

# CoRAL: Collaborative Retrieval-Augmented Large Language Models Improve Long-tail Recommendation

Junda Wu  
juw069@ucsd.edu  
University of California San Diego  
La Jolla, California, USA

Cheng-Chun Chang  
cc4900@columbia.edu  
Columbia University  
New York, New York, USA

Tong Yu  
tyu@adobe.com  
Adobe Research  
San Jose, California, USA

Zhankui He  
zhk004@eng.ucsd.edu  
University of California San Diego  
La Jolla, California, USA

Jianing Wang  
lygwjn@gmail.com  
University of California San Diego  
La Jolla, California, USA

Yupeng Hou  
yphou@ucsd.edu  
University of California San Diego  
La Jolla, California, USA

Julian McAuley  
jmcauley@ucsd.edu  
University of California San Diego  
La Jolla, California, USA

## ABSTRACT

The long-tail recommendation is a challenging task for traditional recommender systems, due to data sparsity and data imbalance issues. The recent development of large language models (LLMs) has shown their abilities in complex reasoning, which can help to deduce users' preferences based on very few previous interactions. However, since most LLM-based systems rely on items' semantic meaning as the sole evidence for reasoning, the collaborative information of user-item interactions is neglected, which can cause the LLM's reasoning to be misaligned with task-specific collaborative information of the dataset. To further align LLMs' reasoning to task-specific user-item interaction knowledge, we introduce collaborative retrieval-augmented LLMs, **CoRAL**, which directly incorporate collaborative evidence into the prompts. Based on the retrieved user-item interactions, the LLM can analyze shared and distinct preferences among users, and summarize the patterns indicating which types of users would be attracted by certain items. The retrieved collaborative evidence prompts the LLM to align its reasoning with the user-item interaction patterns in the dataset. However, since the capacity of the input prompt is limited, finding the minimally-sufficient collaborative information for recommendation tasks can be challenging. We propose to find the optimal interaction set through a sequential decision-making process and develop a retrieval policy learned through a reinforcement learning (RL) framework, **CoRAL**. Our experimental results show that CoRAL can significantly improve LLMs' reasoning abilities on specific recommendation tasks. Our analysis also reveals that CoRAL

can more efficiently explore collaborative information through reinforcement learning.

## KEYWORDS

Large language models, Collaborative Filtering, Long-tail Recommendation

### ACM Reference Format:

Junda Wu, Cheng-Chun Chang, Tong Yu, Zhankui He, Jianing Wang, Yupeng Hou, and Julian McAuley. 2018. CoRAL: Collaborative Retrieval-Augmented Large Language Models Improve Long-tail Recommendation. In *Proceedings of ACM Conference (Conference'17)*. ACM, New York, NY, USA, 11 pages. <https://doi.org/XXXXXXXX.XXXXXXX>

## 1 INTRODUCTION

Recommendation systems are valuable tools for users to explore content that matches their preferences. Traditional data-driven recommendation algorithms (e.g., collaborative filtering) can fail in long-tail recommendation, due to the uneven distribution in user-item interactions [26, 30, 40, 72]. In this paper, we aim to address the challenges associated with long-tail items in collaborative filtering-based recommender systems [70, 72]. In such scenario, long-tail items may have very few associated interactions, such that data-driven algorithms cannot accurately capture user-item interaction patterns [12, 28, 47]. In addition, models trained on such uneven datasets can suffer from selection bias [37, 52], exposure bias [15, 18] and popularity bias [1, 56, 67]. These biases can cause the models to overfit on popular items.

To tackle popularity bias and improve long-tail recommendation performance, data augmentation, and re-balancing methods can be directly applied. Data re-balancing methods [8, 11, 32, 63] try to reduce the distribution discrepancy between popular items and long-tail items in the training stage. However, these methods often obtain sub-optimal solutions due to learning inefficiency problems on long-tail data [70, 72]. This inefficiency leads to knowledge forgetting in the majority of the data, namely the popular items [72]. Since the goal of these methods is to achieve a compromise between the model's attention on popular items and more diversified

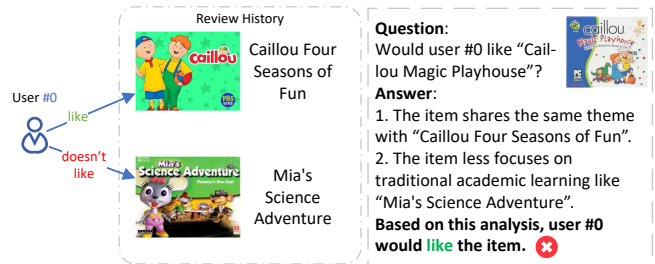
Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than the author(s) must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from [permissions@acm.org](mailto:permissions@acm.org).  
*Conference'17, July 2017, Washington, DC, USA*

© 2018 Copyright held by the owner/author(s). Publication rights licensed to ACM.  
ACM ISBN 978-1-4503-XXXX-X/18/06  
<https://doi.org/XXXXXXXX.XXXXXXX>

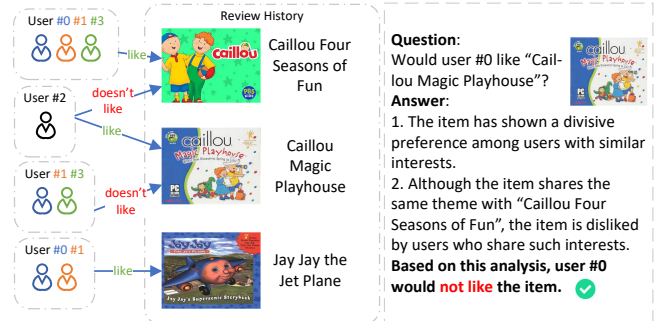
recommendations on long-tail items, achieving accurate recommendations for long-tail items can be challenging. Causal debiasing learning [6, 43, 75] is another line of work that focuses on how to learn the underlying user preferences, instead of simply learning user-item correlation from the data. Since long-tail items only have limited numbers of previous interactions, the model’s fine-grained reasoning abilities become essential for learning user preferences.

Large language models (LLMs) have recently demonstrated great reasoning abilities on very complex reasoning tasks [46, 49, 65], in which fine-grained reasoning paths can be generated to help with obtaining the correct answers [55, 73]. Previous works have also tried to adapt LLMs’ reasoning abilities to recommender systems [53, 54]. One line of previous works tries to use the language description of item content as the reasoning context [16, 42], which can be augmented by the LLM’s internal knowledge [4, 57, 62, 69]. By representing items as natural language, item representation distribution is aligned with the LLM’s knowledge. This alignment allows for a universal semantic representation of items, potentially mitigating the issue of long-tail items. In addition, by aligning recommendation tasks to the reasoning paradigms of language models, LLMs can be leveraged for more fine-grained reasoning based on the semantic contexts of users and items [53, 54]. However, due to several misalignment issues of LLMs in recommendation [31], directly prompting LLMs can be problematic. Specifically, the LLM’s understanding of user preferences over items can be misaligned with real user-item interaction patterns. For example, in Figure 1a, conventional LLM-based methods may simply recommend similar items (e.g., “Caillou Magic Playhouse” is recommended because the user likes “Caillou Four Seasons of Fun”).

In this work, we propose to formulate long-tail recommendation tasks as natural language reasoning problems and use LLMs to enable fine-grained reasoning about user preferences on long-tail items. To further align the reasoning of LLMs with task-specific knowledge of user-item interactions, we introduce collaborative retrieval-augmented LLMs, **CoRAL**, which directly incorporate collaborative evidence into the prompts, via collaborative prompting. For example, in Figure 1b, additional user-item interaction information can be retrieved by an additional lightweight model and included in the prompt. Based on the retrieved collaborative information, the LLM can find that although the items share the same theme (e.g., “Caillou”), the item (e.g., “Caillou Four Seasons of Fun”) is disliked by users who share such interests (e.g., preference on “Caillou Four Seasons of Fun”). However, due to the limited size of the prompts, a large amount of collaborative information cannot be included, and duplicate information may also distract the LLM’s reasoning process. Thus, we develop a retrieval policy to find the minimal-sufficient collaborative information which serves as the supporting evidence of the LLM’s reasoning on the given user-item pair. Since the number of users and items in a recommendation task is significantly larger than the capacity of a prompt input in LLMs, the retrieval policy should learn to explore diversified users and items for potential information gain, as well as exploit the collected collaborative information to maximize prediction accuracy. Based on the necessity to balance between exploration and exploitation in learning the retrieval policy, we propose to formulate the retrieval process as a sequential decision-making problem and employ reinforcement learning methods to maximize the long-term reward.



(a) Conventional item-based [16, 42] LLM reasoning process.



(b) Collaborative Retrieval Augmented LLM reasoning process.

**Figure 1: By text comprehension and extracting information-rich semantic features [41], LLMs in (a) can handle long-tail items [42], but still cannot directly leverage collaborative information. To handle long-tail items in collaborative filtering-based recommender systems, by collaborative prompting, LLMs in (b) can reason the fact that even if the current item shares the same theme with previously liked items, users with similar interests still dislike this item, which provides the rationale to not recommending it.**

To improve the data efficiency and reduce the early exploration time, we also propose to use the data from popular items to provide a warm start for the learning of the retrieval policy. Before the reinforcement learning stage, the user and item representations are learned from the data from popular items by conventional collaborative filtering methods. Then, the retrieval policy network will be initialized by the learned user and item representations, which improves the exploration efficiency and helps to solve the reward-sparsity problem. We summarize our contributions as follows:

- We identify the misalignment problem between the LLM’s reasoning process and long-tail recommendation tasks, which is caused by the lack of collaborative information.
- We propose to retrieve additional user-item interactions as collaborative information for collaborative prompting, which helps to align the LLM’s reasoning process to general recommendation tasks.
- We formulate the retrieval process as a sequential decision-making task and propose an RL framework, in which the retrieval policy learns to find the minimal-sufficient collaborative information specific to a recommendation task.

- To further improve data efficiency, we propose to learn the prior knowledge from more abundant data of popular items and to provide a warm start for the retrieval policy.

## 2 RELATED WORKS

### 2.1 Long-tail Recommendation

Long-tail recommendation plays a crucial role in mitigating the issue of highly skewed distributions of long-tail items in recommendation tasks. Previous works have proposed knowledge transfer learning methods based on novel model architecture designs. Zhang et al. [70] introduces a dual transfer learning framework by leveraging a model-level knowledge transfer and an item-level transfer to link head and tail items through shared features. Zhang et al. [72] propose a Cross Decoupling Network (CDN). This network aims to improve tail item recommendations while simultaneously maintaining overall system efficiency. Wu et al. [58] propose a domain transfer learning method (DACIR) for the sequential recommendation. However, the popularity bias in such long-tail or cold-start recommendation datasets is less discussed in their model designs. To address the issue of bias within the dataset in recommender systems, extensive efforts [6, 27, 59, 60, 75] have been dedicated to designing different training frameworks via causal debiasing. In this work, we propose to handle long-tail items by leveraging LLMs' abilities in fine-grained reasoning on collaborative information and extracting rich semantic features.

### 2.2 LLMs in Recommender Systems

A few studies underscore the expanding role of LLMs in recommender systems. Harte et al. [16] focus on enhancing sequential recommendation models with LLM embeddings, while Sanner et al. [42] explore the use of LLMs in processing language-based user preferences in dialog interfaces. A surge of approaches [4, 57, 62, 69] propose content augmentation method to reduce the cost of re-training or fine-tuning LLMs. To improve LLMs' reasoning capability for recommender systems, [53] proposes RecMind, specifically designed to deliver personalized recommendations. Wang et al. [54], on the other hand, use a retriever-reranking framework to enhance collaborative in-context understanding. However, the LLM's understanding of user preferences over items can be misaligned with real user-item interaction patterns, due to the lack of enough collaborative information.

To further align the LLM's reasoning process to specific recommendation tasks, several works have proposed to inject collaborative knowledge into LLMs by soft-prompt instruction tuning [71, 74]. [23] also introduces a method to transform discrete task-specific prompts into continuous prompt vectors, effectively linking IDs and words while decreasing inference time. However, such methods require an even larger amount of data to achieve good alignment, which would not help in our long-tail recommendation setting. Instead, we propose a lightweight retrieval policy to augment collaborative information in LLMs' reasoning process.

## 3 PROBLEM FORMULATION

In this paper, we focus on collaborative filtering-based recommender systems with long-tail items [70, 72]. We formulate long-tail recommendation as a complex reasoning task in LLMs [20, 22, 48].

When the LLM-empowered recommender system interacts with a user  $u \in \mathcal{U}$ , the task is to predict the user's preference for a long-tail item  $i \in \mathcal{I}$ . Since LLMs have no prior internal knowledge about a certain user's preference as well as the collaborative information of a certain recommendation task, the reasoning process can only be enabled by incorporating supporting evidence. To provide the information for the user  $u$ , items  $\mathcal{I}_u^{supp} = (i_1^u, i_2^u, \dots, i_m^u) \subseteq \mathcal{I}$  from the user's previous interactions are included, which helps the LLM to understand the preference of the user [5, 17, 53]. To help the LLM further understand how the user  $u$  would rate a certain item  $i$ , collaborative information will be retrieved as evidence of user-item interaction patterns. Due to the limitation of the LLM's reasoning context capacity, instead of including all the user-item interaction information, for a certain user-item pair  $z = (u, i)$ , the retrieval policy  $\pi_\theta$  will find a sequence of supporting users  $\mathcal{U}_z^{coll} = (u_1^z, u_2^z, \dots, u_l^z) \subseteq \mathcal{U}$  and a sequence of supporting items  $\mathcal{I}_z^{coll} = (i_1^z, i_2^z, \dots, i_t^z) \subseteq \mathcal{I}$ . At each time step  $t$ , the retrieval policy  $\pi_\theta$  needs to retrieve the next user-item pair  $(u_{t+1}^z, i_{t+1}^z)$  to augment current supporting evidence. In this work, we focus on how to obtain a minimal-sufficient information support for the LLM to deduce the accurate rating of  $z$ .

### 3.1 MDP Formulation for Retrieval Policy

We formulate the sequential retrieval process as a Markov Decision Process (MDP)  $\mathcal{M} = (\mathcal{S}, \mathcal{A}, P, r, \rho, \gamma)$ , where

- $\mathcal{S}$  is a continuous state space that encodes the collaborative information from the retrieved users and items as well as their collaborative preference patterns. The details of the design of state  $s \in \mathcal{S}$  encoding networks are explained in Section 4.2.
- $\mathcal{A}$  is a continuous action space that represents the retrieval queries of the next user and next item. At time step  $t + 1$ , the retrieval queries will try to retrieve the most relevant user and item in their feature spaces, as well as try to explore potentially useful users and items in the under-explored regions. The details of the action  $a \in \mathcal{A}$  prediction are explained in Section 4.2.
- $P : \mathcal{S} \times \mathcal{A} \times \mathcal{S} \rightarrow \mathbb{R}$ , is the state transition probability distribution, which captures the dynamics of the retrieval process.
- $r : \mathcal{S} \times \mathcal{A} \rightarrow \mathbb{R}$ , is the reward function  $r(s, a)$  given current state  $s$  and action  $a$ . The details of the reward function design are explained in Section 3.2.

At each time step  $t$ , the policy will generate the retrieval query  $a_t \in \mathcal{A}$  and retrieve the next user-item pair  $(u_t^z, i_t^z)$ . The learning objective is to find the optimal retrieval policy  $\pi_\theta^* : \mathcal{S} \rightarrow \mathcal{A}$  to achieve the long-term goal of obtaining a minimal-sufficient information support for the LLM, by maximizing the cumulative reward:

$$\pi_\theta^* = \arg \max_{\pi \in \Pi} \mathbb{E} \left[ \sum_{t=0}^T \gamma^t r(s_t, a_t) \right],$$

in which  $\gamma$  is the discount rate of future rewards and  $\Pi$  is the policy search space. When the retrieval policy observes the next user-item pair  $(u_t^z, i_t^z)$ , the state is updated by encoding the user-item pair into the state information  $s_{t+1} = P(\cdot | s_t, [u_t^z, i_t^z])$ .

### 3.2 Reward Function

For each time step  $t$ , the user-item rating prediction  $y_t$  will be prompted from the large language model  $P_\phi$  by the context  $C_t$  constructed from user information  $\mathcal{I}_u^{supp}$ , and previously collected collaborative information  $\mathcal{U}_z^{coll}$  and  $\mathcal{I}_z^{coll}$ ,

$$p_t = P_\phi(y_t | C_t), \quad (1)$$

$$C_t = C(\mathcal{I}_u^{supp}, \mathcal{U}_z^{coll}, \mathcal{I}_z^{coll}), \quad (2)$$

in which  $p_t$  is the prediction likelihood of whether the user  $u$  likes the item  $i$ , and  $C$  is the prompt template (detailed description in Section 4.1) which composes the collected information into a natural language query.

Since the motivation of the retrieval policy is to maximize cumulative information gain, we use the marginal information gain at each time step  $t$  as the reward signal  $r_t$ , which is calculated by the prediction discrepancy,

$$r_t(s_t, (u_t^z, i_t^z)) = \underbrace{|p_{t-1} - y^{gt}|}_{\text{discrepancy at } t-1} - \underbrace{|p_t - y^{gt}|}_{\text{discrepancy at } t}, \quad (3)$$

in which  $y^{gt} = \mathbf{M}(u, i)$  is the ground truth label of the user's preference on the item from the training rating matrix  $M$ . Following [35, 76], we reward those retrieved user-item pairs which will lead the constructed prompt to find a more accurate prediction based on the LLM  $P_\phi$ .

## 4 PROPOSED FRAMEWORK: CORAL

In this section, we first explain the prompting method designed to incorporate collaborative information and collect the LLM's prediction as the recommendation prediction. Then, we introduce the collaborative retrieval policy network as well as the reinforcement learning process, which is illustrated in Algorithm 1. In reinforcement learning, we treat the LLM as part of the environment and thus the LLM is frozen while a lightweight retrieval policy is learnable with significantly fewer learning parameters. To further improve the policy's learning efficiency and accommodate long-tail recommendation, we propose to use collaborative filtering models learned on the short-head data as the model initialization (detailed experimental settings and comparison results are in Section 5.3).

### 4.1 Collaborative Prompting

In this section, we will explain how to construct the prompt  $C_t$  with the retrieved collaborative information in Eq. (2), and how to obtain the prediction probability  $p_t$  in Eq. (1), given the retrieval results from the policy  $\pi_\theta$ . Details about the policy network design will be explained in Section 4.2.

**Collaborative Information.** At time step  $t$ , given the user-item pair  $z = (u, i)$ , the retrieval policy  $\pi_\theta$  obtains the supporting users  $\mathcal{U}_z^{coll}$  and items  $\mathcal{I}_z^{coll}$ . To describe the users and items in natural language and incorporate them into the prompt, we represent each user  $u_i^z \in \mathcal{U}_z^{coll}$  by their user index  $\mathbf{idx}_U(u_i^z)$ , and each item  $i_i^z \in \mathcal{I}_z^{coll}$  by its item index  $\mathbf{idx}_I(i_i^z)$ . We further extract a short text description  $\mathbf{desc}_I(i_i^z)$  for each item from the metadata (detailed descriptions in Section 5.1.1), to assist the LLM's understanding of the item. Based on the rating matrix  $M$  in the training dataset, we

can summarize users' shared preference for each item  $i \in \mathcal{I}_z^{coll}$  in the following format:

$$\begin{aligned} \text{POS}(i, \mathcal{U}_z^{coll}) &= \left\{ \mathbf{M}(i, u) \geq y^{thresh}, u \in \mathcal{U}_z^{coll} \right\}, \\ \text{NEG}(i, \mathcal{U}_z^{coll}) &= \left\{ \mathbf{M}(i, u) < y^{thresh}, u \in \mathcal{U}_z^{coll} \right\}, \end{aligned} \quad (4)$$

in which the rating threshold  $y^{thresh}$  is to determine if the rating is positive or negative. By aggregating the preference of a group of users for each item, the length of the prompt can be significantly reduced, and such descriptions prompt the LLM to focus more on the comparative preference among the users. To construct the first part of the prompt which contains collaborative information, we design the prompt as follows:

- **Role-play:** As a recommender system please solve the following problem.
- **Collaborative Information:** Repeat  $i \in \mathcal{I}_z^{coll}$ 
  - The item  $\mathbf{desc}_I(i)$  is liked by the users  $\text{POS}(i, \mathcal{U}_z^{coll})$ .
  - The item  $\mathbf{desc}_I(i)$  is disliked by the users  $\text{NEG}(i, \mathcal{U}_z^{coll})$ .
- **Summarization:** Try to understand the pattern that the item  $\mathbf{desc}_I(i)$  is typically liked by what kinds of users based on the above information.

Based on our empirical observations, the last **Summarization** instruction is essential to align the LLM's reasoning with the goal of this task.

**User Preference Representation.** To include more information on the user's preference, we follow previous works [4, 57, 62, 69] to include the user's previously interacted items  $\mathcal{I}_z^{supp}$  and their text descriptions. However, different from previous works, we also divide the previous items  $\mathcal{I}_z^{supp}$  of user  $u$  into positive and negative sets, and then query the LLM to deduce the rating for the user-item pair  $z = (u, i)$ ,

$$\begin{aligned} \text{POS}(\mathcal{I}_z^{supp}, u) &= \left\{ \mathbf{M}(i, u) \geq y^{thresh}, i \in \mathcal{I}_z^{coll} \right\}, \\ \text{NEG}(\mathcal{I}_z^{supp}, u) &= \left\{ \mathbf{M}(i, u) < y^{thresh}, i \in \mathcal{I}_z^{coll} \right\}, \end{aligned} \quad (5)$$

Then, we construct the second part of the prompt by including the user's previously interacted items  $\mathcal{I}_z^{supp}$ :

- **User's Positive Preference:** Items the user  $\mathbf{idx}_U(u)$  likes are as follows:  $\text{POS}(\mathcal{I}_z^{supp}, u)$ .
- **User's Negative Preference:** Items the user  $\mathbf{idx}_U(u)$  does not like are as follows:  $\text{NEG}(\mathcal{I}_z^{supp}, u)$ .
- **Query:** For the item described as  $\mathbf{idx}_I(i)$ , would you recommend it to the user  $\mathbf{idx}_U(u)$ ?

With the prompt design (denoted as  $C$ ) described above, we can aggregate the information retrieved at time step  $t$  and transform the information into a natural language prompt  $C_t = C(\mathcal{I}_u^{supp}, \mathcal{U}_z^{coll}, \mathcal{I}_z^{coll})$ . To obtain the LLM's prediction as well as its confidence score, we extract the prediction probability  $p_t = P_\phi(y_t | C_t)$  of the next token generated from the LLM. Specifically, we strictly ask the LLM to answer either "Yes" or "No" without additional text provided, and we take the probability of the LLM generating the token "Yes" as our final score  $p_t$ .

## 4.2 Retrieval Policy Network

In this section, we design the retrieval policy  $\pi_\theta$  to sequentially include additional users and items, which may provide an information gain for the LLM’s reasoning. Since the prompt has only a limited capacity of users and items included, the goal of the retrieval policy is to construct a minimal-sufficient prompt that contains complete information about the current recommendation task of the user-item pair  $z = (u, i)$ . Specifically, the retrieval policy needs to maximize its long-term information gain by maximizing the cumulative reward function.

Instead of learning the action distribution over all the users and items like value-based reinforcement learning methods [33, 44], we choose to directly learn the continuous vector representations of the next user and item based on the DDPG algorithm [25], which helps to learn a low-rank decision space and also makes the solution more scalable even with new users and items included during the inference stage.

**4.2.1 State Representation.** For each user-item pair  $z = (u, i)$ , the retrieval process starts with the user-item embedding  $\mathbf{s}_0 = [\mathbf{u}, \mathbf{i}] \in \mathbb{R}^{2d}$ , in which  $\mathbf{u}$  and  $\mathbf{i}$  are the user and item embeddings in  $d$  dimensions. Notably, the user and item embeddings are randomly initialized by the multivariate normal distribution  $\mathbf{u} \sim \mathcal{N}(\boldsymbol{\mu}, \Sigma)$  and  $\mathbf{i} \sim \mathcal{N}(\boldsymbol{\mu}, \Sigma)$ , in which  $\boldsymbol{\mu}$  is a  $d$ -dimensional zero-vector and  $\Sigma$  is a  $d \times d$  unit matrix.

During the early stage of the reinforcement learning process, when the retrieval policy behaves randomly, similar users and items are likely to be retrieved due to the large-scale user and item spaces, which makes the reward of the policy’s exploration very sparse. To overcome the exploration difficulty, we initialize the policy with the embeddings pre-trained on the portion of the dataset with popular items, which can provide a warm start for the learning of the retrieval policy (detailed comparison results are explained in Section 5.2 and Section 5.3).

**4.2.2 User-item Retrieval.** At each time step  $t$ , based on the current state  $\mathbf{s}_t$ , the retrieval policy  $\pi_\theta$  will find the next user-item pair. Due to the large user and item spaces, direct exploration in the discrete spaces of the users and items can be extremely inefficient. Thus, we employ a continuous action space  $\mathcal{A} \subseteq \mathbb{R}^{2d}$  which also covers the user-item embedding space. The retrieval policy will first generate a user-item query based on the current state  $\mathbf{a}_{t+1} = [\mathbf{a}_{t+1}^u, \mathbf{a}_{t+1}^i] = \pi_\theta(\cdot | \mathbf{s}_t)$ , and try to find the nearest user and item in terms of a distance measurement  $d(\cdot, \cdot)$  defined on the embedding spaces,

$$u_{t+1}^z = \min_{u \in \mathcal{U}} d(\mathbf{u}, \mathbf{a}_{t+1}^u), \quad i_{t+1}^z = \min_{i \in \mathcal{I}} d(\mathbf{i}, \mathbf{a}_{t+1}^i), \quad (6)$$

in which  $\mathbf{u}$  and  $\mathbf{i}$  denote the embeddings of the user  $u$  and item  $i$  respectively, and we use Euclidean distance for  $d$ . The retrieved user and item will be added to the collaborative information  $\mathcal{U}_z^{\text{coll}}$  and  $\mathcal{I}_z^{\text{coll}}$ .

**4.2.3 State Transition.** The state  $\mathbf{s}_t$  encodes the current collaborative information, which will be updated for each time step after the user and item is retrieved. To track the retrieval process and aggregate the collected information, we use the multi-layer perception

model (MLP) for state transition modeling,

$$\mathbf{s}_{t+1} = \text{MLP}(\mathbf{s}_t, [\mathbf{u}_{t+1}^z, \mathbf{i}_{t+1}^z]), \quad (7)$$

in which  $\mathbf{u}_{t+1}^z$  and  $\mathbf{i}_{t+1}^z$  are the embeddings of the retrieved user and item at time step  $t$ .

## 4.3 Minimal-sufficient Collaborative Information via Reinforcement Learning

We follow the standard DDPG [25] reinforcement learning framework to train our retrieval policy with the continuous action space. In the Actor-Critic framework [25], the critic is learning a Q-value function with episodic mini-batch sampled from the replay buffer [34],

$$L(\theta^Q) = \mathbb{E}_{s,a,r,s'} [r + \gamma Q_{\theta^{Q'}}(s', \pi_{\theta^\mu}(\cdot | s')) - Q_{\theta^Q}(s, a)]^2, \quad (8)$$

in which  $\theta^{Q'}$  is the target network [25] of the critic, which is fixed during the act network update. Based on the learning objective in Eq. (8), we can derive the gradient  $\nabla_{\theta^Q} L(\theta^Q)$  to update the act network of the critic. Since the critic provides an approximation of the Q-value function, the optimization step of the actor network can be achieved by policy gradient,

$$\nabla_{\theta^\mu} L(\theta^\mu) = \mathbb{E}_s [\nabla_a Q_{\theta^Q}(s, a) \nabla_{\theta^\mu} \pi_{\theta^\mu}(\cdot | s)], \quad (9)$$

similar to the critic network, the policy gradient only updates the act network of the actor, while the target actor network  $\pi_{\theta^\mu}$  will be synchronized after each update step.

To further enable continuous space exploration, we follow [25] to add exploration  $\mathcal{N}$  noise to the target policy  $\pi_{\theta^\mu}$  to find unexplored but informative users and items,

$$\pi_{\theta^\mu'}(\cdot | s) = \pi_{\theta^\mu}(\cdot | s) + \mathcal{N}, \quad (10)$$

in which we choose the Ornstein–Uhlenbeck [38] random process as the exploration process  $\mathcal{N}$ .

## 5 EXPERIMENTS

In this section, we conduct extensive experiments on multiple datasets to investigate the following research questions (RQs):

- **RQ1:** How does collaborative information help to align the LLM’s reasoning process to general recommendation tasks?
- **RQ2:** Can CoRAL find sufficient collaborative evidence to enhance LLMs’ reasoning?
- **RQ3:** Can CoRAL find minimally-sufficient collaborative evidence to fit the size of prompts?

### 5.1 Experimental Settings

**5.1.1 Datasets.** We evaluate CoRAL and baselines on four Amazon Product [36] tasks, which are used in the evaluations of many collaborative filtering methods [68]:

- **Appliances** refers to a category of home and kitchen devices sold on Amazon. This subset contains 602,777 reviews with 515,650 users and 30,252 products. We use the “title” of the items in the metadata as item descriptions.
- **Gift Cards** on Amazon are prepaid store value cards that can be used as an alternative to cash. This subset contains 147,194 reviews with 128,877 users and 1,548 products. We use the “description” of the items in the metadata.

**Algorithm 1** Training Procedure of CoRAL

---

**Input** episode length  $L$ , Maximum steps in an episode  $T$ .  
**Initialize** actor network  $\theta^\mu$  and critic network  $\theta^Q$   
**Initialize** target networks  $\theta^{\mu'} \leftarrow \theta^\mu$  and  $\theta^{Q'} \leftarrow \theta^Q$   
**Initialize** the replay buffer  $\mathcal{D} = \emptyset$   
**while**  $l \leq L$  **do**  
  Receive a user-item pair  $z = (u, i)$   
  **Initialize**  $\mathcal{I}_u^{supp} = \emptyset$ ,  $\mathcal{U}_z^{coll} = \emptyset$ ,  $\mathcal{I}_z^{coll} = \emptyset$   
  Construct the prompt of user preference  $\mathcal{I}_u^{supp}$  as in Eq. (5)  
  Get the initial prediction  $p_0$  according to Eq. (2)  
  **while**  $t \leq T$  **do**  
    **User-item Retrieval**  
    Generate the current action  $\mathbf{a}_t$  from the policy  $\pi_{\theta^{\mu'}}$   
    Locate the next user-item pair  $(u_t^z, i_t^z)$  as in Eq. (6)  
    Add to the support sets,  $\mathcal{U}_z^{coll} \leftarrow u_t^z$  and  $\mathcal{I}_z^{coll} \leftarrow i_t^z$   
  
    **Collaborative Prompting**  
    Construct the prompt of collaborative information  $\mathcal{U}_z^{coll}$   
    and  $\mathcal{I}_z^{coll}$  according to Eq. (4)  
    Get the current prediction  $p_t$  as in Eq. (1)  
    Calculate the current reward  $r_t$  according to Eq. (3)  
    Observe the next state  $s_{t+1}$  according to Eq. (7)  
    Store the transition quadruple  $(s_t, \mathbf{a}_t, r_t, s_{t+1})$  in  $\mathcal{D}$   
  
    **Networks Update**  
    Sample a minibatch of the quadruple  $(s, \mathbf{a}, r, s')$  from  $\mathcal{D}$   
    Calculate the minibatch loss  $L(\theta^Q)$  for the critic network  
    according to Eq. (8)  
    Update the critic network by the gradient  $\nabla_{\theta^Q} L(\theta^Q)$   
    Update the actor network with the sampled policy gradient  
    according to Eq. (9)  
    Update the target networks:  
     $\theta^{Q'} \leftarrow \tau\theta^Q + (1 - \tau)\theta^{Q'}$   
     $\theta^{\mu'} \leftarrow \tau\theta^\mu + (1 - \tau)\theta^{\mu'}$   
  **end while**  
**end while**

---

- **Prime Pantry** on Amazon refers to a service offering a wide range of everyday household items and groceries. The subset contains 471,614 reviews with 247,659 users and 10,814 products. We use the original “description” in the metadata as the item descriptions.
- **Software** goods on Amazon refer to digital products. The subset contains 459,436 reviews with 375,147 users and 21,663 products. The software product titles in the metadata are used as the item descriptions.

Due to the missing descriptions of items in the metadata of some datasets, we use the item titles as the replacement. To determine the boundary between popular items and long-tail items, we follow the typical 80/20 rule [29, 45, 64, 66], which defines the least 80% items as long-tail items. Due to the sparsity of the datasets, many users and items only have very few entries in the datasets, in which case collaborative information is almost inaccessible. To maintain a certain number of interaction data samples, We follow [9, 19, 24, 71] to filter out users and items with fewer than 5 interactions. For

learning-based baselines and CoRAL, we use 70% of the long-tail data and the remaining data in the dataset as the training data. The remaining 30% of the long-tail data is split equally into the validation and test data. We follow the standard preprocessing method [51, 71] to convert the original 5-score into binary labels by the threshold of 3.

**5.1.2 Metrics.** We follow the metrics *AUC* and *F1* in long-tail recommendation [13, 67] and collaborative filtering [3, 71]. The Area Under the Curve (AUC) metric is a performance measurement for classification models that evaluates the tradeoff between true positive rate and false positive rate across different thresholds, where a higher AUC indicates better model performance. The F1 metric is a statistical measure used in classification tests, combining precision and recall to provide a score that balances both false positives and false negatives, calculated as the mean of precision and recall.

**5.1.3 Baselines.** We introduce baselines in our experiments from three lines of work, collaborative filtering, popularity debiasing, and LLM-based recommendation methods:

Collaborative Filtering:

- **AFM** [61]: A model learns the significance of each feature interaction from data through a neural attention network.
- **DCN** [50]: A deep neural network featuring a cross-structure is designed for enhanced efficiency in learning specific bounded-degree feature interactions.
- **DFM** [14]: A unified neural network architecture for recommender systems is proposed, integrating factorization machines and deep learning.
- **WDL** [10]: A method that integrates wide linear models with deep neural networks is proposed for enhancing recommender systems. This approach synergistically leverages the strengths of both memorization and generalization.

Popularity debiasing baselines directly enhance the collaborative filtering methods via causal debiasing:

- **IPS** [43]: An approach utilizes causal inference techniques to address selection biases in data, ensuring unbiased performance estimation with biased data.
- **CausE** [6]: A domain adaptation technique is developed to train on historical data that captures results from a recommendation system biased by a specific policy and makes predictions for recommendation outcomes under random exposure conditions.

To understand the benefit of collaborative information in LLMs, we consider a LLM-based baseline:

- **LLM-Language** [42]: A LLM prompting method that describes the user’s interacted items and the user’s preferences before asking for the user’s preference on new items.

To understand the behavior of our approach, we consider variants of CoRAL:

- **CoRAL-Method**: The collaborative information augmented LLMs, in which the retrieval policy is initialized by the *Method*. The *Method* in our experiments includes DFM, WDL, AFM, and DCN.

	Software		Prime Pantry		Gift Cards		Appliances		Average	
	AUC	F1	AUC	F1	AUC	F1	AUC	F1	AUC	F1
<b>AFM</b> [61]	75.12	58.39	69.47	52.51	46.93	61.56	76.86	65.52	67.10	59.49
<b>DCN</b> [50]	76.75	66.20	73.30	49.99	55.59	67.07	80.70	71.15	71.59	63.60
<b>DFM</b> [14]	76.04	66.63	72.92	57.86	66.76	60.01	81.83	77.37	74.39	65.47
<b>WDL</b> [10]	78.20	69.25	73.77	56.43	60.81	57.66	73.82	74.56	71.65	64.48
<b>IPS</b> [43]	78.24	71.32	72.24	61.65	64.79	63.95	82.28	75.65	74.39	66.23
<b>CausE</b> [6]	77.78	70.84	73.69	59.80	70.51	65.39	76.86	72.04	74.71	67.02
<b>LLM-Language</b> [42]	73.10	66.32	51.48	41.47	83.52	74.85	74.36	70.52	70.61	63.29
<b>CoRAL-random</b>	77.56	58.60	64.07	50.15	91.30	59.66	77.51	61.35	77.61	57.44
<b>CoRAL-DFM</b>	95.25	88.68	<b>93.32</b>	<b>86.73</b>	96.52	67.51	90.87	86.76	<b>93.99</b>	82.42
<b>CoRAL-WDL</b>	93.97	<b>91.18</b>	87.08	80.52	92.22	70.74	92.55	<b>89.22</b>	91.45	82.92
<b>CoRAL-AFM</b>	<b>93.99</b>	88.41	89.10	86.17	<b>98.99</b>	<b>76.17</b>	<b>92.66</b>	84.55	93.69	<b>83.83</b>
<b>CoRAL-DCN</b>	91.74	87.20	85.75	77.59	97.16	70.63	91.73	86.28	91.59	80.43

Table 1: Experimental results (AUC and F1) on four Amazon Product datasets.

- **CoRAL-random**: The LLM is also augmented by collaborative information. However, the retrieval policy is just a rule-based model which randomly retrieves items.

*5.1.4 Implementation Details.* We implement our retrieval policy network using PyTorch 2.1. For reinforcement learning, the DDPG [25] policy network is implemented using *Stable-Baselines3* [39]. We set the memory buffer size to 1000 and the training batch size to 16 for both the actor and the critic. For the continuous action for the next user and item, we set the dimensions for both as 128, which aligns with the size of the user and item embedding. The Ornstein-Uhlenbeck noise [38] added in each dimension of the continuous action space for exploration is set to zero-mean and the standard deviation as  $\sigma = 0.1$ . The reinforcement learning process starts at the 10-th iteration, which enables a warm start. We use Adam optimizer [21] for all the model learning with the learning rate as 0.001, and we set the maximal learning iterations to 2,000.

We implement the reinforcement learning environment using Gym[7], and we use a GPT-4 [2] model as the backbone large language model to provide reward. During the training stage, we allow up to 10 interactions within a single episode and enable early stop if the absolute value of the discrepancy between the predicted rating and the ground truth rating is less than 0.1. During the evaluation stage, for each data sample, we consistently let the policy retrieve 5 rounds of users and items as collaborative information.

## 5.2 Recommendation Performance (RQ1)

We show the comparison results of CoRAL and various baselines to demonstrate the effectiveness of augmenting LLMs with collaborative information as reasoning evidence.

*5.2.1 Effect of the Retrieval Policy.* In Table 1, we can observe that by adding collaborative information into the LLM’s prompt, even if the retrieved users and items are randomly chosen, **CoRAL-random** consistently outperforms **LLM-Language** in terms of the AUC scores. Such observations may imply that collaborative information is still crucial to specific recommendation tasks, even if the LLM can understand the general semantic meanings of the items and the user’s preference. On the other hand, we also observe

that **CoRAL-random** generally performs worse (except for Prime Pastry) than **LLM-Language** in terms of the F1 scores. One reasonable explanation is that since **CoRAL-random** is not specifically curating its selection of users and items, irrelevant information may bring additional bias into the recommendation process and cause the model to have a poor precision-recall trade-off.

*5.2.2 Effect of Online Reinforcement Learning.* In Table 1, we observe an inconsistency in performance comparison between traditional recommendation baselines and LLM-based baselines, as the LLM-based methods can sometimes (e.g., Prime Pastry and Appliances) perform substantially worse than traditional baselines. This inconsistency in the performance of LLM-based methods suggests the misalignment between the LLM’s reasoning and specific recommendation tasks. The proposed method CoRAL specifically aligns the LLM’s reasoning with the recommendation tasks through reinforcement learning. With the LLM’s reasoning process aligned to the user-item interaction patterns, we can observe significant improvements up to 21.1% and 25.1% for AUC and F1 scores respectively on average.

## 5.3 Sufficient Collaborative Information from Popular Items (RQ2)

We conduct analytical experiments in this section to show how learning from popular items can benefit online reinforcement learning. We choose DFM and WDL as the backbone models in the analytical experiments to show their different learning behaviors.

*5.3.1 Comparison to Randomly Initialized Policy.* In Table 2, we compare our method, which initializes the retrieval policy by learning from the popular items, with the variant that randomly initializes the policy. We show that the policies initialized with models learned from popular items are generally performing better than randomly initialized policies, which suggests the data-efficiency advantage of our method. Since at the early steps of reinforcement learning, the exploration stage may take a long time to navigate and find high-value actions through trial and error, without efficient exploration strategies or some good embedding spaces, the actor networks could easily overfit and fail to explore better actions.

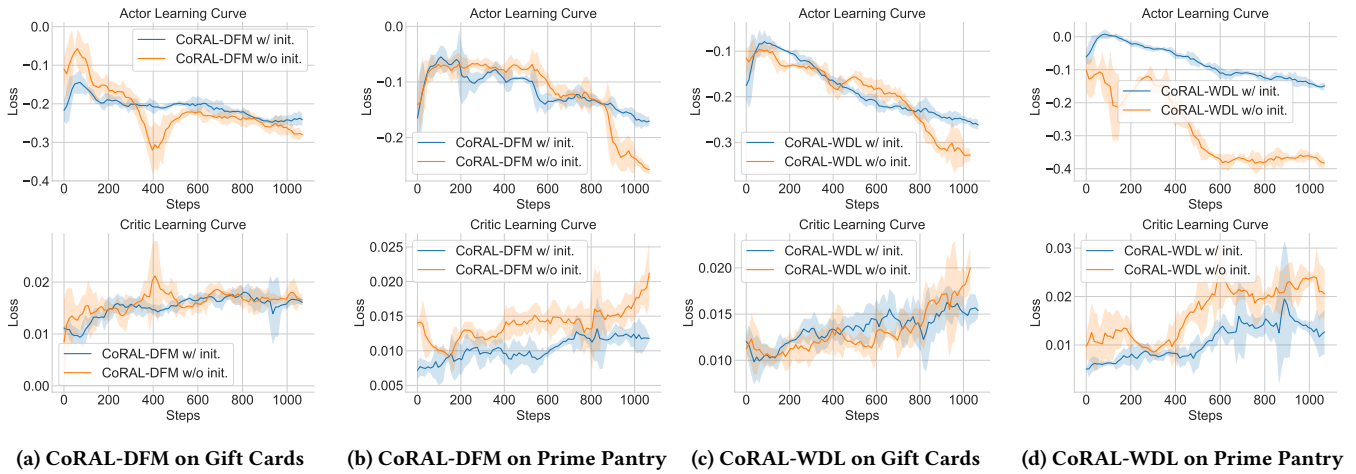


Figure 2: CoRAL’s (DFM and WDL) learning curves on Gift Cards and Prime Pantry datasets.

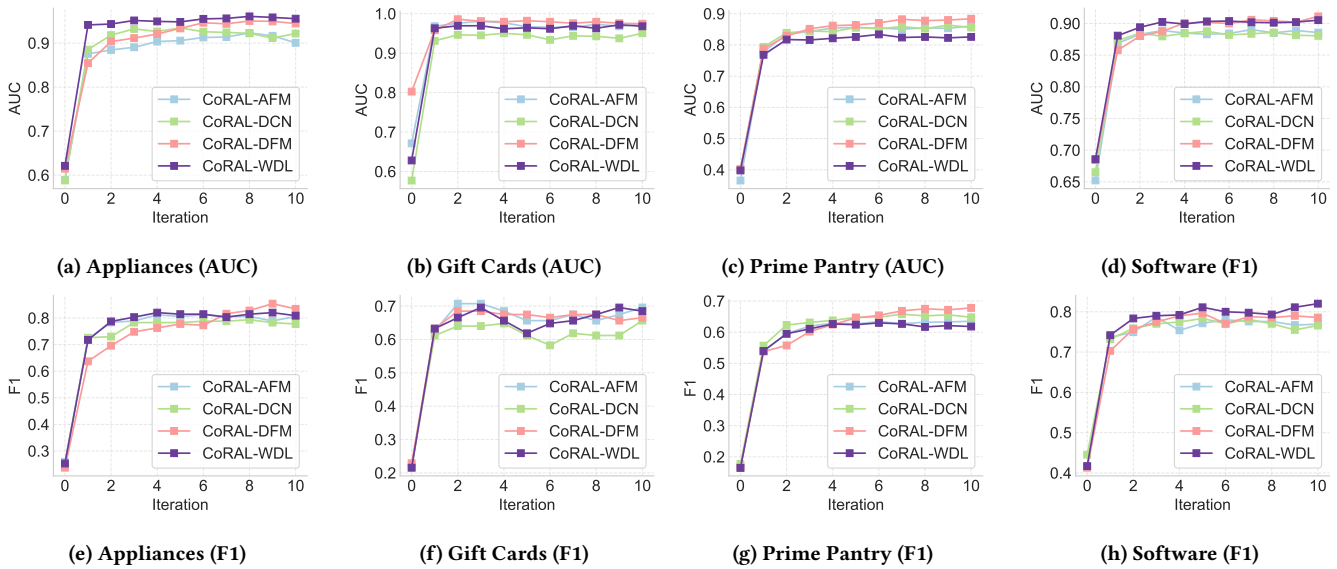


Figure 3: CoRAL’s performance (AUC and F1) w.r.t number of iterations of user-item retrieval on four Amazon Product datasets.

5.3.2 Actor-critic Learning Curves. In Figure 2, we show the learning curves of the actor and critic networks, for policies with and without short-head data initialization. We choose the more challenging datasets, Gift Cards, and Prime Pantry, in terms of methods general F1 performance in Table 1, which suggests that these tasks require better balancing between exploration and exploitation. We can observe a consistent pattern that the actor networks of the policies with random initialization converge faster than policies with short-head data initialization. However, the critic networks of the policies with random initialization have higher learning loss than policies with short-head data initialization. Such an observation suggests that the randomly initialized policies could easily overfit and thus the actor-critic learning process can be done asynchronously.

With the pre-trained user and item embedding spaces on the short-head training dataset, the exploration in the continuous embedding space can be more efficient.

### 5.4 Minimally-sufficient Collaborative Information from Iterative Retrieval (RQ3)

In Figure 3, we show the models’ performance w.r.t the number of rounds of retrieval. We observe that for all the policies of CoRAL, in each iteration, they manage to retrieve informative users and items, while consistently achieving information gain with the information gain margin decreasing. Comparing the backbones DFM and DCN of CoRAL, we find a common exploration-exploitation behavior of these two policies, as the DCN acts more greedy and reaches its upper-bound performance sooner, while the DFM tends



		CoRAL-DFM		CoRAL-WDL	
		w/ init.	w/o init.	w/ init.	w/o init.
Software	AUC	<b>95.25</b>	93.58	<b>93.97</b>	92.35
	F1	<b>88.68</b>	88.36	<b>91.18</b>	88.87
Prime Pantry	AUC	<b>93.32</b>	89.33	87.08	<b>89.49</b>
	F1	<b>86.73</b>	80.76	80.52	<b>81.86</b>
Gift Cards	AUC	<b>96.52</b>	<b>96.52</b>	92.22	<b>96.98</b>
	F1	<b>67.51</b>	64.25	<b>70.74</b>	68.81
Appliances	AUC	90.87	<b>94.48</b>	<b>92.55</b>	91.84
	F1	86.76	<b>88.82</b>	<b>89.22</b>	83.00
Average	AUC	<b>93.99</b>	93.48	91.45	<b>92.66</b>
	F1	<b>82.42</b>	80.55	<b>82.92</b>	80.64

**Table 2: Ablation study of CoRAL’s performance with or without short-head data initialization for DFM and WDL as the collaborative filtering backbones.**

to be more explorative in the early stage and achieves better final performance. Such an observation suggests the importance of the exploration-exploitation trade-off, which can be more efficiently achieved through the proposed reinforcement learning framework.

## 6 CONCLUSION

In this paper, we focus on collaborative filtering-based recommender systems with long-tail items [70, 72]. We introduce **CoRAL**, an approach for enhancing long-tail recommendations in traditional collaborative filtering-based recommender systems, overcoming the limitations of data sparsity and imbalance that hamper collaborative filtering methods. **CoRAL** integrates collaborative retrieval-augmented LLMs to align the model’s reasoning with actual user-item interaction patterns. This alignment is pivotal in addressing the common oversight in LLM-based systems that rely heavily on semantic interpretations, neglecting the collaborative dimensions of user-item interactions. Additionally, CoRAL employs a reinforcement learning framework to develop a retrieval policy, identifying an optimal set of user-item interactions as the supporting evidence for the LLM’s reasoning. This strategy ensures minimal yet sufficient collaborative information is used, enhancing the LLM’s ability to accurately deduce user preferences and interaction dynamics, hence offering a significant improvement on LLM-based recommendation.

## REFERENCES

- [1] Himan Abdollahpouri, Masoud Mansoury, Robin Burke, Bamshad Mobasher, and Edward Malthouse. 2021. User-centered evaluation of popularity bias in recommender systems. In *Proceedings of the 29th ACM Conference on User Modeling, Adaptation and Personalization*. 119–129.
- [2] Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman, Diogo Almeida, Janko Altschmidt, Sam Altman, Shyamal Anadkat, et al. 2023. Gpt-4 technical report. *arXiv preprint arXiv:2303.08774* (2023).
- [3] Vito Walter Anelli, Alejandro Bellogín, Tommaso Di Noia, and Claudio Pomo. 2021. Reenvisioning the comparison between neural collaborative filtering and matrix factorization. In *Proceedings of the 15th ACM Conference on Recommender Systems*. 521–529.
- [4] Jinheon Baek, Nirupama Chandrasekaran, Silviu Cucerzan, Sujay Kumar Jauhar, et al. 2023. Knowledge-Augmented Large Language Models for Personalized

- Contextual Query Suggestion. *arXiv preprint arXiv:2311.06318* (2023).
- [5] Keqin Bao, Jizhi Zhang, Yang Zhang, Wenjie Wang, Fuli Feng, and Xiangnan He. 2023. Tallrec: An effective and efficient tuning framework to align large language model with recommendation. *arXiv preprint arXiv:2305.00447* (2023).
- [6] Stephen Bonner and Flavian Vasile. 2018. Causal embeddings for recommendation. In *Proceedings of the 12th ACM conference on recommender systems*. 104–112.
- [7] Greg Brockman, Vicki Cheung, Ludwig Pettersson, Jonas Schneider, John Schulman, Jie Tang, and Wojciech Zaremba. 2016. Openai gym. *arXiv preprint arXiv:1606.01540* (2016).
- [8] Jonathon Byrd and Zachary Lipton. 2019. What is the effect of importance weighting in deep learning?. In *International conference on machine learning*. PMLR, 872–881.
- [9] Chong Chen, Min Zhang, Yongfeng Zhang, Weizhi Ma, Yiqun Liu, and Shaoping Ma. 2020. Efficient heterogeneous collaborative filtering without negative sampling for recommendation. In *Proceedings of the AAAI Conference on Artificial Intelligence*, Vol. 34. 19–26.
- [10] Heng-Tze Cheng, Levent Koc, Jeremiah Harmsen, Tal Shaked, Tushar Chandra, Hrishu Aradhya, Glen Anderson, Greg Corrado, Wei Chai, Mustafa Ipsir, et al. 2016. Wide & deep learning for recommender systems. In *Proceedings of the 1st workshop on deep learning for recommender systems*. 7–10.
- [11] Yin Cui, Menglin Jia, Tsung-Yi Lin, Yang Song, and Serge Belongie. 2019. Class-balanced loss based on effective number of samples. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*. 9268–9277.
- [12] Zhen Gong, Xin Wu, Lei Chen, Zhenzhe Zheng, Shengjie Wang, Anran Xu, Chong Wang, and Fan Wu. 2023. Full Index Deep Retrieval: End-to-End User and Item Structures for Cold-start and Long-tail Item Recommendation. In *Proceedings of the 17th ACM Conference on Recommender Systems*. 47–57.
- [13] Yulong Gu, Zhuoye Ding, Shuaiqiang Wang, Lixin Zou, Yiding Liu, and Dawei Yin. 2020. Deep multifaceted transformers for multi-objective ranking in large-scale e-commerce recommender systems. In *Proceedings of the 29th ACM International Conference on Information & Knowledge Management*. 2493–2500.
- [14] Huifeng Guo, Ruiming Tang, Yunming Ye, Zhenguo Li, and Xiuqiang He. 2017. DeepFM: a factorization-machine based neural network for CTR prediction. *arXiv preprint arXiv:1703.04247* (2017).
- [15] Shantanu Gupta, Hao Wang, Zachary Lipton, and Yuyang Wang. 2021. Correcting exposure bias for link recommendation. In *International Conference on Machine Learning*. PMLR, 3953–3963.
- [16] Jesse Harte, Wouter Zorgdrager, Panos Louridas, Asterios Katsifodimos, Dietmar Jannach, and Marios Fragkoulis. 2023. Leveraging large language models for sequential recommendation. In *Proceedings of the 17th ACM Conference on Recommender Systems*. 1096–1102.
- [17] Wang-Cheng Kang, Jianmo Ni, Nikhil Mehta, Maheswaran Sathiamoorthy, Lichan Hong, Ed Chi, and Derek Zhiyuan Cheng. 2023. Do LLMs Understand User Preferences? Evaluating LLMs On User Rating Prediction. *arXiv preprint arXiv:2305.06474* (2023).
- [18] Sami Khenissi and Olfa Nasraoui. 2020. Modeling and counteracting exposure bias in recommender systems. *arXiv preprint arXiv:2001.04832* (2020).
- [19] Heung-Nam Kim, Ae-Ttie Ji, Inay Ha, and Geun-Sik Jo. 2010. Collaborative filtering based on collaborative tagging for enhancing the quality of recommendation. *Electronic Commerce Research and Applications* 9, 1 (2010), 73–83.
- [20] Jeonghwan Kim, Giwon Hong, Sung-Hyon Myaeng, and Joyce Whang. 2023. FinePrompt: Unveiling the Role of Finetuned Inductive Bias on Compositional Reasoning in GPT-4. In *Findings of the Association for Computational Linguistics: EMNLP 2023*. 3763–3775.
- [21] Diederik P Kingma and Jimmy Ba. 2014. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980* (2014).
- [22] Huao Li, Yu Quan Chong, Simon Stepputtis, Joseph Campbell, Dana Hughes, Michael Lewis, and Katia Sycara. 2023. Theory of mind for multi-agent collaboration via large language models. *arXiv preprint arXiv:2310.10701* (2023).
- [23] Lei Li, Yongfeng Zhang, and Li Chen. 2023. Prompt distillation for efficient llm-based recommendation. In *Proceedings of the 32nd ACM International Conference on Information and Knowledge Management*. 1348–1357.
- [24] Roger Zhe Li, Julián Urbano, and Alan Hanjalic. 2021. Leave no user behind: Towards improving the utility of recommender systems for non-mainstream users. In *Proceedings of the 14th ACM International Conference on Web Search and Data Mining*. 103–111.
- [25] Timothy P Lillicrap, Jonathan J Hunt, Alexander Pritzel, Nicolas Heess, Tom Erez, Yuval Tassa, David Silver, and Daan Wierstra. 2015. Continuous control with deep reinforcement learning. *arXiv preprint arXiv:1509.02971* (2015).
- [26] Siyi Liu and Yujia Zheng. 2020. Long-tail session-based recommendation. In *Proceedings of the 14th ACM Conference on Recommender Systems*. 509–514.
- [27] Xu Liu, Tong Yu, Kaige Xie, Junda Wu, and Shuai Li. 2024. Interact with the Explanations: Causal Debaised Explainable Recommendation System. In *Proceedings of the 17th ACM International Conference on Web Search and Data Mining*. 472–481.
- [28] Yaokun Liu, Xiaowang Zhang, Minghui Zou, and Zhiyong Feng. 2023. Co-occurrence Embedding Enhancement for Long-tail Problem in Multi-Interest Recommendation. In *Proceedings of the 17th ACM Conference on Recommender*

- Systems. 820–825.
- [29] Andrew Luke, Joseph Johnson, and Yiu-Kai Ng. 2018. Recommending long-tail items using extended tripartite graphs. In *2018 IEEE International Conference on Big Knowledge (ICBK)*. IEEE, 123–130.
- [30] Sichun Luo, Chen Ma, Yuanzhang Xiao, and Linqi Song. 2023. Improving Long-Tail Item Recommendation with Graph Augmentation. In *Proceedings of the 32nd ACM International Conference on Information and Knowledge Management*. 1707–1716.
- [31] Tianhui Ma, Yuan Cheng, Hengshu Zhu, and Hui Xiong. 2023. Large Language Models are Not Stable Recommender Systems. *arXiv preprint arXiv:2312.15746* (2023).
- [32] Aditya Krishna Menon, Sadeep Jayasumana, Ankit Singh Rawat, Himanshu Jain, Andreas Veit, and Sanjiv Kumar. 2020. Long-tail learning via logit adjustment. *arXiv preprint arXiv:2007.07314* (2020).
- [33] Volodymyr Mnih, Koray Kavukcuoglu, David Silver, Alex Graves, Ioannis Antonoglou, Daan Wierstra, and Martin Riedmiller. 2013. Playing atari with deep reinforcement learning. *arXiv preprint arXiv:1312.5602* (2013).
- [34] Volodymyr Mnih, Koray Kavukcuoglu, David Silver, Andrei A Rusu, Joel Veness, Marc G Bellemare, Alex Graves, Martin Riedmiller, Andreas K Fidjeland, Georg Ostrovski, et al. 2015. Human-level control through deep reinforcement learning. *nature* 518, 7540 (2015), 529–533.
- [35] Vishvak Murahari, Prithvijit Chattopadhyay, Dhruv Batra, Devi Parikh, and Abhishek Das. 2019. Improving generative visual dialog by answering diverse questions. *arXiv preprint arXiv:1909.10470* (2019).
- [36] Jianmo Ni, Jiacheng Li, and Julian McAuley. 2019. Justifying recommendations using distantly-labeled reviews and fine-grained aspects. In *Proceedings of the 2019 conference on empirical methods in natural language processing and the 9th international joint conference on natural language processing (EMNLP-IJCNLP)*. 188–197.
- [37] Zohreh Ovaisi, Ragib Ahsan, Yifan Zhang, Kathryn Vasilaky, and Elena Zheleva. 2020. Correcting for selection bias in learning-to-rank systems. In *Proceedings of The Web Conference 2020*. 1863–1873.
- [38] Grigoris A Pavliotis. 2016. *Stochastic processes and applications*. Springer.
- [39] Antonin Raffin, Ashley Hill, Adam Gleave, Anssi Kanervisto, Maximilian Ernestus, and Noah Dormann. 2021. Stable-baselines3: Reliable reinforcement learning implementations. *The Journal of Machine Learning Research* 22, 1 (2021), 12348–12355.
- [40] Hossein A Rahmani, Mohammadmehdi Naghiaci, Mahdi Dehghan, and Mohammad Aliannejadi. 2022. Experiments on generalizability of user-oriented fairness in recommender systems. In *Proceedings of the 45th International ACM SIGIR Conference on Research and Development in Information Retrieval*. 2755–2764.
- [41] Xie Runfeng, Cui Xiangyang, Yan Zhou, Wang Xin, Xuan Zhanwei, Zhang Kai, et al. 2023. Lkpnr: Llm and kg for personalized news recommendation framework. *arXiv preprint arXiv:2308.12028* (2023).
- [42] Scott Sanner, Krisztian Balog, Filip Radlinski, Ben Wedin, and Lucas Dixon. 2023. Large language models are competitive near cold-start recommenders for language-and item-based preferences. In *Proceedings of the 17th ACM conference on recommender systems*. 890–896.
- [43] Tobias Schnabel, Adith Swaminathan, Ashudeep Singh, Navin Chandak, and Thorsten Joachims. 2016. Recommendations as treatments: Debiasing learning and evaluation. In *international conference on machine learning*. PMLR, 1670–1679.
- [44] John Schulman, Filip Wolski, Prafulla Dhariwal, Alec Radford, and Oleg Klimov. 2017. Proximal policy optimization algorithms. *arXiv preprint arXiv:1707.06347* (2017).
- [45] Rama Syamala Sreepada and Bidyut Kr Patra. 2020. Mitigating long tail effect in recommendations using few shot learning technique. *Expert Systems with Applications* 140 (2020), 112887.
- [46] Yiming Tan, Dehai Min, Yu Li, Wenbo Li, Nan Hu, Yongrui Chen, and Guilin Qi. 2023. Can ChatGPT Replace Traditional KBQA Models? An In-Depth Analysis of the Question Answering Performance of the GPT LLM Family. In *International Semantic Web Conference*. Springer, 348–367.
- [47] Shuai Tang and Xiaofeng Zhang. 2021. CADPP: An Effective Approach to Recommend Attentive and Diverse Long-tail Items. In *IEEE/WIC/ACM International Conference on Web Intelligence and Intelligent Agent Technology*. 218–225.
- [48] Jianing Wang, Qiusi Sun, Nuo Chen, Xiang Li, and Ming Gao. 2023. Boosting Language Models Reasoning with Chain-of-Knowledge Prompting. *arXiv preprint arXiv:2306.06427* (2023).
- [49] Jianing Wang, Junda Wu, Yupeng Hou, Yao Liu, Ming Gao, and Julian McAuley. 2024. InstructGraph: Boosting Large Language Models via Graph-centric Instruction Tuning and Preference Alignment. *arXiv preprint arXiv:2402.08785* (2024).
- [50] Ruoxi Wang, Bin Fu, Gang Fu, and Mingliang Wang. 2017. Deep & cross network for ad click predictions. In *Proceedings of the ADKDD'17*. 1–7.
- [51] Wenjie Wang, Yiyang Xu, Fuli Feng, Xinyu Lin, Xiangnan He, and Tat-Seng Chua. 2023. Diffusion Recommender Model. *arXiv preprint arXiv:2304.04971* (2023).
- [52] Xuanhui Wang, Michael Bendersky, Donald Metzler, and Marc Najork. 2016. Learning to rank with selection bias in personal search. In *Proceedings of the 39th International ACM SIGIR conference on Research and Development in Information Retrieval*. 115–124.
- [53] Yancheng Wang, Ziyang Jiang, Zheng Chen, Fan Yang, Yingxue Zhou, Eunah Cho, Xing Fan, Xiaojiang Huang, Yanbin Lu, and Yingzhen Yang. 2023. Recmind: Large language model powered agent for recommendation. *arXiv preprint arXiv:2308.14296* (2023).
- [54] Yu Wang, Zhiwei Liu, Jianguo Zhang, Weiran Yao, Shelby Heinecke, and Philip S Yu. 2023. DRDT: Dynamic Reflection with Divergent Thinking for LLM-based Sequential Recommendation. *arXiv preprint arXiv:2312.11336* (2023).
- [55] Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Fei Xia, Ed Chi, Quoc V Le, Denny Zhou, et al. 2022. Chain-of-thought prompting elicits reasoning in large language models. *Advances in Neural Information Processing Systems* 35 (2022), 24824–24837.
- [56] Tianxin Wei, Fuli Feng, Jiawei Chen, Ziwei Wu, Jinfeng Yi, and Xiangnan He. 2021. Model-agnostic counterfactual reasoning for eliminating popularity bias in recommender system. In *Proceedings of the 27th ACM SIGKDD Conference on Knowledge Discovery & Data Mining*. 1791–1800.
- [57] Wei Wei, Xubin Ren, Jiabin Tang, Qinyong Wang, Lixin Su, Suqi Cheng, Junfeng Wang, Dawei Yin, and Chao Huang. 2023. Llmrec: Large language models with graph augmentation for recommendation. *arXiv preprint arXiv:2311.00423* (2023).
- [58] Junda Wu, Zhihui Xie, Tong Yu, Handong Zhao, Ruiyi Zhang, and Shuai Li. 2022. Dynamics-aware adaptation for reinforcement learning based cross-domain interactive recommendation. In *Proceedings of the 45th International ACM SIGIR Conference on Research and Development in Information Retrieval*. 290–300.
- [59] Junda Wu, Tong Yu, and Shuai Li. 2021. Deconfounded and explainable interactive vision-language retrieval of complex scenes. In *Proceedings of the 29th ACM International Conference on Multimedia*. 2103–2111.
- [60] Yu Xia, Junda Wu, Tong Yu, Sungchul Kim, Ryan A Rossi, and Shuai Li. 2023. User-regulation deconfounded conversational recommender system with bandit feedback. In *Proceedings of the 29th ACM SIGKDD Conference on Knowledge Discovery and Data Mining*. 2694–2704.
- [61] Jun Xiao, Hao Ye, Xiangnan He, Hanwang Zhang, Fei Wu, and Tat-Seng Chua. 2017. Attentional factorization machines: Learning the weight of feature interactions via attention networks. *arXiv preprint arXiv:1708.04617* (2017).
- [62] Jing Yao, Wei Xu, Jianxun Lian, Xiting Wang, Xiaoyuan Yi, and Xing Xie. 2023. Knowledge Plugins: Enhancing Large Language Models for Domain-Specific Recommendations. *arXiv preprint arXiv:2311.10779* (2023).
- [63] Xinyang Yi, Ji Yang, Lichan Hong, Derek Zhiyuan Cheng, Lukasz Heldt, Aditee Kumthekar, Zhe Zhao, Li Wei, and Ed Chi. 2019. Sampling-bias-corrected neural modeling for large corpus item recommendations. In *Proceedings of the 13th ACM Conference on Recommender Systems*. 269–277.
- [64] Hongzhi Yin, Bin Cui, Jing Li, Junjie Yao, and Chen Chen. 2012. Challenging the long tail recommendation. *arXiv preprint arXiv:1205.6700* (2012).
- [65] Junchi Yu, Ran He, and Rex Ying. 2023. Thought propagation: An analogical approach to complex reasoning with large language models. *arXiv preprint arXiv:2310.03965* (2023).
- [66] Arlisa Yuliawati, Hamim Tohari, Rahmad Mahendra, and Indra Budi. 2022. On the Long Tail Products Recommendation using Tripartite Graph. *International Journal of Advanced Computer Science and Applications* 13, 1 (2022).
- [67] Fan Zhang and Qijie Shen. 2023. A Model-Agnostic Popularity Debias Training Framework for Click-Through Rate Prediction in Recommender System. In *Proceedings of the 46th International ACM SIGIR Conference on Research and Development in Information Retrieval*. 1760–1764.
- [68] Kaike Zhang, Qi Cao, Fei Sun, Yunfan Wu, Shuchang Tao, Huawei Shen, and Xueqi Cheng. 2023. Robust Recommender System: A Survey and Future Directions. *arXiv preprint arXiv:2309.02057* (2023).
- [69] Wenxuan Zhang, Hongzhi Liu, Yingpeng Du, Chen Zhu, Yang Song, Hengshu Zhu, and Zhonghai Wu. 2023. Bridging the Information Gap Between Domain-Specific Model and General LLM for Personalized Recommendation. *arXiv preprint arXiv:2311.03778* (2023).
- [70] Yin Zhang, Derek Zhiyuan Cheng, Tiansheng Yao, Xinyang Yi, Lichan Hong, and Ed H Chi. 2021. A model of two tales: Dual transfer learning framework for improved long-tail item recommendation. In *Proceedings of the web conference 2021*. 2220–2231.
- [71] Yang Zhang, Fuli Feng, Jizhi Zhang, Keqin Bao, Qifan Wang, and Xiangnan He. 2023. Collm: Integrating collaborative embeddings into large language models for recommendation. *arXiv preprint arXiv:2310.19488* (2023).
- [72] Yin Zhang, Ruoxi Wang, Derek Zhiyuan Cheng, Tiansheng Yao, Xinyang Yi, Lichan Hong, James Caverlee, and Ed H Chi. 2023. Empowering Long-tail Item Recommendation through Cross Decoupling Network (CDN). In *Proceedings of the 29th ACM SIGKDD Conference on Knowledge Discovery and Data Mining*. 5608–5617.
- [73] Zhuosheng Zhang, Aston Zhang, Mu Li, and Alex Smola. 2022. Automatic chain of thought prompting in large language models. *arXiv preprint arXiv:2210.03493* (2022).
- [74] Bowen Zheng, Yupeng Hou, Hongyu Lu, Yu Chen, Wayne Xin Zhao, and Ji-Rong Wen. 2023. Adapting large language models by integrating collaborative semantics for recommendation. *arXiv preprint arXiv:2311.09049* (2023).

- [75] Yu Zheng, Chen Gao, Xiang Li, Xiangnan He, Yong Li, and Depeng Jin. 2021. Disentangling user interest and conformity for recommendation with causal embedding. In *Proceedings of the Web Conference 2021*. 2980–2991.
- [76] Yong Zhuang, Tong Yu, Junda Wu, Shiqu Wu, and Shuai Li. 2022. Spatial-Temporal Aligned Multi-Agent Learning for Visual Dialog Systems. In *Proceedings of the 30th ACM International Conference on Multimedia*. 482–490.