

# Make Large Language Model a Better Ranker

Anonymous ACL submission

## Abstract

The evolution of Large Language Models (LLMs) has significantly enhanced capabilities across various fields, leading to a paradigm shift in how Recommender Systems (RSs) are conceptualized and developed. However, existing research primarily focuses on point-wise and pair-wise recommendation paradigms. These approaches prove inefficient in LLM-based recommenders due to the high computational cost of utilizing Large Language Models. While some studies have delved into list-wise approaches, they fall short in ranking tasks. This shortfall is attributed to the misalignment between the objectives of ranking and language generation. To this end, this paper introduces the Language Model Framework with Aligned Listwise Ranking Objectives (ALRO). ALRO is designed to bridge the gap between the capabilities of LLMs and the nuanced requirements of ranking tasks within recommender systems. A key feature of ALRO is the introduction of soft lambda loss, an adaptation of lambda loss tailored to suit language generation tasks. Additionally, ALRO incorporates a permutation-sensitive learning mechanism that addresses position bias, a prevalent issue in generative models, without imposing additional computational burdens during inference. Our evaluative studies reveal that ALRO outperforms existing embedding-based recommendation methods and the existing LLM-based recommendation baselines, highlighting its efficacy.

## 1 Introduction

The rapid advancement in Large Language Models (LLMs), known by ChatGPT<sup>1</sup>, has marked a significant milestone in demonstrating their versatility in zero-shot and few-shot learning across various domains. These models, effectively employed in sectors like Question Answering and Information Retrieval, have shown remarkable adaptability and

<sup>1</sup><https://chat.openai.com>

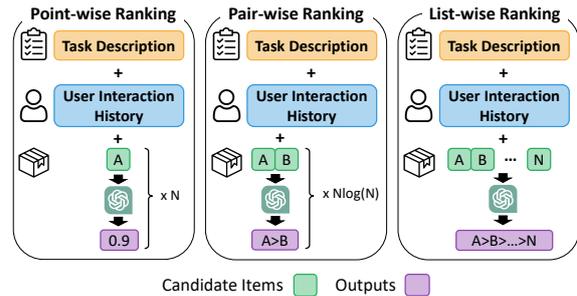


Figure 1: The comparison of point-wise, pair-wise, and list-wise ranking in LLM-based recommendation.

reliability. Their ability to efficiently handle tasks usually requiring extensive domain-specific training has sparked a surge in research aimed at exploring their potential across diverse applications, e.g. Recommender System.

In the context of recommender systems, the application of LLMs has attracted considerable attention. Wu et al. (2023) demonstrates a novel paradigm in using Large Language Models as agents in recommendation systems. This approach leverages their natural language processing strengths for more context-sensitive recommendations. Concurrently, investigations conducted in Bao et al. (2023) and Li et al. (2023) explore the capability of LLM in point-wise recommendation, revealing how language models can be adapted for suggesting products. Qin et al. (2023) investigate the use of pairwise ranking prompts to enhance recommendation systems. Despite these advancements, as depicted in Figure 1, a significant limitation of these methods is their high computational cost, stemming from the iterative call of LLMs to evaluate each candidate item.

In leveraging Large Language Models for recommendation systems, the list-wise ranking method stands out for its computational efficiency. Yet, as Dai et al. (2023) illustrates, executing list-wise ranking effectively is fraught with challenges. The



underestimate the inherent complexities and challenges associated with implementing an efficient list-wise LLM-based recommendation system.

## 2.2 Learning to Rank

Learning to Rank (LTR) constitutes a fundamental component in information retrieval systems, aimed at ordering entities by their relevance. This domain is categorized into three main methodologies according to the design of the loss function: pointwise, pairwise, and listwise approaches. Pointwise methods focus on predicting the absolute relevance of individual items, typically framed as classification or regression tasks (Li et al., 2007; Crammer and Singer, 2001). Pairwise strategies, in contrast, emphasize the relative importance between item pairs, with the goal of accurately determining the more relevant item in a pair (Freund et al., 2003; Burges et al., 2005; Chapelle and Keerthi, 2010). The listwise approaches extend this concept by considering the entire item list as the training unit, aiming to directly optimize the overall item ordering to align with ranking objectives (Xu and Li, 2007; Cao et al., 2007; Taylor et al., 2008; Xia et al., 2008; Burges, 2010). In this paper, we present an innovative adaptation of the lambda loss function (Wang et al., 2018) tailored for natural language generation, leveraging the pairwise approach to enhance the coherence of generated texts. This adaptation underscores the potential of LTR methodologies to extend beyond traditional retrieval tasks.

## 3 Problem Statement

We define the sequential recommendation ranking problem as follows. Let  $\mathcal{U}$  represent the set of users and  $\mathcal{I}$  denote the set of items. For any given user  $u \in \mathcal{U}$ , their historical interactions with items are represented by  $\mathcal{H}_u = \{h_1, h_2, \dots, h_k\}$ , where each  $h_i \in \mathcal{I}$  signifies an item that user  $u$  has previously interacted with. With this notation in place, the ranking problem is formalized as follows:

**Definition 1** For a user  $u$ , consider  $\mathcal{C}_u = \{c_1, c_2, \dots, c_m\}$  as the set of candidate items for recommendation, where each  $c_i \in \mathcal{I}$  and  $m \leq |\mathcal{I}|$ . The goal is to devise a ranking function  $F : \mathcal{H}_u \times \mathcal{C}_u \rightarrow S_m$  that accurately predicts the permutation  $\tau \in S_m$  that best orders the items in  $\mathcal{C}_u$ . The set  $S_m$  is the symmetric group of all  $m$ -element permutations, encapsulating every possible arrangement of the candidate items.

```

### Instruction:
Given the user's interaction history, which reveals their items
preferences, generate a preference-based ranking of the
provided candidate items, your task is to rank a list of new
candidate movies. Your ranking should include all the candidate
movies provided, and it should be based solely on the user's
preferences, without regard to the initial order of the candidates.

### Input:
[User Interaction History]:
title: Independence Day genres: Action|SciFi|War rating: 3
title: Close Encounters of the Third Kind (1977) genres:
Drama|Sci-Fi rating: 4
...
[Candidate Items]:
(A) title: Starman genres: Adventure|Drama|Romance
(B) title: Independence Day genres: Action|SciFi|War
...
### Response:
Given the historical interaction, the ranking result is: B A C ...

```

Figure 2: A template sample for ranking in LLM-based recommendation.

## 4 Methodology

In this section, we elucidate the constraints inherent in prevailing prompting paradigms when addressing list-wise recommendation tasks. Our learning framework is developed with four distinct components: Template Design, Supervised Fine-Tuning, Soft Lambda Loss, and Permutation-Sensitive Learning.

### 4.1 Template Design

Before delving into the specifics of our learning module, we delineate the process of transforming the ranking task into a language generation problem. Drawing inspiration from Alpaca Tuning (Taori et al., 2023), we employ a natural language prompt template, denoted as  $T_{\text{src}}(\mathcal{H}_u, \mathcal{C}_u)$ , which transmutes the input user history  $\mathcal{H}_u$  and context  $\mathcal{C}_u$ , inclusive of item attributes such as names, categories, and descriptions, into a structured format. This transformation additionally aids in creating target text templates  $T_{\text{tgt}}(\tau)$ , representing the permutation that arranges candidate items according to user preferences. An example of the template is illustrated in Figure 2. Leveraging the structured template, we reframe the ranking problem as a language generation problem.

### 4.2 Supervised Fine-Tuning

With the language generation problem that given  $T_{\text{src}}(\mathcal{H}_u, \mathcal{C}_u)$  that aims to predict  $T_{\text{tgt}}(\tau)$ , we implement a supervised fine-tuning paradigm that leverages the Low-Rank Adaptation (LoRA) approach, as introduced by Hu et al. (2022). The core idea

250 behind LoRA is to adapt pre-trained models in a  
 251 parameter-efficient manner, enabling effective fine-  
 252 tuning on downstream tasks with minimal modi-  
 253 fications to the original model parameters. The  
 254 fine-tuning process is formulated by the following  
 255 loss function:

$$\mathcal{L}_{\text{sft}} = - \sum_{t=1}^{|y|} \log(P_{\theta}(y_t|x, y_{<t})), \quad (1)$$

257 where  $\mathcal{L}_{\text{sft}}$  denotes the supervised fine-tuning loss,  
 258 and  $P_{\theta}(y_t|x, y_{<t})$  represents the conditional prob-  
 259 ability of predicting the token  $y_t$  given the input  
 260 tokens  $x$  and the preceding tokens  $y_{<t}$ . In this con-  
 261 text,  $x$  and  $y$  correspond to the tokenized represen-  
 262 tations of  $T_{\text{src}}(\mathcal{H}_u, \mathcal{C}_u)$  and  $T_{\text{tgt}}(\tau)$ , respectively.

263 This supervised fine-tuning process utilizes tar-  
 264 get tokens that correspond to the correctly ranked  
 265 list of candidate answers, which are subsequently  
 266 adjusted to reflect user preferences. The objective  
 267 is for the model to accurately rank a list of items,  
 268 ensuring that the generated responses are not only  
 269 relevant but also personalized to the user’s interests.

### 270 4.3 Soft Lambda Loss (SLL)

271 The widely adopted cross-entropy loss in language  
 272 generation, derived from next-token prediction dur-  
 273 ing supervised fine-tuning, faces a fundamental  
 274 misalignment with the objectives of ranking. Such  
 275 a discrepancy undermines the efficacy of cross-  
 276 entropy loss when applied to the specific demands  
 277 of ranking, leading to suboptimal performance in  
 278 these contexts. Unlike the existing LTR framework  
 279 (Wang et al., 2018), this is not straightforward to  
 280 directly optimize on Normalized Discounted Cumu-  
 281 lative Gain (NDCG) when dealing with language  
 282 models that generate ranked token probabilities in-  
 283 crementally. Traditional ranking losses, such as  
 284 Lambda loss (Wang et al., 2018) or SoftRank (Tay-  
 285 lor et al., 2008), are not directly applicable. The  
 286 Lambda loss, for example, is defined as:

$$\mathcal{L}_{\text{rank}} = \sum_{i=1}^{|\tau|} \sum_{j:\tau_j < \tau_i} \delta_{i,j} |G_i - G_j| \cdot \log_2 \left( 1 + e^{-\sigma(s_i - s_j)} \right), \quad (2)$$

288 where

$$\delta_{ij} = \left| \frac{1}{D_{|i-j|}} - \frac{1}{D_{|i-j|+1}} \right|, \quad (3)$$

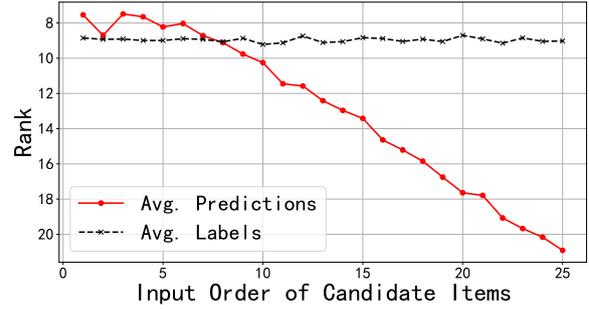


Figure 3: Demonstration of position bias. The figure shows how the placement of candidate items in the input sequence can significantly alter the ranking results produced by a Language Model.

290 with  $G_i$  and  $D_i$  following the definitions from  
 291 NDCG, and  $s_i$  representing the model-derived pre-  
 292 diction score. In large language models, the rank-  
 293 ing order is typically determined by using the  
 294  $argmax$  function on the output probabilities of  
 295 tokens, which is non-differentiable and thus unsuit-  
 296 able for the training process.

297 To overcome this, we propose a method that com-  
 298 bines the soft-argmax function with Lambda loss  
 299 to calculate the deviation of predicted probabilities  
 300 from the ideal ranking order. We define a differ-  
 301 entiable ranking score for the generative model by  
 302 substituting the traditional  $argmax$  function in  $s_i$   
 303 with the soft-argmax, expressed as:

$$s_i = \max_j \frac{e^{\gamma y_{j,i}}}{\sum_k e^{\gamma y_{j,k}}} \cdot j, \quad (4)$$

305 where  $y_{i,j}$  denotes the output probability of the lan-  
 306 guage model for the  $j$ th position and token  $i$ , and  
 307  $\gamma$  represents the scaled value that adjusts the dis-  
 308 tribution of softmax. By making the computation  
 309 of  $s_i$  differentiable with the soft-argmax method,  
 310 we align the objectives of language generation with  
 311 those of the ranking task. Overall, Soft Lambda  
 312 Loss follows the Equation 2, which is derived from  
 313 Wang et al. (2018), by replacing  $s_i$  with Equation 4  
 314 to get a differentiable objective.

### 315 4.4 Permutation-Sensitive Loss (PSL)

316 In list-wise recommendation tasks with Large Lan-  
 317 guage Models (LLMs), position bias emerges as  
 318 a formidable challenge, with the order of the can-  
 319 didate input sequence notably swaying the rank-  
 320 ing outcomes. As depicted in Figure 3, language  
 321 models exhibit a propensity to assign higher rank-  
 322 ings to candidates positioned at the beginning of  
 323 the list. This tendency highlights the significant

influence of candidate positioning on model evaluations, underscoring the imperative for developing methodologies to counteract these biases.

It is worth noted that the observed phenomenon depend exclusively on natural language generation tasks with the sequence of input candidates. This contrasts with embedding-based recommendation systems, where the order of inputs does not influence outcomes by calculating the score of user and item pair separately. The effect of permutation on the output is described by the inequation:

$$F(T(\mathcal{H}_u, \mathcal{C}_u)) \neq F(T(\mathcal{H}_u, \mathcal{C}'_u)), \quad (5)$$

where  $F(\cdot)$  denotes the language model, and  $\mathcal{C}'_u = \{c_{\pi(0)}, c_{\pi(1)}, \dots, c_{\pi(m)}\}$  represents a permuted candidate list from the original candidate list  $\mathcal{C}_u$ , with  $\pi(\cdot)$  as the permutation function that rearranges the candidates. This equation highlights the dependency of the model on the sequence in which inputs are provided, distinguishing it from conventional recommendation approaches where input order is often inconsequential.

Although Hou et al. (2023) proposed the bootstrapping method, which shuffles the candidate items multiple times and takes average scores as the final ranking result, it is inefficient as it repetitively calls language models in the inference stage to get average ranking. To alleviate this issue without burdening the inference in the recommendation, we propose a permutation loss that aims to minimize the output distribution distance between the original candidate list  $\mathcal{C}_u$  and the permuted candidate list  $\mathcal{C}'_u$  within the fine-tuning stage. By adopting cross-entropy loss that measures the distance between two distributions, we empower the model with permutation invariant capability. The loss function could be formulated as:

$$\mathcal{L}_{\text{cont}} = - \sum_{t=1}^{|y|} P_{\theta}(y_t|x, y_{<t}) \log P_{\theta}(y'_t|x', y'_{<t}), \quad (6)$$

where  $x$  and  $x'$  are the prompt derived from  $T(\mathcal{H}_u, \mathcal{C}_u)$  and  $T(\mathcal{H}_u, \mathcal{C}'_u)$  respectively, and  $y$  and  $y'$  are the labels for the given prompts.

#### 4.5 Training Objective

Overall, we provide the soft lambda loss  $\mathcal{L}_{\text{rank}}$  with permutation-sensitive framework  $\mathcal{L}_{\text{cont}}$  to address the issues mentioned above, which goes beyond the naive supervised fine-tuning. The objective function is reformulated as:

$$\mathcal{L} = \mathcal{L}_{\text{sft}} + \alpha \mathcal{L}_{\text{rank}} + \beta \mathcal{L}_{\text{cont}}, \quad (7)$$

where  $\alpha, \beta$  are hyperparameters that adjust the importance of each loss.

## 5 Experiment

In our study, we conducted a comprehensive evaluation of our model across two real-world datasets. This was complemented by an ablation study, robustness tests, and efficiency evaluations. Our experiment was directed by the following pivotal research questions:

- **(RQ1)** Does the proposed framework surpass existing baselines in both embedding-based and LLM-based recommendation models?
- **(RQ2)** What extent does supervised fine-tuning on recommendation-specific corpus enhance Large Language Model performance?
- **(RQ3)** How crucial is the involvement of our proposed module for metrics improvement?
- **(RQ4)** How does permutation-sensitive learning compare to bootstrapping methods in terms of performance and efficiency?
- **(RQ5)** How does the ALRO framework improve performance across different parameter sizes of the backbone language model compared to traditional supervised fine-tuning?

Through these explorations, we aim to elucidate the contributions of domain-specific fine-tuning with our novel modules to the advancements in LLM-based recommendation systems.

### 5.1 Dataset

We selected 3 widely adopted open-source datasets to evaluate the effectiveness of our framework. *Movie* represents the dataset from MovieLens-1M<sup>2</sup>. *Music* is the dataset from the "CDs & Vinyl" subsets of Amazon<sup>3</sup> product reviews dataset. The details of the datasets are shown in Table 1. With the dataset selected user history  $\mathcal{H}_u$  with 20 items and  $\mathcal{C}_u$  with 25 items. The output permutation  $\tau$  is derived from the future candidate rating from the user. The datasets are partitioned into training, validation, and testing subsets with a ratio of 8:1:1.

### 5.2 Baselines and Evaluation Metrics

To evaluate the effectiveness of our framework, we select several state-of-the-art baselines, which could be categorized into Non-Sequential Recommendation, Sequential Recommendation, and Large Language Model-based Recommendation. It

<sup>2</sup><https://grouplens.org/datasets/movielens/1m/>

<sup>3</sup><https://jmcauley.ucsd.edu/data/amazon/>

Table 1: Statistics of datasets.

| Datasets           | Movie     | Music   |
|--------------------|-----------|---------|
| # Users            | 6040      | 1040    |
| # Items            | 3952      | 33191   |
| # Actions          | 1,000,209 | 139,459 |
| Sparsity           | 95.80%    | 99.60%  |
| Avg. # Tokens/Item | 20.76     | 22.82   |

is worth noting that we introduce the BERT-based model as the backbone to extract the textual information of items in both Non-Sequential Recommendation and Sequential Recommendation.

- Non-Sequential Recommendation: **NCF** (He et al., 2017) adopts neural network with collaborative filtering for recommendation, **DIN** (Kang and McAuley, 2018) involves user interest modeling based on user behavior with attention mechanism.
- Sequential Recommendation: **GRU4Rec** (Hidasi et al., 2016) is a session-based recommendation with GRU-based recurrent network. **SASRec** (Hidasi et al., 2016) utilizes the self-attention network with positional embedding to capture the user’s sequential behavior information. **CORE** (Hou et al., 2022) uses the representation consistent framework to unify the session and items representation space. **NARM** (Li et al., 2017) decomposed the user behavior into global and local forms with attention networks for sequential recommendation.
- Large Language Model-based Recommendation: **Zero-shot LLM**, **Few-shot LLM**, **Tall-Rec** (Bao et al., 2023) fine-tunes LLM with instruction tuning for point-wise recommendation, **ES4Rec** (Li et al., 2023) introduces pre-trained item embedding as prompt with an adapter to fine-tune the LLM. We use Llama-7b as the base model for all the LLM-based baselines.

To assess the performance of various models in ranking tasks, we employ Normalized Discounted Cumulative Gain (NDCG) at different cutoff levels as our evaluation metric, specifically NDCG@k with  $k$  values of 3, 10, and 25.

### 5.3 Implementation Details

Our experiments were conducted on a cluster of 12 Linux servers, each equipped with 8 A800 GPUs. For the backbone model, we utilized the Llama-2-

7b<sup>4</sup> with BF16 precision, available on Huggingface. The supervised fine-tuning step was implemented using the PyTorch framework and peft library, applying the LoRA technique with a rank setting of 16. We used the AdamW (Loshchilov and Hutter, 2019) optimizer with a learning rate of 1e-4 and batch size as 1 for SFT, complemented by 32 gradient accumulation steps with a total of 10 training epochs. We also used DeepSpeed (Rasley et al., 2020) with ZeRO stage as 2 for distributed training. For our loss function, we fine-tuned the hyperparameters, setting  $\alpha$  equal to 0.1,  $\beta$  equal to 0.01, and  $\gamma$  equal to 2.

### 5.4 Overall Performance (RQ1)

To validate the performance of our proposed framework, ALRO, we executed comparative analyses against established baseline methods, with the results presented in Table 2. The following observations were made:

- ALRO consistently outperformed the baselines across various metrics and datasets, unequivocally demonstrating its superiority in ranking tasks within recommender systems.
- Large Language Models (LLMs) without fine-tuning fell short against traditional methods, highlighting the crucial role of supervised fine-tuning for LLMs in recommendation contexts.
- The point-wise methods, such as TALLRec and ES4Rec, underperformed in ranking tasks. This is likely due to their inability to account for the interrelationships among candidates, a critical aspect that list-wise approaches inherently address.

These insights confirm the significance of our ALRO framework in enhancing the efficacy of ranking in recommendation systems and underscore the necessity for appropriate fine-tuning of LLMs to fully leverage their potential recommendation.

### 5.5 Effect of Supervised Fine-Tuning (RQ2)

Prompting techniques have showcased the profound ability of language models to interpret and execute tasks with remarkable precision. (Liu et al., 2023) However, the efficacy of these techniques is challenged when applied to specialized domains such as recommendation systems, particularly due to the potential misalignment between the pre-training corpus and the intricate requirements of ranking tasks. As depicted in Table 3, this discrep-

<sup>4</sup><https://huggingface.co/meta-llama/Llama-2-7b-hf>

Table 2: Performance Comparison. Optimal outcomes across all models are emphasized in bold, while second-best performances are distinguished by underlining. Evaluation metrics include NDCG at ranks 3, 10, and 25.

| Dataset       |           | Movie         |               |               | Music         |               |               |
|---------------|-----------|---------------|---------------|---------------|---------------|---------------|---------------|
| NDCG          |           | @3            | @10           | @25           | @3            | @10           | @25           |
| Non-Seq Rec   | NCF       | 0.6235        | 0.6971        | 0.8469        | 0.6791        | 0.7089        | 0.8724        |
|               | DIN       | <u>0.6475</u> | <u>0.7028</u> | <u>0.8609</u> | 0.6249        | 0.6645        | 0.8521        |
| Seq Rec       | GRU4Rec   | 0.6326        | 0.6915        | 0.8551        | 0.6761        | 0.7050        | 0.8701        |
|               | SASRec    | 0.6059        | 0.6737        | 0.8462        | 0.6704        | 0.7029        | 0.8705        |
|               | COREave   | 0.5159        | 0.5984        | 0.8100        | 0.6249        | 0.6645        | 0.8521        |
|               | NARM      | 0.6078        | 0.6732        | 0.8462        | 0.6821        | 0.7043        | 0.8713        |
| LLM-based Rec | Zero-shot | 0.5149        | 0.5958        | 0.8095        | 0.6420        | 0.6636        | 0.8549        |
|               | Few-shot  | 0.5185        | 0.5968        | 0.8104        | 0.6417        | 0.6706        | 0.8572        |
|               | E4SRec    | 0.5393        | 0.6271        | 0.8206        | 0.6054        | 0.6516        | 0.8470        |
|               | TALLRec   | 0.5854        | 0.6511        | 0.8361        | <u>0.7066</u> | <u>0.7262</u> | <u>0.8826</u> |
| Ours          | ALRO      | <b>0.6551</b> | <b>0.7124</b> | <b>0.8835</b> | <b>0.7102</b> | <b>0.7428</b> | <b>0.8915</b> |

Table 3: Comparison of Zero-shot, Few-shot and Supervised Fine-Tuning with Llama-7b backbone.

| Dataset   | Movie         |               |               |
|-----------|---------------|---------------|---------------|
| NDCG      | @3            | @10           | @25           |
| Zero-shot | 0.5149        | 0.5958        | 0.8095        |
| Few-shot  | 0.5185        | 0.5968        | 0.8104        |
| SFT       | <b>0.5843</b> | <b>0.6609</b> | <b>0.8396</b> |

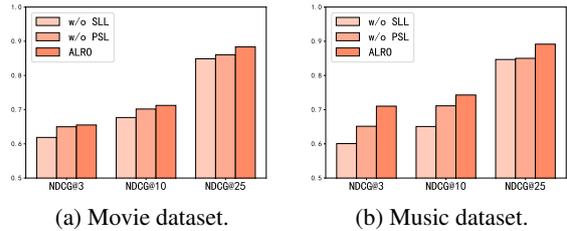


Figure 4: Ablation study on multiple datasets.

506 any is notably pronounced in medium-sized lan-  
 507 guage models like Llama-7b, where simple prompt-  
 508 ing may not suffice to activate the model’s ranking  
 509 capabilities effectively.

510 To address this gap, our study delves into the im-  
 511 pact of supervised fine-tuning on the performance  
 512 of language models in recommendation-related  
 513 tasks. Through a comparative analysis encompass-  
 514 ing zero-shot, few-shot, and supervised fine-tuning  
 515 approaches, we unveil a substantial improvement  
 516 in model performance by supervised fine-tuning,  
 517 with metrics enhancing by over 10%. This im-  
 518 provement is attributed to the fine-tuning process,  
 519 which effectively adjusts the model’s outputs to  
 520 better align with specific task requirements. This  
 521 approach overcomes the shortcomings of conven-  
 522 tional prompting techniques that often yield non-  
 523 parsable outputs, thereby enhancing the model’s  
 524 ability to rank information more accurately.

## 5.6 Ablation Study (RQ3)

525 In our research, we conducted an ablation study  
 526 to distinguish the contributions of distinct com-  
 527 ponents within our proposed framework, system-  
 528 atically omitting each module for comparative  
 529 analysis against the complete model. This in-  
 530 volved evaluating two key variants: Exclusion  
 531 of soft lambda loss (w/o SLL) and Exclusion of  
 532 permutation-sensitive learning (w/o PSL). The ex-  
 533 periment results, as depicted in Figure 4, revealed  
 534 that both components significantly improve the sys-  
 535 tem’s ability to rank candidates. Notably, the ob-  
 536 served reduction in NDCG can be ascribed to the  
 537 exclusion of the soft lambda loss, underscoring  
 538 the crucial role of objective alignment in improv-  
 539 ing the efficacy of language model-based recom-  
 540 mendation systems. Additionally, the performance  
 541 decrement observed upon removing Permutation-  
 542 Sensitive Learning further underscores the critical  
 543 influence of position bias on ranking performance.  
 544

Table 4: Comparative analysis of bootstrapping and permutation-sensitive learning. 'p@i' denotes the number of permutations applied in bootstrapping. The original permutation represented by p@1 is consistent with the prompt in ALRO. TPQ represents the average inference time per query measured in seconds.

| Dataset | Movie         |               |               |        |
|---------|---------------|---------------|---------------|--------|
| NDCG    | @3            | @10           | @25           | TPQ    |
| p@1     | 0.6188        | 0.6770        | 0.8487        | 8.019  |
| p@3     | 0.6486        | 0.7022        | 0.8607        | 25.430 |
| p@5     | <b>0.6566</b> | 0.7112        | 0.8646        | 46.243 |
| ALRO    | 0.6551        | <b>0.7124</b> | <b>0.8835</b> | 8.019  |

### 5.7 Comparison of Bootstrapping and Permutation-Sensitive Learning (RQ4)

Our research introduces a permutation-sensitive learning approach designed to address position bias, which affects the outcomes based on the order of candidate lists. While the bootstrapping method (Hou et al., 2023), offers a solution to this bias, it significantly increases inference time. We evaluated the effectiveness of permutation-sensitive learning compared to bootstrapping, aiming to reduce position bias without burdening the inference stage. Our comparisons included the original model without modifications, and bootstrapping with permutations executed 3 and 5 times. As demonstrated in Table 4, our method achieves comparable outcomes to bootstrapping while reducing inference times by approximately 5-fold. This indicates that our approach effectively mitigates the inference time issue through well-designed learning objectives.

### 5.8 Effect of Model parameter size (RQ5)

In this section of our research paper, we delve into the adaptability and efficacy of our learning framework across several LLM-based recommender systems, spanning various model sizes. Specifically, we selected four distinct models for our analysis: OPT-125M, Pythia 1.4B, Pythia-2.7B, and Llama-7B. By applying our framework to these models, we aim to showcase the consistent and significant performance enhancements it offers compared to traditional supervised fine-tuning approaches. As depicted in Figure 5, there is a clear correlation between model parameter size and performance, which serves to emphasize the capacity of our learning framework to augment the effectiveness of rec-

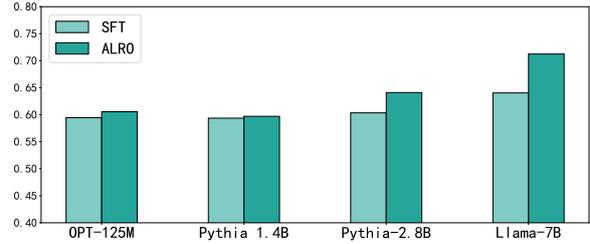


Figure 5: Enhancements achieved by ALRO across various model sizes on Movie dataset, measured using NDCG@10 metric.

ommender systems across a spectrum of language model sizes. Notably, the enhancements provided by our framework are more significant in larger models than in smaller ones, this may be attributed to the innate reasoning capability of language models. Overall, the experiment highlights the versatility and broad applicability of our framework in improving system performance.

## 6 Conclusion

In this research, we tackled the intricacies of employing large language models as ranking agents in recommender systems, with a focus on refining list-wise ranking methods for greater efficiency and accuracy. We proposed a cutting-edge framework that integrates soft lambda loss and permutation-sensitive learning. The integration of soft lambda loss is important as it bridges the objective between LLM's natural language generation and the specific demands of ranking tasks, enhancing the relevance and accuracy of the recommendations. Furthermore, permutation-sensitive learning approaches effectively address the issue of position bias, providing an improvement over traditional bootstrapping methods without imposing additional computational demands during inference. Our comprehensive evaluation across various datasets confirms the success of our method, marking a step forward in the application of LLMs as recommendation agents. This progress not only enhances accuracy but also maintains computational efficiency, striking a balance between the two.

## 7 Limitation

While our framework adeptly aligns the objectives of ranking and language generation, it falls short in fully harnessing the explainability potential inherent in language models. The supervised fine-tuning process, augmented by joint loss optimization, ef-

617 fectively enhances the model’s performance in list-  
618 wise ranking tasks, particularly in recommendation  
619 systems. However, this process inadvertently un-  
620 dermines the model’s proficiency in tasks beyond  
621 recommendation, limiting its versatility. Further-  
622 more, although our method demonstrates efficacy  
623 in ranking a set of 25 items, scalability becomes a  
624 concern as the number of candidates increases sig-  
625 nificantly. This limitation arises due to constraints  
626 such as context limits or the propensity for for-  
627 getting in Large Language Models, compromising  
628 the model’s ability to maintain performance consis-  
629 tency across varying candidate sizes.

## 630 References

631 Keqin Bao, Jizhi Zhang, Yang Zhang, Wenjie Wang,  
632 Fuli Feng, and Xiangnan He. 2023. [Tallrec: An effective and efficient tuning framework to align large language model with recommendation](#). In *Proceedings of the 17th ACM Conference on Recommender Systems, RecSys 2023, Singapore, Singapore, September 18-22, 2023*, pages 1007–1014. ACM.

633  
634  
635  
636  
637

638 Chris Burges, Tal Shaked, Erin Renshaw, Ari Lazier,  
639 Matt Deeds, Nicole Hamilton, and Greg Hullender.  
640 2005. Learning to rank using gradient descent. In  
641 *Proceedings of the 22nd international conference on*  
642 *Machine learning*, pages 89–96.

643  
644  
645

646 Christopher JC Burges. 2010. From ranknet to lamb-  
647 darank to lambdamart: An overview. *Learning*,  
648 11(23-581):81.

649  
650

646 Zhe Cao, Tao Qin, Tie-Yan Liu, Ming-Feng Tsai, and  
647 Hang Li. 2007. Learning to rank: from pairwise  
648 approach to listwise approach. In *Proceedings of the*  
649 *24th international conference on Machine learning*,  
650 pages 129–136.

651  
652  
653

654 Olivier Chapelle and S. Sathiya Keerthi. 2010. [Efficient algorithms for ranking with svms](#). *Inf. Retr.*,  
655 13(3):201–215.

656  
657  
658  
659  
660  
661  
662

663 Yoav Freund, Raj D. Iyer, Robert E. Schapire, and  
664 Yoram Singer. 2003. [An efficient boosting algorithm for combining preferences](#). *J. Mach. Learn. Res.*,  
665 4:933–969.

667 Xiangnan He, Lizi Liao, Hanwang Zhang, Liqiang Nie,  
668 Xia Hu, and Tat-Seng Chua. 2017. [Neural collaborative filtering](#). In *Proceedings of the 26th International Conference on World Wide Web, WWW 2017, Perth, Australia, April 3-7, 2017*, pages 173–182. ACM.

669  
670  
671  
672

673 Balázs Hidasi, Alexandros Karatzoglou, Linas Bal-  
674 trunas, and Domonkos Tikk. 2016. [Session-based recommendations with recurrent neural networks](#). In *4th International Conference on Learning Representations, ICLR 2016, San Juan, Puerto Rico, May 2-4, 2016, Conference Track Proceedings*. 675  
676  
677  
678

679 Yupeng Hou, Binbin Hu, Zhiqiang Zhang, and  
680 Wayne Xin Zhao. 2022. [CORE: simple and effective session-based recommendation within consistent representation space](#). In *SIGIR ’22: The 45th International ACM SIGIR Conference on Research and Development in Information Retrieval, Madrid, Spain, July 11 - 15, 2022*, pages 1796–1801. ACM.

681  
682  
683  
684  
685

686 Yupeng Hou, Junjie Zhang, Zihan Lin, Hongyu Lu,  
687 Ruobing Xie, Julian J. McAuley, and Wayne Xin  
688 Zhao. 2023. [Large language models are zero-shot rankers for recommender systems](#). *CoRR*,  
689 abs/2305.08845. 690

691  
692  
693  
694  
695  
696

697 Edward J. Hu, Yelong Shen, Phillip Wallis, Zeyuan  
698 Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, and  
699 Weizhu Chen. 2022. [Lora: Low-rank adaptation of large language models](#). In *The Tenth International Conference on Learning Representations, ICLR 2022, Virtual Event, April 25-29, 2022*. OpenReview.net.

700  
701

702 Wang-Cheng Kang and Julian J. McAuley. 2018. [Self-attentive sequential recommendation](#). In *IEEE International Conference on Data Mining, ICDM 2018, Singapore, November 17-20, 2018*, pages 197–206. IEEE Computer Society.

703  
704  
705  
706

707 Wang-Cheng Kang, Jianmo Ni, Nikhil Mehta, Mah-  
708 eswaran Sathiamoorthy, Lichan Hong, Ed H. Chi,  
709 and Derek Zhiyuan Cheng. 2023. [Do llms understand user preferences? evaluating llms on user rating prediction](#). *CoRR*, abs/2305.06474.

710  
711  
712

713 Jing Li, Pengjie Ren, Zhumin Chen, Zhaochun Ren, Tao  
714 Lian, and Jun Ma. 2017. [Neural attentive session-based recommendation](#). In *Proceedings of the 2017 ACM on Conference on Information and Knowledge Management, CIKM 2017, Singapore, November 06 - 10, 2017*, pages 1419–1428. ACM.

715  
716  
717  
718  
719  
720

721 Ping Li, Christopher J. C. Burges, and Qiang Wu. 2007. [Mcrank: Learning to rank using multiple classification and gradient boosting](#). In *Advances in Neural Information Processing Systems 20, Proceedings of the Twenty-First Annual Conference on Neural Information Processing Systems, Vancouver, British Columbia, Canada, December 3-6, 2007*, pages 897–904. Curran Associates, Inc.

722  
723  
724  
725  
726  
727  
728  
729  
730  
731  
732

733 Xinhang Li, Chong Chen, Xiangyu Zhao, Yong Zhang,  
734 and Chunxiao Xing. 2023. [E4src: An elegant](#) 721  
722

|     |   |   |     |
|-----|---|---|-----|
| 723 | effective efficient extensible solution of large language models for sequential recommendation. <i>CoRR</i> , abs/2312.02443.   | rank metrics. In <i>Proceedings of the 2008 International Conference on Web Search and Data Mining</i> , pages 77–86.   | 779 |
| 724 |   |   | 780 |
| 725 |   |   | 781 |
| 726 | Pengfei Liu, Weizhe Yuan, Jinlan Fu, Zhengbao Jiang, Hiroaki Hayashi, and Graham Neubig. 2023. Pre-train, prompt, and predict: A systematic survey of prompting methods in natural language processing. <i>ACM Computing Surveys</i> , 55(9):1–35.  | Xuanhui Wang, Cheng Li, Nadav Golbandi, Michael Bendersky, and Marc Najork. 2018. The lambdaloss framework for ranking metric optimization. In <i>Proceedings of the 27th ACM international conference on information and knowledge management</i> , pages 1313–1322.                                       | 782 |
| 727 |   |   | 783 |
| 728 |   |   | 784 |
| 729 |   |   | 785 |
| 730 |   |   | 786 |
| 731 | Ilya Loshchilov and Frank Hutter. 2019. <b>Decoupled weight decay regularization</b> . In <i>7th International Conference on Learning Representations, ICLR 2019, New Orleans, LA, USA, May 6-9, 2019</i> . OpenReview.net.   | Likang Wu, Zhi Zheng, Zhaopeng Qiu, Hao Wang, Hongchao Gu, Tingjia Shen, Chuan Qin, Chen Zhu, Hengshu Zhu, Qi Liu, et al. 2023. A survey on large language models for recommendation. <i>arXiv preprint arXiv:2305.19860</i> .  | 787 |
| 732 |   |   | 788 |
| 733 |   |   | 789 |
| 734 |   |   | 790 |
| 735 |   |   | 791 |
| 736 | Zhen Qin, Rolf Jagerman, Kai Hui, Honglei Zhuang, Junru Wu, Jiaming Shen, Tianqi Liu, Jialu Liu, Donald Metzler, Xuanhui Wang, and Michael Bendersky. 2023. <b>Large language models are effective text rankers with pairwise ranking prompting</b> . <i>CoRR</i> , abs/2306.17563.   | Fen Xia, Tie-Yan Liu, Jue Wang, Wensheng Zhang, and Hang Li. 2008. Listwise approach to learning to rank: theory and algorithm. In <i>Proceedings of the 25th international conference on Machine learning</i> , pages 1192–1199.   | 792 |
| 737 |   |   | 793 |
| 738 |   |   | 794 |
| 739 |   |   | 795 |
| 740 |   |   | 796 |
| 741 |   |   | 797 |
| 742 | Zhaopeng Qiu, Xian Wu, Jingyue Gao, and Wei Fan. 2021. U-bert: Pre-training user representations for improved recommendation. In <i>Proceedings of the AAAI Conference on Artificial Intelligence</i> , volume 35, pages 4320–4327.   | Jun Xu and Hang Li. 2007. <b>Adarank: a boosting algorithm for information retrieval</b> . In <i>SIGIR 2007: Proceedings of the 30th Annual International ACM SIGIR Conference on Research and Development in Information Retrieval, Amsterdam, The Netherlands, July 23-27, 2007</i> , pages 391–398. ACM. | 798 |
| 743 |   |   | 799 |
| 744 |   |   | 800 |
| 745 |   |   | 801 |
| 746 |   |   | 802 |
| 747 | Jeff Rasley, Samyam Rajbhandari, Olatunji Ruwase, and Yuxiong He. 2020. <b>DeepSpeed: System optimizations enable training deep learning models with over 100 billion parameters</b> . In <i>KDD '20: The 26th ACM SIGKDD Conference on Knowledge Discovery and Data Mining, Virtual Event, CA, USA, August 23-27, 2020</i> , pages 3505–3506. ACM.   | Zhi Zheng, Wenshuo Chao, Zhaopeng Qiu, Hengshu Zhu, and Hui Xiong. 2024. Harnessing large language model in text-rich sequential recommendation. In <i>Proceedings of the ACM Web Conference 2024, WWW 2024</i> . ACM.  | 803 |
| 748 |   |   | 804 |
| 749 |   |   | 805 |
| 750 |   |   | 806 |
| 751 |   |   | 807 |
| 752 |   |   | 808 |
| 753 |   |   |     |
| 754 | Devendra Singh Sachan, Mike Lewis, Mandar Joshi, Armen Aghajanyan, Wen-tau Yih, Joelle Pineau, and Luke Zettlemoyer. 2022. <b>Improving passage retrieval with zero-shot question generation</b> . In <i>Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing, EMNLP 2022, Abu Dhabi, United Arab Emirates, December 7-11, 2022</i> , pages 3781–3797. Association for Computational Linguistics. |   |     |
| 755 |   |   |     |
| 756 |   |   |     |
| 757 |   |   |     |
| 758 |   |   |     |
| 759 |   |   |     |
| 760 |   |   |     |
| 761 |   |   |     |
| 762 |   |   |     |
| 763 | Weiwei Sun, Lingyong Yan, Xinyu Ma, Shuaiqiang Wang, Pengjie Ren, Zhumin Chen, Dawei Yin, and Zhaochun Ren. 2023. <b>Is chatgpt good at search? investigating large language models as re-ranking agents</b> . In <i>Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing, EMNLP 2023, Singapore, December 6-10, 2023</i> , pages 14918–14937. Association for Computational Linguistics.         |   |     |
| 764 |   |   |     |
| 765 |   |   |     |
| 766 |   |   |     |
| 767 |   |   |     |
| 768 |   |   |     |
| 769 |   |   |     |
| 770 |   |   |     |
| 771 |   |   |     |
| 772 | Rohan Taori, Ishaan Gulrajani, Tianyi Zhang, Yann Dubois, Xuechen Li, Carlos Guestrin, Percy Liang, and Tatsunori B. Hashimoto. 2023. Stanford alpaca: An instruction-following llama model. <a href="https://github.com/tatsu-lab/stanford_alpaca">https://github.com/tatsu-lab/stanford_alpaca</a> .  |   |     |
| 773 |   |   |     |
| 774 |   |   |     |
| 775 |   |   |     |
| 776 |   |   |     |
| 777 | Michael Taylor, John Guiver, Stephen Robertson, and Tom Minka. 2008. Sofrank: optimizing non-smooth   |   |     |
| 778 |   |   |     |