REWARDRANK: OPTIMIZING TRUE LEARNING-TO-RANK UTILITY

Anonymous authorsPaper under double-blind review

000

001

002003004

006

008 009

010 011

012

013

014

016

017

018

019

021

025

026027028

029

031

033

034

037

038

040

041

042

043

044

045

046

047

048

051

052

ABSTRACT

Traditional ranking systems optimize offline proxy objectives that assume simplistic user behaviors, overlooking factors such as position bias and item diversity. As a result, they fail to improve desired counterfactual ranking utility like click-throughrate or purchase probability computed during online A/B tests. We introduce RewardRank, a data-driven learning-to-rank (LTR) framework for counterfactual utility maximization. RewardRank first trains a reward model to predict the utility of any ranking from logged data, then optimizes a ranker to maximize this reward via a differentiable soft permutation operator. To address the absence of any large-scale reproducible counterfactual LTR benchmarks, we propose two evaluation suites: (i) Parametric Oracle Evaluation (PO-Eval), which uses an open-source click model as the counterfactual oracle on Baidu-ULTR dataset, and (ii) LLM-As-User Evaluation (LAU-Eval), which simulates realistic user interactions using a large language model on Amazon-KDD-Cup dataset. RewardRank achieves the highest counterfactual utility on both suites and sets a new state of the art in relevance performance on Baidu-ULTR with real-click signals, demonstrating the feasibility of directly optimizing ranking policies for counterfactual utility. Our code is available at: https://anonymous.4open.science/r/RewardRank-EE46

1 Introduction

The goal of any ranking system is to model human decision-making in a way that maximizes user engagement and utility. However, real-world user behavior is shaped by subtle, context-dependent cognitive biases that traditional ranking losses fail to capture. Engagement often drops when users are presented with redundant or overly similar items, whereas introducing diversity or strategically positioning items can significantly enhance interest. For example, the decoy effect—where the presence of a less-attractive item increases preference for a similar alternative—has been observed in search interactions and shown to meaningfully influence user choices (Wang et al., 2025a). Other well-documented biases include position bias (Chen et al., 2024; Hager et al., 2024; Zou et al., 2022), brand bias (Li et al., 2025), and similarity aversion (Tversky & Simonson, 2004). In online advertising, the goal is often to maximize the probability that a user clicks on any item in the list, rather than just the top-ranked one. If data shows that users tend to click on the second position, it may be optimal to place the most engaging ad there to improve overall performance. Likewise, in recommendation scenarios, users may prefer a diverse mix of product styles or brands over a cluster of nearly identical, albeit highly relevant, items. Traditional ranking losses, which emphasize relevance at individual positions (typically the top), are ill-suited for modeling such list-level behaviors (Figure 1). They overlook the fact that user utility depends not just on which items are shown, but how they are arranged, highlighting the limitations of handcrafted objectives in capturing the interactive and comparative nature of real user decision-making.

A natural way to model user behavior is by learning preferences over full permutations of items within a query group (i.e., a query and its associated items). The ideal objective is to identify and rank those permutations that are most likely to drive user engagement, which can be formulated as a likelihood maximization problem: maximizing the probability of observing high-engagement permutations while minimizing that of unengaged ones.

However, the combinatorial explosion of the permutation space quickly renders this approach intractable; for instance, ranking 10 items results in 10! (over 3.6 million) possible arrangements. To

055

057

058

060

061

062

063

064

065

066

067

068

069

071

072

073

074

075

076

077

079

081

082

083

084

085

880

090

091

092

096

098

099

100

101

102 103

104

105

107

address this, recent approaches adopt a utility-based framework (Feng et al., 2021; Shi et al., 2023; Xi et al., 2024; Ren et al., 2024; Wang et al., 2025b), where a utility model is trained to score permutations based on user preferences, and a ranker is subsequently optimized to generate item orders that maximize the predicted utility. While this framework reduces the combinatorial burden, it introduces two key challenges. First is the classic exploration—exploitation dilemma: the ranker must leverage known high-utility arrangements while also exploring novel permutations that may yield higher engagement.

Second is the vulnerability to utility model misspecification: if the learned reward model fails to accurately reflect true user preferences, the ranker may be misled, resulting in poor exploration and degraded overall performance.

In operational ranking systems, user interactions are logged for only a small fraction of the total permutation space. For example, in Figure 1, only 3 out of the 120 possible arrangements of 5 items for the query "laptop bag" are observed, constituting the factual space. These observed interactions define the factual/observed space, whereas the vast majority of unexposed, yet potentially high-utility, permutations form the counterfactual/unobserved space. The optimal arrangement that maximizes user engagement may exist anywhere within the full permutation space of 120 arrangements. One of the major challenges in counterfactual ranking lies in reliably evaluating unobserved permutations. Even if the full permutation space is modeled, evaluating ranking strategies under counterfactual settings remains challenging due to the lack of explicit supervision (Agarwal et al., 2019; Gupta et al., 2024b;a; Buchholz et al., 2024). For instance, in Figure 1, 117 out of 120 possible arrangements remain unobserved, making their evaluation inherently counterfactual. Existing approaches, such as offline A/B testing, inverse propensity scoring, or other debiasing techniques, are often costly, statisti-



Figure 1: **Counterfactual ranking with true learning-to-rank utility**. Three arrangements for the query "laptop bag", with item relevance/rating scores (0–5*). The top row ranks purely by relevance but suffers from similarity aversion due to identical color and style, lowering engagement. The middle row improves diversity but surfaces low-relevance items early, which may deter clicks. The bottom row balances diversity and relevance, placing distinct yet relevant items in top positions, leading to higher predicted utility and user engagement. (Figures are generated by GPT-40)

cally unstable, or difficult to scale, making counterfactual evaluation a central bottleneck in listwise utility optimization.

To address these challenges, we propose REWARDRANK, a counterfactual utility maximization framework that models user behavior over full item permutations. Rather than scoring items in isolation, we learn a permutation-aware utility function that captures user preferences at the list level. To enable differentiable optimization over permutations, we employ the *SoftSort* operator (Prillo & Eisenschlos, 2020) to construct soft position embeddings, allowing end-to-end training of the ranking model with respect to utility gradients. To mitigate the effects of reward model misspecification—where the learned utility may diverge from actual user preferences—we introduce a correction term in the ranker's training objective that improves robustness during optimization. For evaluation, we present two scalable, fully automated protocols that assess counterfactual performance without requiring human labels. *Parametric Oracle Evaluation (PO-Eval)* uses a pretrained, position-aware oracle to provide soft supervision and serve as a proxy for user behavior. *LLM-As-User Evaluation (LAU-Eval)* leverages large language models to simulate user preferences and assess ranking quality across unobserved permutations. Together, these methods enable efficient benchmarking of counterfactual ranking strategies and help align learned rankings with actual or simulated user utility.

Our key contributions can be summarized as:

- We introduce REWARDRANK, a framework for counterfactual utility maximization that learns a permutation-aware reward model capturing human list-level preferences and behavioral biases.
- We enable end-to-end ranking optimization using differentiable soft permutation operators, and incorporate a reward correction term to mitigate reward model misspecification.

• We propose two automated evaluation protocols: *PO-Eval* (parametric oracle) and *LAU-Eval* (LLM-as-user), and construct reproducible testbeds for scalable counterfactual ranking evaluation.

2 RELATED WORK

Traditional ranking methods. Traditional learning-to-rank (LTR) methods are typically categorized into three classes: point-wise, pair-wise, and list-wise approaches. Point-wise methods treat ranking as a regression or classification problem by independently assigning relevance scores to each item (Burges et al., 2005a;b). While computationally efficient, they neglect interactions among items in the ranked list. Pair-wise approaches, including RankSVM (Joachims, 2002), RankBoost (Freund et al., 2003), and LambdaMART (Burges, 2006; Wu et al., 2010), aim to learn relative preferences between item pairs, improving over point-wise methods but still failing to capture full list-level dependencies. In contrast, list-wise methods optimize objectives over the entire ranking, such as NDCG (Cao et al., 2007; Xia et al., 2008), offering better alignment with evaluation metrics.

Counterfactual Learning-to-Rank. Prior work in counterfactual learning-to-rank (CLTR) primarily addresses position bias in implicit feedback using methods such as inverse propensity scoring (IPS) (Joachims et al., 2017) and doubly robust estimation (Oosterhuis, 2023). Extensions include modeling trust bias (Agarwal et al., 2019) and jointly correcting for both position and trust biases (Vardasbi et al., 2020). Recent approaches explore policy optimization via proximal updates (Gupta et al., 2024b) and extend this to trust-aware CLTR through proximal ranking objectives (Gupta et al., 2024a). While effective, these methods often focus narrowly on position bias or make strong assumptions, underscoring the need for broader utility-driven ranking frameworks, as pursued in this work.

Utility-oriented counterfactual reranking. Reranking methods enhance an initial ranked list by applying a secondary model to better optimize downstream objectives such as user utility, fairness, or diversity (Xi et al., 2024; Wang et al., 2025b). Recent work in counterfactual ranking predominantly follows a two-stage framework consisting of a generator and an evaluator (Xi et al., 2024; Shi et al., 2023; Ren et al., 2024; Wang et al., 2025b). For example, URCC (Xi et al., 2024) trains a set-aware utility model and employs a context-sensitive pairwise LambdaLoss to guide the ranker. NLGR (Wang et al., 2025b) leverages neighboring lists within a generator-evaluator setup for utility optimization. PRS (Feng et al., 2021) adopts beam search to generate candidate permutations and evaluates them using a permutation-wise scoring model, while PIER (Shi et al., 2023) uses SimHash to select top-K candidates from the full permutation space efficiently.

Reranking methods rely on a strong base ranker trained on factual data and typically explore counterfactuals near these initial permutations, leading to conservative exploration that limits discovery of globally optimal rankings. Moreover, they often overlook correcting for reward misspecification, which is crucial for reliable counterfactual learning.

Differential approximation to ranking. A key challenge in learning-to-rank is the mismatch between evaluation metrics (e.g., NDCG, MAP) and surrogate loss functions amenable to gradient-based optimization, due to the non-differentiable nature of sorting operations. To address this, prior work has either proposed smooth approximations to the rank function (e.g., ApproxNDCG (Qin et al., 2010)) or introduced differentiable approximations to argsort using soft permutation matrices (Grover et al., 2019; Prillo & Eisenschlos, 2020); for instance, PiRank (Swezey et al., 2021) and NeuralNDCG (Pobrotyn & Białobrzeski, 2021) utilize NeuralSort as a temperature-controlled surrogate. Another line of work leverages the Plackett–Luce distribution to model ranking policies in a differentiable manner (Oosterhuis, 2021). Methods like PG-RANK (Gao et al., 2023) use policy gradients to optimize the expected reward over the Plackett–Luce distribution based on REINFORCE, while ListNet (Cao et al., 2007) and ListMLE (Xia et al., 2008) employ the Plackett–Luce framework to derive smooth list-wise objectives.

The Plackett–Luce (PL) model offers efficient, closed-form gradients for ranking, but counterfactual learning with PL requires Monte Carlo sampling, leading to high-variance estimates in large action spaces. While variance reduction helps (Gao et al., 2023), unbiased learning fundamentally depends on stochastic logging, which is incompatible with real-world deterministic rankers designed for stability and trust. Soft permutation relaxations like SoftSort (Prillo & Eisenschlos, 2020) approximate permutations in continuous space, enabling gradient-based optimization without sampling. Though

computationally more expensive, they reduce variance and support end-to-end utility maximization. We pair SoftSort with a learned reward model that generalizes over logged data, enabling scalable training under deterministic logs. This approach trades unbiasedness for stability and practicality in real-world ranking systems.

3 LEARNING-TO-RANK PROBLEM: UTILITY MAXIMIZATION VS SORTING

A data sample of an LTR problem is a *query group* (QG) consisting of a query, q, and a set of L items, $\{x_\ell\}_{\ell=1}^L$, where L may vary. The *query* may represent, for example, a search string, a user profile, or other contextual information like device type and page layout. The *items* are candidate entities like webpages, songs, or products retrieved by an upstream system. We assume that the QGs are drawn i.i.d. from a distribution \mathcal{P} , i.e. $(q, \{x_\ell\}) \sim \mathcal{P}$. When a user is presented with a ranking/arrangement (permutation), $\pi: [L] \to [L]$, of the items of a QG, i.e. $(x_{\pi(1)}, \ldots, x_{\pi(L)})$, they interact with the ranked items, yielding a stochastic utility $U(q, \{x_\ell\}, \pi) \in \mathbb{R}$, which is a hidden function of the QG and the ranking. In typical internet systems, the utility can represent outcomes such as whether a user clicks or purchases any item, or continuous measures such as the total minutes of media consumed. Our objective is to learn a ranking policy, f, mapping the QGs to permutations, that maximizes the expected utility return, i.e.

$$f^* = \operatorname*{argmax}_{f} \mathbb{E}_{(q,\{x_{\ell}\}) \sim \mathcal{P}}[U(q,\{x_{\ell}\}, \pi = f(q,\{x_{\ell}\}))] \tag{1}$$

Based on the choice of the utility, this objective corresponds to business metrics like click-through rate, units sold, or streamed minutes. The main challenge here is that the hidden stochastic utility function U is not directly observable. Instead we are given a training dataset, \mathcal{D} , consisting of N QGs (indexed by i) and their observed utility $\{u_i\}$ under some logged rankings $\{\pi_i\}$, i.e. $\mathcal{D} = \{(q_i, \{x_{i,\ell}\}_{\ell=1}^{L_i}, \pi_i, u_i)\}_{i \in [N]}$. We assume that similar hold-out test and validation dataset are also available. This setting can be viewed as an offline one-step reinforcement-learning problem in which the state space is comprised of all possible QGs in the support of \mathcal{P} , the action space is comprised of all item permutations, and the reward is the observed utility.

In practice, most QGs are unique, so we observe only one out of L! possible rankings for each. Consequently, even if we propose a better alternative ranking for a given QG, the *counterfactual* utility it would have obtained remains unknown. To address this, traditional LTR algorithms optimize heuristic offline ranking metrics like Normalized Discounted Cumulative Gain (NDCG) (Järvelin & Kekäläinen, 2002; Burges, 2006), averaged over a test set. When a user interacts with a ranked QG, we also obtain per-item feedback signals $\{y_\ell \geq 0\}$ (e.g. whether an item was clicked or purchase, or how many minutes it was streamed). Usually, the overall QG-level utility u is some function of these per-item signals. Then, the NDCG of any new ranking $\widehat{\pi}$, on a QG with feedbacks $\{y_\ell\}$ can be defined as

$$NDCG(\widehat{\pi}, \{y_{\ell}\}) \stackrel{\triangle}{=} \frac{DCG(\widehat{\pi}, \{y_{\ell}\})}{DCG(\pi^*, \{y_{\ell}\})} \in [0, 1], \text{ where } DCG(r, \{y_{\ell}\}) \stackrel{\triangle}{=} \sum_{\ell=L} \frac{2^{y_{\ell}} - 1}{\log_2(1 + r^{-1}(\ell))}.$$
(2)

DCG assigns a gain $2^{y\ell}-1$ for the item x_ℓ in a test QG, but its contribution to the metric is discounted by its position $r^{-1}(\ell)$ under the ranking r. Thus NDCG is maximized when the items are ranked in the descending order of their feedback values. Traditional LTR methods (Burges, 2006; Swezey et al., 2021), aim to maximize NDCG by optimizing various continuous relaxations of it. This heuristic of learning to move items with higher feedback signal to the top of list have been highly successful, potentially because (i) items with positive feedback are usually relevant, and (ii) users tend to focus their attention on the top of the list. However, such offline metrics are now well-known to be sub-optimal as they do not perfectly align with the true (hidden) utility (1) we aim to maximize (Wang et al., 2023; Jeunen et al., 2024). In the next section, we propose REWARDRANK, a novel ranking framework designed to maximize the true LTR utility in a purely data-driven manner without the aid of any heuristics or user behavior assumptions.

4 REWARDRANK: DATA-DRIVEN LTR UTILITY MAXIMIZATION

In this section, we present the REWARDRANK framework, which aims to maximize the true (hidden) LTR utility defined in (1). At a high level, REWARDRANK proceeds in two stages. First, using the

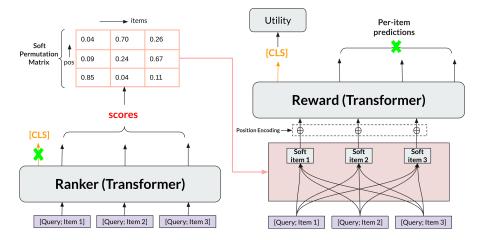


Figure 2: **REWARDRANK.** A ranker scores the items in a query group. These scores are used to compute soft item embeddings via a soft permutation matrix. Position encoded soft item embeddings are passed into a reward to estimate its utility. Finally, the ranker is optimized to maximize the predicted utility.

logged training data \mathcal{D} , it learns a reward model that predicts the counterfactual utility for any QG and permutation. Then it trains a ranker using the reward model's predictions as supervision, so as to maximize the expected counterfactual LTR utility of the ranker's item arrangement (ranking) policy.

4.1 STAGE 1: LEARNING THE UTILITY USING A REWARD MODEL

Let $g(q, \{x_\ell\}, \pi; \phi)$ denote the reward model (parameterized with ϕ) to predict the scalar utility for the QG $(q, \{x_\ell\})$ and a ranking π using. It is trained solely on the logged query groups, rankings and observed utilities in the training dataset \mathcal{D} . When the utility $U \in \{0, 1\}$ is a binary random variable (e.g. click, purchase), we train $g \in [0, 1]$ by minimizing the average binary cross-entropy loss between the observed u_i and the predicted utilities $\widehat{u}_i(\phi) := g(q_i, \{x_{i\ell}\}, \pi_i; \phi)$ over all $i \in [N]$:

$$\min_{\phi} \left[\text{RewardLoss}(\phi) \stackrel{\triangle}{=} -\frac{1}{N} \sum_{i=1}^{N} [u_i \log(\widehat{u}(\phi)) + (1 - u_i) \log(1 - \widehat{u}(\phi))] \right]. \tag{3}$$

When the true utility is a continuous random variable (e.g. minutes a song is streamed) we can use regression losses such mean squared error (MSE) $\min_{\phi}(1/N)\sum_{i=1}^{N}\|u_i-\widehat{u}_i(\phi)\|^2$. In our experiments, the reward model is instantiated with a transformer encoder, Enc, due to its ability to model functions over sequences (ranked list of items). Before passing a QG into Enc each query-item pair $[q,x_\ell]$ is embedded using a text encoder to create the token embedding \mathbf{e}_ℓ . Next, the ranking of the items π is encoded through position encodings $\{\mathbf{p}_k\}$. Then the position encoded tokens are passed to the transformer Enc. Finally, the predicted utility is computed as the sigmoid of a linear function of the <code>[CLS]</code> token output. This can be succinctly represented as:

$$g(q, \{x_{\ell}\}_{\ell=1}^{L}, \pi; \phi) = \sigma \left[\mathbf{v}^{\top} \operatorname{Enc}_{\text{reward}}^{\text{[CLS]}} \left(\left\{ \mathbf{e}_{\pi(k)} + \mathbf{p}_{k} \right\}_{k=1}^{L} \right) \right], \tag{4}$$

where $\mathbf{e}_{\pi(k)}$ is the token embedding of the query and the k-th ranked item. During the second stage of training the ranker, we freeze the reward model parameters ϕ .

Auxiliary per-item predictor: Typically the observed utility u is a byproduct of user's interaction with the items. So, we hypothesize that predicting the per-item feedback signals $\{y_\ell\}$ as an auxiliary task would improve the overall quality of the LTR utility prediction. Thus we include an auxiliary prediction head on the item tokens' outputs to predict the feedback signal observed at each ranked position $k \in [L]$. When $y_\ell \in \{0,1\}$ is binary, the predictions can be instantiated as

$$\widehat{y}_{k}(\phi) \stackrel{\triangle}{=} \sigma \left[\widetilde{\mathbf{v}}^{\top} \operatorname{Enc}_{\text{reward}}^{(k)} \left(\left\{ \mathbf{e}_{\pi(k)} + \mathbf{p}_{k} \right\}_{k=1}^{L} \right) \right], \quad \forall k \in [L],$$
 (5)

where $\operatorname{Enc}_{\operatorname{reward}}^{(k)}$ is the output token at the k-th position. We can learn $\widehat{y}_k(\phi)$ alongside $\widehat{u}(\phi)$ by adding the average cross-entropy loss between \widehat{y}_k and $y_{\pi(k)}$,

ItemLoss
$$(\phi) \stackrel{\triangle}{=} -\frac{1}{\sum_{i} L_{i}} \sum_{i=1}^{N} \sum_{k=1}^{L_{i}} [y_{i\pi(k)} \log(\widehat{y}_{ik}(\phi)) + (1 - y_{i\pi(k)}) \log(1 - \widehat{y}_{ik}(\phi))].$$
 (6)

as an additional regularizer to RewardLoss (3). Note that during the training of the ranker in the next stage, these auxiliary predictions can be discarded. Our ablation study in Section 5.1 shows that the per-item loss provides a moderate boost in performance.

4.2 STAGE 2: RANKER REWARD MAXIMIZATION THROUGH SOFT SORTING

Typically, rankers are modeled as scorers which assigns scores to the items in a QG. Then the items are ordered in the descending order of their scores to obtain the final ranking. We follow the same pattern and define $f(q, \{x_\ell\}; \theta)$ as a scoring-based ranker which maps a QG $(q, \{x_\ell\})$ to a set of item scores $\{s_\ell\}$. Following our reward model design, we instantiate f using the same transformer backbone architecture. Since QG has an unordered set of items, we do not use position encoding. Finally the score is computed as the sigmoid of a linear function of the output item tokens, i.e.

$$s_{\ell} \stackrel{\triangle}{=} f_{\ell}(q, \{x_{\ell}\}; \theta) \stackrel{\triangle}{=} \sigma \left[\mathbf{w}^{\top} \operatorname{Enc}_{\operatorname{ranker}}^{(\ell)} \left(\{ \mathbf{e}_{\ell} \}_{\ell=1}^{L} \right) \right], \quad \forall \ell \in [L],$$

where $\operatorname{Enc}_{\operatorname{ranker}}^{(\ell)}$ is the output token of the ℓ -th item. Our goal is to optimize the effective ranking $\widehat{\pi}$ induced by these scores so that it maximizes the expected counterfactual utility, which is a hidden from us. This is where the reward model comes in handy, as it helps us predict the counterfactual utility as $\widehat{u} := g(q, \{x_\ell\}, \widehat{\pi})$. However, since sorting (of the scores) is a discontinuous operation, it is challenging to optimize the scores to maximize the reward. To enable an end-to-end optimization of the scorer f, we resort to a continuous relaxation of the sorting operation.

Soft Permutation via SoftSort. SoftSort (Prillo & Eisenschlos, 2020) is a continuous relaxation of sorting operation. It defines a *unimodal row-stochastic* matrix (Swezey et al., 2021) as the *soft* permutation matrix $\widehat{\Pi}^{(\tau)} \in [0,1]^{L \times L}$. Row k of this matrix corresponds to a probability distribution of the k-the ranked item over the set of all items. Formally, we define

$$\widehat{\Pi}_{k,\ell}^{(\tau)} \stackrel{\triangle}{=} \frac{\exp\left(-\frac{1}{\tau}\left|s_{\ell} - s_{\widehat{\pi}(k)}\right|\right)}{\sum_{\ell'=1}^{L} \exp\left(-\frac{1}{\tau}\left|s_{\ell'} - s_{\widehat{\pi}(k)}\right|\right)}, \quad \forall \, k, \ell \in [L],$$

$$(8)$$

where τ is a temperature parameter and $\widehat{\pi}(k)$ is the k-th ranked items when (hard) sorting by the scores $\{s_\ell\}$. $\widehat{\Pi}^{(\tau)}$ is a continuous function of the scores $\{s_\ell\}$ and when $\tau \to 0$, $\widehat{\Pi}^{(\tau)}$ tends to the binary hard-permutation matrix $\widehat{\Pi}$, where

$$\lim_{\tau \to 0} \widehat{\Pi}_{k,\ell}^{(\tau)} = \widehat{\Pi}_{k,\ell} \stackrel{\triangle}{=} \mathbb{I}\{k = \widehat{\pi}(\ell)\}, \quad \forall \, k, \ell \in [L],$$
(9)

assuming the scores are unique. Using this soft permutation matrix, we can compute a soft item embedding $\hat{e}_k^{(\tau)}$ at position k as the following convex combination of the true item embeddings

$$\widehat{\mathbf{e}}_{k}^{(\tau)} \stackrel{\triangle}{=} \sum_{\ell \in [L]} \widehat{\Pi}_{k,\ell}^{(\tau)} \mathbf{e}_{\ell} . \tag{10}$$

It is easy to verify that $\widehat{\mathbf{e}}_k^{(\tau)} \to \mathbf{e}_{\widehat{\pi}(k)}$ when $\tau \to 0$. Note that there are alternate soft permutation matrices like NeuralSort (Grover et al., 2019), but we adopt SoftSort for its simplicity and state of the art performance (Prillo & Eisenschlos, 2020). We then compute a soft reward for these soft item embeddings using

$$\widehat{g}(\theta) \stackrel{\triangle}{=} g(q, \{x_{\ell}\}, \widehat{\Pi}^{(\tau)}) \stackrel{\triangle}{=} \sigma \left[\mathbf{v}^{\top} \operatorname{Enc}_{\text{reward}}^{\text{[CLS]}} \left(\left\{ \widehat{\mathbf{e}}_{k}^{(\tau)} + \mathbf{p}_{k} \right\}_{k=1}^{L} \right) \right], \tag{11}$$

This allows us to compute an approximate predicted reward (11) as a continuous function over the ranker scores $\{s_\ell\}$ through the SoftSort matrix. Finally, we optimize the parameters of scorer f to maximize the average approximate reward over the training set in an end-to-end manner:

$$\min_{\theta} \left[\text{RankerLoss}(\theta) \stackrel{\triangle}{=} \frac{1}{N} \sum_{i=1}^{N} \widehat{g}_i(\theta) \right]. \tag{12}$$

Even though REWARDRANK is maximizing the predicted utility of the soft ranking, we hypothesize that it generalizes well and produces rankings with higher expected counterfactual utility than prior LTR methods.

Mitigating reward misspecification. One challenge of reward modeling the hidden counterfactual utility is model misspecification, i.e. a gap between the predicted and the true utilities. A misspecified reward can misguide the ranker into wrong ranking policies. To mitigate this issue we propose a sample reweighting scheme which modifies the ranker loss as

RankerLoss<sup>(
$$\lambda$$
)</sup> $(\theta) \stackrel{\triangle}{=} \frac{1}{N} \sum_{i=1}^{N} w_i \cdot \widehat{g}_i(\theta)$, where $w_i = 1 - \lambda |u_i - \widehat{u}_i(\theta)|$, $\forall i \in [N]$. (13)

Above loss is a pessimistic lowerbound to $\operatorname{RankerLoss}(\theta)$ (12). This reward down-weighting scheme is motivated by a conjecture that when the observed utility u_i for the i-th training QG and the corresponding prediction $\widehat{u}_i(\theta)$ are different, the utility prediction on new ranking of this QG would also be less reliable. Through an ablation in Section 5.1 we show that reward misspecification correction slightly improves the REWARDRANK performance.

5 EXPERIMENTAL RESULTS

Datasets. Public large-scale datasets for learning-to-rank (LTR), especially in counterfactual settings, are scarce. To the best of our knowledge, we propose the first reproducible testbeds for counterfactual ranking evaluation. We utilize two existing large-scale LTR datasets: Baidu-ULTR (Hager et al., 2024; Zou et al., 2022) and Amazon KDD-Cup (Reddy et al., 2022), to construct these testbenches, enabling rigorous evaluation of permutation-aware ranking policies. Baidu-ULTR contains 1.8M query groups (11.7M query-document pairs) and 590K validation/test sessions. Amazon KDD-Cup comprises 130K queries and 2.6M annotated query-product pairs with rich textual metadata. We generate 400K training and 50K validation/test query groups by sampling permutations of products per query. See Appendix B.1 for further details.

Implementation Details and Baselines. Our reward models and rankers are based on a transformer architecture with 12 layers, 768 hidden dimensions, 12 attention heads, and roughly 110M parameters. We set $\tau=0.5$ and $\lambda=0.7$ for all Rewardrank experiments, based on tuning over a held-out set. Ablations with varying values and further implementation details are provided in Appendix B.2. For comparison, we implement two utility-based counterfactual ranking methods: URCC (Xi et al., 2024), which uses a LambdaLoss-based pairwise objective, and PG-rank (Gao et al., 2023), which applies Plackett–Luce modeling with policy gradients. Our variants, URCC* and PG-rank*, replace the original architectures with our transformer-based reward model for improved performance. While PG-rank was originally explored with known NDCG utility, PG-rank* uses the same learned reward model as Rewardrank. Additionally, we train standard LTR baselines: ListNet (Cao et al., 2007), ListMLE (Xia et al., 2008), LambdaRank (Wang et al., 2018), and PiRank (Swezey et al., 2021), all using the same transformer architecture for fair comparison across supervision methods.

5.1 Large-Scale Reproducible Testbenches for Counterfactual LTR

To enable reproducible evaluation of ranking policies without online A/B testing, we introduce two complementary testbeds: PO-Eval, which leverages a parametric click model, and LAU-Eval, which simulates human-like shopping behavior via LLM reasoning. Together, they enable holistic, offline assessment of ranking algorithms under both statistical and behavioral lenses.

Parametric Oracle Evaluation (PO-Eval). To simulate a click-based counterfactual recommendation setting, we build a testbed from the Baidu-ULTR dataset (Hager et al., 2024), employing a pretrained parametric IPS model as the oracle for supervision. This model estimates the click probability at position ℓ as $P(C) = P(E_{\ell}) \cdot P(R_{q,i})$, where $P(E_{\ell})$ is the position-dependent examination probability and $P(R_{q,i})$ is the relevance-conditioned click probability, independent of position. We use this oracle to sample binary clicks for training and later reuse it for counterfactual evaluation of new ranking policies. For each ranked query group (QG), we compute the expected utility as the probability of at least one click and the observed utility as a binary indicator of at least

379380381382

394

395

396

397

398

399

400

401

402

403

404

405

406

407

408

409

410 411

412

413

414

415

416

417

418

419

420

421

422

423

424

425

426

428

429

430

431

Table 1: **Comparison of counterfactual ranking methods under two evaluation settings.** Left partition: IPS-Oracle-based utility and relevance-aware NDCG metrics. Right partition: LLM-Eval purchase prediction metric

Method	PO-Eval		LAU-Eval	
	$\overline{\mathbb{E}[U_{\text{IPS-O}}(q,\{i\}_L)]}$	NDCG _{click}	$\boxed{\mathbb{E}[1\{\text{purchase}=1\}]}$	$NDCG_{purchase}$
Upper-Bound	0.553	1.000	_	_
Policy in data	0.501	_	0.497 ± 0.008	1.00 ± 0.000
ListNet (Cao et al., 2007)	0.523 ± 0.0007	0.376 ± 0.0002	0.521 ± 0.009	0.8611 ± 0.009
ListMLE (Xia et al., 2008)	0.522 ± 0.0007	0.377 ± 0.0002	0.522 ± 0.008	0.8610 ± 0.003
LambdaRank (Wang et al., 2018)	0.524 ± 0.0007	0.378 ± 0.0002	0.523 ± 0.009	0.8610 ± 0.009
PiRank (Swezey et al., 2021)	0.525 ± 0.0007	0.378 ± 0.0002	0.528 ± 0.007	0.8623 ± 0.005
URCC* (Xi et al., 2024)	0.462 ± 0.0005	0.315 ± 0.0004	0.471 ± 0.008	0.8621 ± 0.009
PG-rank* (Gao et al., 2023)	0.501 ± 0.0005	0.327 ± 0.0002	0.489 ± 0.007	0.8630 ± 0.009
REWARDRANK	0.536 ± 0.0007	0.370 ± 0.0002	0.561 ± 0.008	0.8628 ± 0.009

one sampled click. This setup provides a realistic and repeatable framework for evaluating how well learned rankers align with user behavior modeled by the IPS oracle. See Appendix A.1 for details on the parametric model and the derivation of utility metrics. Table 6 reports counterfactual evaluation results using PO-Eval, where we leverage a pre-trained parametric IPS-Oracle to simulate user clicks and assess ranking quality. The IPS-based utility $\mathbb{E}[U_{\text{IPS-O}}(q,\{i\}_L)]$ captures the expected probability of at least one click per ranked list, while NDCGclick measures how high are the originally clicked items in the test dataset ranked. The *Upper-Bound* is computed by ranking items in descending order of P(R), which maximizes utility due to the rearrangement inequality (Day, 1972) (see Appendix A). Among all methods, REWARDRANK achieves the best performance on true utility. Traditional LTR baselines (ListNet, ListMLE, LambdaRank, PiRank), trained on per-item IPS-sampled clicks, fail to fully capture permutation-level utility. URCC* performs poorly due to its reliance on initialization quality, while PG-rank* improves over URCC* but remains suboptimal. Overall, REWARDRANK consistently outperforms existing methods, highlighting the benefit of direct reward optimization and counterfactual supervision. RewardRank attains the highest utility under IPS-Oracle, despite slightly lower NDCG_{click} than some baselines. This reflects a key distinction: proxy metrics like NDCG may not fully capture listwise user utility. By directly optimizing counterfactual reward, RewardRank better aligns with behavioral objectives beyond conventional ranking accuracy.

LLM-based User Simulation (LAU-Eval). While PO-Eval captures position bias via IPS-Oracle supervision, it does not account for broader behavioral patterns such as brand bias, similarity aversion, or irrelevance bias. To complement PO-Eval and more fully assess human-centered ranking behavior, we introduce the LAU-Eval framework. In this setup, a large language model (LLM) is prompted to simulate user shopping behavior given a query and its associated product list from the Amazon KDD-Cup dataset. The prompt incorporates behavioral factors such as position bias, brand bias, irrelevance bias, and color bias (full details are provided in Appendix C.2). The LLM generates a binary purchase decision $D(\text{purchase}) \in \{0,1\}$, which serves as the reward signal for training a reward model and optimizing rankers. For evaluation, the same prompt is used: each ranker's re-ranked outputs are re-scored by the LLM, and performance is reported as the average purchase decision rate on a held-out test set. For LTR methods that do not rely on reward modeling, we instead use the LLM-selected item as the training signal. Higher values indicate stronger alignment with human-centered behavioral criteria. Refer to the Appendix Section C.2 for implementation details. Under LAU-Eval, which evaluates binary purchase decisions produced by the LLM, we observe notable differences across methods. The policy_in_data baseline, representing the original item ordering, achieves an average purchase rate of 0.497. Classical listwise approaches such as ListNet, ListMLE, and LambdaRank provide slight improvements, reaching 0.500–0.513, while PiRank performs marginally better at 0.515, indicating advantages of modeling permutation-level interactions. Interestingly, recent counterfactual methods like URCC* and PG-rank* underperform even the baseline, suggesting that their learning objectives may not fully capture the nuanced user behaviors emphasized in LAU-Eval. In contrast, REWARDRANK achieves the highest purchase rate at 0.561, significantly outperforming all baselines and prior methods. These results establish LAU-Eval as a behaviorally grounded evaluation framework and demonstrate REWARDRANK's effectiveness in

capturing decision-relevant user dynamics beyond conventional ranking criteria. We observe that a high $NDCG_{purchase}$ score does not necessarily translate into high purchase decisions by the LLM. Many baselines achieve similar NDCG values but vary significantly in actual predicted purchases. This highlights that $NDCG_{purchase}$, while ranking-aware, may not faithfully reflect true listwise utility as captured by the LLM. In contrast, RewardRank consistently achieves both high purchase rates and competitive NDCG, demonstrating better alignment with user-centric objectives.

Ablations. We ablate the per-item regularizer and the two parameters of REWARDRANK: the SoftSort temperature τ , which controls the sharpness of the permutation approximation (8), and the misspecification correction strength λ , which down-weights rewards on QGs with high prediction error (13). Removing the auxiliary item-level reward loss (Eqn 6) decreased the final expected counterfactual utility of the learned ranker. This indicates that learning to predict the per-item feedback enhances the reward model's generalization and hence improves downstream ranking performance. We find that $\tau=0.5$ and $\lambda=0.7$ achieve the best trade-off between stability and performance. Full details of these ablations are reported in Appendix D.

5.2 BAIDU-ULTR DATASET WITH REAL USER CLICKS

Table 2: **Baidu-ULTR with real clicks.** RE-WARDRANK achieves SOTA performance. † metrics are taken from Hager et al. (2024)

Method	$DCG_{\rm rel}@5$	$DCG_{rel}@10$
Point IPS [†] (Hager et al., 2024)	4.79	7.43
List IPS [†] (Hager et al., 2024)	5.20	7.88
LambdaRank Hager et al. (2024)	5.45	8.23
ListNet Cao et al. (2007)	5.05	7.64
ListMLE Xia et al. (2008)	5.13	7.88
PiRank Swezey et al. (2021)	5.23	8.01
URCC* Xi et al. (2024)	5.01	7.44
PG-rank* Gao et al. (2023)	5.09	7.62
REWARDRANK	5.83	8.42

While we previously used the Baidu-ULTR dataset within the PO-Eval framework under IPS-Oracle supervision, here we instead rely directly on the real click signals provided in the data. Following the protocol in (Hager et al., 2024), models are trained with a binary click labels for each query groups (u = 1 if any item is clicked, otherwise 0). Since counterfactual evaluation is not feasible here, we follow (Hager et al., 2024) and report Relevance DCG at 5 and 10, computed on the human-assigned relevance labels provided with the test set¹. As shown in Table 2, our method achieves a new state of the art DCG@5 and DCG@10 across all base-

lines. Importantly, these improvements are observed on human-assigned relevance labels that were never used in training by any method. This is particularly noteworthy given that our method is not optimized for relevance DCGThese results highlight both the robustness of our approach and its ability to generalize to real human feedback in large-scale search settings.

6 CONCLUSION

In this work, we introduced REWARDRANK, a general framework for counterfactual learning-to-rank that uses a learned, permutation-aware reward model to optimize the ranker via a differentiable soft permutation operator. To support reproducible counterfactual utility evaluation, we proposed two large-scale evaluation suites: PO-Eval, which leverages a parametric IPS oracle to simulate position-biased click behavior, and LAU-Eval, which uses large language models to capture richer, human-aligned user preferences beyond position. Across these benchmarks REWARDRANK achieves achieved the highest utility across all baselines. Further, it achieved a new state of the art performance on Baidu-ULTR dataset with real clicks. Our findings underscore the value of learning data-driven counterfactual utility models and enabling principled evaluation for building more robust and behaviorally aligned ranking systems.

REFERENCES

Aman Agarwal, Kenta Takatsu, Ivan Zaitsev, and Thorsten Joachims. A general framework for counterfactual learning-to-rank. In *Proceedings of the 42nd International ACM SIGIR Conference on Research and Development in Information Retrieval*, pp. 5–14, 2019.

¹We report DCG rather than NDCG for consistency with (Hager et al., 2024).

491

493

494

495

496

497 498

499

500

501

502

504 505

506 507

508 509

510

511 512

513

514

515

516 517

519

520

521

522 523

524

525

526

527

528

529 530

531

532

534

- 486 Alexander Buchholz, Ben London, Giuseppe Di Benedetto, Jan Malte Lichtenberg, Yannik Stein, and 487 Thorsten Joachims. Counterfactual ranking evaluation with flexible click models. In *Proceedings* 488 of the 47th International ACM SIGIR Conference on Research and Development in Information 489 Retrieval, pp. 1200–1210, 2024.
- Chris Burges. From ranknet to lambdarank to lambdamart: An overview. In Microsoft Research Technical Report, 2006. 492
 - Chris Burges, Tal Ragno, and Quoc Viet Le. Learning to rank using gradient descent. In *ICML*, 2005a.
 - Chris JC Burges, Tal Shaked, Erin Renshaw, Ari Lazier, Matt Deeds, Nicole Hamilton, and Greg Hullender. Ranknet: Gradient descent learning for rank function optimization. In SIGIR, 2005b.
 - Zhe Cao, Tao Qin, Tie-Yan Liu, Ming-Feng Tsai, and Hang Li. Learning to rank: from pairwise approach to listwise approach. In Proceedings of the 24th international conference on Machine learning, pp. 129–136, 2007.
 - Yufei Chen, Wei Li, Jianyang Sun, and Ruiming Wang. Reducing popularity influence by addressing position bias in recommender systems. arXiv preprint arXiv:2412.08780, 2024. URL https: //arxiv.org/abs/2412.08780.
 - Peter W Day. Rearrangement inequalities. Canadian Journal of Mathematics, 24(5):930–943, 1972.
 - Yufei Feng, Yu Gong, Fei Sun, Junfeng Ge, and Wenwu Ou. Revisit recommender system in the permutation prospective. arXiv preprint arXiv:2102.12057, 2021.
 - Yoav Freund, Raj Iyer, Robert E Schapire, and Yoram Singer. An efficient boosting algorithm for combining preferences. In *JMLR*, 2003.
 - Ge Gao, Jonathan D Chang, Claire Cardie, Kianté Brantley, and Thorsten Joachim. Policy-gradient training of language models for ranking. arXiv preprint arXiv:2310.04407, 2023.
 - Aditya Grover, Eric Wang, Aaron Zweig, and Stefano Ermon. Stochastic optimization of sorting networks via continuous relaxations. arXiv preprint arXiv:1903.08850, 2019.
 - Shashank Gupta, Harrie Oosterhuis, and Maarten de Rijke. Practical and robust safety guarantees for advanced counterfactual learning to rank. In Proceedings of the 33rd ACM International Conference on Information and Knowledge Management, pp. 737–747, 2024a.
 - Shashank Gupta, Harrie Oosterhuis, and Maarten de Rijke. Proximal ranking policy optimization for practical safety in counterfactual learning to rank. arXiv preprint arXiv:2409.09881, 2024b.
 - Philipp Hager, Romain Deffayet, Jean-Michel Renders, Onno Zoeter, and Maarten de Rijke. Unbiased learning to rank meets reality: Lessons from baidu's large-scale search dataset. In *Proceedings* of the 47th International ACM SIGIR Conference on Research and Development in Information Retrieval, pp. 1546–1556, 2024.
 - Kalervo Järvelin and Jaana Kekäläinen. Cumulated gain-based evaluation of ir techniques. ACM Transactions on Information Systems (TOIS), 20(4):422–446, 2002.
 - Olivier Jeunen, Ivan Potapov, and Aleksei Ustimenko. On (normalised) discounted cumulative gain as an off-policy evaluation metric for top-n recommendation. In Proceedings of the 30th ACM SIGKDD conference on knowledge discovery and data mining, pp. 1222–1233, 2024.
 - Thorsten Joachims. Optimizing search engines using clickthrough data. In KDD, 2002.
- 535 Thorsten Joachims, Adith Swaminathan, and Tobias Schnabel. Unbiased learning-to-rank with biased 536 feedback. In Proceedings of the tenth ACM international conference on web search and data 537 mining, pp. 781–789, 2017.
 - Wouter Kool, Herke van Hoof, and Max Welling. Buy 4 reinforce samples, get a baseline for free! 2019.

- Ying Li, Yu Wang, and Kevin C.C. Chan. We match! building online brand engagement behaviours through social media influencers. *Journal of Retailing and Consumer Services*, 76:103640, 2025. doi: 10.1016/j.jretconser.2024.103640. URL https://www.sciencedirect.com/science/article/pii/S0969698924004429.
 - Ilya Loshchilov and Frank Hutter. Decoupled weight decay regularization. *arXiv preprint arXiv:1711.05101*, 2017.
 - Harrie Oosterhuis. Computationally efficient optimization of plackett-luce ranking models for relevance and fairness. In *Proceedings of the 44th International ACM SIGIR Conference on Research and Development in Information Retrieval*, pp. 1023–1032, 2021.
 - Harrie Oosterhuis. Doubly robust estimation for correcting position bias in click feedback for unbiased learning to rank. *ACM Transactions on Information Systems*, 41(3):1–33, 2023.
 - Przemysław Pobrotyn and Radosław Białobrzeski. Neuralndcg: Direct optimisation of a ranking metric via differentiable relaxation of sorting. *arXiv preprint arXiv:2102.07831*, 2021.
 - Sebastian Prillo and Julian Eisenschlos. Softsort: A continuous relaxation for the argsort operator. In *International Conference on Machine Learning*, pp. 7793–7802. PMLR, 2020.
 - Tao Qin, Tie-Yan Liu, and Hang Li. A general approximation framework for direct optimization of information retrieval measures. *Information retrieval*, 13:375–397, 2010.
 - Chandan K Reddy, Lluís Màrquez, Fran Valero, Nikhil Rao, Hugo Zaragoza, Sambaran Bandyopadhyay, Arnab Biswas, Anlu Xing, and Karthik Subbian. Shopping queries dataset: A large-scale esci benchmark for improving product search. *arXiv preprint arXiv:2206.06588*, 2022.
 - Yuxin Ren, Qiya Yang, Yichun Wu, Wei Xu, Yalong Wang, and Zhiqiang Zhang. Non-autoregressive generative models for reranking recommendation. In *Proceedings of the 30th ACM SIGKDD Conference on Knowledge Discovery and Data Mining*, pp. 5625–5634, 2024.
 - Xiaowen Shi, Fan Yang, Ze Wang, Xiaoxu Wu, Muzhi Guan, Guogang Liao, Wang Yongkang, Xingxing Wang, and Dong Wang. Pier: Permutation-level interest-based end-to-end re-ranking framework in e-commerce. In *Proceedings of the 29th ACM SIGKDD Conference on Knowledge Discovery and Data Mining*, pp. 4823–4831, 2023.
 - Robin Swezey, Aditya Grover, Bruno Charron, and Stefano Ermon. Pirank: Scalable learning to rank via differentiable sorting. *Advances in Neural Information Processing Systems*, 34:21644–21654, 2021.
 - Amos Tversky and Itamar Simonson. Extremeness aversion and attribute-balance effects in choice. *Journal of Consumer Research*, 31(2):249–257, 2004. doi: 10.1086/172209. URL https://academic.oup.com/jcr/article/31/2/249/1824942.
 - Ali Vardasbi, Harrie Oosterhuis, and Maarten de Rijke. When inverse propensity scoring does not work: Affine corrections for unbiased learning to rank. In *Proceedings of the 29th ACM International Conference on Information & Knowledge Management*, pp. 1475–1484, 2020.
 - Shuheng Wang, Weijie Zhang, Yiqun Xu, and Ji-Rong Wen. Decoy effect in search interaction: Understanding user behavior. In *Proceedings of the ACM Web Conference* 2025. ACM, 2025a. doi: 10.1145/3708884. URL https://dl.acm.org/doi/10.1145/3708884.
 - Shuli Wang, Xue Wei, Senjie Kou, Chi Wang, Wenshuai Chen, Qi Tang, Yinhua Zhu, Xiong Xiao, and Xingxing Wang. Nlgr: Utilizing neighbor lists for generative rerank in personalized recommendation systems. *arXiv preprint arXiv:2502.06097*, 2025b.
 - Xiaojie Wang, Ruoyuan Gao, Anoop Jain, Graham Edge, and Sachin Ahuja. How well do offline metrics predict online performance of product ranking models? In *Proceedings of the 46th International ACM SIGIR conference on Research and Development in Information Retrieval*, pp. 3415–3420, 2023.

- Xuanhui Wang, Cheng Li, Nadav Golbandi, Michael Bendersky, and Marc Najork. The lambdaloss framework for ranking metric optimization. In Proceedings of the 27th ACM international conference on information and knowledge management, pp. 1313–1322, 2018. Qiang Wu, Chris JC Burges, Krysta M Svore, and Jianfeng Gao. Adapting boosting for information retrieval measures. In Information Retrieval, 2010. Yunjia Xi, Weiwen Liu, Xinyi Dai, Ruiming Tang, Qing Liu, Weinan Zhang, and Yong Yu. Utility-oriented reranking with counterfactual context. ACM Transactions on Knowledge Discovery from Data, 18(8):1-22, 2024. Fei Xia, Tie-Yan Liu, Jue Wang, Wensheng Zhang, and Hang Li. Listwise approach to learning to rank: theory and algorithm. In ICML, 2008.
 - Lixin Zou, Haitao Mao, Xiaokai Chu, Jiliang Tang, Wenwen Ye, Shuaiqiang Wang, and Dawei Yin. A large scale search dataset for unbiased learning to rank. *Advances in Neural Information Processing Systems*, 35:1127–1139, 2022.

A PROOFS AND CONCEPTUAL DETAILS

A.1 CLICK-BASED UTILITY FOR PO-EVAL.

The IPS-Oracle models user clicks as a probabilistic function of both examination and item relevance. Specifically, the click probability at position ℓ is modeled as:

$$P(C_{q,i_{\pi_{\ell}},\ell}) = P(E_{\ell}) \cdot \sigma(R_{q,i_{\pi_{\ell}}})$$

where $P(E_\ell)$ is the position-dependent examination probability, and $\sigma(R_{q,i_{\pi_\ell}})$ is the relevance-based click probability estimated by the reward model. Given a query group $(q,\{i\}_L,\pi)$, the click indicator for each item is sampled as:

$$c_{q,i_{\pi_{\ell}},\ell} \sim \text{Bernoulli}(P(C_{q,i_{\pi_{\ell}},\ell}))$$

We define the utility of a ranked list as a binary signal indicating whether at least one item was clicked:

$$U(q, \{i\}_L, \pi) = \begin{cases} 1, & \text{if } \sum_{\ell=1}^L c_{q, i_{\pi_\ell}, \ell} > 0, \\ 0, & \text{otherwise.} \end{cases}$$
 (14)

To obtain a differentiable approximation of this binary utility, we define the expected probability of at least one click as:

$$U_{\text{IPS}}(q, \{i\}_L, \pi) = 1 - \prod_{\ell=1}^{L} \left(1 - P(E_{\ell}) \cdot \sigma\left(R_{q, i_{\pi_{\ell}}}\right) \right)$$
 (15)

This smoothed utility approximates the expected user engagement over the list and serves as a continuous training signal. The reward model is trained to predict the binary group-level utility $U \in \{0,1\}$, while the ranker maximizes the soft expected utility U_{IPS} under permutations. This formulation bridges synthetic click modeling with realistic counterfactual feedback, enabling effective utility-based ranking even in the absence of direct supervision on full permutations.

A.2 IDEAL IPS-ORACLE: REARRANGEMENT INEQUALITY

Theorem 1 (Ideal Ranking Maximizes Utility via Rearrangement Inequality). Let $\mathbf{r} = (r_1, \dots, r_n) \in \mathbb{R}^n_{\geq 0}$ be a vector of predicted relevance scores, and let $\mathbf{e} = (e_1, \dots, e_n) \in \mathbb{R}^n_{\geq 0}$ be a non-increasing sequence of examination probabilities: $e_1 \geq e_2 \geq \dots \geq e_n$. Let π^* be the permutation that sorts \mathbf{r} in descending order: $r_{\pi^*(1)} \geq r_{\pi^*(2)} \geq \dots \geq r_{\pi^*(n)}$. Then, for any permutation $\pi \in \mathcal{S}_n$, we have:

$$\sum_{i=1}^{n} e_i \cdot r_{\pi^*(i)} \ge \sum_{i=1}^{n} e_i \cdot r_{\pi(i)}$$

Proof. This is a direct consequence of the classical rearrangement inequality (Day, 1972). Among all permutations π of the relevance scores, the weighted sum $\sum_i e_i \cdot r_{\pi(i)}$ is maximized when the $r_{\pi(i)}$ are ordered in the same way as the e_i , i.e., both decreasing. Hence, sorting \mathbf{r} in descending order and aligning it with the already sorted \mathbf{e} gives the maximal utility.

Above analysis shows that ideal ranking order under the IPS Oracle is ordering the items such the sorting of item relevance scores and examination probabilities result in the same permutation.

A.3 DETAILS OF BASELINES

A.3.1 PG-RANK*: PG-RANK WITH LEARNED REWARD MODEL.

We extend the PG-Rank framework (Gao et al., 2023) by replacing the handcrafted reward (e.g., NDCG) with a learned reward model $g(q,\{i\}_L,\pi)$ that scores entire permutations based on user utility. The goal is to maximize the expected reward under the Plackett–Luce distribution induced by the ranker's scores:

$$\mathcal{L}_{\text{PG-reward}}(\theta) = \mathbb{E}_{\pi \sim \mathbb{P}_{\theta}} \left[g(q, \{i\}_{L}, \pi) \right]$$
 (16)

where $\mathbb{P}_{\theta}(\pi)$ is the Plackett–Luce distribution over permutations, parameterized by model scores s_1, \ldots, s_L for each item in the query group. To enable backpropagation through the sampled permutations, we adopt the Gumbel-Softmax trick as in the original PG-Rank implementation, which provides a continuous relaxation of the discrete sampling process.

The gradient of this objective is estimated using the REINFORCE trick with a baseline b for variance reduction (adopted from PG-rank (Gao et al., 2023; Kool et al., 2019)):

$$abla_{\theta} \mathcal{L}_{\text{PG-reward}} \approx \frac{1}{K} \sum_{k=1}^{K} \left[\left(g(\pi^{(k)}) - b \right) \cdot abla_{\theta} \log \mathbb{P}_{\theta}(\pi^{(k)}) \right]$$
 (17)

where $\pi^{(k)} \sim \mathbb{P}_{\theta}$ are K Monte Carlo samples drawn from the Plackett–Luce model.

The log-probability of a sampled permutation π under this model is given by:

$$\log \mathbb{P}_{\theta}(\pi) = \sum_{k=1}^{L} \left[s_{\pi(k)} - \log \sum_{j=k}^{L} \exp(s_{\pi(j)}) \right]$$
 (18)

This formulation allows us to train the ranking model directly on learned, utility-aligned reward signals using fully differentiable, sample-based policy gradients.

A.3.2 URCC* WITH LEARNED REWARD MODEL.

URCC (Xi et al., 2024) proposes a two-stage counterfactual reranking framework that jointly learns a set-aware utility function and a context-aware reranker. The utility model in URCC is itself learned from data and used to guide the optimization of the reranker via a pairwise ranking loss over permutations. Since the official implementation of URCC is not publicly available, we reimplemented the method using our own architecture.

In our version of URCC*, we retain the core two-stage structure but implement the utility model $g(q,\{i\}_L,\pi)$ as a Transformer-based encoder trained to predict user utility over full permutations. Given a query q and a set of items $\{i\}_L$, the reward model assigns a scalar score to a permutation π :

$$S_{\pi} = g(q, \{i\}_{L}, \pi) \tag{19}$$

Following URCC, we then train the ranker f_{θ} to maximize this learned reward by optimizing a context-aware pairwise loss. For a pair of permutations (π^+, π^-) such that $g(q, \{i\}_L, \pi^+) > g(q, \{i\}_L, \pi^-)$, we minimize the following objective:

$$\mathcal{L}_{\text{URCC-reward}}(\theta) = \mathbb{E}_{(\pi^+, \pi^-) \sim \mathcal{P}} \left[\log \left(1 + \exp \left(-(S_{\pi^+} - S_{\pi^-}) \right) \right) \right] \tag{20}$$

Here, \mathcal{P} denotes the set of sampled permutation pairs with preference orderings induced by the reward model. Our implementation uses neighborhood-based sampling (e.g., pairwise swaps) to construct π^+ and π^- from the base ranking.

Thus, while our training procedure is structurally consistent with the original URCC framework, we employ a more expressive Transformer-based reward model to capture user behavior better and align optimization with utility-oriented objectives.

A.4 COMPARISON OF TIME COMPLEXITY AND COUNTERFACTUAL SPACE EXPLORATION

Table 3 compares the time complexity of three methods: URCC*, PG-rank*, and REWARDRANK. The per-iteration time complexity is analyzed based on the number of calls to the reward model.

• $URCC^*$: n^2 , where n is the number of items in the list. $URCC^*$ explores the neighborhood of factual permutations, leading to quadratic complexity due to pairwise comparisons.

Table 3: **Comparison of Time Complexity** for URCC*, PG-rank*, and REWARDRANK in term of number of calls to the reward model per iteration on Baidu-ULTR dataset.

Method	Time Complexity	Wall-Clock Time	Description
PiRank	1	\sim 6 hours	No call to the reward model.
URCC*	n^2	\sim 34 hours	Neighborhood search, pessimistic.
PG-rank*	k	\sim 16 hours ($k=10$)	Needs large k for convergence.
RewardRank	1	\sim 7 hours	Full counterfactual space exploration.

URCC* only explores the neighborhood of factual permutations, meaning it performs limited counterfactual exploration. This approach is considered pessimistic because it does not explore the entire space of possible rankings, which could miss potentially better arrangements.

- **PG-rank***: k, where k is the number of Monte Carlo (MC) samples. While k is typically smaller than n, PG-rank* requires large k values and variance reduction baselines to converge. PG-Rank uses Monte Carlo (MC) sampling to explore a broader counterfactual space, but this approach requires large MC samples to converge effectively. To ensure stable and accurate exploration, PG-Rank relies on variance-reduction baselines. However, it still faces challenges in accurately capturing all potential counterfactual configurations without a very large number of samples.
- **REWARDRANK:** 1, as it performs a single call to the reward model. RewardRank explores the entire counterfactual space efficiently and can focus on more certain regions with reward misspecification mitigation.

In Table 3, we also provide the overall wall-clock time to train the model under the above method for the Baidu-ULTR dataset. Each model is trained for 21 epochs.

B EXPERIMENTATION DETAILS

B.1 DATASETS

Baidu-ULTR Reranking Dataset. The Baidu-ULTR dataset (Hager et al., 2024), a large-scale subset of the Baidu-ULTR corpus (Zou et al., 2022), contains user click interactions over web search queries. It includes 1.8M query groups (11.7M query-document pairs) and 590K validation/test sessions (4.8M pairs). The authors of (Hager et al., 2024) provide BERT-based CLS embeddings for each query-document pair.

We use the large-scale reranking dataset introduced by (Hager et al., 2024): publicly available at: https://huggingface.co/datasets/philipphager/baidu-ultr_uva-mlm-ctr, derived from the original Baidu-ULTR corpus (Zou et al., 2022). This dataset is constructed from real-world user interactions on Baidu's production search engine and is designed to support robust evaluation of learning-to-rank models in counterfactual settings.

Each session consists of a user query, a candidate list of documents retrieved by an upstream ranker, the original presented ranking, and user interaction logs (e.g., clicks and dwell time). For each query-document pair, the dataset provides both sparse lexical features (e.g., BM25, TF-IDF, query likelihood) and dense semantic representations.

Table 4: Statistics of the Baidu-ULTR reranking dataset (Hager et al., 2024).

Split	#Query Groups	#Query-Document Pairs
Training Validation/Test	1,857,248 590,612	11,738,489 4,797,378
Total	2,447,860	16,535,867

To generate the dense features, the authors pretrain a BERT-style model, referred to as MonoBERT, from scratch using masked language modeling (MLM) on the full Baidu corpus. This model is trained in a mono-encoder configuration and outputs a [CLS] token embedding for each query-document pair. These CLS embeddings are included in the dataset and serve as fixed, high-quality dense features for downstream reranking. The pretrained MonoBERT model and inference code are publicly available at: https://github.com/philipphager/baidu-bert-model.

Amazon KDD-cup. The KDD-Cup dataset (Reddy et al., 2022) contains 130K queries and 2.6M annotated query-product pairs in English, Japanese, and Spanish. Each query is linked to up to 40 products with rich textual metadata (titles, descriptions, bullet points), making it well-suited for LLM-based evaluation, unlike Baidu-ULTR. Although the presentation order is not recorded, the dataset primarily consists of relevant query-product pairs that were shown to users. For training, validation, and testing, we sample five random permutations of length 8 per query, resulting in 400,000 training and 50,000 validation/test groups. We use the English subset of the product search dataset released as part of the KDD Cup 2022 challenge (Reddy et al., 2022), which contains real-world queries and associated candidate products from Amazon. Each query-product pair is annotated using the ESCI labeling scheme: Exact match, Substitute, Complement, or Irrelevant.

Each query group is identified by a unique query_id and paired with 10–40 product candidates. For each product, the dataset provides structured metadata including:

• product_title,

- product_brand,
- product_color
- product_description,
- product_bullet_point (optional fields)
- product_id,
- product_locale, and
- · ESCI relevance label

To construct our training and evaluation sets, we sample 5 random permutations of length 8 from each query group. Note that we do not use the human-annotated ESCI labels provided in the dataset. Instead, we leverage the LLM's capability for contextual understanding to generate relevance labels automatically. Ideally, the relevance judgments produced by the LLM should align closely with those of human annotators. This yields approximately 392K query groups for training and 20K for validation, and 20K for testing. For a given query group, we encode each query-item pair into sentence embeddings using the all-MiniLM-L6-v2: https://huggingface.co/sentence-transformers/all-MiniLM-L6-v2 model from Sentence Transformers. The input format for the sentence transformer is constructed as:

```
{query} [SEP] {product_title} Brand:{brand}
Color:{color}
```

Table 5: Statistics of the Amazon KDD Cup (ESCI) dataset (English subset).

Split	#Query groups	#Query-Product Pairs
Training	78,447	627,576
Validation	4,000	32,000
Test	4,000	32,000
Total	86,447	691,576
Total (including 5 random permutations)	392,235	3,137,880

B.2 IMPLEMENTATION DETAILS

We use a transformer architecture for both the reward model and the ranker across all methods to ensure a consistent architectural backbone. The model contains 12 transformer layers, 768 hidden

 dimensions, 12 attention heads, and approximately 110M parameters. All models are trained with a learning rate of 2×10^{-5} using the AdamW optimizer (Loshchilov & Hutter, 2017) with a weight decay of 10^{-2} . We use a batch size of 512 and train for 21 epochs, applying a learning rate decay at epoch 12 via a step-based learning rate scheduler. All experiments are conducted using 2 NVIDIA A100 GPUs (40GB each).

For our method, REWARDRANK, we use a soft permutation temperature $\tau=0.5$ and reward correction term $\lambda=0.7$. In the PG-rank* baseline, which replaces the handcrafted NDCG utility with our learned reward model, we apply Gumbel-Softmax sampling with temperature 0.1 to approximate permutation sampling from the Plackett–Luce distribution. We report PG-rank* results for different Monte Carlo samples (MC =1,5,10) to evaluate variance in reward estimation.

In our URCC* implementation, we follow the original two-stage design: a set-aware utility model and a pairwise ranker. The utility model is trained with a binary cross-entropy loss computed over per-item logits derived from the transformer encoder outputs. Specifically, for each item in the permutation, we pool its embedding from the encoder, apply dropout, and project it through a shared per-item classifier. The per-item predictions are matched to click labels, and their aggregated loss forms the utility supervision.

As an additional baseline, we include a Naive-ranker trained with a relaxed NDCG objective following the PiRank formulation (Swezey et al., 2021), allowing listwise supervision using soft permutation matrices. All baselines are trained using the same reward data and input embeddings to isolate the impact of the learning objective.

Representative code for our implementations of REWARDRANK, PG-rank*, URCC* baselines, and evaluation procedures is included in the supplementary material.

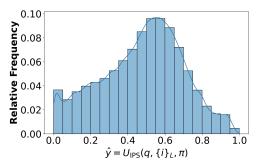
C COUNTERFACTUAL EVALUATION PROTOCOLS

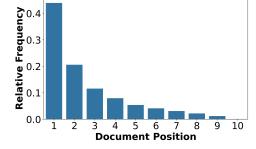
C.1 PO-EVAL DETAILS

PO-Eval provides a click-based framework for counterfactual evaluation of ranking models. Using the pre-trained Inverse Propensity Scoring model (IPS-Oracle) (Hager et al., 2024)² on the Baidu-ULTR dataset, it generates soft click probabilities for items in a ranked query group. These probabilities serve as counterfactual labels, enabling the evaluation of how effectively a ranker can model user engagement patterns reflected in clicks.

As the Baidu-ULTR dataset is derived from user interaction logs, click activity is heavily concentrated in the top-ranked positions, reflecting strong position bias (see Figure 3b). In contrast, the distilled soft utility (\hat{y}) generated by the IPS-Oracle exhibits a more uniform distribution across positions (Figure 3a), indicating that the oracle has successfully learned to correct for position bias. Under the

²https://github.com/philipphager/baidu-bert-model





(a) Soft-utility distribution on Baidu-ULTR generated by IPS-Oracle computed using Eqn 15.

(b) Position distribution of clicks in Baidu-ULTR.

Figure 3: Distributions extracted from IPS-Oracle analysis on Baidu-ULTR.

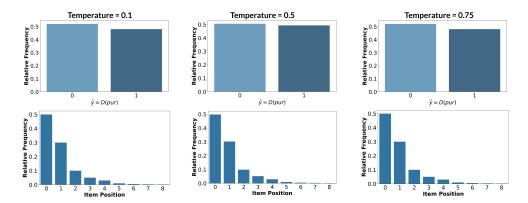


Figure 4: Effect of Sampling Temperature on LLM-Simulated Behavior in LAU-Eval. We visualize the distribution of binary purchase decisions (top) and item positions (bottom) generated by Claude Sonnet 3.5 v2 under three sampling temperatures: 0.1, 0.5, and 0.75. Each sample corresponds to a ranked list generated during LAU-Eval. As temperature increases, the purchase signal slightly diversifies, while positional biases remain consistent across settings. These results suggest that LAU-Eval is robust to moderate sampling variability, with LLMs producing stable user-like behavior under soft prompting.

PO-Eval protocol, ranking methods aim to implicitly learn position debiasing from the IPS-Oracle's soft utility, as indicated by high $U_{\text{IPS-O}}(q, \{i\}_L, \hat{\pi})$.

Training and evaluating ranking schemes. fix ndcg Using the learned reward model, any ranker f can be optimized via the reward maximization objective defined in Eqn 12. To evaluate its performance under the IPS-Oracle, we define the following metric: Given a query group $(q,\{i\}_L)$ and predicted relevance scores $s=[s_1,\ldots,s_L]$, the induced permutation is $\hat{\pi}=\operatorname{argsort}(s)$. For each position $\hat{\pi}_\ell$, the examination probability is $P(E_{\hat{\pi}_\ell})$, and the associated relevance score R_{q,i_ℓ} is provided by the IPS-Oracle. The overall utility is computed as the probability of at least one click: $U_{\text{IPS}}(q,\{i\}_L,\hat{\pi})$, which serves as the primary evaluation metric. It reflects how well f aligns with the user behavior modeled by the IPS-Oracle; higher values indicating better alignment. Additionally, we report NDCG_{rel}@10, which measures how much the predicted ranking respects the relevance scores $R_{q,i}$.

We incorporate the examination probabilities from (Hager et al., 2024), which are defined as:

$$P(E) = \{1: 1.0000, 2: 0.6738, 3: 0.4145, 4: 0.2932, 5: 0.2079, 6: 0.1714, 7: 0.1363, 8: 0.1166\}$$

C.2 LAU-EVAL DETAILS

We use Claude 3.5 Sonnet v2 with a temperature of 0.5 and a context window of 5,000 tokens. The LLM is prompted using a consistent instruction template, as illustrated in Figure 5. To evaluate a ranker with LAU-Eval, its predicted scores are converted into item positions, which are then used to reorder the input list. This reordered list is passed to the LLM alongside the original query, and the LLM outputs a binary decision regarding purchase. We include representative query groups and the corresponding LLM responses to demonstrate this pipeline.

To assess the robustness of LAU-Eval under different sampling conditions, we examine how varying the temperature of the LLM decoding process affects its outputs. Figure 4 shows the distributions of LLM-simulated purchase decisions and selected item positions at temperatures 0.1, 0.5, and 0.75. While purchase rates exhibit slight variation, the LLM consistently favors top-ranked items—reflecting realistic user behavior in shopping scenarios.

Instruction prompt for LLM. We design the LLM-Eval instruction to incorporate behavioral biases such as position bias, brand preference, irrelevance filtering, similarity aversion, and color bias, guiding the LLM to consider both relevance and context-dependent preferences. Given a query and

 an ordered product list, the LLM estimates (i) the probability of purchasing at least one item and (ii) the selected item, without explicit relevance constraints. We illustrate the instruction prompt using an example from the Amazon KDD-Cup dataset (Reddy et al., 2022), as shown in Figure 5.

Ranking Evaluations. We present the LLM's response to the initial list in Figure 6, including the full reasoning behind the response. It is noteworthy how the LLM is able to reason about the biases present in the query groups effectively. For each initial list, we also show the LLM's response to the rearranged list generated by Claude, depicted in Figure 7. As seen, the initial arrangements in Figure 6 lead to a no purchase decision, whereas REWARDRANK generates arrangements that increase the likelihood of purchase according to the LLM. Furthermore, the LLM's response enhances the interpretability of LLM-Eval, demonstrating how REWARDRANK's ranking capabilities align with the LLM's reasoning process.

Initially, we experimented with smaller language models such as Llama-3.1-8B: meta-llama/Llama-3.1-8B-Instruct and DeepSeek-R1-Distill: deepseek-ai/DeepSeek-R1-Distill-Llama-8B. However, these models were unable to generate appropriate responses to the instructions. Our experiments revealed that larger models were better at understanding the context.

It is important to note that LAU-Eval is used to simulate user behavior dynamics that may influence user decisions. Our selection of biases and instruction prompt serves as a proof-of-concept demonstrating that an LLM can be used as a proxy user to study counterfactual ranking strategies. We acknowledge that there are likely many variants of instruction prompts that could be designed to simulate user behavior. This area of exploration could be a direction for future work.

LLM Prompt (Probability Estimation Task)

You are shopping for a given query. Your task is to estimate the likelihood of purchasing any item in a provided list. Please answer yes or no, indicating whether you wish to purchase any item from the given list. Consider the relative relevance of items in the list when making your decisions. Be frugal, as a typical human user would be—most users buy when the list is highly relevant, and often make no purchase when following behavioral criteria are not met. You enter a 'query' into the shopping system, and it returns some items mentioned in the 'products'. The items are presented in the given order, with 1st item shown at the top of the list and the last item shown at the bottom.

Your query-products shopping list:

Query: 11 iphone pro screen protector Products:

```
{"B09CGJ8RW1":"title":"JETech Screen Protector and Camera Lens",
             "brand": "JETech", "color": "Transparent",
"B07515P7PT": "title": "JETech Screen Protector for iPhone 11 Pro",
             "brand": "JETech", "color": "Clear",
"B075S8V728":"title":"Ailun for Apple iPhone 11 Pro/iPhone",
             "brand": "Ailun", "color": NA,
"B07STC633H": "title": "UNBREAKcable Screen Protector for iPhone 11",
             "brand": "UNBREAKcable", "color": NA,
"B07D6XR7FM": "title": "TETHYS Glass Screen Protector for iPhone 11",
             "brand": "TETHYS", "color": "Transparent",
"B073DLZWX7": "title": "Maxboost Screen Protector for Apple iPhone",
             "brand": "Maxboost", "color": "Clear",
"B07FP41MC5": "title": "Trianium (3 Packs) Screen Protector",
             brand: "Trianium", color: "Clear",
"B09BQRWG15": "title": "YRMJK Screen Protector Compatible iPhone",
             brand: "YRMJK", color: NA}
```

Relevance Score: The relevance score shows how relevant the item is given the query. For every query-item pair, it is a numerical value between 0 and 1. You should consider the following criteria:

- 1. **Position bias**: where the items appearing near the top are more likely to be clicked. The position score decreases based on the following examination probabilities: position_scores = { 1: 1.0000, 2: 0.6738, 3: 0.4145, 4: 0.2932, 5: 0.2079, 6: 0.1714, 7: 0.1363, 8: 0.1166 } If the relevant item is not near the top, it will reduce the probability of purchase irrespective of its relevance.
- 2. **Brand bias**: If items from the same brand are placed adjacent to each other, the user is less likely to make a purchase. High brand bias means adjacent items are from the same brand.
- 3. **Irrelevance bias**: Multiple irrelevant items near the top reduce the chance of purchasing any item. This measures contextual dissimilarity at top positions.
- 4. **Color bias**: A lack of visual diversity (e.g., similar-colored items grouped together) discourages purchases.

Note that high brand-bias, irrelevance-bias, or color-bias harms the user experience and should lower the chance of purchasing. Act like a frugal user: only purchase if all criteria are satisfied; otherwise, do not purchase.

Task: Given a query and its candidate product list, decide whether any item should be purchased and report the decision (no exact totals needed). Use a strict policy: purchase only if all criteria are satisfied; otherwise, output no purchase. **Decision rule**: Determine the Final decision (yes/no) by jointly considering relevance score, position bias, brand bias, irrelevance bias, and color bias. Proceed to purchase only if all criteria are satisfactorily met. If not, do not purchase. Be especially cautious about brand and color biases—purchase only when the biases are very low. If the decision is yes, select exactly one item to purchase.

Output: D(purchase) = <answer>
If decision = yes, then Item to be purchased = <answer>

Figure 5: LLM instructions for a sample taken from the Amazon KDD-Cup dataset (Reddy et al., 2022). A total of 8 products are displayed for the query "11 iphone pro screen protector". For visual clarity, we have shortened the titles and product descriptions.

LLM's response to the initial item list

Let me analyze this shopping scenario step by step:

1. **Relevance Analysis:** Most products are screen protectors for iPhone 11 Pro/X/XS (5.8-inch), which is highly relevant to the query - Top items (JETech, Ailun) are particularly relevant as they specifically mention iPhone 11 Pro.

Position Bias:- First few positions have higher probability of being selected- Top 3 positions (JETech, JETech, Ailun) have the highest position scores.

Brand Bias:- First two positions are JETech brand (high brand bias at top)- This creates a negative shopping experience- Having same brand adjacent to each other reduces purchase likelihood.

Irrelevance Bias:- Most items are relevant to the query- However, some products like Trianium (4th position) is for iPhone 11 Pro Max (6.5"), which is incorrect size- Overall irrelevance bias is low as most products match the query.

Color Bias:- Many items are listed as "Clear" or "Transparent"- High color similarity among products- This creates a moderate color bias.

Final Analysis: Despite high relevance of products, there are several concerns: 1. High brand bias at top positions (JETech repeated)2. High color bias (many transparent/clear options)3. The position bias favors top items, but those have brand repetition 4. Some size mismatches in the list. Given these factors, especially the high brand bias at top positions and color similarity, a frugal user would likely be hesitant to make a purchase from this list.

D(**purchase**) = no; The brand repetition at top positions and lack of diversity in colors significantly harm the shopping experience, making it less likely for a careful shopper to make a purchase, even though the products themselves are relevant to the query.

Figure 6: LLM response to the initial list of items shown in Figure 5.

LLM's response to the ranked list provided by REWARDRANK

Let me analyze this shopping scenario carefully:

- 1. **Relevance Analysis:** All products are screen protectors for iPhone 11 Pro/X/XS, which matches the query perfectly- Most products offer multiple packs (2-3 pieces)- All are tempered glass protectors- High relevance across items.
- 2. **Position Bias**:- Top positions have higher visibility- First three items (UNBREAKcable, TETHYS, Ailun) have strongest position advantage.
- 3. **Brand Bias**:- JETech appears twice (positions 5 and 6) creates brand bias- Other brands are well distributed- Moderate brand bias due to JETech repetition.
- 4. **Irrelevance Bias**:- Low irrelevance bias as all products are highly relevant to query- All products are specifically designed for the requested phone model- Core features (tempered glass, screen protection) are consistent.
- 5. **Color Bias:** Most items are clear/transparent- Very similar color schemes across products-High color bias due to lack of variety.

Final Decision Analysis:- While relevance is high and position bias favors top items- The high color bias (all transparent) and moderate brand bias (JETech repetition) are concerning- However, given this is a screen protector where transparency is expected/desired, color bias should be weighted less heavily- The first three positions offer different brands and good specifications-TETHYS (2nd position) offers 3-pack with edge-to-edge coverage and installation frame.

D(purchase) = yes; Item to be purchased = TETHYS Glass Screen Protector (B07D6XR7FM)Reasoning: Selected TETHYS due to good position (2nd), differentiated brand, comprehensive features, and 3-pack value proposition. The high relevance and acceptable biases make this a reasonable purchase despite some concerns.

Figure 7: LLM response to the ranked list generated by REWARDRANK.

Table 6: Ablation studies on Baidu-ULTR dataset under PO-Eval. This table presents ablation results for RewardRank under different configurations, including variations in SoftSort temperature τ , misspecification correction regularization λ , and the addition of the auxiliary reward loss term from Eqn 6. We also report results for PG-rank* using different numbers of Monte Carlo samples. The counterfactual evaluation metric is $\mathbb{E}[U_{\text{IPS-O}}(q, \{i\}_L)]$.

Method	$\mathbb{E}[U_{\text{IPS-O}}(q,\{i\}_L)]$	
Upper-Bound	0.553	
ListNet Cao et al. (2007)	0.523 ± 0.0007	
ListMLE Xia et al. (2008)	0.522 ± 0.0007	
LambdaRank Wang et al. (2018)	0.524 ± 0.0007	
PiRank Swezey et al. (2021)	0.525 ± 0.0007	
URCC*	0.462 ± 0.0005	
PG-rank* (mc=1)	0.481 ± 0.0006	
PG-rank* (mc=5)	0.495 ± 0.0005	
PG-rank* (mc=10)	0.501 ± 0.0005	
SoftSort Temperature	τ	
RewardRank ($\tau = 0.1, \lambda = 0.0$)	0.532 ± 0.0005	
RewardRank ($\tau = 0.2, \lambda = 0.0$)	0.533 ± 0.0005	
RewardRank ($\tau = 0.7, \lambda = 0.0$)	0.531 ± 0.0005	
RewardRank ($\tau = 1.0, \lambda = 0.0$)	0.530 ± 0.0005	
Misspecification Correcti	ion λ	
RewardRank ($\tau = 0.5, \lambda = 0.1$)	0.532 ± 0.0005	
RewardRank ($\tau = 0.5, \lambda = 0.3$)	0.534 ± 0.0007	
RewardRank ($\tau = 0.5, \lambda = 0.7$)	$\textbf{0.536} \pm \textbf{0.0007}$	
RewardRank ($\tau = 0.5, \lambda = 1.0$)	0.533 ± 0.0007	
Auxiliary Per-Item Regularizer Eqn 6		
RewardRank (reward loss = Eqn 3)	0.528 ± 0.0005	
RewardRank (reward loss = Eqn $3 + \text{Eqn } 6$)		

D FURTHER ABLATION STUDIES

We use the Baidu-ULTR dataset to study how the performance of REWARDRANK varies with two key hyperparameters: the temperature τ of the SoftSort operator, which controls permutation sharpness, and the regularization strength λ for reward misspecification correction introduced in Eqn 13. Varying $\tau \in \{0.1, 0.2, 0.7, 1.0\}$ shows that moderate temperature ($\tau = 0.2 - 0.5$) achieves the best utility and relevance alignment. Too low a temperature leads to unstable gradients due to near-hard permutations, while higher values oversmooth rankings, diluting learning signals. Fixing $\tau = 0.5$, we ablate the correction term with $\lambda \in \{0.0, 0.1, 0.3, 0.7, 1.0\}$. As shown in Table 6 and visualized in Figure 8, moderate correction ($\lambda = 0.5 - 0.7$) yields the best trade-off, by down-weighting unreliable samples without discarding informative ones. This results in higher IPS utility, confirming the benefit of explicitly mitigating reward misspecification.

We explore the impact of incorporating an auxiliary item-level reward loss (Eqn 6) into the training objective of the reward model. As shown in Table 6, adding this auxiliary loss to the list-level cross-entropy objective (Eqn 3) improves expected utility from 0.528 to 0.536. This indicates that learning to predict the per-item feedback as an auxiliary task enhances the reward model's generalization and improves the downstream utility-optimized ranking.

As shown in Figure 8, increasing λ progressively reduces the influence of unreliable reward estimates by lowering their instance weights, leading to more stable learning. For illustration, we display soft utility scores from Eqn 15; however, all experiments use binary utility signals as defined in Eqn 14.

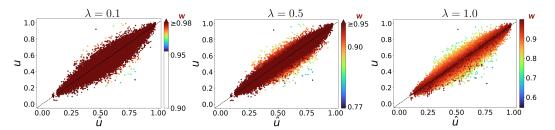


Figure 8: **Effect of reward misspecification correction on PO-Eval testbed.** Scatter plots show the relationship between predicted utility \hat{u} and true utility u for different values of the correction strength $\lambda \in \{0.1, 0.5, 1.0\}$. Each point corresponds to a ranked list, colored by its instance weight $w = 1 - \lambda |y - \hat{y}|$ based on reward model error. As λ increases, the influence of low-confidence predictions (lower w) diminishes, effectively down-weighting misspecified instances. This correction improves stability by emphasizing samples with well-aligned predicted rewards.

E INFERENCE COST AND LIMITATIONS

Inference Cost. The main inference cost in our work arises from using large language models (LLMs) for ranking and purchase probability estimation. These models require significant computational resources, especially for large datasets and permutations of items. Optimizations like batch processing and multi-GPU use help manage costs, but scalability remains a challenge. Caching frequently accessed queries can further reduce repeated computation costs.

Limitations. While both PO-Eval and LAU-Eval provide valuable insights into ranking quality and user preferences, there are inherent limitations in each approach. These limitations arise from their reliance on specific biases and the quality of input data, which may affect their performance in diverse real-world scenarios. Below, we outline the key limitations of each method:

- **PO-Eval Limitations:** While PO-Eval provides a robust baseline for position-debiasing, it is limited in behavioral scope. It primarily focuses on mitigating position bias without considering other nuanced user preferences, such as brand bias or contextual relevance, which can lead to suboptimal performance in more complex scenarios.
- LAU-Eval Limitations: LAU-Eval captures richer heuristics and offers more context-aware ranking, but it depends heavily on the quality and stability of the LLM outputs. Inconsistent or noisy outputs from the LLM can negatively affect the reliability of the evaluation, as the method assumes that the LLM accurately reflects user preferences in all scenarios.

These limitations highlight areas for future improvement, such as incorporating additional user behavior modeling and enhancing the robustness of the LLM outputs.