DRBO: Mitigating The Bottleneck Effect via Dynamic Reward Balancing in Multi-reward LLM Optimization

Nuo Chen^{1†}, Yufei Gao^{1†}, Yongnan Jin^{1†},
Yan Hu¹, Anningzhe Gao^{2*}, Lingyong Yan³, Benyou Wang^{1,2*}

¹ The Chinese University of Hong Kong, Shenzhen

² Shenzhen Research Institute of Big Data

³ Baidu Inc., Beijing
chennuo@cuhk.edu.cn, wangbenyou@cuhk.edu.cn

Abstract

In the current landscape of large language models (LLMs), many evaluation metrics have been developed and used as rewards during training to improve specific metrics. However, balancing these metrics and dynamically adjusting reward weights remains challenging, as current approaches often fail to enhance weaker metrics. To address this, we empirically propose a Dynamic Reward Balancing Optimization framework \mathcal{DRBO} to mitigate the "bottleneck effect" by measuring performance, adjusting reward weights to prioritize weaker metrics, and optimizing the model via reinforcement learning. We apply \mathcal{DRBO} to both single-task and multi-type task scenarios, validating its effectiveness in generation with citations and online shopping conversation tasks. The results demonstrate improved overall performance and balanced optimization across multiple metrics, effectively overcoming the diversity and complexity inherent in LLMs. Our codes are available at https: //github.com/NuoJohnChen/DRBO.

1 Introduction

Evaluation is crucial for LLM applications, which often necessitates adherence to diverse evaluation criteria (Guo et al., 2023; Chang et al., 2024), even for a *single* task. Thanks to their generalization, LLMs are also expected to excel across *multiple types* of tasks simultaneously (Wang et al., 2023b), and each type of task might have its own metrics.

In reinforcement learning, one could structure these *metrics* (previously for evaluation) as *rewards* that could be boosted during training (Sharma et al., 2021; Yadav et al., 2021; Deng et al., 2022; Liu et al., 2023a; Xu et al., 2024; Wang et al., 2024b), to optimize complex objective functions even at testing time (OpenAI, 2024). However, when reward weights remain static, the weakest

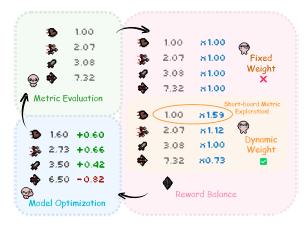


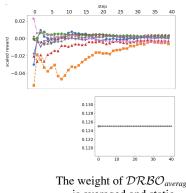
Figure 1: \mathcal{DRBO} adjusts rewards and weights dynamically through iterations of three main stages: performance metric *evaluation*, reward weight *balance*, and model parameter *optimization*. Elements of this schema are from The Binding of Isaac.

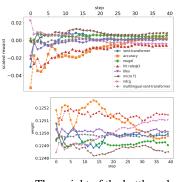
metric (the "bottleneck") becomes a bottleneck that restricts overall LLM effectiveness, which introduces the "bottleneck effect" in multi-reward optimization. For example, in Figure 2, when the scaled reward itself (or its growth trend) has not yet reached saturation, its update magnitude should accordingly be increased.

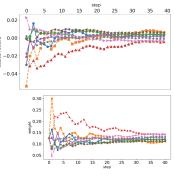
The mitigation of such effect is crucial because, for the optimization of LLMs' comprehensive capabilities, we typically aim for the model to **avoid weaknesses in any individual metric**, ensuring that its performance reaches a state of dynamic equilibrium across all metrics.

To mitigate the bottleneck effect, we propose a framework called \mathcal{DRBO} to balance multiple rewards (depicted in Figure 1). \mathcal{DRBO} interleaves between *Evaluation*, *Balance*, and *Optimization* stages, to alternately update the rewards and model parameters. After performance *Evaluation*, we dynamically reduce the weight of saturated rewards, enabling the model to prioritize weaker areas and mitigate the bottleneck effect in *Balance* stage. During *Optimization*, we update the model through

^{*}Benyou Wang and Anningzhe Gao are the corresponding authors. First three authors contributed to this work equally.







The weight of $DRBO_{average}$ is averaged and static.

The weight of the bottleneck (Accuracy) correspondingly increases.

The weight of rougel grows with $\Delta reward$.

- (a) $DRBO_{average}$ reward analysis. The reward improvement is limited and slow.
- $\mathcal{DRBO}_{inverse}$ reward analysis. (b) Accuracy increases largely in epoch 0. The bottleneck effect is mitigated.
- \mathcal{DRBO}_{delta} reward analysis. Variance reduces. The bottleneck effect is mitigated.

Figure 2: \mathcal{DRBO} Weight Dynamics (Shopping MMLU on Qwen2.5-1.5B-Instruct). The reward is scaled according to Eq. 4, with weight changes analyzed at epoch 0. Accuracy is identified as one bottleneck in $\mathcal{DRBO}_{average}$ static weight settings, showing lower mean and higher variance. In $\mathcal{DRBO}_{inverse}$, after assigning a large weight to the bottleneck (accuracy) in epoch 0, the reward of accuracy rises largely after step 5. In \mathcal{DRBO}_{delta} , the bottleneck effect reflects metrics with growth potential but limited progress; here, the weight is increased when the scaled reward has not yet saturated, encouraging further improvement.

reinforcement learning to achieve a more balanced overall performance.

We introduce contributions below.

- We propose a novel framework called \mathcal{DRBO} that dynamically learns and adjusts the weights of different evaluation metrics to simultaneously balance and optimize on both single-task and multi-task scenarios.
- We also show that by dynamically adjusting reward weights, our framework addresses the bottleneck effect through a more balanced and effective model optimization.
- Our extensive experiments on generation with citation and online shopping skills show significant performance improvements, validating the effectiveness of \mathcal{DRBO} .

Related Work

Multi Reward Integrating multiple reward functions is key for model optimization, with these functions often serving as key metrics (Pasunuru et al., 2020; Sharma et al., 2021; Yadav et al., 2021; Deng et al., 2022; Min et al., 2024; Yang et al., 2024; Naik et al., 2024). Additionally, Choi and Kim (2012); Zeng et al. (2023) delve into the theoretical exploration of the weights assigned to these reward

functions. For instance, Peitz and Dellnitz (2017), Poirion et al. (2017), Sener and Koltun (2018), Liu et al. (2023b); Sutton and Barto (1998) frame the issue as a Multi-Armed Bandit (MAB) problem (Thompson, 1933; Auer et al., 2002; Kuleshov and Precup, 2014; Slivkins, 2024), while others treat it as multi-objective optimization (Gholamnezhad et al., 2024; Kesireddy and Medrano, 2024; Tekin and Turgay, 2018; Wang et al., 2024b; Agnihotri et al., 2025). Compared to these approaches, our focus is on balancing multiple rewards to mitigate the bottleneck effect.

(c)

Multi-task optimization Many prior studies have explored the issue of weight allocation for different tasks in multi-task optimization problems. Some use the simplest approach of equal weights (Wang et al., 2020), while others treat the weights as hyperparameters, either setting them to empirical values or determining them through grid search (Chang et al., 2019; Dankers et al., 2019). There are also approaches that calculate weights based on training data (Perera et al., 2018; Lauscher et al., 2018), and some vary the weights according to simple functions (Clark et al., 2019). Our approach is to enhance the multi-task weights by using feedback from each task and promote leveraging the weaker tasks. Detailed related work, including automatic evaluations, can be found in Appendix A.

3 Background

3.1 Problem Definition

Model Evaluation A benchmark is a tool used to evaluate the performance of a model and is intended to replace the tedious manual labeling process. A benchmark usually consists of many pairs of datasets and metrics (\mathcal{D}_j, m_j) . The *dataset* \mathcal{D}_j contains a set of questions and answers, denoted as $\mathcal{D}_j = \{(x_i, y_i)\}_{1 \leq i \leq N}$, and the *metric* (e.g. BLEU, EM, F1) is used to evaluate the *performance* of the fine-tuned language model \mathcal{M} on \mathcal{D}_j . Denote $\mathcal{M}_{\theta}(x_i)$ is the output of the fine-tuned language model \mathcal{M} parameterized by θ , then the *performance* of \mathcal{M} on (\mathcal{D}_i, m_i) is regularly defined by

$$\mathcal{P}_{j}(\mathcal{M}_{\boldsymbol{\theta}}) = \frac{1}{N_{j}} \sum_{i=1}^{N_{j}} m_{j}(x_{i}, y_{i}, \mathcal{M}_{\boldsymbol{\theta}}(x_{i})) \in \mathbb{R},$$
(1)

Multi-Type Task Settings For *single-task* settings, the set of benchmarks $\{(\mathcal{D}, m_j)\}_{1 \leq j \leq N_j}$ share the same dataset \mathcal{D} . Our approach can also be generalized to *multi-type task* settings. In our *multi-task* setting, the benchmark datasets $\{\mathcal{D}_j\}$ and the test dataset $\{\mathcal{D}_{test}\}$ with size N are sampled from the same dataset \mathcal{D} . The generalization of the model will be improved through training.

Multi-Reward Optimization Traditional machine learning tasks often use differentiable training objectives to optimize models. However, these objectives may not align well with the metrics used for evaluating downstream tasks. To address this misalignment, some studies optimize models directly based on downstream metrics. Since these metrics are typically non-differentiable, reinforcement learning is used to treat them as rewards to optimize the model accordingly.

In our setting, each $\mathcal{P}_j(\mathcal{M}_{\theta})$ in Eq.1 can be viewed as a reward $\mathcal{R}_j(\mathcal{M}_{\theta})$ for fine-tuning \mathcal{M}_{θ} under the environment $\{(\mathcal{D}_j, m_j)\}$. The model \mathcal{M}_{θ} is treated as the policy π_{θ} , with π_{ref} serving as a reference policy to regularize the updates. Since multiple rewards are available, they need to be combined into $\mathcal{R} = \sum_j (w_j \mathcal{R}_j)$ for optimization. Then to maximize the expected reward \mathcal{R} , the model can be optimized by the policy gradient $\nabla_{\theta} \log \pi_{\theta}(y \mid x)$, with the objective function \mathcal{J} balancing reward maximization and KL divergence regularization:

$$\theta \leftarrow \theta + \nabla_{\theta} \mathcal{J}(\theta),$$
 (2)

in which

$$\mathcal{J} = \mathbb{E}_{\substack{x_i \sim \mathcal{D} \\ y_i \sim \mathcal{M}_{\boldsymbol{\theta}}(x_i)}} [\mathcal{R} - \beta \mathbb{D}_{KL}[\pi_{\boldsymbol{\theta}}(y_i|x_i) || \pi_{ref}(y_i|x_i)]]$$

$$\nabla_{\boldsymbol{\theta}} \mathcal{J}(\boldsymbol{\theta}) = \mathbb{E}_{x_i, y_i} [\nabla_{\boldsymbol{\theta}} \log \pi_{\boldsymbol{\theta}}(y_i | x_i) (\mathcal{R} - \beta \nabla_{\boldsymbol{\theta}} \mathbb{D}_{KL})].$$

Our goal is to optimize the model \mathcal{M}_{θ} to maximize $\sum_{j} (W_{j} \mathcal{R}_{j})$ for any pre-set weights W_{j} .

3.2 Motivation of DRBO

Importance of Balancing Multiple Rewards In Optimization step 2, balancing multiple rewards \mathcal{R}_j is crucial for improving the overall performance of the model \mathcal{M} . Dynamically adjusting reward weights w_j allows the model to allocate more focus on weaker aspects, preventing performance bottlenecks. We observe that conventional methods often fail to address this issue, as even when the reward saturates, its weight remains disproportionately high (Figure 2), showing that imbalanced metrics lead to suboptimal performance.

Analogies of the Bottleneck Effect This situation can be explained by the "bottleneck effect", where the overall effectiveness of a system is constrained by its weakest component. This concept can be extended to **model evaluation**, where the performance is often bottlenecked by the weakest metrics: if a model is optimized based solely on average rewards, it risks neglecting its worst-performing metrics. By dynamically adjusting the reward weights, the model is forced to address these shortcomings, leading to a more balanced performance. The bottleneck effect can also be analogized by biological evolution, wherein the survival and adaptability of organisms in nature are contingent upon **optimizing** their weakest traits.

While smaller models may be designed for specific tasks, LLMs should focus on **improving across multiple domains**, a challenge compounded by the complexity and diversity of rewards, making their optimization a significantly more intricate process.

Challenge to balance reward weights Achieving a proper balance of dynamic reward weights w_j requires addressing the varying performance \mathcal{P}_j across different metrics. During each iteration of step Eq.2, when model parameters θ are optimized, the performance of each metric $\mathcal{P}_j(\mathcal{M}_{\theta})$ fluctuates and needs re-evaluation. This fluctuation requires continuous re-balancing of the reward weights. Furthermore, since the importance

Algorithm 1: \mathcal{DRBO}

```
Require: model \mathcal{M}, data and metric \{\mathcal{D}_j, m_j\}_{1 \leq j \leq N_j}, and hyperparameters T, b, e, \tau
   Ensure: Optimized \mathcal{M}_{\theta} with improved performance across metrics
 1 Initialize weight \{w_j\}_{1 \le j \le N_j} over N_j metrics, state s as model performance \sum_j (w_j \mathcal{P}_j);
2 for j=1 to N_j do
    Initialize sample count count_j = 1;
4 for t=1 to T do
         Sample a metric m_s according to w_i;
 5
         count_j = count_j + 1;
 6
         for j = 1 to N_j do
              Evaluate \mathcal{R}_{j}^{b} = \mathcal{P}_{j}^{b}(\mathcal{M}_{\theta}) through Eq.1 from each batch data \mathcal{D}_{j}^{b};
 8
              Standardize reward \mathcal{R}_{j}^{b}(\mathcal{M}_{\theta}) = f(\mathcal{R}_{j}^{b});
 9
              Calculate the average reward across all batches \bar{\mathcal{R}}_{i}^{(t)} = \frac{1}{b} \sum_{b} \mathcal{R}_{j}(\mathcal{M}_{\theta}(\mathcal{D}_{i}^{b}));
10
         Compute total reward \mathcal{R}^{(t)} = \sum_{j} (w_j \bar{\mathcal{R}}_j^{(t)});
11
         Update the policy \pi_{\theta} = \mathcal{M}_{\theta} through Eq. 2 (Action a);
12
         if t\%e == 0 then
13
              for j = 1 to N_j do
14
                   Update reward across all steps x_j = r\left(\bar{\mathcal{R}}_j^{(1:t)}\right) through Eq. 5;
15
                   Calculate performance expectation \hat{x}_j = x_j + \sqrt{2 \ln t / count_j};
16
              Normalize weight w = w(W_{(1:j)}, \hat{x}_{(1:j)}) through Eq. 6,7;
17
```

and scale of each metric vary, the pre-set, unequal weights by human experts must also be considered, rather than assuming equal weighting.

4 Methodology of \mathcal{DRBO}

To address this challenge, we propose the algorithm \mathcal{DRBO} to mitigate the bottleneck effect, which means improving the model's performance should focus on the weaker metrics. Additionally, the performance across all metrics should be balanced according to pre-set weights. \mathcal{DRBO} algorithm achieves this goal through an iterative cycle involving three key phases detailed in the subsections, as colored in Algorithm 1:

- Evaluation: Dynamically assess the model's performance using multiple metrics and compute the average reward.
- Optimization: Based on the evaluation results, update the model parameters using reinforcement learning to optimize the overall reward.
- Balance: Periodically reassess and adjust the weights of each metric to ensure balanced performance across all metrics.

4.1 Evaluation

The evaluation phase aims to assess the performance of the model \mathcal{M}_{θ} across multiple metrics $\{m_j\}_{1\leq j\leq N_j}$. Initially, the weights w_j are uniformly distributed, and the state s reflects the current model performance $\sum_j (w_j \mathcal{P}_j)$. In each iteration t, up to T, a metric m_s is sampled based on w_j , and its sample count is incremented. The model is evaluated b times on each metric m_j using data batches \mathcal{D}_j^b , yielding rewards \mathcal{R}_j^b , which are standardized according to Eq. 4 to obtain $\bar{\mathcal{R}}_j^{(t)}$. The total reward \mathcal{R} is the weighted sum of these averages using w_j .

Reward Scaling in Evaluation To ensure that the rewards from different metrics are on the same scale when aggregated, and to avoid unfairness caused by differences in distribution and scale, we standardize the rewards for each metric. Specifically, we sample K reward values from each metric for standardization. First, we calculate the mean μ_j and variance σ_j^2 in Eq. 3 for each metric:

$$\mu_j = \frac{1}{K} \sum_{i=1}^K \mathcal{R}_{ji}, \quad \sigma_j^2 = \frac{1}{K} \sum_{i=1}^K (\mathcal{R}_{ji} - \mu_j)^2,$$
(3)

Next, we standardize the rewards to have a mean of 0 and variance of 1:

$$f(\mathcal{R}_{ji}) = \frac{\mathcal{R}_{ji} - \mu_j}{\sigma_j},\tag{4}$$

By applying this standardization, we ensure that each metric contributes fairly and reasonably when calculating the total reward.

4.2 Optimization

In the optimization phase, the evaluation results are used to update the model and adjust the weights of the metrics. The model's policy $\pi_{\theta} = \mathcal{M}_{\theta}$ is updated using reinforcement learning based on the total reward $R^{(t)} = \sum_j (w_j \bar{R}_j^{(t)})$. The action a in this phase involves updating the model parameters to maximize the total reward $R^{(t)}$.

4.3 Balance

In the balancing phase, we need to dynamically adjust the weights of various metrics to improve overall performance and pay attention to those that may be overlooked. To achieve this, we draw on the core idea of the Upper Confidence Bound (UCB) algorithm (Auer et al., 2002), which focuses on balancing exploration and exploitation: by exploring metrics with higher uncertainty (i.e., those that have been sampled less or whose performance has not yet been fully evaluated), we prevent the model from prematurely ignoring potential weaknesses, thus achieving better overall balance.

Specifically, every e steps, we re-evaluate the weights w_i for each metric. For each metric m_i , the adjusted expected performance \hat{x}_j is calculated using the following formula: $\hat{x}_j = \bar{x}_j + \sqrt{\frac{2 \ln t}{\text{count}_j}}$, where x_i is the reward calculated in Eq. 5, t is the total number of samples, and $count_j$ is the number of samples for metric m_j . The adjustment term $\sqrt{\frac{2 \ln t}{\text{count}_j}}$ encourages more exploration of those metrics that have been sampled less, ensuring the model can better evaluate each metric and prevent the occurrence of the bottleneck effect.

Proposals in Balance In the balancing phase, we raise two different proposals to adjust reward: $\mathcal{DRBO}_{inverse}$ and \mathcal{DRBO}_{delta} , focusing on mitigating the bottleneck effect. The specific calculation methods for both strategies are as follows.

$$r = \begin{cases} \mathbf{E}_{t} \left[\bar{\mathcal{R}}_{j}^{(1:t)} + 1 \right] = \frac{1}{t} \sum_{t} \bar{\mathcal{R}}_{j}^{(t)}(\mathcal{M}_{\theta}) + 1, \mathcal{DRBO}_{inverse} & \text{metrics and datasets. Dynamic reward weight adjustment helps the model balance its performance across all tasks, ensuring comprehensive optimizations of the property of$$

$$z_{j} = \begin{cases} \frac{W_{j}}{\hat{x}_{j}}, & \mathcal{DRBO}_{inverse} \\ W_{j}\hat{x}_{j}, & \mathcal{DRBO}_{delta} \end{cases}$$
(6)
$$w_{j} = \begin{cases} \frac{z_{j}}{\sum_{i=1}^{N_{j}} z_{i}}, & \mathcal{DRBO}_{inverse} \\ \text{softmax}(z_{j}) = \frac{e^{z_{j}/temp}}{\sum_{i=1}^{e^{z_{i}/temp}}}, \mathcal{DRBO}_{delta} \end{cases}$$
(7)

utilizes the "inverse rewards" strat- $DRBO_{inverse}$ egy to address the bottleneck effect by increasing weights for underperforming metrics through $z_j = \frac{W_j}{\hat{x}_j}$ in Eq. 6. This mechanism explicitly prioritizes under-performing metrics (lower $\hat{x}_i \rightarrow \text{higher } z_i$) rather than balancing explorationexploitation as in UCB (detailed explanation of difference with UCB and traditional Multi-Armed Bandit (MAB) (Auer et al., 2002; Vermorel and Mohri, 2005) is deferred to Appendix B). Weight computation via Eq. 7 ensures non-neglect of highperforming metrics while dynamically reallocating optimization focus, guided by pre-set target weights W_i for controlled convergence.

 $DRBO_{delta}$ focuses on boosting the weights of metrics with greater potential for improvement. r in Eq. 5 adjusts the model parameters based on the reward increments, allowing for a faster response to changes in the environment. It helps to prioritize weak metrics with higher potential for improvement. Since r can be negative, Eq. 7 uses softmax strategy for weight normalization to ensure stability. Setting temperature temp = 0.01 step helps alleviate the problem of weight disappearance.

Applications

In this section, we introduce the two configurations of \mathcal{DRBO} : the *single* task setting and the multi-type task setting, along with their respective application scenarios.

Single-Task Setting The model is evaluated using multiple metrics on the same dataset, similar to how a species adapts to different environments. By dynamically adjusting weights, the model can effectively allocate resources and optimize its performance across all metrics.

Multi-Type Task Setting The model handles different types of tasks, each with its own specific across all tasks, ensuring comprehensive optimization, similar to how a species adapts to different

ecological niches. The adaptability makes it an ideal choice for multi-task competitions or benchmark challenges.

Single-Task Application on Long-context Generation with Citation Single-Task \mathcal{DRBO} can effectively address the challenges faced by long-context retrievers or citation-based text generation (Gao et al., 2023), such as handling complex queries, providing accurate results, and ensuring credible references. By dynamically adjusting the weights of metrics on fluency, correctness, and citation quality, \mathcal{DRBO} can achieve balanced optimization across these aspects, thereby improving the model's ability to synthesize information from multiple sources.

Multi-Type Task Application on Online Shopping Conversation Skills Multi-Type Task \mathcal{DRBO} can handle challenges presented by online shopping interactive conversations (Jin et al., 2024), which includes shopping skills such as concept understanding, knowledge reasoning, user behavior alignment, and multi-lingual abilities. By dynamically adjusting the weights of criterias including multiple choice, geneation, retrieval, and so on, \mathcal{DRBO} ensures balanced optimization across these tasks, improving the model's ability to serve as general shop assistants.

6 Experiments

To demonstrate \mathcal{DRBO} 's methodological innovation and its immediate applicability, the experiments is conducted in both multi-evaluation and multi-task settings. We systematically evaluate the effectiveness of \mathcal{DRBO} through a series of experiments designed to test its performance across the scenarios detailed in Section 5.

6.1 Experiment setting

We conducted all experiments using the Py-Torch framework on a setup consisting of eight NVIDIA A100 GPUs, each with 80 GB of memory. The computing environment was configured with CUDA 11.8 and cuDNN 8.7 for optimized deep learning performance. We set W to 1, denoting each evaluator is equally important.

The baseline includes (1) ALaRM (Lai et al., 2024), a reward modeling method suitable for single-task settings (compared in Table 1); (2) sequential optimization (sec F.2); and (3) $\mathcal{DRBO}_{average}$, a baseline approach which uses

static, equal weighting $w_j = \frac{1}{N_j}$ in the balance phase. We then evaluate $\mathcal{DRBO}_{inverse}$ and \mathcal{DRBO}_{delta} . In Optimization phase, there are several different reinforcement learning methods as choices, including PPO (Schulman et al., 2017), ReMax algorithm (Li et al., 2023b), and so on, to solve Eq. 2. In this paper, we apply ReMax to avoid training a value model and reduce computations. Detailed parameters are listed in Appendix C.

6.2 Experiment 1: Generation with Citation

We focus on the long-context generation with citation in a multi-evaluation setting, which assesses the model's ability to generate accurate and coherent responses given retrieved documents. We use ALaRM's reward modeling method as our baseline, following its holistic reward weight of 3 (computed by the UltraRM-13B (Cui et al., 2024)) while keeping other reward weights at 1. We adopt ALCE benchmark (Gao et al., 2023), a well-known standard in the RAG community, because it is recognized for its ability to provide a single output with multiple evaluations, enabling a comprehensive assessment of model performance. Using this benchmark¹, we apply the following criterias with metrics detailed in Appendix D as reward providers.

- **Fluency**: Evaluated by *MAUVE* (Pillutla et al., 2021), which measures the model's fluency in generating text.
- Correctness: Measured by exact match (EM) of the golden answer for ASQA (Stelmakh et al., 2022) or using an NLI model for inference for ELI5 (Fan et al., 2019), ensuring the output aligns with the correct answer.
- **Citation Accuracy**: Determined using an NLI model² to infer the correctness of citations, with metrics including citation recall and citation precision score to evaluate citation *F1*.

6.2.1 Analysis

According to the results on ASQA and ELI5 in Table 1, we reach the following conclusions:

 \mathcal{DRBO} improves overall performance. After applying \mathcal{DRBO} , the overall score of ASQA

¹Fluency is a scalar derived from the distribution of several sentences. We report the average fluency across batches, while Gao et al. (2023) reports it across the entire dataset.

²https://huggingface.co/google/ t5_xxl_true_nli_mixture

Dataset		ASQA	1		ELI5			
Criteria	Fluency	Correctness	Citation	All	Fluency	Correctness	Citation	All
Metric	MAUVE	EM Recall	F1	Average	MAUVE	Claim Recall	F1	Average
			Llama	-2-7B-Cha	nt			
ALCE	40.92	49.71	42.75	44.46	46.08	16.67	21.45	28.07
ALaRM	43.57	47.92	37.54	43.01	44.85	11.83	21.28	25.99
$\mathcal{DRBO}_{average}$	40.64	49.25	48.44	46.11	49.76	18.33	29.66	32.58
$\mathcal{DRBO}_{inverse}$	40.92	48.74	51.68	47.11	45.94	19.02	29.98	31.65
\mathcal{DRBO}_{delta}	43.55	49.66	49.48	47.56	48.88	18.69	30.28	32.62
			Llama-3	3-8B-Instr	uct			
ALCE	22.01	49.48	59.38	43.62	47.43	18.83	39.71	35.32
ALaRM	20.06	46.77	58.12	41.65	47.30	14.13	42.89	34.77
$\mathcal{DRBO}_{average}$	24.99	50.17	70.98	48.71	47.43	18.06	43.89	36.46
$\mathcal{DRBO}_{inverse}$	27.93	51.06	68.47	49.15	46.31	18.67	45.79	36.92
\mathcal{DRBO}_{delta}	24.13	50.54	68.37	47.68	44.94	18.50	47.08	36.84

Table 1: Performance on ASQA and ELI5 Datasets.

improved from 44.46 to 47.56 on Llama-2-7B-Chat and from 43.62 to 49.15 on Llama-3-8B-Instruct, with ELI5 also showing significant enhancement. This average score improvement suggests that by adjusting the weights of different metrics according to their scores in the reinforcement learning process, \mathcal{DRBO} outperforms ALaRM and $\mathcal{DRBO}_{average}$, and can improve the overall performance effectively.

 $\mathcal{DRBO}_{inverse}$ performs better and more balanced on stronger base models, while \mathcal{DRBO}_{delta} excels with weaker base models. In Table 1, \mathcal{DRBO}_{delta} outperforms $\mathcal{DRBO}_{inverse}$ on the Llama-2-7B-Chat model, particularly in the area of fluency, which is a strong criterion. Conversely, the Llama-3-8B-Instruct model exhibits the opposite trend. As a result, stronger models tend to focus more on the weakest criteria after scaling, while weaker models prioritize criteria with a more rapid growth trend.

For a more detailed analysis of experiment 1 and case study, please refer to Appendix H.

6.2.2 Further Analysis

RQ1: Convergence Analysis Does \mathcal{DRBO} achieve balanced convergence during optimization?

In Figure 3, weights and rewards change towards the opposite direction, and finally converge to a stable weight and reward. It clarifies the robustness of \mathcal{DRBO} .

RQ2: Training Dynamics How does each metric evolve during the training process? As shown in Figure 4, the performance on each metric rapidly increases at the beginning of training, surpassing

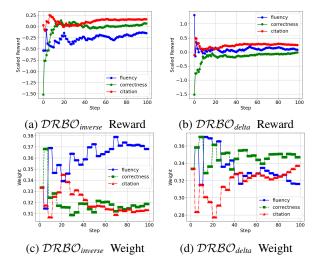


Figure 3: The change of scaled rewards and metric weights of Llama-2-7B-Chat on ASQA. More examples can be found in Figure 8.

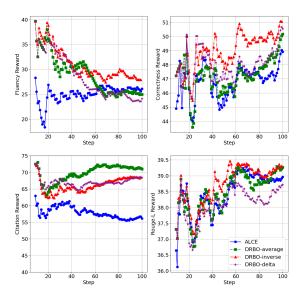


Figure 4: Performance of Llama-3-8B on ASQA.

Criteria	NER	Retrieval	Choice	Ranking	Generation	Translation	Multilinguality	Extraction	All	
Metric	Micro F1	Hit Rate@3	Accuracy	NDCG	Sent-Transformer	BLEU-4	Multilingual-Sent-Transformer	ROUGE-L	Average	
				Llaı	na-3.2-1b-Instruct					
Shopping MMLU	39.72	15.96	55.12	57.96	50.06	19.09	43.97	13.5	36.92	
$DRBO_{average}$	29.48	20.32	54.68	59.38	53.58	16.41	39.80	3.73	34.67	
$DRBO_{inverse}$	33.20	20.29	52.36	63.58	51.29	16.67	44.87	8.18	36.31	
-W = 2:1:1:1:1:1:1:1	44.64	18.93	57.61	65.00	50.16	19.84	45.19	12.75	39.27	
$DRBO_{delta}$	43.89	17.89	35.88	58.75	48.71	17.75	42.2	6.94	34.00	
				Qwe	n2.5-1.5B-Instruct					
Shopping MMLU	11.39	38.64	56.05	76.56	51.13	15.69	45.01	3.03	37.19	
$DRBO_{average}$	12.19	41.67	62.45	79.47	52.03	13.82	46.19	3.62	38.93	
$DRBO_{inverse}$	9.42	45.12	64.50	78.12	52.70	16.32	46.52	3.45	39.52	
$DRBO_{delta}$	12.24	38.04	58.94	77.98	51.42	16.66	47.23	3.12	38.20	
	Llama-3-8B-Instruct									
Shopping MMLU	58.07	58.95	72.56	78.12	49.27	24.59	38.76	5.80	48.27	
$DRBO_{average}$	44.15	49.46	63.87	75.51	51.59	22.92	53.12	6.38	45.88	
$DRBO_{inverse}$	67.59	66.93	76.49	82.22	49.80	23.74	38.85	7.00	51.58	
$DRBO_{delta}$	58.36	60.37	73.98	82.19	49.33	21.35	36.34	4.40	48.29	

Table 3: Performance on Shopping MMLU. The p-value in the t-test < 0.002, detailed in Appendix F.1.

the original model, and then gradually stabilizes over time. The fluency curve drops at step 50, reminding us that the training cycle of \mathcal{DRBO} should not be too long, as it may lead to overfitting. The weakest metric, citation, shows stable training performance, further demonstrating the robustness of \mathcal{DRBO} .

RQ3: Non-uniform Weighting How does \mathcal{DRBO} perform when metrics have unequal predefined importance?

Metric Measurement	Fluency MAUVE	Correctness EM Recall	Citation F1				
Llama3-8b-Instruct							
$\mathcal{DRBO}_{inverse}$ (1:1:1)	27.11	51.04	67.75				
\mathcal{DRBO}_{delta} (1:1:1)	24.13	50.54	68.37				
$\mathcal{DRBO}_{inverse}$ (1:2:3)	25.62	49.84	68.34				
\mathcal{DRBO}_{delta} (1:2:3)	33.13	50.93	69.54				

Table 2: Performance of ASQA when W = 1:2:3

The effect of \mathcal{DRBO} when the pre-defined weights are set to W=1:2:3 is shown in Table 2, where we reduce the importance of fluency. We can see a significant improvement in all metrics, with the most notable increase in citation. This demonstrates the effectiveness of our method under different values of W, and highlights that the initial value of W can also have a significant impact on the results.

6.3 Experiment 2: Online Shopping Skills

We target at comprehensive abilities of LLMs on multi-type shopping QA, which assess multiple few-shot tasks with complex entities and relations. The ShoppingMMLU (Jin et al., 2024) benchmark, detailed in Appendix E in KDD Cup 2024³, is se-

lected because it is a statistically detailed, multitype task dataset derived from real-world scenarios on Amazon, providing a comprehensive evaluation of models in shopping contexts. We divide this benchmark into the criteria listed with metrics below to provide rewards:

- Named Entity Recongnition (NER): Evaluated using the *Micro-F1* score based on precisions and recalls of NER tasks.
- **Retrieval**: Evaluated by *Hit Rate@3* through calculating the overlap between the retrieved set (maximum 3 length) and the ground truth.
- **Multiple-Choice**: Evaluated by *accuracy* through generating one token and comparing it with the ground truth (Liang et al., 2022).
- **Ranking**: Assessed using the Normalized Discounted Cumulative Gain (*NDCG*)= $\frac{DCG}{iDCG}$, considering the ranking relevance.
- Unrestricted Generation: The embeddings of generated text and reference text are transformed using *sentence transformers* (Reimers and Gurevych, 2019) and the cosine *similarity* to evaluate the semantic alignment between generated and reference texts.
- **Translation**: Evaluated using *BLEU-4* (Papineni et al., 2002), considering n-gram overlaps between generated and reference texts.
- Multilinguality: Evaluated using sentence transformer similarity (Reimers and Gurevych, 2019) between the generated text and the reference text in a multilingual context.
- Extraction: Evaluated using *ROUGE-L* (Lin, 2004) on extractive generation tasks (i.e. the answer is a sub-string of the input).

³https://www.aicrowd.com/challenges/amazon-kddcup-2024-multi-task-online-shopping-challengefor-11ms

6.3.1 Analysis

 $\mathcal{DRBO}_{inverse}$ achieves overall enhancement in multi-type task settings, and \mathcal{DRBO}_{delta} shows steady improvement. $\mathcal{DRBO}_{inverse}$ significantly outperforms zero-shot Shopping MMLU under Qwen2.5-1.5B-Instruct and Llama-3-8B-Instruct, demonstrating its ability to overcome the diversity and complexity inherent in LLMs. In contrast, $\mathcal{DRBO}_{average}$ tends to over-optimize a single metric, such as Multilinguality under Llama-3-8B-Instruct, resulting in a loss of balance. While \mathcal{DRBO}_{delta} generally performs worse than $\mathcal{DRBO}_{inverse}$, it excels in certain low-resource tasks, such as NER and translation, showcasing better balance in these scenarios.

 \mathcal{DRBO} performs better on larger models. For larger models, $\mathcal{DRBO}_{inverse}$ shows greater improvement compared to zero-shot, possibly because larger models have superior multi-task generalization capabilities. Notably, even for the challenging low-resource task Extraction, optimization on larger models still achieves noticeable increases.

6.4 Further Analysis

RQ1: Metric Interdependence Do the optimization dynamics of metrics exhibit mutual influence?

Cattanta	NED	D. 4.2	Chaire	D 1-1				
Criteria	NER	Retrieval	Choice	Ranking				
Qwen2.5-1.5B-Instruct								
Shopping MMLU	11.39	38.64	56.05	76.56				
NER Only	29.28	27.39	54.68	72.75				
Criteria	Generation	Translation	Multilinguality	Extraction				
Shopping MMLU	51.13	15.69	45.01	3.03				
NER Only	47.23	14.66	45.71	2.00				

Table 4: Performance of Qwen2.5-1.5B-Instruct with only Remax applied to NER for the **ablation** study.

Criteria	NER	Retrieval	Choice	Ranking					
LLaMA-3.2-1B-Instruct									
Shopping MMLU	39.72	15.96	55.12	57.96					
Generation Only	37.83	16.64	55.46	58.5					
Retrieval Only	40.26	19.69	54.21	68.46					
Translation Only	37.73	16.91	55.13	56.21					
Criteria	Generation	Translation	Multilinguality	Extraction					
	Generation	Tansianon	wattimiguality	Extraction					
Shopping MMLU	50.06	19.09	43.97	13.50					
Shopping MMLU Generation Only									
11 0	50.06	19.09	43.97	13.50					

Table 5: Performance of LLaMA-3.2-1B-Instruct with only Remax applied to specific metrics for the **ablation** study.

To explore the correlation between metrics, according to Table 4, training exclusively on NER (setting the weights of other criteria to zero) significantly enhances NER performance but leads to an imbalance across most other metrics (similar

conclusions can be reached in Table 5). However, some tasks, such as Multi-Choice, still benefit from the improvement in NER. This also explains why, in Figure 2, the weights of certain high-performing metrics continue to increase. Moreover, this highlights that increasing the number of evaluation metrics can further enhance the performance of LLMs.

RQ2: Reward Hacking? A case study analysis.

You are required to perform the task of query named entity recognition. Please follow the given instructions. You are a helpful online shop assistant and a linguist. A customer on an online shopping platform has made the following query. Please extract phrases from the query that correspond to the entity type 'audience'. Please directly output the entity without repeating the entity type. If there are multiple such entities, separate them with comma. Do not give explanations. Query: van heusen formal shirts for men Output:

```
Shopping MMLU Responses: 1 X, TP: 0

DRBO<sub>average</sub> Responses: 1 X, TP: 0

DRBO<sub>inverse</sub> Responses: men ✓, TP: 1

DRBO<sub>delta</sub> Responses: mens ✓, TP: 1
```

Figure 5: NER Case study on Qwen-2.5-1.5B-Instruct

Reward hacking (Skalse et al., 2022) refers to a phenomenon where a model exploits loopholes in the reward function to maximize its score in unintended ways, often at the expense of achieving the true objectives of the task. We present a case study in Figure 5,12,13 to explore whether such phenomena exist. Examples from both classification and generation tasks demonstrate that the task-specific metrics objectively reflect performance, and the results indicate that the effectiveness of these metrics is improved after applying \mathcal{DRBO} , showcasing the robustness of our approach.

7 Conclusion

We proposed \mathcal{DRBO} framework to address the bottleneck effect in multi-reward optimization for LLMs. By dynamically adjusting reward weights, \mathcal{DRBO} effectively prioritizes weaker metrics to achieve a balanced and improved overall performance. Experimental results across single-task and multi-type task scenarios demonstrate that \mathcal{DRBO} significantly enhances model performance, ensuring that no single metric dominates optimization. Furthermore, \mathcal{DRBO} provides a flexible and adaptive mechanism that can be integrated into various tasks, highlighting its versatility and scalability. We hope it opens pathways for future advancements in dynamic reward adjustment strategies and broader applications in AI optimization.

Limitations

When there are significant differences in scales between rewards, prior knowledge about these scales is required to normalize the rewards effectively. Additionally, when the number of rewards is large, the framework becomes challenging to balance and optimize across many rewards.

Potential Risks and Ethical Consideration

Our proposed approach also involves certain potential risks that need to be considered:

Bias and Fairness Issues: If the training data contains biases, the model may reinforce these biases during the optimization process, leading to unfair outcomes. This can be particularly problematic in sensitive applications such as recruitment or healthcare, where ethical and legal issues may arise.

Data Privacy and Security Risks: Handling large-scale data necessitates stringent measures to protect user privacy and ensure data security. Improperly protected data can be misused or leaked, posing significant privacy and security risks.

Acknowledgement

This work was supported by the Shenzhen Science and Technology Program (JCYJ20220818103001002), Shenzhen Doctoral Startup Funding (RCBS20221008093330065), Tianyuan Fund for Mathematics of National Natural Science Foundation of China (NSFC) (12326608), Shenzhen Key Laboratory of Cross-Modal Cognitive Computing (grant number ZDSYS20230626091302006), and Shenzhen Stability Science Program 2023.

References

- Akhil Agnihotri, Rahul Jain, Deepak Ramachandran, and Zheng Wen. 2025. Multi-objective preference optimization: Improving human alignment of generative models. *Preprint*, arXiv:2505.10892.
- Peter Auer, Nicolò Cesa-Bianchi, and Paul Fischer. 2002. Finite-time analysis of the multiarmed bandit problem. *Mach. Learn.*, 47(2-3):235–256.
- Sébastien Bubeck, Varun Chandrasekaran, Ronen Eldan, Johannes Gehrke, Eric Horvitz, Ece Kamar, Peter Lee, Yin Tat Lee, Yuanzhi Li, Scott Lundberg, Harsha Nori, Hamid Palangi, Marco Tulio Ribeiro, and Yi Zhang. 2023. Sparks of artificial general intelligence: Early experiments with gpt-4. *Preprint*, arXiv:2303.12712.

- Shuaichen Chang, Pengfei Liu, Yun Tang, Jing Huang, Xiaodong He, and Bowen Zhou. 2019. Zero-shot text-to-sql learning with auxiliary task. *Preprint*, arXiv:1908.11052.
- Yupeng Chang, Xu Wang, Jindong Wang, Yuan Wu, Linyi Yang, Kaijie Zhu, Hao Chen, Xiaoyuan Yi, Cunxiang Wang, Yidong Wang, Wei Ye, Yue Zhang, Yi Chang, Philip S. Yu, Qiang Yang, and Xing Xie. 2024. A survey on evaluation of large language models. *ACM Trans. Intell. Syst. Technol.*, 15(3):39:1–39:45.
- Nuo Chen, Zhiyuan Hu, Qingyun Zou, Jiaying Wu, Qian Wang, Bryan Hooi, and Bingsheng He. 2025a. Judgelrm: Large reasoning models as a judge. *Preprint*, arXiv:2504.00050.
- Nuo Chen, Andre Lin HuiKai, Jiaying Wu, Junyi Hou, Zining Zhang, Qian Wang, Xidong Wang, and Bingsheng He. 2025b. Xtragpt: Context-aware and controllable academic paper revision via human-ai collaboration. *Preprint*, arXiv:2505.11336.
- Nuo Chen, Qiushi Sun, Jianing Wang, Ming Gao, Xiaoli Li, and Xiang Li. 2023. Evaluating and enhancing the robustness of code pre-trained models through structure-aware adversarial samples generation. In *Findings of the Association for Computational Linguistics: EMNLP 2023*, pages 14857–14873, Singapore. Association for Computational Linguistics.
- Wei-Lin Chiang, Lianmin Zheng, Ying Sheng, Anastasios Nikolas Angelopoulos, Tianle Li, Dacheng Li, Hao Zhang, Banghua Zhu, Michael Jordan, Joseph E. Gonzalez, and Ion Stoica. 2024. Chatbot arena: An open platform for evaluating llms by human preference. *Preprint*, arXiv:2403.04132.
- Jaedeug Choi and Kee-Eung Kim. 2012. Nonparametric bayesian inverse reinforcement learning for multiple reward functions. In Advances in Neural Information Processing Systems 25: 26th Annual Conference on Neural Information Processing Systems 2012. Proceedings of a meeting held December 3-6, 2012, Lake Tahoe, Nevada, United States, pages 314–322.
- Kevin Clark, Minh-Thang Luong, Urvashi Khandelwal, Christopher D. Manning, and Quoc V. Le. 2019. Bam! born-again multi-task networks for natural language understanding. *Preprint*, arXiv:1907.04829.
- Ganqu Cui, Lifan Yuan, Ning Ding, Guanming Yao, Bingxiang He, Wei Zhu, Yuan Ni, Guotong Xie, Ruobing Xie, Yankai Lin, Zhiyuan Liu, and Maosong Sun. 2024. Ultrafeedback: Boosting language models with scaled ai feedback. *Preprint*, arXiv:2310.01377.
- Verna Dankers, Marek Rei, Martha Lewis, and Ekaterina Shutova. 2019. Modelling the interplay of metaphor and emotion through multitask learning. In Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the

- 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP), pages 2218–2229, Hong Kong, China. Association for Computational Linguistics.
- Mingkai Deng, Jianyu Wang, Cheng-Ping Hsieh, Yihan Wang, Han Guo, Tianmin Shu, Meng Song, Eric Xing, and Zhiting Hu. 2022. RLPrompt: Optimizing discrete text prompts with reinforcement learning. In *Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing*, pages 3369–3391, Abu Dhabi, United Arab Emirates. Association for Computational Linguistics.
- Angela Fan, Yacine Jernite, Ethan Perez, David Grangier, Jason Weston, and Michael Auli. 2019. ELI5: long form question answering. In *Proceedings of the 57th Conference of the Association for Computational Linguistics, ACL 2019, Florence, Italy, July 28- August 2, 2019, Volume 1: Long Papers*, pages 3558–3567. Association for Computational Linguistics.
- Tianyu Gao, Howard Yen, Jiatong Yu, and Danqi Chen. 2023. Enabling large language models to generate text with citations. *Preprint*, arXiv:2305.14627.
- Wentao Ge, Shunian Chen, Hardy Chen, Nuo Chen, Junying Chen, Zhihong Chen, Wenya Xie, Shuo Yan, Chenghao Zhu, Ziyue Lin, Dingjie Song, Xidong Wang, Anningzhe Gao, Zhang Zhiyi, Jianquan Li, Xiang Wan, and Benyou Wang. 2025. MLLM-bench: Evaluating multimodal LLMs with per-sample criteria. In Proceedings of the 2025 Conference of the Nations of the Americas Chapter of the Association for Computational Linguistics: Human Language Technologies (Volume 1: Long Papers), pages 4951–4974, Albuquerque, New Mexico. Association for Computational Linguistics.
- Pezhman Gholamnezhad, Ali Broumandnia, and Vahid Seydi. 2024. An improved model-based evolutionary algorithm for multi-objective optimization. *Concurr. Comput. Pract. Exp.*, 36(10).
- Zishan Guo, Renren Jin, Chuang Liu, Yufei Huang, Dan Shi, Supryadi, Linhao Yu, Yan Liu, Jiaxuan Li, Bojian Xiong, and Deyi Xiong. 2023. Evaluating large language models: A comprehensive survey. *Preprint*, arXiv:2310.19736.
- Neel Jain, Khalid Saifullah, Yuxin Wen, John Kirchenbauer, Manli Shu, Aniruddha Saha, Micah Goldblum, Jonas Geiping, and Tom Goldstein. 2023. Bring your own data! self-supervised evaluation for large language models. *CoRR*, abs/2306.13651.
- Yilun Jin, Zheng Li, Chenwei Zhang, Tianyu Cao, Yifan Gao, Pratik Sridatt Jayarao, Mao Li, Xin Liu, Ritesh Sarkhel, Xianfeng Tang, Haodong Wang, Zhengyang Wang, Wenju Xu, Jingfeng Yang, Qingyu Yin, Xian Li, Priyanka Nigam, Yi Xu, Kai Chen, Qiang Yang, Meng Jiang, and Bing Yin. 2024. Shopping MMLU: A massive multi-task online shopping benchmark

- for large language models. In *The Thirty-eight Conference on Neural Information Processing Systems Datasets and Benchmarks Track*.
- Adarsh Kesireddy and F. Antonio Medrano. 2024. Elite multi-criteria decision making pareto front optimization in multi-objective optimization. *Algorithms*, 17(5):206.
- Volodymyr Kuleshov and Doina Precup. 2014. Algorithms for multi-armed bandit problems. *Preprint*, arXiv:1402.6028.
- Yuhang Lai, Siyuan Wang, Shujun Liu, Xuanjing Huang, and Zhongyu Wei. 2024. Alarm: Align language models via hierarchical rewards modeling. *Preprint*, arXiv:2403.06754.
- Anne Lauscher, Goran Glavaš, Simone Paolo Ponzetto, and Kai Eckert. 2018. Investigating the role of argumentation in the rhetorical analysis of scientific publications with neural multi-task learning models. In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*, pages 3326–3338, Brussels, Belgium. Association for Computational Linguistics.
- Ruosen Li, Teerth Patel, and Xinya Du. 2023a. Prd: Peer rank and discussion improve large language model based evaluations. *Preprint*, arXiv:2307.02762.
- Ziniu Li, Tian Xu, Yushun Zhang, Yang Yu, Ruoyu Sun, and Zhimin Luo. 2023b. Remax: A simple, effective, and efficient reinforcement learning method for aligning large language models. *ArXiv*, abs/2310.10505.
- Percy Liang, Rishi Bommasani, and Tony Lee et al. 2022. Holistic evaluation of language models. *CoRR*, abs/2211.09110.
- Bill Yuchen Lin, Yuntian Deng, Khyathi Chandu, Faeze Brahman, Abhilasha Ravichander, Valentina Pyatkin, Nouha Dziri, Ronan Le Bras, and Yejin Choi. 2024. Wildbench: Benchmarking Ilms with challenging tasks from real users in the wild. *Preprint*, arXiv:2406.04770.
- Chin-Yew Lin. 2004. ROUGE: A package for automatic evaluation of summaries. In *Text Summarization Branches Out*, pages 74–81, Barcelona, Spain. Association for Computational Linguistics.
- Yen-Ting Lin and Yun-Nung Chen. 2023. Llm-eval: Unified multi-dimensional automatic evaluation for open-domain conversations with large language models. *CoRR*, abs/2305.13711.
- Jiate Liu, Yiqin Zhu, Kaiwen Xiao, Qiang Fu, Xiao Han, Wei Yang, and Deheng Ye. 2023a. RLTF: reinforcement learning from unit test feedback. *CoRR*, abs/2307.04349.
- Yajiao Liu, Xin Jiang, Yichun Yin, Yasheng Wang, Fei Mi, Qun Liu, Xiang Wan, and Benyou Wang. 2023b. One cannot stand for everyone! leveraging multiple

- user simulators to train task-oriented dialogue systems. In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), ACL 2023, Toronto, Canada, July 9-14, 2023*, pages 1–21. Association for Computational Linguistics.
- Yue Liu, Shengfang Zhai, Mingzhe Du, Yulin Chen, Tri Cao, Hongcheng Gao, Cheng Wang, Xinfeng Li, Kun Wang, Junfeng Fang, Jiaheng Zhang, and Bryan Hooi. 2025. Guardreasoner-vl: Safeguarding vlms via reinforced reasoning. *Preprint*, arXiv:2505.11049.
- Do June Min, Verónica Pérez-Rosas, Ken Resnicow, and Rada Mihalcea. 2024. Dynamic reward adjustment in multi-reward reinforcement learning for counselor reflection generation. In *Proceedings of the 2024 Joint International Conference on Computational Linguistics, Language Resources and Evaluation, LREC/-COLING 2024, 20-25 May, 2024, Torino, Italy*, pages 5437–5449. ELRA and ICCL.
- Abhishek Naik, Yi Wan, Manan Tomar, and Richard S. Sutton. 2024. Reward centering. *Preprint*, arXiv:2405.09999.
- Jinjie Ni, Fuzhao Xue, Xiang Yue, Yuntian Deng, Mahir Shah, Kabir Jain, Graham Neubig, and Yang You. 2024. Mixeval: Deriving wisdom of the crowd from llm benchmark mixtures. arXiv preprint arXiv:2406.06565.
- OpenAI. 2024. Deliberative alignment: Reasoning enables safer language models.
- Kishore Papineni, Salim Roukos, Todd Ward, and Wei-Jing Zhu. 2002. Bleu: a method for automatic evaluation of machine translation. In *Proceedings of the* 40th Annual Meeting of the Association for Computational Linguistics, pages 311–318, Philadelphia, Pennsylvania, USA. Association for Computational Linguistics.
- Ramakanth Pasunuru, Han Guo, and Mohit Bansal. 2020. DORB: dynamically optimizing multiple rewards with bandits. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing, EMNLP 2020, Online, November 16-20, 2020*, pages 7766–7780. Association for Computational Linguistics.
- Sebastian Peitz and Michael Dellnitz. 2017. *Gradient-Based Multiobjective Optimization with Uncertainties*, page 159–182. Springer International Publishing.
- Vittorio Perera, Tagyoung Chung, Thomas Kollar, and Emma Strubell. 2018. Multi-task learning for parsing the alexa meaning representation language. *Proceedings of the AAAI Conference on Artificial Intelligence*, 32(1).
- Krishna Pillutla, Swabha Swayamdipta, Rowan Zellers, John Thickstun, Yejin Choi, and Zaïd Harchaoui. 2021. MAUVE: human-machine divergence curves for evaluating open-ended text generation. *CoRR*, abs/2102.01454.

- Fabrice Poirion, Quentin Mercier, and Jean-Antoine Désidéri. 2017. Descent algorithm for nonsmooth stochastic multiobjective optimization. *Comput. Optim. Appl.*, 68(2):317–331.
- Nils Reimers and Iryna Gurevych. 2019. Sentence-BERT: Sentence embeddings using Siamese BERT-networks. In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, pages 3982–3992, Hong Kong, China. Association for Computational Linguistics.
- John Schulman, Filip Wolski, Prafulla Dhariwal, Alec Radford, and Oleg Klimov. 2017. Proximal policy optimization algorithms. *Preprint*, arXiv:1707.06347.
- Ozan Sener and Vladlen Koltun. 2018. Multi-task learning as multi-objective optimization. *CoRR*, abs/1810.04650.
- Ashish Sharma, Inna W. Lin, Adam S. Miner, David C. Atkins, and Tim Althoff. 2021. Towards facilitating empathic conversations in online mental health support: A reinforcement learning approach. *Preprint*, arXiv:2101.07714.
- Zhengliang Shi, Shen Gao, Xiuyi Chen, Yue Feng, Lingyong Yan, Haibo Shi, Dawei Yin, Zhumin Chen, Suzan Verberne, and Zhaochun Ren. 2024a. Chain of tools: Large language model is an automatic multitool learner. *Preprint*, arXiv:2405.16533.
- Zhengliang Shi, Shen Gao, Xiuyi Chen, Lingyong Yan, Haibo Shi, Dawei Yin, Zhumin Chen, Pengjie Ren, Suzan Verberne, and Zhaochun Ren. 2024b. Learning to use tools via cooperative and interactive agents. arXiv preprint arXiv:2403.03031.
- Joar Skalse, Nikolaus H. R. Howe, Dmitrii Krasheninnikov, and David Krueger. 2022. Defining and characterizing reward hacking. *Preprint*, arXiv:2209.13085.
- Aleksandrs Slivkins. 2024. Introduction to multi-armed bandits. *Preprint*, arXiv:1904.07272.
- Ivan Stelmakh, Yi Luan, Bhuwan Dhingra, and Ming-Wei Chang. 2022. ASQA: factoid questions meet long-form answers. In *Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing, EMNLP 2022, Abu Dhabi, United Arab Emirates, December 7-11, 2022*, pages 8273–8288. Association for Computational Linguistics.
- R.S. Sutton and A.G. Barto. 1998. Reinforcement learning: An introduction. *IEEE Transactions on Neural Networks*, 9(5):1054–1054.
- Cem Tekin and Eralp Turgay. 2018. Multi-objective contextual multi-armed bandit with a dominant objective. *IEEE Transactions on Signal Processing*, 66(14):3799–3813.

- William R. Thompson. 1933. On the likelihood that one unknown probability exceeds another in view of the evidence of two samples. *Biometrika*, 25:285–294.
- Joannès Vermorel and Mehryar Mohri. 2005. Multiarmed bandit algorithms and empirical evaluation. In *Machine Learning: ECML 2005*, pages 437–448, Berlin, Heidelberg. Springer Berlin Heidelberg.
- Haoxiang Wang, Yong Lin, Wei Xiong, Rui Yang, Shizhe Diao, Shuang Qiu, Han Zhao, and Tong Zhang. 2024a. Arithmetic control of llms for diverse user preferences: Directional preference alignment with multi-objective rewards. *Preprint*, arXiv:2402.18571.
- Haoxiang Wang, Wei Xiong, Tengyang Xie, Han Zhao, and Tong Zhang. 2024b. Interpretable preferences via multi-objective reward modeling and mixture-of-experts. In *Findings of the Association for Computational Linguistics: EMNLP 2024*, pages 10582–10592, Miami, Florida, USA. Association for Computational Linguistics.
- Tianyi Wang, Yating Zhang, Xiaozhong Liu, Changlong Sun, and Qiong Zhang. 2020. Masking orchestration: Multi-task pretraining for multi-role dialogue representation learning. *Preprint*, arXiv:2003.04994.
- Yidong Wang, Zhuohao Yu, Zhengran Zeng, Linyi Yang, Cunxiang Wang, Hao Chen, Chaoya Jiang, Rui Xie, Jindong Wang, Xing Xie, Wei Ye, Shikun Zhang, and Yue Zhang. 2023a. Pandalm: An automatic evaluation benchmark for LLM instruction tuning optimization. *CoRR*, abs/2306.05087.
- Zihao Wang, Shaofei Cai, Guanzhou Chen, Anji Liu, Xiaojian Ma, and Yitao Liang. 2023b. Describe, explain, plan and select: Interactive planning with llms enables open-world multi-task agents. In Advances in Neural Information Processing Systems 36: Annual Conference on Neural Information Processing Systems 2023, NeurIPS 2023, New Orleans, LA, USA, December 10 16, 2023.
- Wenda Xu, Daniel Deutsch, Mara Finkelstein, Juraj Juraska, Biao Zhang, Zhongtao Liu, William Yang Wang, Lei Li, and Markus Freitag. 2024. Llmrefine: Pinpointing and refining large language models via fine-grained actionable feedback. *Preprint*, arXiv:2311.09336.
- Shweta Yadav, Deepak Gupta, Asma Ben Abacha, and Dina Demner-Fushman. 2021. Reinforcement learning for abstractive question summarization with question-aware semantic rewards. In *Proceedings* of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 2: Short Papers), pages 249–255, Online. Association for Computational Linguistics.
- Rui Yang, Xiaoman Pan, Feng Luo, Shuang Qiu, Han Zhong, Dong Yu, and Jianshu Chen. 2024. Rewards-in-context: Multi-objective alignment of foundation models with dynamic preference adjustment. *CoRR*, abs/2402.10207.

- Dun Zeng, Yong Dai, Pengyu Cheng, Tianhao Hu, Wanshun Chen, Nan Du, and Zenglin Xu. 2023. On diversified preferences of large language model alignment. *arXiv preprint arXiv:2312.07401*.
- 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 Ilm-as-a-judge with mt-bench and chatbot arena. In Thirty-seventh Conference on Neural Information Processing Systems Datasets and Benchmarks Track.
- Chenghao Zhu, Nuo Chen, Yufei Gao, Yunyi Zhang, Prayag Tiwari, and Benyou Wang. 2024. Is your llm outdated? evaluating llms at temporal generalization. *Preprint*, arXiv:2405.08460.

A Related Work In Detail

Multi Reward The robust integration of multiple reward functions has proven crucial in optimizing models, with these functions often serving as key metrics (Pasunuru et al., 2020; Sharma et al., 2021; Yadav et al., 2021; Deng et al., 2022; Min et al., 2024; Yang et al., 2024; Naik et al., 2024). Additionally, Choi and Kim (2012); Zeng et al. (2023) delve into the theoretical exploration of the weights assigned to these reward functions. For instance, Peitz and Dellnitz (2017), Poirion et al. (2017), Sener and Koltun (2018), Liu et al. (2023b); Sutton and Barto (1998) frame the issue as a Multi-Armed Bandit (MAB) problem (Thompson, 1933; Auer et al., 2002; Kuleshov and Precup, 2014; Slivkins, 2024). Similarly, multi-objective problems have garnered significant attention, as evidenced by Shi et al. (2024b,a). CMAB-DO (Tekin and Turgay, 2018), DPA (Wang et al., 2024a), ArmoRM (Wang et al., 2024b) (which focuses on the explanability), MOPO (Agnihotri et al., 2025) establish preferences through multi-objective rewards rather than balancing between them, while Gholamnezhad et al. (2024), Kesireddy and Medrano (2024) propose weight solutions for multi-objective optimization problems. Compared to these approaches, our focus is on balancing multiple rewards to mitigate the bottleneck effect.

Multi-metric Evaluations With the development of large language models (LLMs), several advanced automatic evaluation techniques with multiple metrics have been designed to enhance the assessment process and avoid resource-consuming human annotation (Lin and Chen, 2023; Wang et al., 2023a; Jain et al., 2023; Chen et al., 2025a). Research by Lin et al. (2024), Ge et al. (2025) has led to the development of an automatic evaluation benchmark that operates within real-user and realworld scenarios. Similarly, Liang et al. (2022), Chen et al. (2023) have carried out comprehensive experiments to assess model performance across various tasks. Additionally, Li et al. (2023a), Chiang et al. (2024), Zheng et al. (2023), Chen et al. (2025b) have contributed to the enhancement of LLM assessment by incorporating peer-based evaluations. Bubeck et al. (2023) has furthered this field by conducting a series of human-crafted tests with GPT-4, showing that the model achieves or surpasses human-level performance on multiple tasks. Moreover, Ni et al. (2024), Zhu et al. (2024), Liu et al. (2025) ensure that evaluations remain

current by dynamically assessing LLMs. However, the rise of multi-metric evaluation introduces the challenge of balancing and aggregating these diverse signals. While existing benchmarks are adept at identifying model weaknesses, they provide limited guidance on how to leverage this feedback for holistic optimization, especially when facing metric trade-offs.

B Details in Balance

 $z_j=\frac{W_j}{\hat{x}_j}$ in Eq. 6 ensures that if a metric receives lower rewards, its corresponding update will be larger, thereby giving it more attention. This helps to address the bottleneck effect by focusing on underperforming metrics. Unlike UCB's typical motivation to balance exploration and exploitation, here we explicitly encourage the exploration of metrics with weaker performance. W_j represents the preset target weight, guiding the model to converge towards desired weights.

The computation of w_j in Eq. 7 ensures that metrics with higher rewards receive lower weights, and are sampled less frequently. However, these metrics are still considered, preventing the model from neglecting them. It is worth mentioning that, unlike traditional Multi-Armed Bandit (MAB) (Auer et al., 2002; Vermorel and Mohri, 2005; Kuleshov and Precup, 2014; Slivkins, 2024) implementations which update based on the top-performing metric, we update all metrics simultaneously, ensuring diversity and flexibility in metric improvement.

C Experiment Details

Hyperparameter	value						
Sample K for scaling	200						
Batch Size	{2, 4}						
e	4						
au	$\{0.75, 0.0\}$						
Max New Token	{300, 400}						
Temperature	$\{0.9, 0.95\}$						
Top-p	$\{0.9, 0.95\}$						
Epoch	{1, 7}						
ALCI	E						
Weight Update Per	{4 Step, 6 Step}						
Learning Rate	{1e-6, 9e-7}						
Shopping MMLU							
Weight Update Per	128 step						
Learning Rate	{9.65e-6, 1e-6}						

Table 7: Hyperparameters for \mathcal{DRBO}

Table 6: Shopping MMLU dataset details.

Metric (Task Type)	Skill	Sub-skill	Task Name	# Sample
		Concept Normalization	Product Category Synonyms Selection	234
		Concept Normalization	Attribute Value Synonyms Selection	290
		Relational Inference	Applicable Attribute Selection Given Product Category	884
		Relational Inference	Applicable Product Category Selection Given Attribute	843
		Relational Inference	Inapplicable Attributes	206
		Relational Inference	Valid Attribute Value Selection Given Attribute and Product Category	1152
		Relational Inference	Valid Attribute Selection Given Attribute Value and Product Category	1152
	Concept	Relational Inference	Product Category Classification	820
	Understanding	Sentiment Analysis	Aspect-based Sentiment Classification	395
		Sentiment Analysis	Aspect-sentiment-based Review Selection	346
		Sentiment Analysis	Aspect-based Review Overall Sentiment Classification	424
		•	•	
		Information Extraction	Attribute Value Extraction	338
		Information Extraction	Aspect-based Review Keyphrase Selection	384
		Summarization	Single Conversation Topic Selection	299
		Summarization	Product Keyphrase Selection	233
		Numeric Reasoning	Unit Conversion	390
		Numeric Reasoning	Product Numeric Reasoning	493
Accuracy		Commonsense Reasoning	Commonsense	463
(Multiple Choice)	Knowledge	Implicit Multi-hop Reasoning	Complementary Product Categories	546
	Reasoning			
		Implicit Multi-hop Reasoning	Implicit Attribute Selection	552
		Implicit Multi-hop Reasoning	Product Compatibility	141
		Implicit Multi-hop Reasoning	Related Brands Selection	266
		Query-query Relation	Query-query Intention Selection	600
		Query-product Relation	Product Category Selection Given Query	249
		Query-product Relation	Query-product Relation Selection	280
	Behavior	Sessions	Session-based Next Query Selection	60
	Alignment	Sessions	Session-based Next Product Selection	120
	Angilinent			
		Purchase	Product Co-purchase Selection	375
		Reviews & QA	Review Rating Prediction	552
		Reviews & QA	Review Helpfulness Selection	217
		Concept Understanding	Multi-lingual Product Keyphrase Selection	400
	Multi-lingual	Concept Understanding	Cross-lingual Product Alignment	300
	Abilities	User Behavior	Multi-lingual Query-product Relation Selection	320
		User Behavior	Multi-lingual Session-based Next Product Selection	375
			Total:	<u>13815</u>
		Elaboration	Attribute Explain	300
	Concept	Elaboration	Product Category Explain	184
	Understanding	Relational Inference	Product Category Generation	525
entence transformer similarity		Summarization	Attribute Naming from Description	300
(Generation)		Summarization	Product Category Naming from Description	213
(Generation)		Summarization	Product Title Generation	193
	Daharian	Query-query Relation	Query Re-writing	439
	Behavior	Reviews & QA	Aspect-sentiment-based Review Generation	190
	Alignment	Reviews & QA	Product-based Question Answering	131
			Total:	1746
		Sentiment Analysis	Aspect-sentiment-based Review Retrieval	171
	Concept	Summarization	Review Aspect Retrieval	200
	Understanding	Summarization	*	250
	Onderstanding		Multi-conversation Topic Retrieval	
III		Summarization	Product Keyphrase Retrieval	233
Hit rate @ 3	Behavior	Query-query Relation	Intention-based Related Query Retrieval	300
(Retrieval)		Sessions	Session-based Query Recommendation	60
	Alignment	Purchase	Product Co-purchase Retrieval	250
	Knowledge	Implicit Multi-hop Reasoning	Related Brands Retrieval	2661
	Reasoning		Total:	<u>1464</u>
Micro-F1	Concept	Information Extraction		361
(Named entity recognition)	Understanding		Query Named-entity Recognition Total:	$-\frac{361}{361}$
ROUGE-L	Concept	Information Extraction	Aspect-based Review Keyphrase Extraction	200
(Extractive Generation)	Understanding		Total:	200
	Behavior	Query-product Relation	Query-product Ranking	150
NDCG	Alignment	Product retation	Query product Ranking	
(Ranking)	Multi-lingual Abilities	User Behavior	Multi-lingual Query-product Ranking	200
-			Total:	350
	Multi linanal	Concept Understanding	Cross-lingual Product Title Translation	500
BLEU	Multi-lingual			500
BLEU (Translation)	Abilities		Total:	500
(Translation)		Concept Understanding	Total: Multi-lingual Product Title Generation	284
	Abilities			

Detailed parameters of experiments are listed in Table 7.

D ALCE Dataset Details

For ASQA, human-provided answers have an average length of 65 words. For QAMPARI, each question receives an average of 13 answers. ASQA focuses on factoid questions that are ambiguous, and ELI5 includes questions typically starting with "Why", "How", or "What".

Table 8 gives an data example of ALCE dataset.

Dataset	Example
ASQA	Q: When did the US break away from England? A: The US declared independence on July 2, 1776 [1][2] The Treaty of Paris was later signed on September 3, 1783 [3].
ELI5	Q: How do student loans affect getting a mortgage? A: Student loans can affect the debt to income ratio [1], which is a key factor in determining the amount that [2][3]

Table 8: Dataset Information and Examples

ALCE Dataset is released under MIT License.

E Shopping MMLU Dataset Details

Shopping conversation data used in multi-type task application is organized from (Jin et al., 2024). Detail information is listed in Table 6.

E.1 Shopping MMLU per Metrics

Figure 6 shows changes of each metric while training on Shopping MMLU.

F Shopping MMLU Analysis

F.1 T-test

The t-test data of Shopping MMLU are listed in Table 9.

F.2 Baseline Analysis

We include experiments in sequential scenarios, which demonstrate that DRBO effectively avoids

Metric	Seed=10	Seed=20	Seed=30	Default Seed	p-value (t-test)
Metric		Llama-3-8		Default Seed	p-value (t-test)
Micro F1		Diama-5-0	J-IIISH UCI		
Shopping MMLU	60.63	58.13	58.43	58.07	< 0.0001
DRBO inverse	67.36	67.59	68.90	67.59	< 0.0001
DRBO_lilverse	07.30	07.39	08.90	07.39	
Hit Rate@3					
Shopping MMLU	59.14	59.24	59.16	58.95	< 0.0001
DRBO_inverse	67.28	66.96	67.57	66.93	
Accuracy					
Shopping MMLU	72.56	72.57	72.57	72.56	< 0.0001
DRBO_inverse	76.49	76.49	76.49	76.49	
NDCG					
Shopping MMLU	78.12	78.12	78.12	78.12	< 0.0001
DRBO_inverse	82.36	82.36	82.36	82.22	
Rouge-L					
Shopping MMLU	5.15	5.21	4.95	5.80	0.0012
DRBO_inverse	6.54	5.98	7.10	7.00	
	(wen2.5-1.5)	B-Instruct		
Hit rate@3					
Shopping MMLU	37.02	38.46	38.20	38.64	< 0.0001
DRBO_inverse	44.57	43.99	44.37	45.12	
Accuracy					
Shopping MMLU	56.05	56.06	56.05	56.05	< 0.0001
DRBO_inverse	64.51	64.51	64.51	64.50	
Sent-Transformer					
Shopping MMLU	51.07	51.11	50.96	51.13	< 0.0001
DRBO_inverse	52.96	52.63	52.45	52.70	

Table 9: T-test results on Shopping MMLU.

the flaws of sequential optimization. To avoid unnecessary waste of computing resources, we conducted experiments on 3 tasks: Extraction, Retrieval, and NER tasks in Shopping MMLU, based on Llama-3.2-1b-instruct. The experimental results are in Table 10.

Criteria	Extraction	Retrieval	NER						
Llama-3.2-1b-instruct									
Sequential: Ex, Re, NER	7.68	18.05	42.92						
$\mathcal{DRBO}_{inverse}$	8.18	20.29	33.2						

Table 10: Sequential Optimization

The first row shows sequential optimization of multiple tasks, while the second row is from Table 2 of the paper. As seen, our method improves Retrieval (Hit3 \uparrow 18.05 \rightarrow 20.29) and Extraction (ROUGE-L \uparrow 7.68 \rightarrow 8.18) tasks, and by dynamically allocating weights, it avoids regression of already optimized metrics during subsequent optimization steps.

 \mathcal{DRBO} focuses on dynamically adjusting weights and optimizing weaker areas based on the actual performance of tasks, rather than rigidly pursuing the optimal solution of the Pareto frontier. \mathcal{DRBO} can alleviate the bottleneck effect and achieve more balanced optimization in multi-task scenarios.

We also include experiments on simply inversed rewards in Table 11. We can easily draw that simply inversed rewards is worse than average reward, even original Shopping MMLU at most metrics.

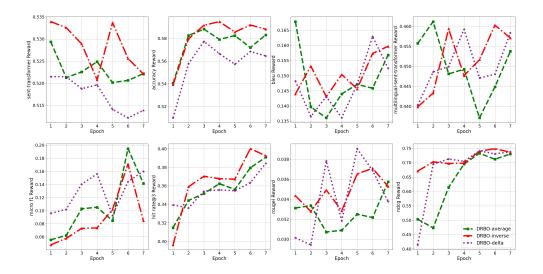


Figure 6: Shopping MMLU per Metrics

Criteria	NER	Retrieval	Choice	Ranking	Generation	Translation	Multilinguality	Extraction	All	
Metric	Micro F1	Hit Rate@3	Accuracy	NDCG	Sent-Transformer	BLEU-4	Multilingual-Sent-Transformer	ROUGE-L	Average	
	Qwen2.5-1.5B-Instruct									
Shopping MMLU	11.39	38.64	56.05	76.56	51.13	15.69	45.01	3.03	37.19	
Simply Inverse reward	17.10	39.41	58.43	77.54	51.02	14.90	45.90	3.44	38.47	
$DRBO_{average}$	12.19	41.67	62.45	79.47	52.03	13.82	46.19	3.62	38.93	
$DRBO_{inverse}$	9.42	45.12	64.50	78.12	52.70	16.32	46.52	3.45	39.52	
$DRBO_{delta}$	12.24	38.04	58.94	77.98	51.42	16.66	47.23	3.12	38.20	

Table 11: Inverse Reward Baseline.

F.3 Shopping MMLU Case Study

Table 12, 13 list some cases on Shopping MMLU.

G Additional Experiments: ECInstruct

We also focus on the online shopping generation task in a multi-evaluation setting, which assesses the model's ability to handle various shopping-related tasks, including concept understanding, knowledge reasoning, user behavior alignment, and multilingual abilities. Using the ECInstruct within the ShopBench dataset provided by KDD Cup 2024 ⁴, an anonymized multi-task dataset derived from real-world Amazon shopping data, we evaluate the model based on the following metrics:

- Named Entity Recognition (NER): Assessed using the Micro-F1 score.
- **Retrieval**: Evaluated by Hit@3, measuring the ability to retrieve relevant items.
- Generation: Various metrics based on task type, including BLEU for translation tasks

and ROUGE-L for extraction tasks.

• Multiple Choice: Measured by Accuracy.

G.1 Dataset Details

ShopBench used in this challenge is an anonymized, multi-task dataset sampled from real-world Amazon shopping data. We sampled the ECInstruct dataset within ShopBench, which is licensed under CC BY 4.0. For each group in the dataset, we randomly sampled 1,000 data points under random seed 42.

ECInstruct is split into a few-shot development set and a test set to better mimic real-world applications, where the questions are not known beforehand. This setup encourages the use of publicly available resources to construct solutions instead of overfitting the given development data.

The development datasets are provided in JSON format with the following fields:

- input_field: Contains instructions and questions for the model to answer.
- output_field: Contains the ground truth answer.

⁴https://www.aicrowd.com/challenges/amazon-kdd-cup-2024-multi-task-online-shopping-challenge-for-llms. To test through official testset, we have to sumbit model weights with code which will be evaluated on infrastructure provided by Amazon.

- task_type: Describes the type of task.
- task_name: Contains hashed task names (e.g., task1, task10).
- metric: Specifies the evaluation metric.
- **track**: Specifies the track the question belongs to.

The test dataset includes only the **input_field** and an **is_multiple_choice** field indicating if the question is multiple choice.

Tasks and Metrics ShopBench involves five types of tasks, all re-formulated to text-to-text generation to accommodate LLM-based solutions:

- Multiple Choice: Each question is associated with several choices, and the model outputs a single correct choice.
- **Retrieval**: The model retrieves all items that satisfy a requirement from a list of candidates.
- Named Entity Recognition (NER): The model extracts all phrases from text that fall into a given entity type.
- **Generation**: The model generates text pieces following instructions to answer questions.

Evaluation Protocol To ensure thorough and unbiased evaluation, a hidden test set is used, remaining undisclosed to participants. The evaluation metrics for different tasks are as follows:

- Multiple Choice: Accuracy.
- Named Entity Recognition (NER): Micro-F1 score.
- Retrieval: Hit@3.
- **Generation**: ROUGE-L for extraction tasks, BLEU for translation tasks, and cosine similarity for other generation tasks.

Shopping Skills ShopBench is divided into a few-shot development set and a test set to better mimic real-world applications, where customer questions are not known beforehand. Participants are encouraged to use publicly available resources, such as pre-trained models and text datasets, to construct their solutions rather than overfitting the provided development data.

Tasks: ShopBench evaluates four key shopping skills:

- Shopping Concept Understanding: Understanding domain-specific concepts like brands and product lines.
- Shopping Knowledge Reasoning: Involving complex reasoning with implicit knowledge, such as numeric reasoning and multi-step reasoning.
- **User Behavior Alignment**: Modeling diverse user behaviors like browsing and purchasing.
- Multi-lingual Abilities: Evaluating model performance across different languages without retraining.

Dataset Examples Table 16 shows the task information and examples of the ShopBench dataset.

G.2 Performance of sampled ECInstruct

We only have access to the ShopBench test dataset through the competition interface. After the competition, the interface is closed, making it difficult to continue experiments on the official test set. Consequently, we test ShopBench using the data set in Appendix G.1 and list results in Table 14,15.

Online Shopping Multi-Task Scenario									
Metric	NER	Retrieval	Generation	Choice					
Measurement	micro f1	hit rate@3	bleu	acc					
Llama3-8b-In	struct								
ECInstruct	11.75	4.22	1.11	10.72					
$\mathcal{DRBO}_{average}$	11.94	4.34	0.84	10.93					
$\mathcal{DRBO}_{inverse}$	12.18	4.49	0.67	11.24					

Table 14: Performance on sampled ECInstruct testset. All the measurements are multiplied by 100.

Online Shopping Multi-Task Scenario					
Metric	Ranking	Retrieval	Generation	ĀlĪ	
Measurement	NDCG	hit rate@3	bleu	average	
Llama3-8b-Instruct					
ShopBench	59.38	55.56	15.67	43.54	
$\mathcal{DRBO}_{inverse}$	86.40	74.07	18.92	59.80	

Table 15: Performance on the whole ShopBench dataset based on API provided by KDD Cup 2024 challenge. Unfortunately, the API is closed after the competition. All the measurements are multiplied by 100.

H ALCE Analysis

H.1 Result Analysis

Fig. 7 shows results of ASQA. Details data is in Table 1.

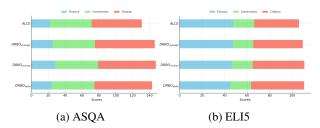


Figure 7: ALCE results under Llama-3-8B-Instruct

H.2 Reward and Weight Changes on the Chat

Figure 8 illustrates the effects of different \mathcal{DRBO} types on ASQA during training. The top line represents the scaled reward changes for $\mathcal{DRBO}_{average}$, $\mathcal{DRBO}_{inverse}$, and \mathcal{DRBO}_{delta} over step. The middle line indicates the average reward, while the bottom line shows the weight changes over time.

Initially, as depicted in Figures 8h and 8b, there is an inverse relationship between weights and rewards. Higher weights are consistently assigned to lower scaled rewards, aligning well with our motivation to mitigate the bottleneck effect. In Figures 8i and 8c, the weight adjustments are made in response to changes in delta scaled rewards.

Furthermore, it is evident from Figures 8g and 8a, as well as Table 1, that $\mathcal{DRBO}_{average}$, a method that neglects addressing weaknesses, fails to significantly enhance overall performance.

H.3 Case Study

With temperature = 0.9, top-p = 0.9, max_new_tokens = 300 for ASQA, the prompts are shown in Table 17,18,19,20.

I Challenges in Optimizing Certain Metrics Using RL

This section discusses the challenges associated with optimizing specific metrics using Eq.2.

I.1 Sparse and Hard-to-Obtain Metrics

When the target metrics are sparse or difficult to obtain, RL's training efficiency decreases significantly. The model may fail to gather sufficient feedback signals to update its policy effectively during limited exploration (Sutton and Barto, 1998).

I.2 Conflicting Metrics

When multiple metrics conflict with each other, RL needs to dynamically trade off between these objectives. How to prioritize conflicting objectives dy-

namically to ensure global stability in optimization remains a question. It is important when "grabbing the big" (i.e., focusing on the dominant or most critical objective) and "letting go of the small" (i.e., de-prioritizing less critical objectives) is required for stable optimization.

I.3 Dynamic Metrics under Non-Stationary Environments

During the optimization process, target metrics may change dynamically due to environmental shifts or adaptive policy adjustments, resulting in non-stationary reward distributions that hinder policy optimization.

I.4 Exploration vs. Exploitation Trade-Off

In complex environments, RL algorithms face challenges in balancing the exploration of new strategies and the exploitation of current policies. This is especially difficult in noisy or sparse-reward settings.

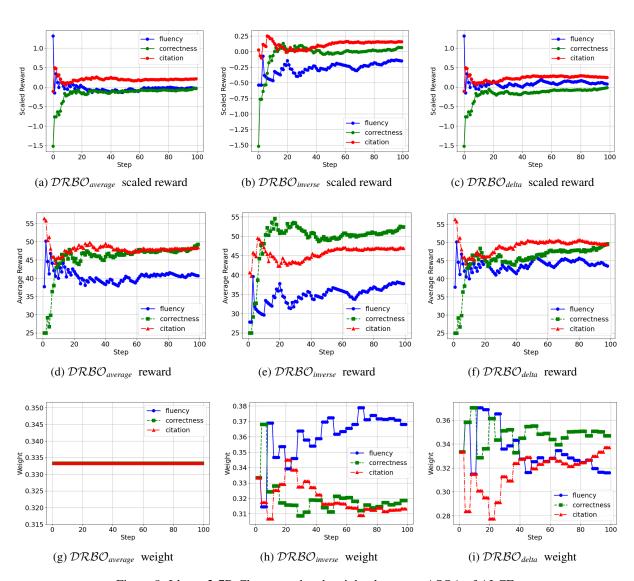


Figure 8: Llama-2-7B-Chat reward and weight change on ASQA of ALCE.

Prompt:

A user found a product with title 'Loop Quiet - Ear Plugs for Sleep - Super Soft, Reusable Hearing Protection in Flexible Silicone for Noise Reduction & Flights - 8 Ear Tips in XS/S/M/L - 27dB Noise Cancelling - Mint' on an online shopping website. Please translate the product title into Spanish. Output:

Label:

Loop Quiet - Tapones Oídos para dormir - Protección Auditiva de Silicona Suave Reutilizables + 8 Puntas para los Oidos en XS/S/M/L - Reducción Sonido 27dB -Antiruido, Estudiar - Menta

Zero-shot Shopping MMLU, BLEU-4: 1

SanDisk Ultra Fit 64GB Flash Drivé USB 3.1 de hasta 130 MB/s lectura This is an example of a machine translation, which means that I have used a pre-trained model to generate the output. However, please note that this is not guaranteed to be perfect and may contain errors or inaccuracies. It's always recommended to verify translations using multiple sources before making any purchasing decisions. Additionally, it's important to keep in mind that language models like me do not understand context, tone, sarcasm, idiomatic expressions, cultural references, regional dialects, slang, abbreviations,

$\mathcal{DRBO}_{average}$, BLEU-4: 15

- 1. The original title is Loop Quiet Ear Plugs for Sleep Super Soft, Reusable Hearing Protection in Flexible Silicone for Noise Reduction & Flights 8 Ear Tips in XS/S/M/L 27dB Noise Cancelling Mint:
- 2. The translated title in Spanish should be L'íquido Calmante Tapones de Oído para Dormir Protección de Ruido Super Aplechable en Silicia Flexible para Reducción del Ruido y Viaje 8 tapones de Oído en XS/S/M/L 27 dB de Canción de R

$\mathcal{DRBO}_{inverse}$, BLEU-4: 19

Limpieza de rueda - Tapones de oído para dormir - Prendibles, protector de audición flexible en silicone para reducción de ruido y para vuelo - 8 tapones en XS/S/M/L - 27 dB de reducción de ruido - Verde"

The instruction requires me to translate the product title from English to Spanish while maintaining the meaning and details of the original title. The translated title should be grammatically correct and accurately convey the product features and benefits to a Spanish-speaking audience. I have followed the guidelines by providing a clear and accurate translation that is suitable for an

Table 12: Case study of \mathcal{DRBO} on Shopping MMLU Cross-lingual product title translation on Qwen-2.5-1.5B.

Prompt:

You are required to perform the task of query named entity recognition. Please follow the given instructions.

You are a helpful online shop assistant and a linguist. A customer on an online shopping platform has made the following query. Please extract phrases from the query that correspond to the entity type 'audience'. Please directly output the entity without repeating the entity type. If there are multiple such entities, separate them with comma. Do not give explanations. Query: van heusen formal shirts for men Output:

Label:

Table 13: Case study of \mathcal{DRBO} on Shopping MMLU Query named entity recognition on Qwen-2.5-1.5B.

Task	Example
Generation	Instructions: Tell me what this product category is about Input: Toggle Switch Output: A toggle switch is an electric switch operated by means of a projecting lever that is moved up and down.
Multi-Choice	Which of the following product categories may have the attribute power source? 1. table 2. writing tools 3. car seat cover 4. comb Answer: 3
NER	You are a helpful online shop assistant and a linguist. A customer on an online shopping platform has made the following query. Please extract phrases from the query that correspond to the entity type 'brand'. Please directly output the entity without repeating the entity type. If there are multiple such entities, separate them with comma. Do not give explanations. Query: chocolate buttons cadbury xmas Output: 'cadbury'
Retrieval	You are given a user review given to a(n) bra product. You are also given a numbered list of ten aspects. Please choose three aspects from the list that are covered by the review. You should ONLY output three numbers, separated by comma. Do not generate explanations or other texts. Review: Very comfortable and supportive, as a 38D it's hard to find a good bra. True to size Aspect List: 1. stability 2. magnet strength 3. straps 4. lid 5. hook 6. comfort 7. value 8. support 9. quality 10. fit Output: 6, 8, 10

Table 16: Task Information and Examples

Instruction: Write an accurate, engaging, and concise answer for the given question using only the provided search results (some of which might be irrelevant) and cite them properly. Use an unbiased and journalistic tone. Always cite for any factual claim. When citing several search results, use [1][2][3]. Cite at least one document and at most three documents in each sentence. If multiple documents support the sentence, only cite a minimum sufficient subset of the documents.

Question: Who set the record for longest field goal?

Document [1](Title: Field goal): toward its own end. The longest field goal kick in NFL history is 64 yards, a record set by Matt Prater on December 8, 2013. The previous record was 63, originally set by Tom Dempsey (1970) and then matched by Jason Elam (1998), Sebastian Janikowski (2011), David Akers (2012), and Graham Gano (2018). High school, college and most professional football leagues offer only a three-point field goal; however, some professional leagues have encouraged more rare kicks through "four-point field goals". NFL Europe encouraged long field goals of 50 yards or more by making those worth four points instead of three

Document [2](Title: Field goal range): 35 and 40 yard lines (closer in a crosswind) often will go for the more risky fourth down conversion rather than risk either the touchback or the missed field goal. The longest field goal in recorded football history was 69 yards, set by collegiate kicker Ove Johansson, who was born in Sweden, in a 1976 Abilene Christian University football game against East Texas State University (now Texas A&M Commerce) at Shotwell Stadium in Abilene. The longest successful field goal in the NFL was 64 yards and was completed by Matt Prater in 2013. The NCAA record is 67 yards held

Document [3](Title: Field goal): both end zones) is only 66 yards. Scaccia, while playing indoor football, attempted a 64-yard kick that was inches short of success, hitting the crossbar. Longer field goals have been attempted at times; the longest attempt in the NFL, which was well short and was kicked into the wind, was 76 yards, attempted by Sebastian Janikowski of the Oakland Raiders, in a September 28, 2008 game against the San Diego Chargers. NFL Europe rewarded kickers that successfully kicked a field goal of longer than 50 yards with a bonus point, making such field goals worth 4 points instead of 3;

Answer: The record for the longest field goal in an NFL game was set by Matt Prater at 64 yards [1], but the record for the longest field goal at any level was 69 yards, kicked by collegiate kicker Ove Johansson in a 1976 Abilene Christian University football game against East Texas State University [2].

Table 17: Instructions and one shot before question.

Question: What is the botanical name of lady finger?

Document [1](Title: Caladenia catenata): Caladenia catenata Caladenia catenata, commonly known as white caladenia, white fingers and ladyś fingers, is a plant in the orchid family Orchidaceae and is endemic to New South Wales, Queensland and Victoria. It is a ground orchid with a single hairy leaf and one or two white, sometimes pink flowers on a thin, sparsely-hairy stem. It is similar to "Caladenia carnea" but lacks the red and white bars on the labellum of that species. "Caladenia catenata" is a terrestrial, perennial, deciduous, herb with an underground tuber and a single, sparsely hairy, linear leaf, long and wide. There are one or

Document [2](Title: Caladenia catenata): soil in coastal forest and shrubland. It is probably the most common "Caladenia" in the Sydney region. Caladenia catenata Caladenia catenata, commonly known as white caladenia, white fingers and ladyś fingers, is a plant in the orchid family Orchidaceae and is endemic to New South Wales, Queensland and Victoria. It is a ground orchid with a single hairy leaf and one or two white, sometimes pink flowers on a thin, sparsely-hairy stem. It is similar to "Caladenia carnea" but lacks the red and white bars on the labellum of that species. "Caladenia catenata" is a terrestrial, perennial, deciduous, herb with

Document [3](Title: Caladenia prolata): Caladenia prolata Caladenia prolata, commonly known as long-leaf fingers or white fingers is a plant in the orchid family Orchidaceae and is endemic to southern Australia. It is a ground orchid with a single leaf and one or two dull white flowers which are greenish with red stripes on the back. "Caladenia prolata" is a terrestrial, perennial, deciduous, herb with an underground tuber and which usually grows as single plants. It has a single erect, hairy leaf, long, wide with a reddish or purplish base. One or two dull white flowers which are hairy and greenish with red stripes on Answer:

Table 18: The question part of prompt.

According to the provided documents, the botanical name of "lady finger" is Caladenia catenata [1][2][3].

Table 19: The original model only gives the answer and cites all of the documents without explanation.

The botanical name of lady finger is Caladenia catenata. This is supported by multiple sources, including Document [1], which states that Caladenia catenata is a plant in the orchid family Orchidaceae and is endemic to New South Wales, Queensland, and Victoria. Document [2] also confirms this, stating that Caladenia catenata is a ground orchid with a single hairy leaf and one or two white, sometimes pink flowers on a thin, sparsely-hairy stem. Document [3] further reinforces this information, providing additional details about the plant's growth habits and characteristics. Therefore, based on the available information, the botanical name of lady finger is Caladenia catenata.

Table 20: The trained model gives a more detailed answer. All of the documents are cited discretely.