating information flow filtering and enhancing language perception through guidance from a frozen language model in the value model. This approach effectively corrects reward biases, leading to more accurate advantage estimation. Experiments on math, science and dialogue tasks show that VRPO outperforms baseline methods such as PPO and GRPO, highlighting the untapped potential of the value model in robust RL training.

579 Limitations

580 In this section, we discuss the potential threats
581 to the validity of our method. VRPO does not
582 solve all forms of reward noise, but offers a prin-
583 cipled way to make the value model more robust
584 and semantically aware in the face of noisy su-
585 pervision. And due to limited computational re-
586 sources, the experiments with rule-based rewards
587 were only conducted on mathematics and science
588 tasks, where each model was trained and evalu-
589 ated on a single dataset. Moreover, the ablation
590 studies were primarily performed during test-time
591 optimization under rule-based reward settings. To
592 mitigate this potential threat, we further conducted
593 experiments across a wide range of datasets and
594 tasks, and performed ablation studies under rule-
595 based reward conditions during training-time opti-
596 mization as well. Although current results indicate
597 no significant differences between the two settings,
598 this broader evaluation helps ensure the robustness
599 of our conclusions.

600 References

601 Marc G. Bellemare, Salvatore Candido, Pablo Samuel
602 Castro, Jun Gong, Marlos C. Machado, Subhodeep
603 Moitra, Sameera S. Ponda, and Ziyun Wang. 2020.
604 [Autonomous navigation of stratospheric balloons us-
605 ing reinforcement learning.](#) *Nature*, 588:77 – 82.

606 Angelo Caregnato-Neto, Luciano Cavalcante Siebert,
607 Arkady Zgonnikov, Marcos Ricardo Omena de Al-
608 buquerque Maximo, and Rubens Junqueira Magal-
609 hães Afonso. 2024. [ARMCHAIR: integrated inverse
610 reinforcement learning and model predictive control
611 for human-robot collaboration.](#) *arXiv e-prints*,
612 arXiv:2402.19128.

613 Stephen Casper, Xander Davies, Claudia Shi, Thomas
614 Krendl Gilbert, Jérémy Scheurer, Javier Rando,
615 Rachel Freedman, Tomasz Korbak, David Lindner,
616 Pedro Freire, Tony Wang, Samuel Marks, Charbel-
617 Raphaël Segerie, Micah Carroll, Andi Peng, Phillip
618 Christoffersen, Mehul Damani, Stewart Slocum, Us-
619 man Anwar, and 13 others. 2023. [Open Prob-
620 lems and Fundamental Limitations of Reinforcement
621 Learning from Human Feedback.](#) *arXiv e-prints*,
622 arXiv:2307.15217.

623 Shreyas Chaudhari, Pranjal Aggarwal, Vishvak Mura-
624 hari, Tanmay Rajpurohit, Ashwin Kalyan, Karthik
625 Narasimhan, Ameet Deshpande, and Bruno Castro
626 da Silva. 2024. [Rlhf deciphered: A critical analysis
627 of reinforcement learning from human feedback for
628 llms.](#) *Preprint*, arXiv:2404.08555.

629 Yuanpei Chen, Tianhao Wu, Shengjie Wang, Xidong
630 Feng, Jiechuang Jiang, Stephen Marcus McAleer, Yi-
631 ran Geng, Hao Dong, Zongqing Lu, Song-Chun Zhu,

and Yaodong Yang. 2022. [Towards Human-Level Bi-
manual Dexterous Manipulation with Reinforcement
Learning.](#) *arXiv e-prints*, arXiv:2206.08686.

Jie Cheng, Gang Xiong, Xingyuan Dai, Qinghai
Miao, Yisheng Lv, and Fei-Yue Wang. 2024.
[RIME: Robust Preference-based Reinforcement
Learning with Noisy Preferences.](#) *arXiv e-prints*,
arXiv:2402.17257.

Ruoxi Cheng, Haoxuan Ma, Weixin Wang, Zhiqiang
Wang, Xiaoshuang Jia, Simeng Qin, Xiaochun Cao,
Yang Liu, and Xiaojun Jia. 2025. [Inverse Reinforce-
ment Learning with Dynamic Reward Scaling for
LLM Alignment.](#) *arXiv e-prints*, arXiv:2503.18991.

Bin Dai, Chen Zhu, and David Wipf. 2018. [Compress-
ing Neural Networks using the Variational Informa-
tion Bottleneck.](#) *arXiv e-prints*, arXiv:1802.10399.

Leo Gao, John Schulman, and Jacob Hilton. 2022. [Scal-
ing Laws for Reward Model Overoptimization.](#) *arXiv
e-prints*, arXiv:2210.10760.

Anirudh Goyal, Riashat Islam, Daniel Strouse, Zafar-
ali Ahmed, Matthew Botvinick, Hugo Larochelle,
Yoshua Bengio, and Sergey Levine. 2019. [InfoBot:
Transfer and Exploration via the Information Bottle-
neck.](#) *arXiv e-prints*, arXiv:1901.10902.

Danijar Hafner, Timothy Lillicrap, Jimmy Ba, and Mo-
hammad Norouzi. 2019. [Dream to Control: Learning
Behaviors by Latent Imagination.](#) *arXiv e-prints*,
arXiv:1912.01603.

Dan Hendrycks, Collin Burns, Saurav Kadavath, Akul
Arora, Steven Basart, Eric Tang, Dawn Song, and
Jacob Steinhardt. 2021. [Measuring mathematical
problem solving with the math dataset.](#) *Preprint*,
arXiv:2103.03874.

Yu-Heng Hung, Ping-Chun Hsieh, Akshay Mete, and
P. R. Kumar. 2023. [Value-biased maximum like-
lihood estimation for model-based reinforcement
learning in discounted linear mdps.](#) *Preprint*,
arXiv:2310.11515.

Borja Ibarz, Jan Leike, Tobias Pohlen, Geoffrey Irving,
Shane Legg, and Dario Amodei. 2018. [Reward learn-
ing from human preferences and demonstrations in
Atari.](#) *arXiv e-prints*, arXiv:1811.06521.

Cassidy Laidlaw, Shivam Singhal, and Anca Dragan.
2025. [Correlated proxies: A new definition and
improved mitigation for reward hacking.](#) *Preprint*,
arXiv:2403.03185.

Aitor Lewkowycz, Anders Andreassen, David Dohan,
Ethan Dyer, Henryk Michalewski, Vinay Ramasesh,
Ambrose Slone, Cem Anil, Imanol Schlag, Theo
Gutman-Solo, Yuhuai Wu, Behnam Neyshabur, Guy
Gur-Ari, and Vedant Misra. 2022. [Solving quan-
titative reasoning problems with language models.](#)
Preprint, arXiv:2206.14858.

685	Xize Liang, Chao Chen, Shuang Qiu, Jie Wang, Yue Wu, Zhihang Fu, Zhihao Shi, Feng Wu, and Jieping Ye. 2024. ROPO: Robust Preference Optimization for Large Language Models . <i>arXiv e-prints</i> , arXiv:2404.04102.	741
686		742
687		743
688		744
689		745
690	K. Michael Martini, Eslam Abdelaleem, and Ilya Nemenman. 2024. Deep Variational Multivariate Information Bottleneck. In <i>APS March Meeting Abstracts</i> , volume 2024 of <i>APS Meeting Abstracts</i> , page T28.009.	746
691		747
692		748
693		749
694		750
695	Yuchun Miao, Sen Zhang, Liang Ding, Rong Bao, Lefei Zhang, and Dacheng Tao. 2024. InfoRM: Mitigating Reward Hacking in RLHF via Information-Theoretic Reward Modeling . <i>arXiv e-prints</i> , arXiv:2402.09345.	751
696		752
697		753
698		754
699		755
700	Niklas Muennighoff, Zitong Yang, Weijia Shi, Xiang Lisa Li, Li Fei-Fei, Hannaneh Hajishirzi, Luke Zettlemoyer, Percy Liang, Emmanuel Candès, and Tatsunori Hashimoto. 2025. s1: Simple test-time scaling . <i>Preprint</i> , arXiv:2501.19393.	756
701		757
702		758
703		759
704		760
705	Julien Perolat, Bart De Vylder, Daniel Hennes, Eugene Tarassov, Florian Strub, Vincent de Boer, Paul Muller, Jerome T. Connor, Neil Burch, Thomas Anthony, Stephen McAleer, Romuald Elie, Sarah H. Cen, Zhe Wang, Audrunas Gruslys, Aleksandra Malysheva, Mina Khan, Sherjil Ozair, Finbarr Timbers, and 15 others. 2022. Mastering the game of Stratego with model-free multiagent reinforcement learning . <i>Science</i> , 378(6623):990–996.	761
706		762
707		763
708		764
709		765
710		766
711		767
712		768
713		769
714	Long Phan, Alice Gatti, Ziwen Han, Nathaniel Li, Josephina Hu, Hugh Zhang, Chen Bo Calvin Zhang, Mohamed Shaaban, John Ling, Sean Shi, Michael Choi, Anish Agrawal, Arnav Chopra, Adam Khoja, Ryan Kim, Richard Ren, Jason Hausenloy, Oliver Zhang, Mantas Mazeika, and 1090 others. 2025. Humanity’s last exam . <i>Preprint</i> , arXiv:2501.14249.	770
715		771
716		772
717		773
718		774
719		775
720		776
721	Riccardo Poiani, Gabriele Curti, Alberto Maria Metelli, and Marcello Restelli. 2024. Inverse Reinforcement Learning with Sub-optimal Experts . <i>arXiv e-prints</i> , arXiv:2401.03857.	777
722		778
723		779
724		780
725	Qwen, :, An Yang, Baosong Yang, Beichen Zhang, Binyuan Hui, Bo Zheng, Bowen Yu, Chengyuan Li, Dayiheng Liu, Fei Huang, Haoran Wei, Huan Lin, Jian Yang, Jianhong Tu, Jianwei Zhang, Jianxin Yang, Jiaxi Yang, Jingren Zhou, and 25 others. 2025. Qwen2.5 technical report . <i>Preprint</i> , arXiv:2412.15115.	781
726		782
727		783
728		784
729		785
730		786
731		787
732	John Schulman, Philipp Moritz, Sergey Levine, Michael Jordan, and Pieter Abbeel. 2018. High-dimensional continuous control using generalized advantage estimation . <i>Preprint</i> , arXiv:1506.02438.	788
733		789
734		790
735		791
736	Wei Sun, Wen Yang, Pu Jian, Qianlong Du, Fuwei Cui, Shuo Ren, and Jiajun Zhang. 2025. Ktae: A model-free algorithm to key-tokens advantage estimation in mathematical reasoning . <i>Preprint</i> , arXiv:2505.16826.	792
737		793
738		794
739		795
740		796
		797
	M-A-P Team, Xinrun Du, Yifan Yao, Kaijing Ma, Bingli Wang, Tianyu Zheng, Kang Zhu, Minghao Liu, Yiming Liang, Xiaolong Jin, Zhenlin Wei, Chujie Zheng, Kaixing Deng, Shuyue Guo, Shian Jia, Sichao Jiang, Yiyan Liao, Rui Li, Qinrui Li, and 76 others. 2025. Supergpqa: Scaling llm evaluation across 285 graduate disciplines . <i>Preprint</i> , arXiv:2502.14739.	798
		799
		800
		801
		802
		803
		804
		805
		806
		807
		808
		809
		810
		811
		812
		813
		814
		815
		816
		817
		818
		819
		820
		821
		822
		823
		824
		825
		826
		827
		828
		829
		830
		831
		832
		833
		834
		835
		836
		837
		838
		839
		840
		841
		842
		843
		844
		845
		846
		847
		848
		849
		850
		851
		852
		853
		854
		855
		856
		857
		858
		859
		860
		861
		862
		863
		864
		865
		866
		867
		868
		869
		870
		871
		872
		873
		874
		875
		876
		877
		878
		879
		880
		881
		882
		883
		884
		885
		886
		887
		888
		889
		890
		891
		892
		893
		894
		895
		896
		897
		898
		899
		900
		901
		902
		903
		904
		905
		906
		907
		908
		909
		910
		911
		912
		913
		914
		915
		916
		917
		918
		919
		920
		921
		922
		923
		924
		925
		926
		927
		928
		929
		930
		931
		932
		933
		934
		935
		936
		937
		938
		939
		940
		941
		942
		943
		944
		945
		946
		947
		948
		949
		950
		951
		952
		953
		954
		955
		956
		957
		958
		959
		960
		961
		962
		963
		964
		965
		966
		967
		968
		969
		970
		971
		972
		973
		974
		975
		976
		977
		978
		979
		980
		981
		982
		983
		984
		985
		986
		987
		988
		989
		990
		991
		992
		993
		994
		995
		996
		997
		998
		999
		1000

798 Ge, Haoran Wei, Huan Lin, Jialong Tang, and 41
 799 others. 2025. [Qwen3 technical report](#). *Preprint*,
 800 arXiv:2505.09388.

801 Tongtian Yue, Jie Cheng, Longteng Guo, Xingyuan
 802 Dai, Zijia Zhao, Xingjian He, Gang Xiong, Yisheng
 803 Lv, and Jing Liu. 2024. [SC-Tune: Unleashing
 804 Self-Consistent Referential Comprehension in
 805 Large Vision Language Models](#). *arXiv e-prints*,
 806 arXiv:2403.13263.

807 Jiazheng Zhang, Wenqing Jing, Zizhuo Zhang, Zhiheng
 808 Xi, Shihan Dou, Rongxiang Weng, Jiahuan Li, Jin-
 809 gang Wang, Mingxu Chai, Shibo Hong, Tao Gui,
 810 and Qi Zhang. 2025. [Two Minds Better Than One:
 811 Collaborative Reward Modeling for LLM Alignment](#).
 812 *arXiv e-prints*, arXiv:2505.10597.

813 Ming Zhou, Jun Luo, Julian Villella, Yaodong Yang,
 814 David Rusu, Jiayu Miao, Weinan Zhang, Mont-
 815 gomery Alban, Iman Fadarar, Zheng Chen, Aurora
 816 Chongxi Huang, Ying Wen, Kimia Hassanzadeh,
 817 Daniel Graves, Dong Chen, Zhengbang Zhu, Nhat
 818 Nguyen, Mohamed Elsayed, Kun Shao, and 18 oth-
 819 ers. 2020. [SMARTS: Scalable Multi-Agent Rein-
 820 forcement Learning Training School for Autonomous
 821 Driving](#). *arXiv e-prints*, arXiv:2010.09776.

822 A Additional Details for VRPO

823 A.1 Pseudocode

824 The full algorithm of VRPO is detailed in Algo-
 825 rithm 1.

826 B Additional Experimental Details

827 B.1 Setup

828 B.1.1 Dataset Construction

829 **Rule-based Setting.** For **test-time augmenta-**
 830 **tion**, we use four math benchmarks: AIME 2024,
 831 AMC 2023, Minerva-Math, and MATH-500, as
 832 well as three QA datasets: GPQA-Extended, Sim-
 833 pleQA, and Humanity’s Last Exam (HLE).

834 For **training-time augmentation**, as an illus-
 835 trative example based on the Qwen3-8B model,
 836 pseudo-labels for mathematical tasks are gener-
 837 ated from 39,000 samples in the Light-R1 (Wen
 838 et al., 2025) dataset. After 5 rounds of majority vot-
 839 ing, samples with at least 3 identical responses are
 840 retained, yielding 31,209 instances for reinforc-
 841 ement learning training. For scientific tasks, pseudo-
 842 labels are generated from 26,529 samples in the Su-
 843 perGPQA dataset. After the same 5-round majority
 844 voting procedure, 10,075 samples with at least 2
 845 consistent responses are preserved for RL train-
 846 ing. Evaluation is conducted on non-overlapping
 847 mathematical benchmarks to ensure fairness.

Algorithm 1: VRPO training process

Require: Dataset $\mathcal{D} = \{(x_t, a_t, r_t)\}_{t=1}^T$, policy
 π_θ , decoder model q_ψ , encoder model f_ϕ ,

bottleneck prior $r(z)$; coefficients $\beta, \lambda_{\text{ent}}, \lambda_{\text{ppl}}$

Ensure: Actor model π_θ , value model f_ϕ, q_ψ

for each VRPO iteration **do**

Sample trajectories using current policy π_θ

Compute rewards $\{r_t\}$ and actions $\{a_t\}$

Compute policy ratio: $\rho_t = \frac{\pi_\theta(a_t|x_t)}{\pi_{\theta_{\text{old}}}(a_t|x_t)}$

for each timestep t **do**

Encode latent: sample $\epsilon \sim \mathcal{N}(0, I)$,

compute $z_t = f_\phi^\mu(x_t) + f_\phi^\sigma(x_t) \cdot \epsilon$

Predict value: $\hat{V}_t = q_\psi(z_t)$

end for

Compute GAE advantage:

$A_t = \text{GAE}(r_t, \hat{V}_t, \hat{V}_{t+1})$

Compute returns: $R_t = A_t + \hat{V}_t$

Value Model Update

// Value Loss with IB Structure

Compute value loss:

$\mathcal{L}_{\text{MSE}} = \frac{1}{T} \sum_t (q_\psi(f_\phi(x_t)) - R_t)^2$

Compute KL loss:

$\mathcal{L}_{\text{KL}} = \frac{1}{T} \sum_t \text{KL}(p_\phi(z_t|x_t) || r(z))$

// Semantic Alignment via Frozen LM Head

for each timestep t **do**

Get token distribution $P_V(y_t|x_t)$ from
frozen LM head

Identify high-uncertainty tokens:

$T_{\text{ent}} = \{t : H[P_V(y_t|x_t)] > \hat{T}_{\text{entropy}}\}$

$T_{\text{ppl}} = \{t : -\log P_V(y_t = y_t^*|x_t) >$

$\hat{T}_{\text{perplexity}}\}$

Compute entropy loss:

$L_{\text{ent}} = \sum_{t \in T_{\text{ent}}} H[P_V(y_t|x_t)]$

Compute perplexity loss:

$L_{\text{ppl}} = \sum_{t \in T_{\text{ppl}}} -\log P_V(y_t = y_t^*|x_t)$

end for

Compute semantic loss:

$\mathcal{L}_{\text{sem}} = \lambda_{\text{ent}} \cdot L_{\text{ent}} + \lambda_{\text{ppl}} \cdot L_{\text{ppl}}$

Combine value loss:

$\mathcal{L}_{\text{value}} = \mathcal{L}_{\text{MSE}} + \beta \cdot \mathcal{L}_{\text{KL}} + \mathcal{L}_{\text{sem}}$

Update value model $\{\phi, \psi\}$ using $\nabla \mathcal{L}_{\text{value}}$

Actor Model Update

Compute clipped PPO objective: $\mathcal{L}_{\text{PPO}} =$

$\mathbf{E}_t [\min(\rho_t A_t, \text{clip}(\rho_t, 1 - \epsilon, 1 + \epsilon) A_t)]$

Update actor model π_θ , using $\nabla \mathcal{L}_{\text{PPO}}$

end for

Scenario	Wealth	Rental	Insurance	Food	Express	Promotion	Loan	Housing	Service	Product	General	Avg
Dialogue Count	87	99	138	120	215	66	70	87	67	69	92	94.73
Avg Turns	5.40	4.22	5.22	3.37	3.76	4.55	4.33	4.68	5.46	4.44	3.46	4.44

Table 4: Dataset statistics across 11 real-world service scenarios. The validation set covers diverse interaction types with varying dialogue lengths, enabling reliable evaluation of robustness and generalization under realistic settings.

Category	Human Annotators				GPT Models		Human AVG	Model AVG
Score	4.595	4.689	4.490	4.330	4.510	4.501	4.526	4.505

Table 5: Comparison of goal completion rate scores between human annotators and GPT-4o. Human annotators achieve an average score of 4.526, while GPT-4o models reach 4.505, showing a negligible performance gap.

Model-based Setting. A real-world dialogue dataset (Honor-Dialogue) is used. The reward model is trained on 36,000 labeled samples, with 3,000 for validation. The policy model is fine-tuned and trained on 80,000 additional conversations, with 8,000 for evaluation. This dataset contains various dialogue tasks in real-life scenarios. It truly reflects the actual conditions of real situations. Such multi-category real-scenario dialogue task data is not available in other datasets. All real information in the data has undergone anonymization, and the use of the data has been authorized and reviewed by the ethics committee.

Honor-Dialogue dataset is constructed based on a goal-driven scenario-oriented design paradigm. For the data within each domain, we construct realistic caller inputs that include scenario-matched latest messages and conversation history, as well as corresponding standardized outputs that comply with the requirements of the target goal. We explicitly mark the dialogue state, response content, and matched target to ensure the quality of the supervised training data. A representative example of the Honor-Dialogue dataset is presented below. (Figure 9) The notation xxx denotes the masking of sensitive information such as user ID and contact details, which is implemented to comply with data privacy regulations.

B.1.2 Baseline Initialization

Rule-based Setting. Test-time augmentation: Qwen2.5-7B-Base (math task) and Qwen2.5-7B-Instruct (QA task), fine-tuned on OpenR1-Math and S1.1 (Muennighoff et al., 2025), respectively. **Training-time augmentation:** Qwen2.5-7B-Cold Start, Qwen3-1.7B, Qwen3-8B and LLaMA 3.1-8B-Instruct are used as baselines. The Qwen2.5-7B-Cold Start model has the same model settings as those in the Test-time augmentation.

Model-based Setting. Both the policy and reward models are initialized from Qwen3-8B, fine-tuned on the Honor-Dialogue dataset.

B.1.3 Training Configuration

Training parameters. In the majority voting during testing, 32 samples are used, while 5 samples are used for voting during training. The weights for the auxiliary losses (entropy and perplexity) are set to 0.5, making 80% of the labels effective. The RL training runs for 1 iteration and the RL setup adopts a training batch size of 128 and a rollout batch size of 128, with actor and critic learning rates of 5×10^{-7} and 5×10^{-6} . During both training and testing, the model generates dialogue tasks of length 4096 for reasoning tasks and 1024 for dialogue tasks. The experiments are conducted on 8 * NVIDIA A100 80GB GPUs.

B.1.4 Dialogue Task Evaluation

Our dialogue evaluation leverages GPT-4o, but the model is not utilized as a free-form judge. Instead, it implements a rigorous rubric-based evaluation protocol designed to mitigate subjectivity and enhance reproducibility. This rubric explicitly defines evaluation metrics, adopts a 1–5 scoring scale with clear criteria, includes two scoring examples and a standardized output format, and the prompt incorporates a well-structured process along with comprehensive dialogue content and contextual information. GPT-4o functions solely as an automated evaluator applying this fixed rubric, rather than an unconstrained scorer.

To further ensure reliability, we conducted cross-validation through two key steps: **1)** multiple sampling runs to verify the stability of rubric execution, and **2)** spot-checked human evaluations that demonstrated high agreement with rubric-based scores. Specifically, for the comparative experiment between human annotators and the model,

we validated the protocol using 1,110 data samples covering 10 scenarios. The statistical details of the validation dataset are presented in Table 4. Four independent annotators scored these samples strictly in accordance with the rubric. They are professional data annotators in the company. Table 5 presents the experimental results for goal completion rate scoring (full score: 5) in dialogue tasks. As shown in Table 5, the mean score assigned by human annotators is extremely close to that of the model, with a difference of only 0.021 points. These findings confirm strong consistency between human and model scores, validating the effectiveness of our rubric. We acknowledge that fully establishing external validity necessitates additional independent annotators. We are in the process of releasing our rubric, evaluation prompts, and real-dialogue dataset to facilitate replication. To illustrate the evaluation logic for dialogue performance, the core prompt section regarding the assessment of conversation logicity is provided as follows in Figure 10, which is a part of goal completion rate.

B.2 Additional experimental details in the Dialogue Task

To further analyze the behavior of different methods in the dialogue task, Table 6 presents the full training trajectory across various steps. A notable observation is that both PPO and GRPO experience severe performance degradation during training, likely due to reward over-optimization or instability, resulting in significantly lower final scores than both the initial model and our method. Specifically, their final average scores drop to 40.13% and 36.17%, respectively.

In contrast, our proposed method (VRPO) maintains high robustness throughout training. While the final performance slightly decreases compared to the cold-start model (from 85.87% to 83.80%), our method avoids collapse and achieves high task performance in early training stages. For example, at step 50 and 100, VRPO reaches 75.90% task completion rate (TCR) and an overall average (AVG) score of 88.13%, outperforming all baselines at these stages.

These results suggest that VRPO effectively preserves model capabilities under heavy noisy supervision and raises the lower bound of post-training performance. This stabilizing effect makes it significantly less prone to reward hacking or collapse compared to standard PPO approaches.

Step	TCR	ACR	GCR	AVG
<i>Cold Start</i>				
0	72.10%	94.30%	91.20%	85.87%
<i>GRPO</i>				
50	69.40%	91.90%	89.60%	83.63%
100	35.80%	50.00%	58.50%	48.10%
150	40.90%	35.00%	41.40%	39.10%
200	41.70%	30.20%	36.90%	36.27%
Last	40.70%	30.10%	37.70%	36.17%
<i>PPO</i>				
50	75.70%	95.40%	93.40%	88.17%
100	65.90%	90.00%	88.30%	81.40%
150	36.00%	58.90%	62.30%	52.40%
200	37.50%	50.90%	55.80%	48.07%
Last	39.10%	38.60%	42.70%	40.13%
<i>Reward model improvement method base PPO</i>				
50	71.90%	92.80%	91.70%	85.47%
100	65.30%	88.80%	88.90%	81.00%
150	37.50%	64.30%	64.80%	55.53%
200	35.70%	56.20%	57.40%	49.77%
Last	34.80%	54.30%	53.90%	47.67%
<i>OURS</i>				
50	75.60%	94.70%	94.10%	88.13%
100	75.90%	94.50%	93.80%	88.07%
150	72.40%	91.70%	91.90%	85.33%
200	70.00%	89.40%	91.90%	83.77%
Last	68.90%	91.30%	91.20%	83.80%

Table 6: Trajectory-level generalization performance across different methods and training steps. Our method maintains high performance even in last stages, beyond GRPO, PPO, and the reward model improvement method based on PPO.

B.3 Additional experimental details in the Math and Science Task

Comparison with Recent Robust Advantage-Based RL Methods We further compare VRPO with recent and conceptually related robust RL methods that explicitly target advantage shaping or robustness, including KTAE and λ -GRPO. While all these approaches aim to mitigate noisy or biased advantages, they differ fundamentally in how token-level credit assignment is addressed.

KTAE adopts an explicit keyword-oriented strategy, directly identifying salient tokens and reweighting advantages based on rule-based key-token selection. In contrast, VRPO does not rely on predefined or externally extracted keywords. Instead, it leverages semantic guidance through entropy- and perplexity-based losses applied to the value model, encouraging advantage estimation to implicitly attend to semantically informative tokens via representation learning. This enables VRPO to capture task-relevant signals in a softer, model-driven manner rather than through explicit token

Domain	Math				Factuality	Science	Knowledge	ALL
Method	MATH500	AIME24	Minerva-Math	AMC23	SampleQA	GPQA	HLE	AVG
<i>Qwen2.5-7B-Cold Start(Weak Model)</i>								
Base	71.60%	6.67%	18.75%	52.50%	2.36%	1.45%	3.29%	22.37%
KTAE	68.60%	6.67%	19.12%	48.33%	1.73%	2.36%	3.24%	21.44%
λ -GRPO	74.00%	10.00%	18.38%	56.67%	2.94%	3.26%	3.34%	24.08%
Ours	72.20%	23.33%	19.85%	50.00%	2.82%	3.44%	3.94%	25.08%
<i>Qwen3-8B(Strong Model)</i>								
Base	87.40%	41.67%	28.68%	75.83%	2.89%	3.10%	2.89%	34.64%
KTAE	86.00%	38.33%	29.41%	78.33%	3.17%	3.99%	2.64%	34.55%
λ -GRPO	90.00%	50.00%	31.62%	80.00%	2.77%	4.35%	3.57%	37.47%
Ours	90.20%	46.67%	31.99%	86.67%	3.21%	4.35%	3.89%	38.14%

Table 7: Comparative accuracy (%) of train-time optimization across mathematical and scientific reasoning benchmarks with recent advantage-based robust RL methods.

masking or selection.

Similarly, λ -GRPO introduces robustness by explicitly penalizing response length through a length-dependent weighting in the softmax computation, thereby discouraging degenerate long responses. VRPO complements this approach by strengthening the value model itself: through the information bottleneck and semantic regularization, the value model learns to better perceive and encode the semantic structure of generated text, allowing advantage estimates to adaptively reflect the quality and relevance of content rather than relying on a length-based prior.

Empirically, these conceptual differences translate into consistent performance gains. Under the same train-time optimization setting described in Section 4.1, VRPO achieves the highest average accuracy across both Qwen-2.5-7B Cold-Start and Qwen-3-8B models (Table 7). In the Qwen-2.5-7B setting, VRPO improves the overall average accuracy to 25.08%, outperforming KTAE (21.44%) and λ -GRPO (24.08%), with gains observed across both mathematical and scientific benchmarks. Similar trends persist for the stronger Qwen-3-8B model, where VRPO again attains the best overall performance (38.14%), surpassing both KTAE and λ -GRPO.

Taken together, these results suggest that VRPO offers a complementary robustness mechanism: rather than explicitly constraining advantages via explicit token rules or length penalties, it enhances the semantic sensitivity and noise-filtering capacity of the value model itself, leading to more reliable advantage estimation under noisy supervision.

VRPO consistently improves reasoning performance under rule-based noisy rewards To fur-

ther evaluate the effectiveness of our approach in rule-based reward scenarios, we conducted a series of experiments covering both mathematical and factual reasoning domains in test-time inference settings. Specifically, we trained and evaluated the model using seven datasets. For each dataset, we performed multiple inferences with the original model and conducted voting-based result filtering, then trained on the filtered data before evaluating the model on the respective datasets. The datasets include four math-focused tasks: MATH500 (Hendrycks et al., 2021), AIME24, Minerva-Math (Lewkowycz et al., 2022), AMC23, and three factual knowledge tasks: SampleQA (Wei et al., 2024), GPQA, and HLE (Humanity’s Last Exam). Table 8 presents accuracy after reasoning-enhanced training using different RL methods.

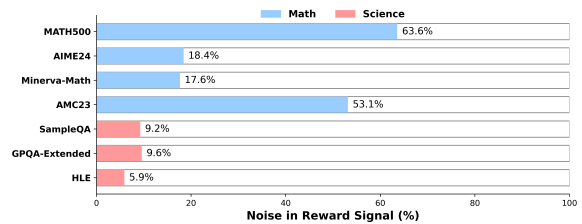


Figure 6: Noise statistics in the math and science tasks with rule-based reward supervision during test-time optimization.

Despite noisy feedback, our value model learns effectively. As shown in Figure 7, the explained variance of our method increases steadily during training, indicating closer alignment between predicted and GAE-derived returns. This trend is accompanied by a consistent drop in reward prediction error and reduction in high-variance outliers. As shown in Table 8, our method consis-

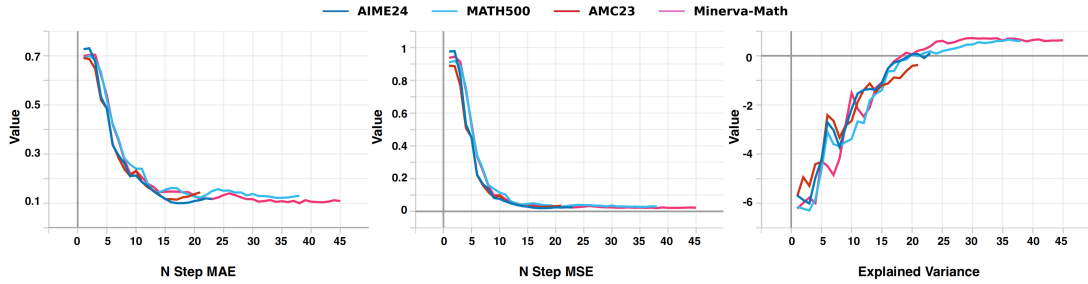


Figure 7: Prediction error of the value model across training steps

Domain	Math				Factuality	Science	Knowledge	ALL
	MATH500	AIME24	Minerva-Math	AMC23	SampleQA	GPQA	HLE	AVG
Cold Start	71.60%	6.67%	18.75%	52.50%	1.62%	2.53%	2.55%	22.32%
GRPO	73.40%	11.67%	19.49%	56.67%	2.47%	3.08%	3.80%	24.37%
PPO	73.60%	10.00%	19.49%	55.00%	2.43%	2.54%	3.56%	23.80%
Ours	74.40%	13.33%	20.22%	61.67%	2.61%	4.17%	4.45%	25.84%

Table 8: Accuracy (%) on test-time optimization with rule-based rewards. Our method achieves gains across various datasets and domains.

tently outperforms other baseline methods in various domains by leveraging a strong value model with semantic-aware filtering. This improvement is evident in both mathematical and knowledge reasoning tasks. For instance, accuracy on AIME24 increases from 6.67% to 13.33%, and on HLE from 2.55% to 4.45%. These results highlight the model’s ability to extract relevant information from noisy or ambiguous feedback and demonstrate its robust generalization across diverse domains.

B.4 Additional Ablation Experiment

Method	MATH 500	AIME 24	Minerva Math	AMC 23	AVG
Base	49.80%	10.00%	17.28%	24.17%	25.31%
GRPO	49.20%	6.67%	16.18%	27.50%	24.89%
PPO	46.20%	6.67%	18.01%	25.00%	23.97%
Ours	49.40%	10.00%	17.65%	27.50%	26.14%

Table 9: Test-time optimization results on mathematical reasoning benchmarks using LLaMA 3.1-8B-Instruct under noisy rule-based reward supervision. Our method is the only one that improves upon the base model, demonstrating superior robustness under noisy feedback.

VRPO effectively mitigates noise in mathematical tasks when applied to the LLaMA model series Recent studies(Wu et al., 2025) suggest that Qwen models may be less reliable for mathematical reasoning. Therefore, we design experiments using the LLaMA 3.1-8B-Instruct model. Tables 9

Method	MATH 500	AIME 24	Minerva Math	AMC 23	AVG
Base	49.80%	10.00%	17.28%	24.17%	25.31%
GRPO	42.40%	3.33%	17.28%	32.50%	23.88%
PPO	41.00%	6.67%	17.65%	25.83%	22.79%
Ours	41.40%	10.00%	18.01%	30.00%	24.85%

Table 10: Training-time optimization results on the math task using LLaMA 3.1-8B-Instruct under noisy reward conditions. Due to weak base performance and high noise, all methods show a performance drop, but ours degrades the least and even improves results on Minerva-Math and AMC23.

and 10 present its performance on four mathematical reasoning benchmarks under rule-based reward supervision. Specifically, Table 9 reports results in the test-time optimization setting, while Table 10 shows training-time optimization outcomes.

In the test-time optimization setting, only our proposed method achieves improvements over the base model, whereas PPO and GRPO both suffer performance degradation, underscoring VRPO’s robustness under noisy supervision. In the training-time optimization setting, due to limited model capacity and the prevalence of noisy samples (with overall data accuracy of only 42.33%), all methods experience performance drops. Nevertheless, our method exhibits the smallest decline and even achieves gains on Minerva-Math and AMC23, demonstrating stronger resilience in noisy training environments. These findings further highlight the generalizability of VRPO across different model

1095 architectures, confirming its effectiveness beyond
1096 Qwen models.

1097 **Effect of Variational Bottleneck on Noise Filter-** 1098 **ing**

1099 As shown in Table 11, we compare different
1100 information bottleneck designs across tasks and ob-
1101 tain three takeaways: (1) A single-layer bottleneck
1102 delivers the best overall performance, with 4.17%
1103 and 4.45% accuracy on GPQA and HLE, indicating
1104 strong noise-filtering capability. (2) A two-layer
1105 bottleneck excels on high-complexity math tasks
1106 but underperforms on others, likely due to overfit-
1107 ting to noise. However, for language models like
1108 Qwen, which are pretrained on rich mathematical
1109 data, deeper value structures may facilitate deeper
1110 understanding, yielding benefits in math domains.
1111 (3) No bottleneck yields modest gains on low-noise
1112 datasets but collapses on noisy ones (e.g., AMC23
1113 drops to 50.83%), showing high susceptibility to
1114 noisy or uncertain rewards. Overall, moderate ar-
1115 chitectural bottlenecks substantially improve the
1116 robustness and cross-task generalization of value
modeling under noisy RL Training.

1117 **VRPO’s two components interact synergistically**
1118 **rather than redundantly** To examine whether
1119 the two core components of VRPO contribute
1120 redundantly or synergistically, we conducted a
1121 component-level ablation study. The results are
1122 reported in Table 13 and Table 11, covering three
1123 representative settings: semantic-only training, IB-
1124 only training, and the full VRPO method that inte-
1125 grates both components.

1126 Under the zero-bottleneck configuration, the
1127 model disables the information bottleneck and re-
1128 tains only the entropy/perplexity-based semantic
1129 loss, forming a semantic-only condition. This vari-
1130 ant already outperforms the cold-start baseline on
1131 most metrics, indicating that semantic guidance
1132 alone can improve robustness under noisy supervi-
1133 sion. Conversely, when using a 0% token activation
1134 ratio, the semantic loss is effectively removed while
1135 the information bottleneck remains active, result-
1136 ing in an IB-only condition. This setting also yields
1137 consistent improvements over the base model.

1138 The fact that both isolated variants surpass the
1139 cold-start baseline demonstrates that each com-
1140 ponent independently enhances reasoning robust-
1141 ness. More importantly, the full VRPO configura-
1142 tion, which jointly combines semantic regulariza-
1143 tion with the information bottleneck, consistently
1144 achieves the best overall performance across both
1145 mathematical and scientific benchmarks. These

1146 results indicate that the two components interact
1147 synergistically rather than redundantly, jointly con-
1148 tributing to more stable and reliable advantage esti-
1149 mation.

1150 **VRPO remains effective across models of dif-** 1151 **ferent scales**

1152 Table 12 presents the results of the
1153 Qwen3-1.7B model on mathematical and scientific
1154 tasks under rule-based reward supervision during
1155 training-time optimization. The results demon-
1156 strate that, under noisy supervision, our method
1157 consistently outperforms baseline approaches such
1158 as PPO and GRPO, achieving an average accuracy
1159 of 58.37% on mathematical tasks and 2.81% on
1160 scientific tasks. These findings confirm the effec-
1161 tiveness of VRPO across model scales, and further
1162 highlight its robustness and generalization ability
even in small-model settings.

1163 **Partial Token Activation Ratio Analysis and** 1164 **Self-Coupling Concern Verification**

1165 To study
1166 how partial token activation influences semantic
1167 learning, we conducted an extensive activation per-
1168 centage sweep on mathematical reasoning tasks.
1169 The results are summarized in Table 13.

1170 We observe that optimal performance does not
1171 come from a single sharply tuned threshold, but
1172 from a broader 80–85% activation interval, indicat-
1173 ing the phenomenon is a robust region rather than a
1174 thresholding artifact. Performance varies smoothly
1175 outside this interval and remains relatively stable.

1176 To further address potential concerns regarding
1177 self-coupling, where entropy or perplexity is esti-
1178 mated by models at different stages of the same
1179 training system, we conducted an additional ex-
1180 periment using an external language model (Qwen-
1181 2.5-7B-Instruct) to compute token-level uncertainty
1182 scores. As shown in Table 14, the same activation
1183 pattern emerges: activating approximately 80% to-
1184 kens again achieves the best average performance.
1185 The consistency of this trend across both internal
1186 and external uncertainty estimators supports the in-
1187 terpretation that the observed improvement reflects
1188 a generalizable property of partial token activation,
1189 rather than a model-specific artifact or self-induced
coupling effect.

1190 **Improved Key Words Attribution in Advantage** 1191 **Estimation with VRPO**

1192 To move beyond quali-
1193 tative visualizations and provide a more quantita-
1194 tive assessment of token-level credit assignment,
1195 we conduct an additional analysis comparing how
different methods identify key tokens when esti-

Domain	Math				Factuality	Science	Knowledge
	MATH500	AIME24	Minerva-Math	AMC23	SampleQA	GPQA	HLE
Zero bottleneck(semantic-only)	70.20%	16.67%	19.49%	50.83%	2.47%	2.17%	3.80%
Two bottleneck	73.20%	16.67%	21.32%	60.00%	2.45%	3.44%	3.62%
One bottleneck(Ours)	74.40%	13.33%	20.22%	61.67%	2.61%	4.17%	4.45%

Table 11: Performance comparison across domains and datasets using different bottleneck strategies. Semantic-only indicates that the experiment uses only the configuration with semantic loss improvements.

Domain	Math				Factuality	Science	Knowledge	ALL
	MATH500	AIME24	Minerva-Math	AMC23	SampleQA	GPQA	HLE	AVG
<i>Qwen3-1.7B</i>								
Base	84.80%	35.00%	19.49%	69.17%	1.78%	2.36%	2.87%	30.78%
GRPO	84.20%	33.33%	22.79%	79.17%	1.60%	1.99%	3.94%	32.43%
PPO	83.80%	35.00%	22.06%	78.33%	1.91%	2.54%	3.94%	32.51%
Reinforce++	84.40%	41.67%	23.53%	75.83%	1.48%	2.36%	3.52%	33.26%
Dr.GRPO	83.40%	36.67%	22.43%	75.83%	1.50%	2.36%	3.10%	32.19%
KTAE	64.40%	18.33%	18.01%	56.67%	1.64%	2.90%	4.36%	23.76%
λ -GRPO	82.80%	40.00%	21.69%	75.83%	1.62%	2.90%	3.89%	32.68%
Ours	85.40%	46.67%	23.90%	77.50%	1.50%	2.72%	4.22%	34.56%
<i>Qwen3-8B</i>								
Base	87.40%	41.67%	28.68%	75.83%	2.89%	3.10%	2.89%	34.64%
GRPO	89.20%	35.00%	28.68%	84.17%	3.03%	2.98%	3.24%	35.19%
PPO	86.40%	36.67%	30.51%	80.00%	2.82%	2.17%	3.29%	34.55%
Reinforce++	89.00%	45.00%	27.21%	84.17%	3.19%	4.35%	3.29%	36.60%
Dr.GRPO	87.80%	45.00%	27.94%	85.83%	2.50%	3.99%	3.10%	36.59%
KTAE	86.00%	38.33%	29.41%	78.33%	3.17%	3.99%	2.64%	34.55%
λ -GRPO	90.00%	53.33%	31.62%	80.00%	2.77%	4.35%	3.57%	37.95%
Ours	90.20%	46.67%	31.99%	86.67%	3.21%	4.35%	3.89%	38.14%

Table 12: Accuracy (%) on train-time optimization with rule-based rewards across models of different scales. Our method achieves consistent improvements under both Qwen3-1.7B and Qwen3-8B models.

Activ- ation	MATH 500	AIME 24	Minerva Math	AMC 23	AVG
0%	72.20%	10.00%	20.22%	51.67%	38.52%
20%	75.00%	6.67%	18.38%	54.17%	38.56%
50%	73.40%	13.33%	19.49%	55.83%	40.51%
70%	72.80%	10.00%	20.22%	58.33%	40.33%
75%	73.20%	13.33%	19.49%	59.17%	41.29%
80% (Ours)	74.40%	13.33%	20.22%	61.67%	42.41%
85%	73.60%	13.33%	23.90%	58.33%	42.29%
90%	71.00%	13.33%	21.69%	59.17%	41.29%
100%	73.80%	16.67%	19.12%	55.00%	41.15%

Table 13: Impact of partial token activation on math reasoning performance using entropy/perplexity-based semantic loss. Activating approximately 80-85% of tokens yields the best overall accuracy.

mating advantages under noisy training conditions. Our goal is to evaluate whether VRPO’s advantage estimates better align with semantically meaningful tokens in the output text.

As an approximate human-interpretable reference, we employ GPT-4o to extract salient keywords from each generated answer, which serve as a proxy for ground-truth key tokens. For each model, we then identify the top-20 tokens with the highest estimated advantages and compute their

Activ- ation	MATH 500	AIME 24	Minerva Math	AMC 23	AVG
50%	70.40%	10.00%	19.85%	53.33%	38.40%
80%	72.80%	16.67%	20.96%	55.00%	41.36%
100%	73.40%	6.67%	18.01%	61.67%	39.94%

Table 14: Impact of partial token activation when entropy/perplexity scores are computed by an external model (Qwen-2.5-7B-Instruct). The optimal activation region remains consistent with the original setting.

overlap with the extracted keywords. We report the **Key Token Overlap (KTO)** score, defined as:

$$\text{KTO} = \frac{|T_{\text{model}} \cap T_{\text{gold}}|}{|T_{\text{gold}}|}, \quad (21)$$

where T_{model} denotes the set of top-advantage tokens produced by the model and T_{gold} denotes the keyword set extracted by GPT-4o.

We evaluate this metric on both mathematical and scientific reasoning benchmarks. Specifically, we sample 50 instances per dataset, except for AIME24 (30 samples) and AMC23 (40 samples), where all available data are used. Due to computational constraints, we focus on a direct comparison

Domain	Math				Factuality	Science	Knowledge	ALL
Method	MATH500	AIME24	Minerva-Math	AMC23	SampleQA	GPQA	HLE	AVG
PPO	21.10%	22.83%	27.70%	19.13%	30.80%	24.55%	22.00%	24.02%
Ours	27.30%	30.33%	30.80%	26.25%	26.20%	26.36%	20.20%	26.78%

Table 15: Top 20 Key Token Overlap (KTO) comparison between PPO and VRPO under noisy test-time optimization. Higher values indicate stronger alignment between high-advantage tokens and semantically meaningful keywords.

1218 between PPO and VRPO under noisy train-time
1219 optimization settings.

1220 The results in Table 15 show that VRPO achieves
1221 higher KTO scores than PPO, yielding an average
1222 improvement of over 2.76 points. This indicates
1223 that VRPO’s advantage estimates are more strongly
1224 aligned with semantically meaningful tokens, sug-
1225 gesting more faithful credit assignment. We also
1226 observe that high-advantage tokens sometimes in-
1227 clude punctuation (e.g., “?”, “!”, “.”), likely be-
1228 cause such tokens carry high attention weights and
1229 mark semantic boundaries in transformer models.
1230 Since these tokens are not included in GPT-4o’s
1231 keyword lists, they tend to lower absolute KTO
1232 values for all methods.

1233 Overall, despite the heuristic nature of the refer-
1234 ence signal, this quantitative analysis supports our
1235 core claim: by integrating semantic regularization
1236 with an information bottleneck in the value model,
1237 VRPO captures key tokens more effectively and
1238 produces more robust advantage signals, which in
1239 turn leads to improved training stability and per-
1240 formance under noisy supervision.

1241 B.5 Additional Sensitivity Analysis

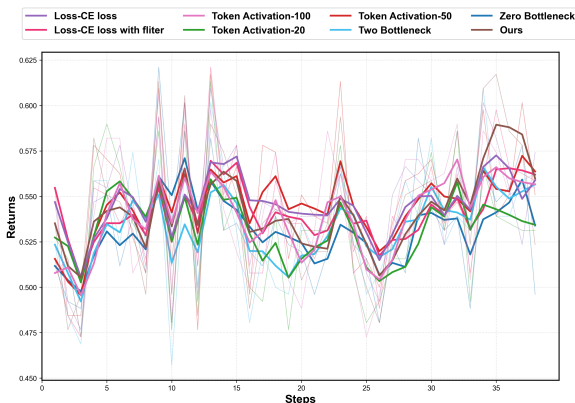


Figure 8: The evolution of training returns under test-time optimization settings on math datasets across different token-selection thresholds, learning strategies, and bottleneck configurations.

1242 We further analyze the training sensitivity of

1243 VRPO. Due to resource constraints, we take the
1244 experiments conducted under the settings of test-
1245 time optimization on math datasets with varying
1246 token-selection thresholds, learning strategies, and
1247 bottleneck configurations as an example for analy-
1248 sis, as illustrated in Figure 8.

1249 Across all examined settings, the training curves
1250 exhibit smooth and low-variance behavior during
1251 the early stages, accompanied by a consistent up-
1252 ward trend. This indicates stable optimization dy-
1253 namics and suggests that moderate variations in
1254 hyperparameters do not disrupt training. While the
1255 performance differences across configurations re-
1256 main relatively small, our proposed setting achieves
1257 slightly higher returns throughout the training pro-
1258 cess with notably high returns observed particularly
1259 in the late stages of training.

1260 These observations confirm that VRPO is not
1261 overly sensitive to hyperparameter tuning. Instead,
1262 its core components interact in a stable and well-
1263 behaved manner, providing mild yet consistent ad-
1264 vantages across different configurations. This ro-
1265 bustness further supports the practicality of VRPO
1266 in noisy reinforcement learning scenarios.

1267 C Additional Visualization

1268 Figures 12, 13, 14, and 15 present a comparison
1269 of advantage estimations between PPO and our
1270 method (VRPO) on the same question outputs re-
1271 gardless of semantic quality. Furthermore, in the
1272 incorrect answer case (Figure 11), under both cor-
1273 rect and incorrect answer conditions.

1274 Our method demonstrates a clear ability to cap-
1275 ture key semantic components. In the correct
1276 answer case (Figure 12), it effectively identifies mean-
1277 ingful concepts such as "massive" and "initial mass
1278 M", while also showing a higher advantage esti-
1279 mate for the correct answer (Figure 12) compared
1280 to the incorrect one (Figure 13). The model not
1281 only captures answer structures like box-based out-
1282 puts but also makes more accurate judgments based
1283 on correctness.

1284 In contrast, PPO fails to capture such distinc-

1285 tions. As seen in Figure 14, the advantage esti-
1286 mation for the correct answer collapses, with later
1287 tokens receiving increasing advantage values. This
1288 pattern can lead to length hacking, where the model
1289 learns to generate longer outputs regardless of se-
1290 mantic quality. Furthermore, in the incorrect an-
1291 swer case (Figure 15), PPO assigns high advan-
1292 tages to both relevant and irrelevant content, indi-
1293 cating a lack of focus on key semantic cues and
1294 ultimately degrading training performance.

System Prompt

Role: You are the AI call assistant of xxx Inc., capable of answering calls clearly and politely on behalf of the user and conversing with the caller. Based on the user’s customized dialogue goals, you should first determine which goal applies to the conversation and respond accordingly. You do not know any personal information about the user beyond the task description and must inform the caller that any information they provide will be passed on to the user.

Dialogue Goals (Customized):

- **0: General Domain** – Used if the call doesn’t match any specific domain. Be polite and professional. Try to answer or guide the caller and clarify their intent.
- **1: E-commerce/Retail** – If caller mentions down jacket performance, mention fast delivery, good warmth at -5°C, ask for cleaning tips, other styles, hotline, or exchange policy, inquire about washing temperature.
- **2: Technical Support** – If license issues arise for WX3510/WX5540, answer related questions, confirm email and phone, send manuals, ask for model and associated phone.
- **3: Financial Services – Pension** – If retirement certification reminder, say not remembered, will handle via WeChat, thank them, ask about past communication and WeChat face verification method.
- **4: Life Services** – If laundry fee update, respond with satisfaction, hope for discount, mention night service needs, ask about pricing explanation, training quality, and monthly discounts.
- **5: Financial Services – Trade Notice** – If transaction confirmation, say user unavailable. If medium-risk, ask for analyst opinion, arrange follow-up next day, inquire about transaction fees.

Caller History and New Input

Conversation History: (*Empty*)

Caller New Input: Hello, this is the Forex Trading Customer Service Center. We are confirming that the currency pair you traded is USD/SGD, and the rate information is ready.

Model Response

Target Domain: *5: Financial Services – Trade Notice*

Dialogue State: *Ongoing*

Response: Hello, this call is currently being handled by the smart assistant. The user is temporarily unavailable. If this is a medium-risk transaction, could we get advice from a market analyst? Please contact us again tomorrow for confirmation.

Figure 9: A multi-turn example from the Honor-Dialogue dataset. The dataset features realistic, task-oriented, multi-domain conversations, in which each model response includes structured states. This example represents the situation of the first conversation.

Core Prompt for Dialogue Logicality

When scoring the intelligent call assistant's responses, please follow these steps:

1. Read and understand the entire call content.

2. Score from the dimension of dialogue logic rationality. The definition of dialogue logic rationality is: whether all responses provided by the assistant during the entire dialogue are clear, coherent, and effectively convey relevant information. A reasonable response should ensure the consistency, coherence, and relevance of information. Use a 1-5 scoring scale with specific criteria as follows:

- **1 point:** The response exhibits severe logical inconsistencies and violates fundamental reasoning principles. The intent of the assistant is unclear, and the response fails to convey meaningful or usable information.
- **2 points:** The response contains multiple inconsistencies or conflicting statements. Although partially related to the query, it lacks sufficient contextual grounding and coherent reasoning, resulting in poor interpretability and weak informational value.
- **3 points:** The response is generally coherent but contains minor logical gaps or discontinuities that affect fluency. While relevant to the query, the explanation lacks depth or clarity in parts, limiting the overall effectiveness of communication.
- **4 points:** The response is logically consistent and well-structured, with clear and relevant content. Information is conveyed effectively and aligns well with the user's intent, enabling smooth and coherent interaction.
- **5 points:** The response demonstrates strong logical consistency, clear reasoning, and precise information delivery. It not only addresses the query accurately but also facilitates deeper understanding, effectively guiding the conversation and enhancing overall interaction quality.

3. Based on the above scoring criteria and combined with the entire dialogue content, give a reasonable score, along with the scoring thinking process, deduction points, and modification suggestions to improve the score without changing the original meaning of the text.

Figure 10: A core prompt for dialogue logicality assessment from the constructed rubric evaluation method. The rubric features clear scoring steps, explicit definition of logical rationality, and detailed 1-5 point end-point criteria, in which each evaluation requires supplementary scoring reasoning, deduction explanations, and optimization suggestions.

Physics Fission Problem

Problem: An atomic nucleus of mass M is at rest. A spontaneous fission occurs in which the nucleus splits into two fragments (and nothing else), such that one fragment is 2 times more massive than the other (in terms of their rest-masses). The sum of rest-masses of the two fragments is 99% of the initial mass M . If the rest-mass energy of the initial nucleus (before fission) was 300 GeV, what is the kinetic energy of the more massive fragment?

Assumptions: Velocities of the fragments are significantly smaller than the speed of light. Ignore the electrons.

Figure 11: Example problem from the GPQA dataset

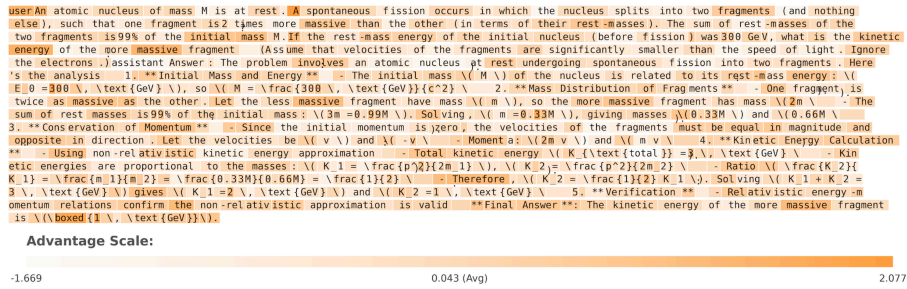


Figure 12: Our method on correct answer. Key information such as "initial mass M" is successfully captured.

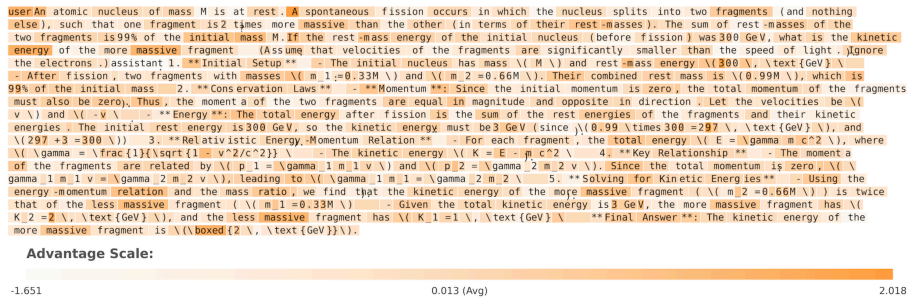


Figure 13: Our method on incorrect answer. The wrong prediction is appropriately assigned low advantage.

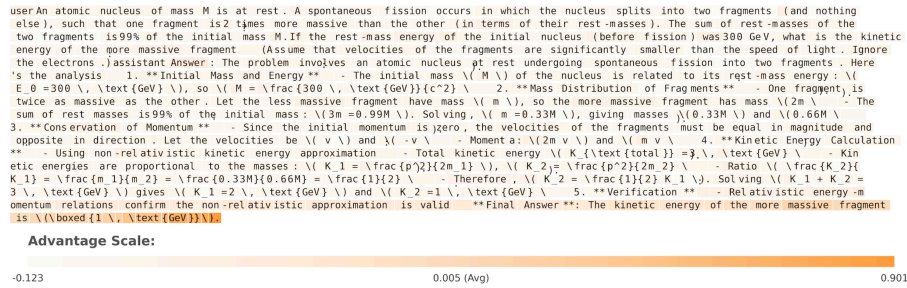


Figure 14: PPO on correct answer. The model fails to capture the semantic core and assigns low or inconsistent advantage.

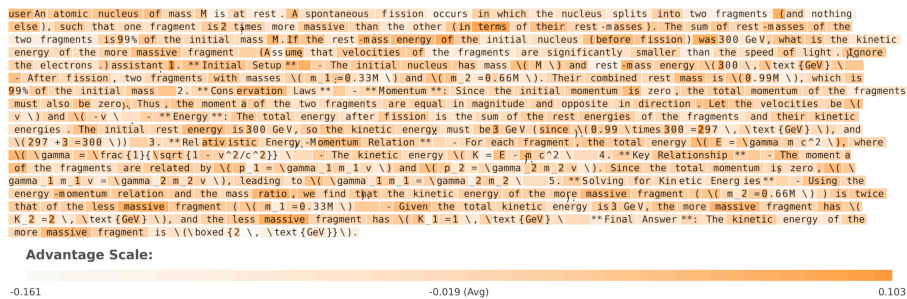


Figure 15: PPO on incorrect answer. High advantage is mistakenly assigned to both key and irrelevant content, weakening semantic discrimination.