

REFINING BIAS AND REWARD IN LLM RECOMMENDER AGENTS THROUGH META-CONTROLLED TOOL INVOCATION

Anonymous authors

Paper under double-blind review

ABSTRACT

Large language model (LLM) agents have recently been brought to recommender systems given their flexible capability of tool use. Although existing approaches adopt the reasoning and acting paradigms for profiling, planning, and memory augmentation, they remain ad hoc and overlook core recommendation challenges in agent-environment interactions, including debiasing and reward estimation in offline learning scenarios. In this paper, we introduce BARO (Bias And Reward Optimization), a meta-controlled, tool-augmented LLM agent framework that explicitly addresses these challenges. BARO employs a two-stage recommendation process: a coarse recommender generates a candidate slate based on user history, and a meta-controller adaptively invokes three specialized tools to refine the recommendation results: a bias detector assesses and mitigates bias in the candidate set, a reward estimator calibrates noisy offline rewards, and an action grounder selects final recommendations from the candidate pool. This design injects bias correction and reward refinement directly into the agent’s decision loop in the recommendations. Empirical results on two benchmark datasets demonstrate that BARO achieves consistent improvements over state-of-the-art methods in metrics such as accuracy, diversity, and fairness. The code will be made publicly available upon acceptance.

1 INTRODUCTION

Recommender systems (RecSys) are the cornerstone of modern digital platforms, such as e-commerce, streaming services, and short video applications (Bobadilla et al., 2013; Burke, 2002), where a key task is to align recommendations with user preferences and improve long-term user satisfaction (Wang et al., 2024a). In modern RecSys, semantic information about users and items has been a key factor in effective recommendations, for example, the profile of a user or the genre of a movie (Harper & Konstan, 2015). However, how to meaningfully leverage this information remains a challenge.

Recently, large language models (LLMs), with their remarkable contextual understanding and reasoning abilities (Touvron et al., 2023; Guo et al., 2025), have been applied to RecSys to capture complex patterns between users, items, and their interactions. These methods can be grouped into three paradigms (Wu et al., 2024). (1) In-context learning (ICL) approaches (Dong et al., 2024) treat LLMs as powerful general purpose functions, extracting user preferences via prompt engineering without modifying model parameters (Liu et al., 2023). Although convenient, ICL methods inherit hallucinations and biases from general-purpose LLMs, since these models are not built for recommendation-specific tasks (Gao et al., 2025a; Cai et al., 2025a). (2) Post-training approaches (Gao et al., 2025b; Hu et al., 2022) fine-tune LLMs in RecSys tasks to alleviate hallucination and bias issues. However, post-training is computationally costly, and the resulting models often lack extensibility to new domains or objectives (Mohammadi et al., 2025; Cai et al., 2025b; Huang et al., 2025). (3) In contrast, LLM agents (Wang et al., 2025b; Ning et al., 2024; Wang et al., 2024b) further extend flexibility by enabling tool use, allowing recommendation pipelines to integrate auxiliary

047 modules beyond static LLM predictions. This flexibility of LLM agents has motivated a surge of effective
048 agent frameworks for RecSys. One branch of work (Zhang et al., 2024c; Shi et al., 2024b) optimizes a
049 memory module to capture long-term engagement, while others (Zhang et al., 2024a; Gong et al., 2024)
050 strengthen planning capabilities to enhance reasoning before micro-level recommendations.

051 Despite the promise of LLM agents in RecSys, most of these agentic frameworks were not designed specif-
052 ically for RecSys, leading to overlooking the interaction bias and the user interaction offline reward issues.
053 Generally, a recommendation process is framed as a sequential decision-making problem, where user prefer-
054 ence can be viewed as implicit rewards (Chen et al., 2023b;a; Xue et al., 2023). Yet current approaches focus
055 mostly on agent memory and planning via prompting the LLMs, in which **biases** and **reward inaccuracies**
056 in offline environments remain largely unaddressed. These issues are critical: biases such as popularity bias
057 or filter bubbles (Gao et al., 2023a;b) can harm fairness and diversity, while inaccurate reward estimation
058 can mislead optimization. Moreover, integrating LLM agents into RecSys introduces additional challenges,
059 including sycophantic behaviors (Pan et al., 2024) and hallucination risks, which exacerbate the inherent
060 bias and reliability concerns.

061 To tackle these challenges, we propose **BARO** (*Bias And Reward Optimization*), a meta-controlled, tool-
062 augmented LLM agent specifically designed for recommendations. BARO adopts a two-stage recommen-
063 dation pipeline. First, an in-context learning-based coarse recommender produces candidate slates from
064 user interaction histories. Then, a **Meta-Controller** governs tool use to refine the slate, invoking three
065 specialized modules: (i) a **Bias Detector**, which assesses the bias level of the candidate set and triggers
066 regeneration when biases exceed thresholds; (ii) a **Reward Estimator**, which evaluates uncertainty and cor-
067 rects noisy offline rewards; and (iii) an **Action Grounder**, which integrates information from all modules
068 to make fine-grained final recommendations. By explicitly embedding debiasing and reward refinement into
069 the decision-making loop, BARO provides a RecSys-specific design that moves beyond generic reasoning
070 and acting agents, yielding state-of-the-art performance on reliable, fair, and accurate recommendations with
071 evaluations on two widely used benchmarks. Our contributions are threefold:

- 072 • We highlight the limitations of existing ReAct-style LLM agents in recommendation and formulate
073 the overlooked challenges of biases and reward inaccuracies in offline environments, motivating a
074 RecSys-specific agent framework.
- 075 • We propose BARO (*Bias And Reward Optimization*), a meta-controlled, tool-augmented LLM agent.
076 BARO integrates a Bias Detector, Reward Estimator, and Action Grounder under a meta-controller,
077 enabling explicit debiasing and reward calibration within a two-stage pipeline.
- 078 • We conduct experiments on two benchmark datasets, comparing BARO against both traditional se-
079 quential recommenders and SOTA LLM-based methods. Results show consistent improvements
080 across diverse metrics, demonstrating the effectiveness and generality of our approach.

081 2 RELATED WORK

083 **Large Language Models for Recommendations.** Recent LLMs have been integrated into recommenda-
084 tion systems and show remarkable potential in downstream tasks (Chen et al., 2025; Bao et al., 2023a; Li
085 et al., 2024). Early studies explore in-context learning to leverage LLMs for recommendation tasks without
086 parameter updates (Bao et al., 2025; Liu et al., 2023; Bao et al., 2023a; Shi et al., 2024b; Gao et al., 2023c; Ye
087 et al., 2024; Shi et al., 2024a; Bao et al., 2023b; Liao et al., 2023; Lin et al., 2024). Yet ICL-based methods
088 suffer from hallucination and bias from their pretraining knowledge (Gao et al., 2025a; Cai et al., 2025a), and
089 fine-tuning methods lack flexibility in adapting to new tasks (Mohammadi et al., 2025). BARO alleviates
090 these issues by integrating LLM agents with a meta-controller, enabling more accurate recommendations
091 under biased and implicit feedback.

092 **Large Language Model Agents.** Recent LLM-based agents with autonomous decision-making processes
093 have been used to tackle complex tasks (Wang et al., 2025b; Ning et al., 2024; Cai et al., 2025a). LLM agents

can dynamically explore complex user preferences and perform multi-step planning to generate personalized recommendations (Li et al., 2024; Bao et al., 2023a; Ning et al., 2024; Huang et al., 2025; Zhao et al., 2024; Wang et al., 2025b; 2024b; Zhang et al., 2024b; Shu et al., 2024; Zhang et al., 2024a;c). Yet these existing methods often overlook key challenges such as bias and reward inaccuracies. Our proposed BARO introduces a RecSys-specific agent design with dedicated Bias Detector and Reward Estimator modules, providing accuracy, fairness, and robustness in recommendations.

3 PRELIMINARIES

Model-based Offline Reinforcement Learning. Reinforcement learning aims to learn policies for solving sequential decision-making problems in Markov Decision Processes (MDPs) (Sutton & Barto, 2018). An MDP is formulated as a 5-tuple $\langle \mathbf{S}, \mathbf{A}, T, R, \gamma \rangle$, where \mathbf{S} denotes the state space and \mathbf{A} represents the set of feasible actions. At timestep t , after taking an action \mathbf{a}_t in state \mathbf{s}_t , the system transitions to the next state \mathbf{s}_{t+1} and receives a reward r_t . These dynamics are captured by the transition function T , while the reward is determined by the reward function R . The objective of RL is to maximize the discounted cumulative reward $G = \sum_{t=0}^{\infty} \gamma^t r_t$, where $\gamma \leq 1$ is the discount factor that balances immediate and future rewards.

While RL typically learns policies through extensive interactions with environments, offline RL instead derives policies solely from pre-collected logs, mitigating the issue of costly interactions in applications such as healthcare and recommender systems (Yu et al., 2024; Chen et al., 2023a). Denoting the offline logs as a set of transitions $\mathcal{D} = \{(s_t, a_t, s_{t+1}, r_t)\}$, model-based offline RL leverages \mathcal{D} to learn a world model that serves as a proxy environment, allowing policies to be trained through interactions with this learned model. A world model usually consists of an estimated transition function, *i.e.*, \hat{T} , and an estimated reward function, *i.e.*, \hat{R} . Although model-based offline RL has demonstrated high sample efficiency and strong performance, a key challenge remains in the inaccuracy of the learned world models.

Problem Formulation. Due to the interactive nature and sequential decision-making paradigm, recommendation tasks can be naturally formulated as MDPs and addressed with RL (Sutton & Barto, 2018; Yu et al., 2024). A state $\mathbf{s} \in \mathbf{S}$ represents the state of the system prior to recommendation, typically including user information, (*e.g.*, recent interactions and side features) as well as item attributes. An action $\mathbf{a} \in \mathbf{A}$ corresponds to a recommended item, and the reward r is derived from user feedback—for example, a movie rating on streaming platforms or the view-time ratio of a short video (Gao et al., 2022a;b). The transition function is often parameterized by sequential models such as Transformers (Vaswani et al., 2017), which encode a fixed-length history of past interactions together with the current interaction to produce the next state. Further, long-term user engagement is commonly captured by the discounted cumulative reward.

4 METHOD

We propose **BARO** (**B**ias **A**nd **R**eward **O**ptimization), a meta-controlled and tool-augmented LLM agent for recommendation. BARO comprises six components: (i) a world model for simulation, (ii) an ICL-based coarse recommender, (iii) a meta-controller that orchestrates tool use, (iv) a bias detector, (v) a reward estimator, and (vi) a fine-tuned action grounder. We next describe how these modules cooperate and learn.

4.1 WORLD MODEL SIMULATION

We adopt a model-based offline RL setup, where a learned reward function \hat{R} serves as the world model to evaluate candidate items without interacting with real users. Following prior works (Gao et al., 2023b;a; Zhang et al., 2024d), \hat{R} is trained using offline logs and a supervised recommendation model, such as DeepFM (Guo et al., 2017). The reward function \hat{R} is then frozen and used to provide immediate feedback for new recommendations. We do not model the transition function \hat{T} explicitly, as the current interaction and history are already fed into LLM agents for reasoning.

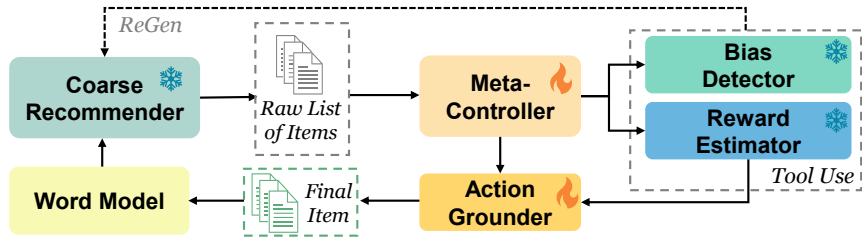




Figure 1: The overall architecture of our proposed **BARO** framework. The process begins with a frozen Coarse Recommender generating an initial *Raw List of Items*. The core of our agent, the Meta-Controller, evaluates this list and decides whether to (1) invoke specialized tools—the Bias Detector or the Reward Estimator—for in-depth analysis, or (2) directly pass the list to the Action Grounder. If significant issues are detected by the tools, a regeneration signal (*ReGen*) prompts the Coarse Recommender to refine the list. Finally, the Action Grounder selects the *Final Item* from the satisfactory candidate list. Components marked with  are trainable agents optimized via reinforcement learning, while those with  are frozen.

After training the world model, following ROLeR (Zhang et al., 2024d), we can obtain a refined reward function \tilde{R} . We estimate uncertainty using the absolute residual between the two rewards:

$$u(s, a) = |\tilde{R}(s, a) - \hat{R}(s, a)|, \quad (1)$$

which captures the disagreement between the original world model and the shaped reward, serving as a proxy for reward uncertainty.

4.2 COARSE RECOMMENDER

The coarse recommender provides a candidate item set as a preliminary recommendation. To keep it simple, we choose an in-context learning design that employs a frozen LLM. The inputs of the coarse recommender consist of (i) recent interaction histories $h_t = [i_{t-w}, \dots, i_{t-1}]$, where w is the window size and t represents the current time step; (ii) the number of regenerations ($\#ReGen$); and (iii) the guidance, g , from the bias detector or the reward estimator. The output of the coarse recommender is a candidate set \mathcal{I}_t :

$$\mathcal{I}_t = f_{CR}(h_t, \#ReGen, g). \quad (2)$$

4.3 META-CONTROLLER

The meta-controller serves as the core orchestrator in BARO. It comprehensively evaluates the quality of the candidate set—including overall bias, reward estimation, and user satisfaction—and then calls the appropriate tool to explicitly mitigate the identified issue. Formally, the meta-controller takes the candidate set derived from the coarse recommender, interaction history, the number of regenerations, and some statistical information, denoted as \mathcal{M} , as its state and makes decisions based on that state:

$$a_{\text{meta}} = \arg \max_{a \in \mathcal{A}_{\text{meta}}} P(a | \mathcal{I}_t, h_t, \#ReGen, \mathcal{M}), \quad (3)$$

where $\mathcal{A}_{\text{meta}} = \{\text{call BD, call RE, call AG}\}$, referring to calling the bias detector, reward estimator, and action grounder, respectively. Then, we elaborate on the details of \mathcal{M} . Aligning with our motivation of debiasing and reward refinement, popularity bias, exposure bias, uncertainty of the reward estimation, and immediate user satisfaction are measured. To be specific, items are categorized as popular if they fall within the top 20% of the rating distribution in the offline logs. Meanwhile, for a user’s interaction history, the tags of the interacted items, e.g., the movie genres, are ranked, and the top 20% are labeled as popular tags. The uncertainty estimation is calculated by Eq. 1. The immediate user satisfaction is measured by the rewards in the candidate set. To sum up, \mathcal{M} is composed of the ratio of popular items $\mathcal{P}_{\text{item}}$, the ratio of items with popular tags \mathcal{P}_{tag} , the mean uncertainty $\bar{u}_{\mathcal{I}_t}$, and the mean reward $\bar{r}_{\mathcal{I}_t}$ of the candidate set.

188 If the meta-controller suspects that the candidate set is biased, it invokes the bias detector for further exami-
 189 nation; a similar procedure is followed when reward estimation is required. Conversely, if the candidate set
 190 is deemed satisfactory, the meta-controller calls the action grounder to generate the final recommendation.
 191 In addition, the meta-controller will summarize the quality of \mathcal{I}_t as $\mathcal{S}_{\text{meta}}$ for downstream modules.

$$192 \mathcal{S}_{\text{meta}} = f_{\text{meta}}(\mathcal{I}_t, h_t, \#\text{ReGen}, \mathcal{M}) \quad (4)$$

193 The meta-controller is implemented by a tunable LLM. Its training will be detailed later.
 194

195 4.4 BIAS DETECTOR

196 The bias detector acts as a specialized expert for assessing the overall bias in the candidate set, *i.e.*, \mathcal{I}_t .
 197 Compared to the meta-controller, the bias detector has a more in-depth evaluation of the bias status. When
 198 a candidate set is suspected of bias, the bias detector is invoked. Only if severe bias is detected does the
 199 bias detector trigger the coarse recommender to regenerate a new candidate set. Its inputs include the candi-
 200 date set, the summary text from the meta-controller, history interactions, interaction history, the number of
 201 regenerations, and the statistical information.

$$202 a_{\text{bd}} = \arg \max_{a \in \mathcal{A}} P(a | \mathcal{S}_{\text{meta}}, \mathcal{I}_t, h_t, \#\text{ReGen}, \mathcal{M}), \quad (5)$$

203 where $\mathcal{A} = \{\text{call AG, ReGen}\}$ and ReGen stands for regeneration. If the bias detector determines that the
 204 overall bias has a negligible impact on fine-grained recommendation, it proceeds by invoking the action
 205 grounder. Otherwise, it triggers the coarse recommender to regenerate a new candidate set. In both cases, a
 206 diagnostic summary \mathcal{S}_{bd} is produced to guide the subsequent action.
 207

$$208 \mathcal{S}_{\text{bd}} = f_{\text{meta}}(\mathcal{S}_{\text{meta}}, \mathcal{I}_t, h_t, \#\text{ReGen}, \mathcal{M}). \quad (6)$$

209 The bias detector is implemented with a frozen LLM. Although a tunable LLM could potentially enhance
 210 its performance, we adopt the frozen version for efficiency considerations.
 211

212 4.5 REWARD ESTIMATOR

213 The reward estimator operates in parallel with the bias detector but serves a complementary role. Specifi-
 214 cally, it diagnoses the overall reward status of the candidate set, focusing on both the uncertainty of reward
 215 estimates and their scale. Formally, its inputs are the same as those of the bias detector, but the summary
 216 text from the meta-controller is different.
 217

$$218 a_{\text{re}} = \arg \max_{a \in \mathcal{A}} P(a | \mathcal{S}_{\text{meta}}, \mathcal{I}_t, h_t, \#\text{ReGen}, \mathcal{M}). \quad (7)$$

219 If the reward estimator concludes that the uncertainty and scale of the reward signals are within an acceptable
 220 range, it proceeds by invoking the action grounder. Otherwise, it calls the coarse recommender to regenerate
 221 a new candidate set. Similar to the bias detector, the reward estimator also produces a diagnostic summary
 222 \mathcal{S}_{re} to guide subsequent decisions.
 223

$$224 \mathcal{S}_{\text{re}} = f_{\text{meta}}(\mathcal{S}_{\text{meta}}, \mathcal{I}_t, h_t, \#\text{ReGen}, \mathcal{M}). \quad (8)$$

225 Similarly, the reward estimator is implemented by a frozen LLM.
 226

227 4.6 ACTION GROUNDER

228 The action grounder is a fine-grained recommender responsible for selecting the final recommendation from
 229 a limited candidate set. Constraining decisions within candidate sets simplifies the recommendation process
 230 and reduces the risk of hallucination. As the final step of the recommendation, the action grounder can
 231 comprehensively consider all information from the former modules:

$$232 e_i = f_{\text{ag}}(\mathcal{S}, \mathcal{I}_t, h_t, \#\text{ReGen}), \quad (9)$$

233 where $\mathcal{S} \in \{\mathcal{S}_{\text{meta}}, \mathcal{S}_{\text{bd}}, \mathcal{S}_{\text{re}}\}$. The action grounder is implemented by a tunable LLM, and its learning is
 234 described in the next part.

4.7 LEARNING PIPELINE

Training. The training of BARO follows a two-phase pipeline. Phase I performs supervised fine-tuning warm-up for the action grounder with reward distillation from the world model under a frozen meta-controller. Phase II applies reinforcement learning on a simulated environment built from the world model.

Phase I: Supervised Fine-Tuning (SFT). Given a state s_t (history h_t and summaries \mathcal{S}) and the retrieved candidate set \mathcal{I}_t , we query the world model to obtain shaped rewards $\{\tilde{R}(s_t, i)\}_{i \in \mathcal{I}_t}$ and uncertainties $\{u(s_t, i)\}_{i \in \mathcal{I}_t}$ (Eq. 1). We form a soft label distribution over \mathcal{I}_t by temperature-scaled normalization:

$$y_i = \frac{\exp(\tilde{R}(s_t, i)/\tau)}{\sum_{j \in \mathcal{I}_t} \exp(\tilde{R}(s_t, j)/\tau)}, \quad i \in \mathcal{I}_t, \quad (10)$$

where $\tau > 0$ controls the sharpness. Denoting the policy of the action grounder as $\pi_\theta(\cdot | s_t, \mathcal{I}_t)$, we minimize a KL-loss with uncertainty weighting:

$$\mathcal{L}_{\text{KD}}(\theta) = \mathbb{E}_{(s_t, \mathcal{I}_t)} \left[w_t \cdot \sum_{i \in \mathcal{I}_t} y_i \left(-\log \pi_\theta(i | s_t, \mathcal{I}_t) \right) \right], \quad (11)$$

where $w_t = \text{stopgrad}(\exp(-\bar{u}_{\mathcal{I}_t}/\tau_u))$ downweights high-uncertainty batches. $\bar{u}_{\mathcal{I}_t}$ is the mean uncertainty on \mathcal{I}_t , and $\tau_u > 0$ is a temperature. This label-free SFT distills the world model into the action grounder on the retrieved candidate sets. During Phase I, both the bias detector and reward estimator use fixed thresholds. As for the meta-controller, we utilize a rule-based policy according to the statistical information \mathcal{M} and keep it static. This yields a stable initialization and avoids early-stage credit assignment issues.

Phase II: Reinforcement Learning (RL). After SFT, we fine-tune the agent with RL on episodes that consist of up to K tool invocations followed by a final selection. K is set to limit the number of ReGen. At each decision step, the action grounder samples $a_t \in \mathcal{I}_t$ from π_θ , and we obtain a shaped reward from the world model:

$$r_t^{\text{ag}} = \tilde{r}(s_t, a_t) - \lambda_{\text{item}} \cdot \mathbb{I}\{a_t \text{ is popular}\} - \lambda_{\text{tag}} \cdot \mathbb{I}\{a_t \text{'s tag is popular}\} - \lambda_u u(s_t, a_t), \quad (12)$$

where the three coefficients control the degree of the corresponding penalties. We optimize π_θ to maximize the expected discounted return $\mathbb{E}[\sum_t \gamma^t r_t^{\text{ag}}]$ with PPO (Schulman et al., 2017):

$$\max_{\theta} \mathbb{E} \left[\min(\rho_t(\theta) \hat{A}_t, \text{clip}(\rho_t(\theta), 1 - \epsilon, 1 + \epsilon) \hat{A}_t) + \beta_{\text{ent}} \mathcal{H}(\pi_\theta) \right], \quad (13)$$

where $\rho_t(\theta) = \frac{\pi_\theta(a_t | s_t, \mathcal{I}_t)}{\pi_{\theta_{\text{old}}}(a_t | s_t, \mathcal{I}_t)}$. \hat{A}_t is the GAE advantage under r_t^{ag} . ϵ is the clipping parameter. $\mathcal{H}(\pi_\theta)$ and β_{ent} are the entropy penalty and its coefficient, respectively.

Training the meta-controller with sparse rewards may be insufficient, as the credit assignment problem hampers effective cooperation between the meta-controller and the action grounder. We therefore adopt a dense reward design that scores improvements in candidate quality when the meta-controller calls a tool, either BD, RE, or AG. Let $\mathcal{M}_t = (\mathcal{P}_{\text{item}}^t, \mathcal{P}_{\text{tag}}^t, \bar{u}_{\mathcal{I}_t}, \bar{r}_{\mathcal{I}_t})$ be the diagnostics computed on \mathcal{I}_t (ratios of popular items/tags, mean uncertainty from Eq. 1, and mean reward). After the meta action transforms $\mathcal{I}_t \rightarrow \mathcal{I}_{t+1}$, we define

$$r_t^{\text{meta}} = w_r (\bar{r}_{\mathcal{I}_{t+1}} - \bar{r}_{\mathcal{I}_t}) - w_p [(\mathcal{P}_{\text{item}}^{t+1} - \mathcal{P}_{\text{item}}^t) + (\mathcal{P}_{\text{tag}}^{t+1} - \mathcal{P}_{\text{tag}}^t)] - w_u (\bar{u}_{\mathcal{I}_{t+1}} - \bar{u}_{\mathcal{I}_t}) - c_{\text{regen}} \cdot \mathbb{I}\{\text{ReGen}\}, \quad (14)$$

where $w_r, w_p, w_u \geq 0$ weight improvements in mean reward, reductions in popularity bias, and reductions in uncertainty, and $c_{\text{regen}} > 0$ penalizes unnecessary regenerations. The meta-controller policy $\pi_\psi^{\text{meta}}(a | \mathcal{I}_t, h_t, \#\text{ReGen}, \mathcal{M}_t)$ (parameterized by ψ) is also optimized by PPO, similar to the optimization in Eq. 12.

In practice, we alternate updates: fix π_ψ^{meta} and optimize π_θ for a few epochs, then fix π_θ and optimize π_ψ^{meta} , and repeat this process until the predefined number of epochs is reached. We adopt parameter-efficient fine-

tuning via LoRA (Hu et al., 2022). During both SFT and RL, only the low-rank adapters are updated, while the base LLM remains frozen. Since BARO is evaluated on proactive recommendation tasks (Zhu et al., 2023), the reshaped rewards in Eq. 12 and Eq. 14 are further refined by incorporating the consecutive reward difference and the L_2 distance to the target items as in T-PRA (Wang et al., 2025b).

Inference. At inference time, BARO performs a single forward pass with a budget of at most K tool calls. Given h_t , the coarse recommender produces \mathcal{I}_t . The meta-controller inspects \mathcal{M}_t and either (i) invokes BD/RE to refine \mathcal{I}_t and updates the summary, or (ii) directly calls the action grounder. The action grounder then selects the final item from the (refined) \mathcal{I}_t . All summaries ($\mathcal{S}_{\text{meta}}, \mathcal{S}_{\text{bd}}, \mathcal{S}_{\text{re}}$) are appended to prompts to ensure consistency and interpretability.

5 EXPERIMENTS

Datasets and Tasks. Two datasets, specifically Steam (Kang & McAuley, 2018) and Amazon-book (Ni et al., 2019), are used in our experiments. The statistical details of these datasets are presented in Table 1. The reward scale for both datasets is structured range from 1 to 5. The Steam dataset, sourced from a gaming platform, measures reward signals based on playtime, whereas the Amazon-Book dataset originates from a book rating platform, with user ratings functioning as reward signals. In accordance with the methodology established by T-PRA (Wang et al., 2025b) to ensure a fair comparison, each dataset is partitioned into 80% for training and 20% for testing.

Table 1: Statistical overview of the Steam and Amazon-book datasets.

Dataset	Usage	# Users	# Items	Density
Steam	Train	6012	190365	0.145%
	Test	6012	190365	0.084%
Amazon	Train	3109	13864	0.788%
	Test	3109	13864	0.320%

The proactive recommendation task, proposed by (Zhu et al., 2023), assigns each user a target item. The policy must gradually guide the user toward ultimately selecting this target item.

Baselines. The baseline methods fall into two categories: (i) one naive baseline and four sequential recommendation baselines, including **POP**: straightforwardly recommends the most popular items at each step; **Caser** (Tang & Wang, 2018): utilizes convolutional filters to extract sequential patterns as local features; **GRU4Rec** (Jannach & Ludewig, 2017) and **SASRec** (Kang & McAuley, 2018): exploit GRU and Transformer to encode sequential features, respectively. (ii) LLM-based recommendation baselines: **IRS** (Zhu et al., 2023): combines Transformer with impressionability mask to implement personalized proactive recommendation; **BiLLP** (Shi et al., 2024b): adopts a two-level decision-making paradigm, forming a ReACT-like framework; **LLM-IPP/LLM-IPP (CoT)/LLM-IPP (ToT)** (Wang et al., 2025a): makes use of in-context learning to handle proactive recommendation. The *CoT* and *ToT* variants are empowered by Chain-of-Thought (Wei et al., 2022) and Tree-of-Thought (Yao et al., 2023), respectively; T-PRA (Wang et al., 2025b): an LLM agent designed for proactive recommendation. It adopts DPO (Rafailov et al., 2023) to tune the adapters for enhanced performance and achieves state-of-the-art (SOTA).

Metrics and Evaluation. We follow the standard evaluation protocol in proactive recommendation following (Zhu et al., 2023; Wang et al., 2025a) and report two groups of metrics.

Simulator-based metrics. An independently trained recommender is used as a user simulator. Based on this simulator, we compute: (i) *Increment of Interest (IoI)* — the increase of the predicted score for the target item after the whole trajectory; (ii) *Increment of Rank (IoR)* — the improvement of the target item’s ranking position; (iii) *Accuracy (Acc)* — the proportion of recommended items that fall within the top- K positions of the simulator ranking where K is a threshold and set as 50.

LLM-based metrics. A frozen LLM evaluator is employed to assess the recommendation sequence. Two aspects are measured: (i) *Acceptance*, the average probability that the user would accept each recommended item; (ii) *Coherence*, the average consistency between successive items in the trajectory.

All compared methods are evaluated with the same simulator and evaluator settings to ensure fairness.

Table 2: Performance comparison on Steam and Amazon Book datasets. Best results are in bold.

Method	Steam					Amazon Book				
	IoI	IoR	Acc.	Acce.	Coh.	IoI	IoR	Acc.	Acce.	Coh.
POP	-0.426	41.8	0.097	0.381	0.482	0.494	-205.7	0.085	0.428	0.595
Caser	0.269	89.9	0.977	0.505	0.239	0.263	317.1	0.969	0.446	<u>0.629</u>
GRU4Rec	0.220	60.8	<u>0.969</u>	0.511	0.240	0.672	725.6	<u>0.956</u>	0.443	0.645
SASRec	0.354	-75.14	0.484	0.259	0.257	1.060	436.8	0.563	0.438	0.506
IRS	0.164	218.7	0.934	0.381	0.249	0.097	166.4	0.883	0.470	0.481
BiLLP	0.427	308.7	0.855	0.523	0.477	1.200	655.3	0.490	0.583	0.604
LLM-IPP	0.259	340.6	0.912	<u>0.651</u>	<u>0.597</u>	1.436	845.6	0.844	0.595	0.557
LLM-IPP (CoT)	0.264	303.0	0.906	0.629	0.580	1.277	803.5	0.836	<u>0.601</u>	0.538
LLM-IPP (ToT)	0.282	244.1	0.861	0.571	0.509	0.944	513.0	0.765	0.522	0.533
T-PRA	<u>0.584</u>	<u>432.8</u>	0.894	0.588	0.403	<u>1.783</u>	<u>1276.5</u>	0.773	0.589	<u>0.629</u>
BARO (Ours)	0.601	469.4	0.886	0.664	0.601	1.976	1324.2	0.802	0.611	0.617

Implementation Details. In BARO, we employ Llama-3.1-8B-Instruct (Dubey et al., 2024) as the backbone for all LLM modules to ensure a fair comparison with T-PRA. Supervised fine-tuning is performed for 3 epochs, followed by reinforcement learning for another 5 epochs, where each epoch alternates between updating the meta-controller and the action grounder. We adopt the AdamW optimizer with a learning rate of 5×10^{-5} . A single run of BARO can be reproduced on two NVIDIA RTX™ A6000 GPUs (48 GB GDDR6) and requires approximately 30 GPU hours.

5.1 OVERALL

According to Table 2, our proposed method **BARO significantly outperforms both traditional and LLM-based baselines** in terms of the core metrics IoI and IoR compared to the strongest baseline T-PRA. Specifically, BARO obtains an IoI of 0.601 and IoR of 469.4 on Steam, and 1.976 and 1324.2 on Amazon Book. This demonstrates the superiority of our Bias and Reward Optimization framework in enhancing long-term user interest and item exposure. Moreover, BARO achieves the highest Acceptance score across both datasets (0.664 on Steam and 0.611 on Amazon Book), showing the importance of using reward estimation and bias mitigation to align recommendations with user preferences and expectations. For the metric of Coherence, it also obtains the best Coherence score on Steam (0.601), and a competitive score on Amazon Book (0.617), comparable to BiLLP (0.604) and T-PRA (0.629), indicating that the two-stage optimization framework enables BARO to maintain semantic consistency across multiple recommendation turns, while aligning the generated content with both user intent and long-term engagement goals.

Table 3: Ablation study results on Steam and Amazon Book datasets.

Method	Steam					Amazon Book				
	IoI	IoR	Acc.	Acce.	Coh.	IoI	IoR	Acc.	Acce.	Coh.
T-PRA	0.584	432.8	0.894	0.588	0.403	1.783	1276.5	0.773	0.589	0.629
-BD	0.492	379.1	0.727	0.499	0.383	0.840	1182.0	0.453	0.423	0.345
-RE	0.599	399.2	0.778	0.640	0.560	1.481	1321.3	0.776	0.594	0.651
-BD, RE	0.360	311.9	0.664	0.308	0.305	0.686	738.5	0.381	0.354	0.217
-SFT	0.540	415.3	0.716	0.451	0.583	1.463	892.9	0.295	0.653	0.639
BARO (ours)	0.601	469.4	0.886	0.664	0.601	1.976	1324.2	0.802	0.611	0.617

5.2 ABLATION STUDY

This section presents the ablation study. -BD refers to the version of BARO without the bias detector. Similarly, -RE refers to BARO without the reward estimator. -BD, RE represents that both the bias detector and the reward estimator are removed, and the meta-controller directly decides whether to call ReGen or the action grounder. -SFT denotes the variant that the action grounder is not trained with SFT. The corresponding results are shown in Table 3. First, removing the bias detector (-BD) consistently hurts performance across almost all metrics, particularly Accuracy (Acc.) on both datasets, confirming that explicitly detecting and correcting popularity bias is crucial for improving recommendation quality. Second, dropping the reward estimator also leads to significant drops, especially in Coherence and Acceptance, indicating that

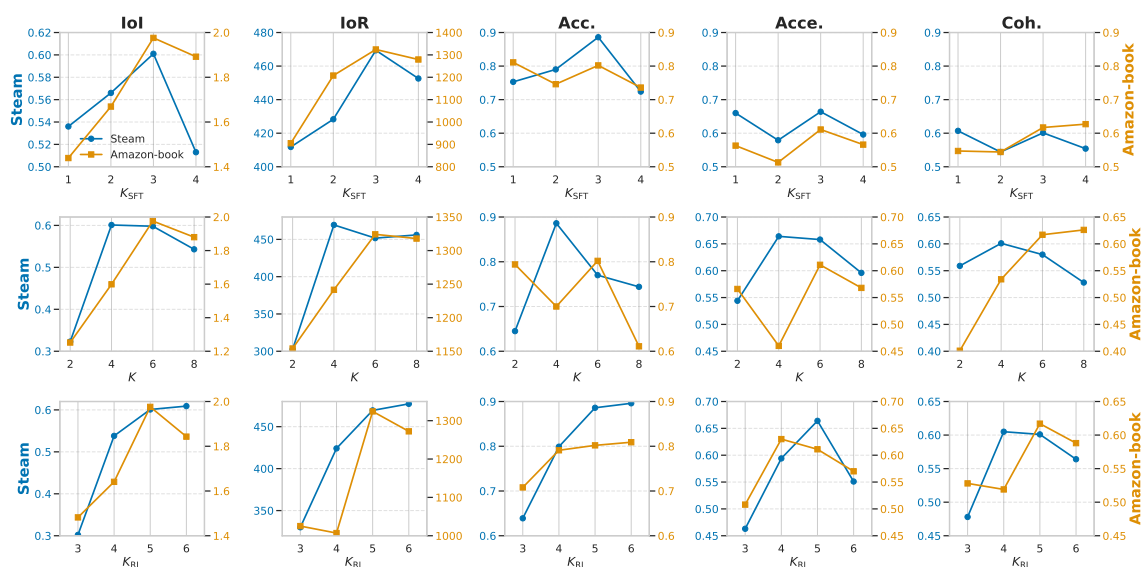


Figure 2: Hyperparameter sensitivity analysis for the number of SFT epochs (K_{SFT}), maximum tool calls (K), and RL epochs (K_{RL}) on the Steam and Amazon-book datasets.

refining noisy world-model rewards improves not only accuracy but also the trajectory consistency perceived by users. Notably, the joint removal of both tools ($-\text{BD}$, RE) yields the most severe degradation, highlighting their complementary roles in addressing distinct yet coupled challenges of bias and reward in offline recommendation. Third, eliminating the SFT initialization ($-\text{SFT}$) reduces overall performance, with a particularly sharp decline in Accuracy on Amazon Book, suggesting that the reward-distilled warm-up is essential for stabilizing subsequent RL training. Finally, the full BARO model achieves the best overall results, outperforming the T-PRA baseline in both IoI and IoR, while simultaneously improving Acceptance and Coherence. These results **validate the effective contributions of each component in BARO**.

5.3 ROBUSTNESS OF HYPERPARAMETERS

In BARO, we investigate the impact of three key factors: the number of SFT epochs K_{SFT} , the number of RL epochs K_{RL} , and the maximum number of tool calls K . The results in Figure 2 demonstrate clear trends. For SFT epochs (top row), we observe that a moderate number (around 3 epochs) yields the best balance across metrics, as too few epochs underfit while excessive tuning causes overfitting. For RL epochs (bottom row), performance steadily improves up to about 5 epochs and then saturates, confirming that policy refinement benefits from additional interactions but plateaus once the agent stabilizes. Regarding the tool budget K (middle row), we find that allowing 4–6 calls provides consistent gains by enabling sufficient refinement, while an overly large K leads to unnecessary regenerations and degraded efficiency. These trends highlight the importance of **training and inference budgets to maximize BARO’s benefits without introducing instability or redundancy**.

6 CONCLUSION

In this work, we presented BARO, a meta-controlled and tool-augmented LLM agent tailored for recommender systems. Unlike prior ReAct-style agents that primarily enhance memory and planning, BARO explicitly addresses two fundamental challenges in offline recommendation: **bias mitigation** and **reward calibration**. By integrating a bias detector, a reward estimator, and an action grounder under the governance of a meta-controller, BARO injects debiasing and reward refinement directly into the decision-making loop.

REFERENCES

- Keqin Bao, Jizhi Zhang, Yang Zhang, Wenjie Wang, Fuli Feng, and Xiangnan He. Large language models for recommendation: Progresses and future directions. In *SIGIR-AP*, 2023a.
- Keqin Bao, Jizhi Zhang, Yang Zhang, Wenjie Wang, Fuli Feng, and Xiangnan He. Tallrec: An effective and efficient tuning framework to align large language model with recommendation. In *RecSys*, 2023b.
- Keqin Bao, Ming Yan, Yang Zhang, Jizhi Zhang, Wenjie Wang, Fuli Feng, and Xiangnan He. Customizing in-context learning for dynamic interest adaption in llm-based recommendation. In *ACL*, 2025.
- Jesús Bobadilla, Fernando Ortega, Antonio Hernando, and Abraham Gutiérrez. Recommender systems survey. *KBS*, 2013.
- Robin Burke. Hybrid recommender systems: Survey and experiments. *User Modeling and User-Adapted Interaction*, 2002.
- Hongru Cai, Yongqi Li, Wenjie Wang, Fengbin Zhu, Xiaoyu Shen, Wenjie Li, and Tat-Seng Chua. Large language models empowered personalized web agents. In *WWW*, 2025a.
- Shihao Cai, Jizhi Zhang, Keqin Bao, Chongming Gao, Qifan Wang, Fuli Feng, and Xiangnan He. Agentic feedback loop modeling improves recommendation and user simulation. In *SIGIR*, 2025b.
- Jiaju Chen, Chongming Gao, Shuai Yuan, Shuchang Liu, Qingpeng Cai, and Peng Jiang. Dlrec: A novel approach for managing diversity in llm-based recommender systems. In *WSDM*, 2025.
- Xiaocong Chen, Siyu Wang, Julian McAuley, Dietmar Jannach, and Lina Yao. On the opportunities and challenges of offline reinforcement learning for recommender systems. *TOIS*, 2023a.
- Xiaocong Chen, Lina Yao, Julian McAuley, Guanglin Zhou, and Xianzhi Wang. Deep reinforcement learning in recommender systems: A survey and new perspectives. *KBS*, 2023b.
- Qingxiu Dong, Lei Li, Damai Dai, Ce Zheng, Jingyuan Ma, Rui Li, Heming Xia, Jingjing Xu, Zhiyong Wu, Baobao Chang, et al. A survey on in-context learning. In *EMNLP*, 2024.
- Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha Letman, Akhil Mathur, Alan Schelten, Amy Yang, Angela Fan, et al. The llama 3 herd of models. *arXiv*, 2024.
- Chongming Gao, Shijun Li, Wenqiang Lei, Jiawei Chen, Biao Li, Peng Jiang, Xiangnan He, Jiaxin Mao, and Tat-Seng Chua. Kuairac: A fully-observed dataset and insights for evaluating recommender systems. In *CIKM*, 2022a.
- Chongming Gao, Shijun Li, Yuan Zhang, Jiawei Chen, Biao Li, Wenqiang Lei, Peng Jiang, and Xiangnan He. Kuairand: An unbiased sequential recommendation dataset with randomly exposed videos. In *CIKM*, 2022b.
- Chongming Gao, Kexin Huang, Jiawei Chen, Yuan Zhang, Biao Li, Peng Jiang, Shiqi Wang, Zhong Zhang, and Xiangnan He. Alleviating matthew effect of offline reinforcement learning in interactive recommendation. In *SIGIR*, 2023a.
- Chongming Gao, Shiqi Wang, Shijun Li, Jiawei Chen, Xiangnan He, Wenqiang Lei, Biao Li, Yuan Zhang, and Peng Jiang. Cirs: Bursting filter bubbles by counterfactual interactive recommender system. *TOIS*, 2023b.

- 470 Chongming Gao, Ruijun Chen, Shuai Yuan, Kexin Huang, Yuanqing Yu, and Xiangnan He. Sprec: Leverag-
471 ing self-play to debias preference alignment for large language model-based recommendations. In *SIGIR*,
472 2025a.
- 473 Chongming Gao, Mengyao Gao, Chenxiao Fan, Shuai Yuan, Wentao Shi, and Xiangnan He. Process-
474 supervised llm recommenders via flow-guided tuning. In *SIGIR*, 2025b.
- 476 Yunfan Gao, Tao Sheng, Youlin Xiang, Yun Xiong, Haofen Wang, and Jiawei Zhang. Chat-rec: Towards
477 interactive and explainable llms-augmented recommender system. *arXiv*, 2023c.
- 478 Peiyuan Gong, Jiamian Li, and Jiaxin Mao. Cosearchagent: a lightweight collaborative search agent with
479 large language models. In *SIGIR*, 2024.
- 481 Daya Guo, Dejian Yang, Haowei Zhang, Junxiao Song, Ruoyu Zhang, Runxin Xu, Qihao Zhu, Shirong Ma,
482 Peiyi Wang, Xiao Bi, et al. Deepseek-r1: Incentivizing reasoning capability in llms via reinforcement
483 learning. *arXiv*, 2025.
- 485 Huifeng Guo, Ruiming Tang, Yunming Ye, Zhenguo Li, and Xiuqiang He. Deepfm: a factorization-machine
486 based neural network for ctr prediction. In *IJCAI*, 2017.
- 487 F Maxwell Harper and Joseph A Konstan. The movielens datasets: History and context. *TIIS*, 2015.
- 488 Edward J Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, Weizhu
489 Chen, et al. Lora: Low-rank adaptation of large language models. *ICLR*, 2022.
- 490 Xu Huang, Jianxun Lian, Yuxuan Lei, Jing Yao, Defu Lian, and Xing Xie. Recommender ai agent: Integrat-
491 ing large language models for interactive recommendations. *TOIS*, 2025.
- 492 Dietmar Jannach and Malte Ludewig. When recurrent neural networks meet the neighborhood for session-
493 based recommendation. In *RecSys*, 2017.
- 494 Wang-Cheng Kang and Julian McAuley. Self-attentive sequential recommendation. In *ICDM*, 2018.
- 495 Yongqi Li, Xinyu Lin, Wenjie Wang, Fuli Feng, Liang Pang, Wenjie Li, Liqiang Nie, Xiangnan He, and
496 Tat-Seng Chua. A survey of generative search and recommendation in the era of large language models.
497 *arXiv*, 2024.
- 498 Jiayi Liao, Sihang Li, Zhengyi Yang, Jiancan Wu, Yancheng Yuan, and Xiang Wang. Llara: Aligning large
499 language models with sequential recommenders. *arXiv*, 2023.
- 500 Xinyu Lin, Wenjie Wang, Yongqi Li, Shuo Yang, Fuli Feng, Yinwei Wei, and Tat-Seng Chua. Data-efficient
501 fine-tuning for llm-based recommendation. In *SIGIR*, 2024.
- 502 Junling Liu, Chao Liu, Peilin Zhou, Renjie Lv, Kang Zhou, and Yan Zhang. Is chatgpt a good recommender?
503 a preliminary study. *arXiv*, 2023.
- 504 Mahmoud Mohammadi, Yipeng Li, Jane Lo, and Wendy Yip. Evaluation and benchmarking of llm agents:
505 A survey. In *SIGKDD*, 2025.
- 506 Jianmo Ni, Jiacheng Li, and Julian McAuley. Justifying recommendations using distantly-labeled reviews
507 and fine-grained aspects. In *EMNLP*, 2019.
- 508 Liang-bo Ning, Shijie Wang, Wenqi Fan, Qing Li, Xin Xu, Hao Chen, and Feiran Huang. Cheatagent:
509 Attacking llm-empowered recommender systems via llm agent. In *SIGKDD*, 2024.

- 517 Alexander Pan, Erik Jones, Meena Jagadeesan, and Jacob Steinhardt. Feedback loops with language models
518 drive in-context reward hacking. In *ICML*, 2024.
- 519
- 520 Rafael Rafailov, Archit Sharma, Eric Mitchell, Christopher D Manning, Stefano Ermon, and Chelsea Finn.
521 Direct preference optimization: Your language model is secretly a reward model. *NIPS*, 2023.
- 522
- 523 John Schulman, Filip Wolski, Prafulla Dhariwal, Alec Radford, and Oleg Klimov. Proximal policy opti-
524 mization algorithms. *CoRR*, abs/1707.06347, 2017.
- 525 Tianhao Shi, Yang Zhang, Zhijian Xu, Chong Chen, Fuli Feng, Xiangnan He, and Qi Tian. Preliminary
526 study on incremental learning for large language model-based recommender systems. In *CIKM*, 2024a.
- 527
- 528 Wentao Shi, Xiangnan He, Yang Zhang, Chongming Gao, Xinyue Li, Jizhi Zhang, Qifan Wang, and Fuli
529 Feng. Large language models are learnable planners for long-term recommendation. In *SIGIR*, 2024b.
- 530 Yubo Shu, Haonan Zhang, Hansu Gu, Peng Zhang, Tun Lu, Dongsheng Li, and Ning Gu. Rah! recsys-
531 assistant-human: A human-centered recommendation framework with llm agents. *IEEE TCSS*, 2024.
- 532
- 533 Richard S Sutton and Andrew G Barto. *Reinforcement learning: An introduction*. MIT press, 2018.
- 534 Jiayi Tang and Ke Wang. Personalized top-n sequential recommendation via convolutional sequence em-
535 bedding. In *WSDM*, 2018.
- 536
- 537 Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay Bash-
538 lykov, Soumya Batra, Prajjwal Bhargava, Shruti Bhosale, et al. Llama 2: Open foundation and fine-tuned
539 chat models. *CoRR*, abs/2307.09288, 2023.
- 540 Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser,
541 and Illia Polosukhin. Attention is all you need. In *NeurIPS*, 2017.
- 542
- 543 Jie Wang, Alexandros Karatzoglou, Ioannis Arapakis, and Joemon M Jose. Reinforcement learning-based
544 recommender systems with large language models for state reward and action modeling. In *SIGIR*, 2024a.
- 545
- 546 Mingze Wang, Shuxian Bi, Wenjie Wang, Chongming Gao, Yangyang Li, and Fuli Feng. Leveraging llms
547 for influence path planning in proactive recommendation. In *WWW Companion*, 2025a.
- 548
- 549 Mingze Wang, Chongming Gao, Wenjie Wang, Yangyang Li, and Fuli Feng. Tunable llm-based proactive
550 recommendation agent. In *ACL*, 2025b.
- 551
- 552 Yancheng Wang, Ziyang Jiang, Zheng Chen, Fan Yang, Yingxue Zhou, Eunah Cho, Xing Fan, Xiaojiang
553 Huang, Yanbin Lu, and Yingzhen Yang. Recmind: Large language model powered agent for recommen-
554 dation. In *Findings of NAACL*, 2024b.
- 555
- 556 Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Fei Xia, Ed Chi, Quoc V Le, Denny Zhou,
557 et al. Chain-of-thought prompting elicits reasoning in large language models. *NIPS*, 2022.
- 558
- 559 Likang Wu, Zhi Zheng, Zhaopeng Qiu, Hao Wang, Hongchao Gu, Tingjia Shen, Chuan Qin, Chen Zhu,
560 Hengshu Zhu, Qi Liu, et al. A survey on large language models for recommendation. *WWW*, 2024.
- 561
- 562 Wanqi Xue, Qingpeng Cai, Zhenghai Xue, Shuo Sun, Shuchang Liu, Dong Zheng, Peng Jiang, Kun Gai,
563 and Bo An. Prefrec: Recommender systems with human preferences for reinforcing long-term user
engagement. In *SIGKDD*, 2023.
- 564
- 565 Shunyu Yao, Dian Yu, Jeffrey Zhao, Izhak Shafran, Tom Griffiths, Yuan Cao, and Karthik Narasimhan. Tree
of thoughts: Deliberate problem solving with large language models. *NIPS*, 2023.

564 Ziang Ye, Zhenru Zhang, Yang Zhang, Jianxin Ma, Junyang Lin, and Fuli Feng. Disentangling reasoning
565 tokens and boilerplate tokens for language model fine-tuning. *arXiv preprint arXiv:2412.14780*, 2024.
566

567 Yuanqing Yu, Chongming Gao, Jiawei Chen, Heng Tang, Yuefeng Sun, Qian Chen, Weizhi Ma, and Min
568 Zhang. Easyrl4rec: An easy-to-use library for reinforcement learning based recommender systems. In
569 *SIGIR*, 2024.

570 An Zhang, Yuxin Chen, Leheng Sheng, Xiang Wang, and Tat-Seng Chua. On generative agents in recom-
571 mendation. In *SIGIR*, 2024a.
572

573 Jizhi Zhang, Keqin Bao, Wenjie Wang, Yang Zhang, Wentao Shi, Wanhong Xu, Fuli Feng, and Tat-Seng
574 Chua. Prospect personalized recommendation on large language model-based agent platform. *arXiv*,
575 2024b.

576 Junjie Zhang, Yupeng Hou, Ruobing Xie, Wenqi Sun, Julian McAuley, Wayne Xin Zhao, Leyu Lin, and Ji-
577 Rong Wen. Agentcf: Collaborative learning with autonomous language agents for recommender systems.
578 In *WWW*, 2024c.
579

580 Yi Zhang, Ruihong Qiu, Jiajun Liu, and Sen Wang. Roler: Effective reward shaping in offline reinforcement
581 learning for recommender systems. In *CIKM*, 2024d.

582 Yuyue Zhao, Jiancan Wu, Xiang Wang, Wei Tang, Dingxian Wang, and Maarten De Rijke. Let me do it for
583 you: Towards llm empowered recommendation via tool learning. In *SIGIR*, 2024.
584

585 Haoren Zhu, Hao Ge, Xiaodong Gu, Pengfei Zhao, and Dik Lun Lee. Influential recommender system. In
586 *ICDE*, 2023.
587
588
589
590
591
592
593
594
595
596
597
598
599
600
601
602
603
604
605
606
607
608
609
610

611 A APPENDIX

612
613 A.1 STATEMENT ON LLM USAGE

614
615 Other than the LLM backbones leveraged in the experiments, the only LLM we used was ChatGPT-5, solely
616 for manuscript polishing and grammar correction.

617
618 A.2 DEMONSTRATION OF PROMPTS

619 620 Prompt template for Coarse Recommender	
621 622 Instruction:	You are a recommendation expert. You should rely on the user’s interaction history, the current 623 regeneration round, and guidance from other modules. The goal is to recommend one item that 624 serves as the most representative item you want to recommend at the current step.
625 626 Input:	The user's interaction history is $\{h_t\}$ (END OF HISTORY). 627 It is your $\{\#ReGen\}$ times to generate an item in the current step. The guidance information from upstream modules is: $\{g\}$.
628 629 Response:	Recommend the [item].

630
631 Figure 3: Prompts to generate an item with Coarse Recommender.

632 633 Prompt template for Meta-controller	
634 635 Instruction:	You are the meta-controller of the recommendation pipeline. Your task is to evaluate the candidate 636 set based on three key factors: popularity ratio, reward scale, and reward uncertainty. - If the set appears strongly biased, call the Bias Detector for further examination. 637 - If rewards are very low or uncertainty is high, call the Reward Estimator for further examination. - If the set quality is acceptable, call the Action Grounder to make the final recommendation. 638 In making decisions, also consider the user’s interaction history and the current regeneration round. 639 Finally, provide a concise summary of the candidate quality to guide downstream modules.
640 641 Input:	Current candidate set is: $\{\text{candidate set}\}$ (END OF CANDIDATE SET). 642 The corresponding statistical information of the candidate set is: $\{\mathcal{M}\}$ (END OF STATISTICAL INFORMATION). 643 The user's interaction history is $\{h_t\}$ (END OF HISTORY). 644 It is the coarse recommender's $\{\#ReGen\}$ times to generate a candidate set in the current step.
645 646 Response:	Call Bias Detector/Reward Estimator/Action Grounder, [summary text: \mathcal{S}_{meta}].

647
648 Figure 4: Prompts for tool-use and generating summary texts with Meta-controller.

Prompt template for Bias Detector	
Instruction:	You are the Bias Detector. Your role is to carefully evaluate whether the candidate set contains excessive popularity bias. Use the provided candidate set, the user’s interaction history, the regeneration count, and the diagnostic summary from the meta-controller to make your decision. - If the bias is severe, call ReGen and provide a short summary to guide the coarse recommender. - If the bias is acceptable, call the Action Grounder and provide a short summary to guide it.
Input:	The candidate set is {candidate set}(END OF CANDIDATE SET). The corresponding statistical information of the candidate set is: $\{\mathcal{M}\}$ (END OF STATISTICAL INFORMATION). The user’s interaction history is $\{h_t\}$ (END OF HISTORY). It is the coarse recommender's $\{\#ReGen\}$ times to generate a candidate set in the current step. The diagnostic summary from the meta-controller is $\{\mathcal{S}_{meta}\}$.
Response:	- Either: call ReGen [diagnostic summary: \mathcal{S}_{bd}]. - Or: call Action Grounder [diagnostic summary: \mathcal{S}_{bd}].

Figure 5: Prompts to evaluate the candidate set’s bias with Bias Detector.

Prompt template for Reward Estimator	
Instruction:	You are the Reward Estimator. Your role is to check whether the candidate set suffers from low rewards or high uncertainty. Use the provided candidate set, the user’s interaction history, the regeneration count, and the diagnostic summary from the meta-controller to make your decision. - If the reward signals are unreliable, call ReGen and provide a short summary to guide the coarse recommender. - If the signals are acceptable, call the Action Grounder and provide a short summary to guide it.
Input:	The candidate set is {candidate set}(END OF CANDIDATE SET). The corresponding statistical information of the candidate set is: $\{\mathcal{M}\}$ (END OF STATISTICAL INFORMATION). The user’s interaction history is $\{h_t\}$ (END OF HISTORY). It is the coarse recommender's $\{\#ReGen\}$ times to generate a candidate set in the current step. The diagnostic summary from the meta-controller is $\{\mathcal{S}_{meta}\}$.
Response:	- Either: call ReGen [diagnostic summary: \mathcal{S}_{re}]. - Or: call Action Grounder [diagnostic summary: \mathcal{S}_{re}].

Figure 6: Prompts to assess the candidate set’s rewards with Reward Estimator.

Prompt template for Action Grounder	
Instruction:	You are the Action Grounder. Your task is to make the final recommendation from the candidate set. Use the user’s interaction history, the regeneration count, and the diagnostic summaries from upstream modules to select the single best item. The goal is to maximize user satisfaction while avoiding hallucination.
Input:	The candidate set is {candidate set}(END OF CANDIDATE SET). The user’s interaction history is $\{h_t\}$ (END OF HISTORY). It is the coarse recommender's $\{\#ReGen\}$ times to generate a candidate set in the current step. The diagnostic summary is from the Meta-Controller/Bias Detector/Reward Estimator, and the corresponding (diagnostic) summary is $\{\mathcal{S}_{meta}/\mathcal{S}_{bd}/\mathcal{S}_{re}\}$.
Response:	Recommend the [item].

Figure 7: Prompts for final recommendation with Action Grounder.

705 A.3 PSEUDO CODE

706 The pseudo-code of BARO is provided here for reproducibility.

709 **Algorithm 1:** BARO’s Training and Inference

710 **Input:** Offline logs \mathcal{D} , world model \hat{R} , coarse recommender f_{CR} , bias detector (BD), reward estimator
 711 (RE), action grounder π_θ , meta-controller π_ψ^{meta}

712 **Output:** Trained BARO agent

713 **Phase I: Supervised Fine-Tuning (SFT)**

714 **for** each batch (s_t, \mathcal{I}_t) sampled from \mathcal{D} **do**

715 Query \hat{R} to obtain $\{\hat{R}(s_t, i)\}, \{u(s_t, i)\}$

716 Construct soft label distribution y_i via Eq. 10

717 Compute loss $\mathcal{L}_{\text{KD}}(\theta)$ with uncertainty weighting (Eq. 11)

718 Update π_θ with LoRA adapters; keep π_ψ^{meta} , BD, and RE frozen

719 **Phase II: Reinforcement Learning (RL)**

720 **for** episode = 1 to N **do**

721 Initialize history h_t , candidate set $\mathcal{I}_t \leftarrow f_{CR}(h_t, \# \text{ReGen}, g)$

722 **for** step = 1 to K **do**

723 Meta-controller selects $a_{\text{meta}} \in \{\text{BD}, \text{RE}, \text{AG}\}$

724 **if** $a_{\text{meta}} = \text{BD}$ **then**

725 BD evaluates \mathcal{I}_t and outputs \mathcal{S}_{bd}

726 If severe bias \rightarrow regenerate \mathcal{I}_t via f_{CR}

727 **else if** $a_{\text{meta}} = \text{RE}$ **then**

728 RE evaluates reward status and outputs \mathcal{S}_{re}

729 If reward unacceptable \rightarrow regenerate \mathcal{I}_t via f_{CR}

730 **else**

731 Invoke action grounder π_θ to select final item a_t

732 Break

733 Compute meta reward r_t^{meta} via Eq. 14

734 Update π_ψ^{meta} with PPO

735 Compute action grounder reward r_t^{ag} via Eq. 12

736 Update π_θ with PPO

737 Alternate updates between π_θ and π_ψ^{meta}

738 **Inference**

739 Given h_t , coarse recommender produces \mathcal{I}_t

740 Meta-controller inspects \mathcal{M}_t and chooses BD, RE, or AG

741 If BD/RE \rightarrow refine \mathcal{I}_t and update summaries

742 Action grounder selects final recommendation a_t from \mathcal{I}_t

743 Append $\mathcal{S}_{\text{meta}}, \mathcal{S}_{\text{bd}}, \mathcal{S}_{\text{re}}$ to ensure interpretability
