

CORONA: A Coarse-to-Fine Framework for Graph-based Recommendation with Large Language Models

Junze Chen
Beijing University of Posts and
Telecommunications
Beijing, China
junze@bupt.edu.cn

Xinjie Yang
Beijing University of Posts and
Telecommunications
Beijing, China
yangxinjie@bupt.edu.cn

Cheng Yang*
Beijing University of Posts and
Telecommunications
Beijing, China
yangcheng@bupt.edu.cn

Junfei Bao
Beijing University of Posts and
Telecommunications
Beijing, China
bupt_baojunfei@bupt.edu.cn

Zeyuan Guo
Beijing University of Posts and
Telecommunications
Beijing, China
1154459434@bupt.edu.cn

Yawen Li
Beijing University of Posts and
Telecommunications
Beijing, China
warmly0716@126.com

Chuan Shi
Beijing University of Posts and
Telecommunications
Beijing, China
shichuan@bupt.edu.cn

Abstract

Recommender systems (RSs) are designed to retrieve candidate items a user might be interested in from a large pool, with a typical approach being the use of graph neural networks (GNNs) to capture high-order interaction relationships. As large language models (LLMs) have demonstrated remarkable success across various domains, researchers are exploring ways to apply their capabilities for improving recommendation performance. However, existing work limits the use of LLMs to either re-ranking recommendation results of traditional RSs or pre-processing the datasets as data augmenters. Both lines of work failed to explore LLMs' capabilities during the filtering process of candidate items, which may lead to suboptimal performance. Instead, we propose to leverage LLMs' reasoning abilities during the candidate filtering process, and introduce *Chain Of Retrieval ON graphs (CORONA)* to progressively narrow down the range of candidate items on interaction graphs with the help of LLMs: (1) First, LLM performs preference reasoning based on user profiles, with the response serving as a query to extract relevant users and items from the interaction graph as *preference-assisted retrieval*; (2) Then, using the information retrieved in the previous step along with the purchase history of target user, LLM conducts intent reasoning to help refine an even smaller interaction subgraph as *intent-assisted retrieval*; (3)

Finally, we employ a GNN to capture high-order collaborative filtering information from the extracted subgraph, performing *GNN-enhanced retrieval* to generate the final recommendation results. The proposed framework leverages the reasoning capabilities of LLMs during the retrieval process, while seamlessly integrating GNNs to enhance overall recommendation performance. Extensive experiments on various datasets and settings demonstrate that our proposed CORONA achieves state-of-the-art (SOTA) performance with an 18.6% relative improvement in recall and an 18.4% relative improvement in NDCG on average. Our code is available on GitHub at <https://github.com/BUPT-GAMMA/CORONA>.

CCS Concepts

• Computing methodologies → Machine learning; • Information systems → Recommender systems.

Keywords

Large Language Models, Graph Neural Networks, Recommendation

ACM Reference Format:

Junze Chen, Xinjie Yang, Cheng Yang, Junfei Bao, Zeyuan Guo, Yawen Li, and Chuan Shi. 2025. CORONA: A Coarse-to-Fine Framework for Graph-based Recommendation with Large Language Models. In *Proceedings of the 48th International ACM SIGIR Conference on Research and Development in Information Retrieval (SIGIR '25)*, July 13–18, 2025, Padua, Italy. ACM, New York, NY, USA, 11 pages. <https://doi.org/10.1145/3726302.3729937>

1 Introduction

Recommender systems aim to filter out candidate items that users might purchase [8, 9, 54]. For effective predictions, it is crucial to capture collaborative filtering (CF) relationships from massive user-item interactions. To this end, graph-based approaches [15, 47, 52, 56] typically construct a bipartite graph with historical user-item interactions, and employ graph neural networks to model

*Corresponding author.

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.
SIGIR '25, Padua, Italy.

© 2025 Copyright held by the owner/author(s). Publication rights licensed to ACM.
ACM ISBN 979-8-4007-1592-1/25/07
<https://doi.org/10.1145/3726302.3729937>

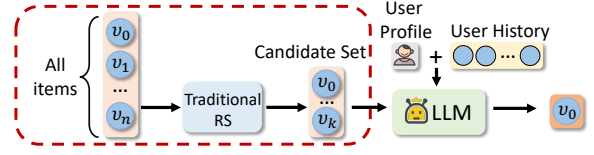
high-order relationships, showing satisfied recommendation performance [10, 57].

The rapid advancements of LLMs have showcased remarkable capabilities in generating, reasoning, and modeling world knowledge. Recommender systems are also anticipated to reap significant benefits from the development of LLMs. Recent efforts have leveraged LLMs for recommendation and achieved promising performance across various scenarios [1, 7, 11, 17, 17, 21, 26, 38, 46, 48, 60, 62, 66, 69]. Existing approaches can be categorized into two types: (1) *Applying LLMs after candidate filtering*. These methods typically integrate relevant information (e.g., item candidates) into natural language-based prompts, and derive recommendations from LLM-generated responses. Due to the context length limitations of LLMs, these methods need traditional RSs to generate candidate item sets as inputs. Early attempts in this category relied on in-context learning (ICL) with frozen LLMs [11, 17], while later methods began to explore different fine-tuning strategies for specialization [1, 13, 64]. (2) *Applying LLMs before candidate filtering*. These methods employ LLMs to enrich user/item attributes along with interactions for pre-processing, and then the augmented dataset will be fed into traditional RSs for candidate filtering [37, 38, 51]. However, as shown in Figure 1, both types of methods failed to explore LLMs' capabilities during the candidate filtering process itself, which may lead to suboptimal performance.

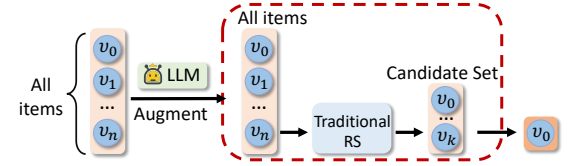
In this paper, we propose to leverage LLMs' reasoning abilities during the candidate filtering process, and introduce a coarse-to-fine framework named Chain Of Retrieval ON grAphs (CORONA) for recommendation. As shown in Figure 1c, CORONA progressively narrows the search space on the user-item interaction graph to improve recommendation accuracy with the help of LLMs. Note that the strength of LLMs lies in their extensive general knowledge and world understanding [3, 35, 53, 65]. Therefore, we leverage LLMs for coarse-grained preference and intent reasoning, rather than relying on them for fine-grained selection from similar candidate items.

Specifically, CORONA includes three stages of retrieval at different granularities: (1) We first use the user profile as input to an LLM for preference reasoning, generating a query embedding to perform *preference-assisted retrieval*, which extracts a subgraph aligned with the user's general preferences from the interaction graph. (2) We then combine the user's purchase history with statistical information from the previous subgraph, prompting the LLM for intent reasoning. The resulting query embedding supports *intent-assisted retrieval*, refining the subgraph to reflect more personalized and short-term user intent. (3) Finally, we apply *GNN-enhanced retrieval*, where the GNN processes the subgraph to capture valuable relationships, producing the final recommendation results. In this way, we can leverage the strong reasoning capabilities of LLMs in the item retrieval process, narrowing down relevant items in the entire dataset. Additionally, our framework seamlessly integrates with traditional GNN models to efficiently capture collaborative filtering information, enhancing overall recommendation performance. Extensive experiments on three datasets show that, on average, our model achieves an 18.6% relative improvement in recall and an 18.4% relative improvement in NDCG compared to the best baseline, highlighting the effectiveness of this framework.

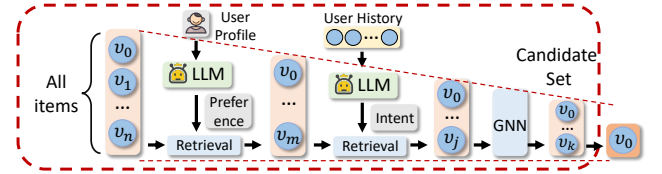
The contribution of this work are three-fold:



(a) Previous methods that apply LLMs after candidate filtering. These methods typically rely on traditional RSs to obtain candidate items, and use LLMs as the final predictors (e.g., TALLRec [1]).



(b) Previous methods that apply LLMs before candidate filtering. These methods typically use LLMs for data augmentation, and then rely on traditional RSs to obtain candidate items (e.g., LLMRec [51]).



(c) Our proposed framework using LLMs for coarse-grained retrieval and GNNs for fine-grained recommendation, directly leveraging LLMs to assist candidate filtering.

Figure 1: Comparisons between previous methods and our proposed coarse-to-fine framework, with the candidate filtering process highlighted by the red dashed box. We leverage LLMs for coarse-grained preference and intent reasoning instead of data augmentation or fine-grained item selection. Our framework progressively narrows the range of candidate items and integrates GNNs for improved performance.

- We introduce a novel framework for graph-based recommendation, utilizing LLMs to assist coarse-grained retrieval, followed by traditional CF methods for fine-grained recommendation. This allows LLMs to capitalize on their strengths in reasoning and directly involve in the candidate filtering process.

- We propose CORONA, a carefully designed three-stage retrieval framework that progressively refines the retrieval process at different levels of granularity. CORONA leverages both the strong capability of LLMs for preference and intent reasoning, as well as typical GNNs for efficient recommendation.

- Extensive experiments on three publicly available datasets demonstrate the effectiveness of our approach for recommendation, with ablation studies validating the necessity of each module.

2 Related Work

2.1 Graph-based Recommendation

Collaborative Filtering is a foundational approach in recommendation systems and has been the subject of extensive research [4]. A

growing trend in the field is to model user-item interactions as a bipartite graph and apply GNNs to capture high-order collaborative relationships, such as NGCF [47], LightGCN [15], GraphPro [56], GRCN [52], PUP [68] and IRGPR [31]. Specifically, NGCF models the connections between users and items by propagating node embeddings across the graph. LightGCN improves efficiency by eliminating redundant components from the message-passing process in graph-based recommendations. GraphPro incorporates both a temporal prompt mechanism and graph-structural prompt learning into its pre-trained GNN architecture. GRCN dynamically refines the interaction graph structure in response to the model's training progress. PUP designs an encoder with GCN on a predefined heterogeneous graph to capture price awareness, investigating the influence of price feature in ranking stage and enhancing performance. IRGPR proposes a heterogeneous graph to fuse the two information sources, one item relation graph to capture multiple item relationships and one user-item scoring graph to include the initial ranking scores, accomplishing personalized re-ranking with the help of GNN. In this work, we extend these approaches with LLMs' world knowledge and reasoning ability, thereby improving recommendation performance.

2.2 LLM-enhanced Recommendation

Recently, using LLMs to enhance recommendation systems has gained significant attention [1, 5, 12, 17, 18, 21, 23, 26–30, 32, 33, 38–42, 45, 49, 53, 55, 59–61, 67, 69]. Existing works can be categorized into the following two types.

2.2.1 Applying LLMs after candidate filtering. Some works (e.g. [11, 17, 21, 48, 69]) align recommendation tasks with natural language, and directly prompt LLMs to generate final recommendation results based on a candidate set generated by traditional RSs. Initially, several approaches use pre-trained LLMs and leverage in-context learning capabilities to perform recommendation. For example, Chat-REC [11] constructs a conversational recommender by converting user profiles and interactions into prompts, allowing LLMs to generate recommendations. LLMRank [17] provides demonstration examples by enhancing the input interaction sequence directly. Other methods (e.g. [1, 60, 63, 64]) tailor LLMs for recommendation by tuning them with recommendation data. For example, P5 [13] transforms user interaction data into textual prompts based on item indexes, which are then utilized for language model training. InstructRec [60] and TALLRec [1] use instructional designs to define recommendation tasks and fine-tune LLMs to follow these instructions for generating recommendations. However, these methods heavily depend on the quality of candidate sets generated by traditional approaches [51].

2.2.2 Applying LLMs before candidate filtering. Some recent work leverages LLMs for data augmentation as a pre-processing step. Specifically, RLMRec [37] proposes a model-agnostic approach to enhance existing recommenders with LLM generated user profiles. LLMRec [51] uses rich online content, including image and text, to augment the interaction graph. In fact, these methods are compatible with our framework for data pre-processing.

We would like to highlight the key difference between these methods and our work: (1) These methods employ LLMs for data

augmentation rather than in the retrieval process, making them incapable of handling dynamically changing user intents with LLMs. (2) In contrast, our approach integrates LLM reasoning into the candidate filtering process, allowing instructions to be easily modified to guide recommendations under different contexts.

3 Methodology

3.1 Problem Statement

3.1.1 Notations. We focus on graph-based recommendation with textual information [37], where the interaction graph consists of two types of nodes, *i.e.*, users and items. Each node is associated with both textual descriptions and numerical features, and the edges represent purchase records. Formally, we denote the interaction graph as $\mathcal{G} = (\mathcal{U}, \mathcal{V}, \mathbf{A})$, where \mathcal{U} is the set of users, \mathcal{V} is the set of items and $\mathbf{A} \in \{0, 1\}^{(|\mathcal{U}|+|\mathcal{V}|) \times (|\mathcal{U}|+|\mathcal{V}|)}$ represents the adjacency matrix. Each user $u \in \mathcal{U}$ has an interaction history L_u as a list of items sorted by interaction time. Both users and items have corresponding textual information, represented as $P_{\mathcal{U}}$ and $T_{\mathcal{V}}$, respectively. In addition, feature vectors of users and items can be respectively derived as $F_{\mathcal{U}}$ and $M_{\mathcal{V}}$.

3.1.2 Problem Definition. We follow the settings of recent advances in LLM-augmented recommendation [37, 51] and focus on Top-N Recommendation. Formally, given the interaction graph \mathcal{G} and associate information (including $P_{\mathcal{U}}, T_{\mathcal{V}}, F_{\mathcal{U}}$ and $M_{\mathcal{V}}$), the goal is to predict the next item v that user u will purchase based on interaction history $L_u = \{v_1, v_2, \dots, v_n\}$. Training set $\mathcal{Z}^{\text{train}}$ and test set $\mathcal{Z}^{\text{test}}$ are both sets of interactions. The adjacency matrix \mathbf{A} ensures that the interactions selected as training and test samples are excluded.

3.2 Framework Overview

The framework of the proposed CORONA is shown in Figure 2. We conduct three stages of retrieval from different granularities. In the first stage, we take the target user's profile as input, and employ an LLM to infer the user's potential preferences. These preferences are then transformed into query embeddings, and help extract the subgraph from the interaction graph as preference-assisted retrieval. The resulting subgraph is further summarized into natural language. In the second stage, the summary is combined with the user's purchase history for LLM-based intent reasoning. Similarly, we will narrow the scope based on the previous retrieval by extracting a smaller subgraph as intent-assisted retrieval. In the third stage, the final subgraph will be passed to a GNN for message passing, and items are scored by the inner products of user/item embeddings for recommendation.

3.3 Preference-assisted Retrieval

In the first stage, we use an LLM to reason about user preferences based on profiles, helping to extract items and users that align with the target user's interests. The process of preference-assisted retrieval can be broken down into three steps.

3.3.1 Preference Reasoning. First, we use an LLM to generate preference reasoning content based on the user's profile, which will further serves as the query embedding. Specifically, for a target user u , the user's profile information P_u includes details such as age,

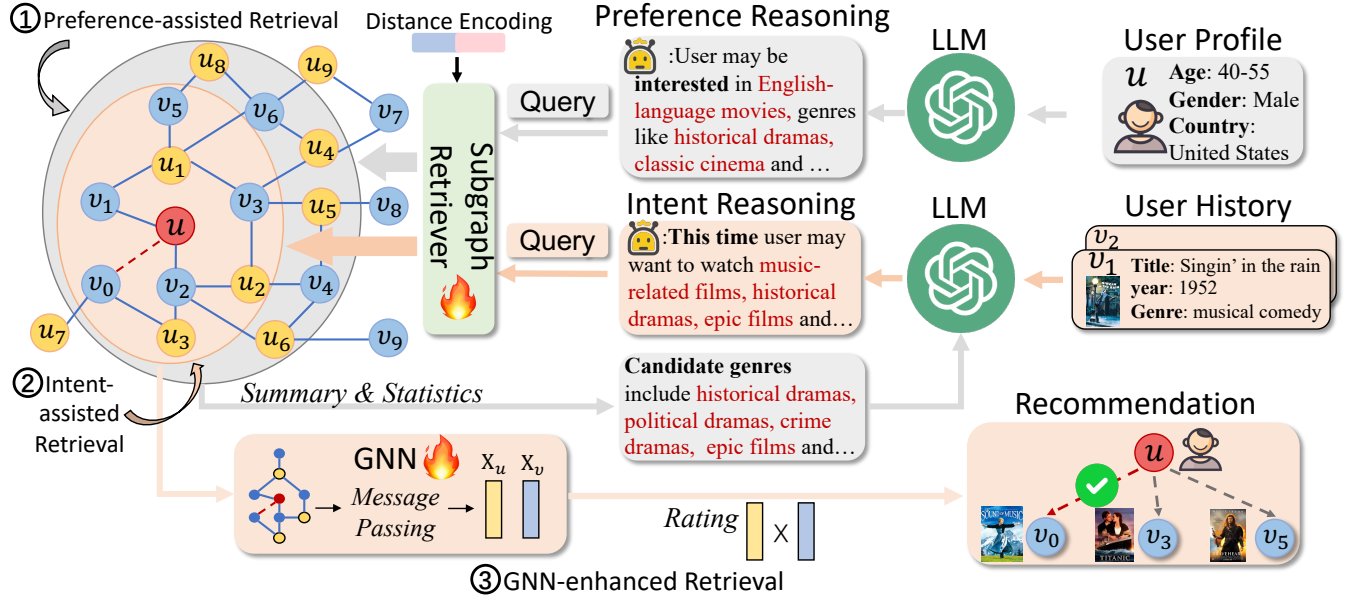


Figure 2: The overall framework of our proposed CORONA, which includes three stages of retrieval from different granularities. Here an LLM is employed to perform coarse-grained preference/intent reasoning based on user profile/history in the first two stages, helping progressively extract a subgraph of relevant users and items from the interaction graph. The final subgraph will be encoded by a GNN for fine-grained recommendation.

gender, country, and other relevant attributes. We assemble these information into a natural language instruction to input into the LLM. The LLM used for preference reasoning can be represented as $LLM_{PR}(\cdot)$, it is expected to utilize its reasoning ability and stored knowledge to infer the user's preference information and output it in natural language form. For instance, in a movie recommendation scenario, if the user's profile P_u indicates she is from France, the LLM should infer that the user might prefer French-language movies in the item set. The instruction for preference reasoning is shown below:

Module: Preference Reasoning

Instruction: Please infer the user's movie preferences based on the provided user profile. This may include preferred categories, styles, origins, or years of movies. Begin by summarizing the relevant information, and then provide your preference reasoning predictions using the user profile details along with your general knowledge.

User profile: Age: 25; Gender: Male; Country: US; Language: English; Occupation: Master student...

Besides the inferred preference, the LLM's responses may include explanations and reasoning steps, which we believe also reflect the user's preferences in some way. Therefore, we treat all generated content as the preference reasoning result $Q_1 = LLM_{PR}(P_u)$. Next, Q_1 will be encoded to form the query embedding $E_{Q_1} = \text{Encode}(Q_1)$ for subgraph retrieval.

3.3.2 Subgraph Retriever. To retrieve the most relevant items based on the query, we design a simple subgraph retriever similar to the

attention mechanism. Following the classic collaborative filtering principle [24], we believe that identifying users with similar interests to the target user will help retrieve more relevant information. While in the interaction graph, users who are closer to each other tend to have more similar interests due to the similarity in their interaction histories. Therefore, we base the retrieval process on similar users, and focus on those closer to the target user u .

Specifically, we first concatenate the distance encoding $\{e_j\}_{j=1}^3$ to the users' features based on their distances from the target user u . We use the number of hops between target user u and user u' to represent the distance $\text{dist}(u', u)$ between them. If users u and u' have interacted with the same item, they are considered one-hop neighbors, and $\text{dist}(u', u) = 1$. Similarly, a user who is a one-hop neighbor of a one-hop neighbor of u is a two-hop neighbor, with $\text{dist}(u', u) = 2$, and so on.

Our data analysis revealed that very few users in $\mathcal{Z}_{\text{train}}$ pay attention to items interacted by users beyond two-hop neighbors: only 3% of users on Netflix [2] and 0.4% on MovieLens [14]. Therefore, users beyond two hops are uniformly assigned e_3 , aiming to reduce the attention given to these distant users. A linear layer Linear_θ with trainable parameters θ is then employed to fuse user embeddings with the distance encoding information and reduce dimensionality. Specifically, with target user u , each $u' \in \mathcal{U}$ can get user embedding $X_{u'}$ as follows:

$$X_{u'} = \begin{cases} \text{Linear}_\theta(\text{CONCAT}(F_{u'}, e_1)) & \text{if } \text{dist}(u', u) = 1 \\ \text{Linear}_\theta(\text{CONCAT}(F_{u'}, e_2)) & \text{if } \text{dist}(u', u) = 2 \\ \text{Linear}_\theta(\text{CONCAT}(F_{u'}, e_3)) & \text{otherwise} \end{cases} \quad (1)$$

Finally, we have user embeddings with encoded distance information $X_{\mathcal{U} \setminus \{u\}}$ for all users \mathcal{U} except target user u .

Based on the query generated by the LLM, we use cosine similarity to find the top- k most similar users \mathcal{U}'_1 . Here the target user u is also included. This process is formalized as:

$$\mathcal{U}'_1 = \text{argtopk}_{\{u' \in \mathcal{U} \setminus \{u\}\}} \cos(E_{Q_1}, X_{u'}) \cup \{u\}. \quad (2)$$

To ensure the extracted subgraph includes the most relevant items, we add all items \mathcal{V}'_1 connected to \mathcal{U}'_1 into the subgraph.

3.3.3 Summary & Statistics. To assist the intent reasoning in next stage, we summarize the information of the retrieved items \mathcal{V}'_1 into natural language as auxiliary information. Specifically, we collect the textual information of items in \mathcal{V}'_1 and then conduct statistics. We include the top 20 most frequent attributes (e.g., movie genres and categories) in the summary as $\text{Summary}(T_{\mathcal{V}'_1})$. The example of text template for summarization and statistics is as follows:

Module: Summary & Statistics

Candidate Genres: The candidate items belong to the following genres: *[historical drama, political drama, epic film, music-related film, crime drama...]*

Candidate Categories: The items can be categorized into: *[independent film, documentary, animation ...]*

Years of release: The movies were released in the following periods: *[1930-1950, 1960-1970, ...]*

3.4 Intent-assisted Retrieval

User intent can be viewed as a more fine-grained and short-term interest compared to preference [25]. Therefore, we use the target user's interaction history $L_u = \{v_1, v_2, \dots, v_n\}$ as the basis for inferring her intent, and perform the second round of subgraph retrieval.

Specifically, given a target user u , we extract the textual information T_{L_u} of the items in u 's interaction history, and input them together with $\text{Summary}(T_{\mathcal{V}'_1})$ into the LLM for intent reasoning.

The intent reasoner can be denoted as $\text{LLM}_{IR}(\cdot, \cdot)$, and LLM's response $Q_2 = \text{LLM}_{IR}(\text{Summary}(T_{\mathcal{V}'_1}), T_{L_u})$. Similar to the previous step, we encode the response content as query embedding $E_{Q_2} = \text{Encode}(Q_2)$ for intent-assisted retrieval. We use cosine similarity to select the top- $\frac{k}{2}$ most relevant users as \mathcal{U}'_2 , and include the target user in the user set. This process can be formalized as:

$$\mathcal{U}'_2 = \text{argtop}_{\frac{k}{2}} \cos(E_{Q_2}, X_{u'}) \cup \{u\}. \quad (3)$$

Along with the items \mathcal{V}'_2 connected to \mathcal{U}'_2 , we have the final retrieved subgraph $\mathcal{G}' = (\mathcal{U}'_2, \mathcal{V}'_2, A[\mathcal{U}'_2, \mathcal{V}'_2])$. The instruction example used for intent reasoning is shown below:

Module: Intent Reasoning

Instruction: Please infer the user's watching intent based on user history and candidate information. The watching intent represents the user's intention for this viewing, and you should make your selection from within the candidate range. Begin by summarizing the relevant information, then provide your intent reasoning predictions using the user history details, candidate sets, and your general knowledge.

User history: No.1: Title: The Last Emperor; Year: 1987; Genre: History drama; No.2: Singin' in the rain...

Candidate summary: Candidate Genres: ...; Candidate Categories: ...; Years of release: ...

To learn the parameters in subgraph retriever (i.e., θ and $\{e_j\}_{j=1}^3$), we regard the users interacted with ground truth item v as true users, and use the following training loss:

$$\mathcal{L}_1 = - \sum_{u' \in \mathcal{N}_v} \left(\frac{\exp(E_{Q_1}^\top \cdot X_{u'})}{\sum_{u'' \in \mathcal{U}} \exp(E_{Q_1}^\top \cdot X_{u''})} + \frac{\exp(E_{Q_2}^\top \cdot X_{u'})}{\sum_{u'' \in \mathcal{U}} \exp(E_{Q_2}^\top \cdot X_{u''})} \right), \quad (4)$$

where \mathcal{N}_v is the set of users connected with item v .

3.5 GNN-enhanced Retrieval

Following previous graph-based recommendation methods [15, 51], we further introduce a GNN to capture the high-order relations within the retrieved subgraph \mathcal{G}' . Formally, a GNN with parameters ϕ encodes the target user u based on the extracted subgraph \mathcal{G}' as H_u . Then we score each item in \mathcal{G}' by the inner product between user embeddings and item features:

$$\text{score}(u, v) = H_u^\top \cdot M_v. \quad (5)$$

Following previous work [15], we use the classical Bayesian Personalized Ranking (BPR) loss for training the GNN. The BPR loss can be calculate by:

$$\mathcal{L}_2 = - \sum_{v' \in \mathcal{V}_{\text{neg}}} \log \sigma(\text{score}(u, v) - \text{score}(u, v')), \quad (6)$$

where \mathcal{V}_{neg} is the set of negative items randomly chosen from \mathcal{V}'_2 . Finally, we select items with the highest scores for top- n recommendation task.

4 Experiments

To validate the effectiveness of our proposed CORONA, we conduct extensive experiments to answer the following research questions (**RQs**): **RQ1:** How effective is our proposed CORONA compared to the state-of-the-art baselines? **RQ2:** Is the proposed CORONA also effective for recommending items in cold-start setting? **RQ3:** Has each component of our framework played its role effectively? **RQ4:** How do different values of key parameters influence the method's performance? **RQ5:** How efficient is CORONA compared with previous methods? **RQ6:** Does the LLM-based reasoning process in CORONA offer some interpretability?

4.1 Experimental Setup

4.1.1 Datasets. We perform experiments on three publicly available datasets, *i.e.*, Netflix¹, MovieLens² and Amazon-book³. For baselines that cannot directly utilize textual information, such as LightGCN [15], we use text encodings as node features. For the Netflix and MovieLens datasets, we use the same split as in LLM-Rec [51]; for the Amazon-book dataset, we follow the split from RLMRec [37].

- Netflix dataset [2] is sourced from the Kaggle website. We use BERT [36] to encode the textual information of users and items, obtaining user features F_U and item features M_V , respectively.

- MovieLens dataset [14] is sourced from ML-10M. We encode these textual information using BERT [36] as features F_U and M_V .

- Amazon-book dataset [34] contains book review records from 2000 to 2014. Information are encoded by BERT [36] to obtain features F_U and M_V .

4.1.2 Evaluation Protocols. Following previous work [15, 50, 51], we evaluate our approach in the top- K item recommendation task using two common metrics: Recall ($R@K$) and Normalized Discounted Cumulative Gain ($N@K$), where K is set to 10, 20, and 50. We employ the all-ranking strategy, and report averaged results from five independent runs.

4.1.3 Methods for Comparison. To fully demonstrate the effectiveness of our proposed CORONA, we compare a number of baselines from three groups. (1) Graph-based Collaborative Filtering Methods: These approaches leverage GNNs to capture the structural relationships in the interaction graph, including NGCF [47], LightGCN [15], GraphPro [56] and GRCN [52]. (2) LLM for Recommendation: These methods apply LLMs to recommendation tasks to improve performance metrics, including TALLRec [1], BinLLM [63], RLMRec [37], LLMRec [51]. We also compare CORONA with an alternative LLM retrieval method: G-retriever [16], following its original approach described in the paper to perform subgraph construction and then recommend items based on similarity scoring.

For the experiments under item cold-start setting, we select several methods specifically designed for cold-start scenarios along with the “LLM for Recommendation” methods as baselines. The selected baselines include DropoutNet [44], ALDI [19], TALLRec [1], BinLLM [63], RLMRec [37], LLMRec [51], LLM-Ins [20].

4.1.4 Implementation Details. We follow existing methods [51] to obtain the textual attribute of users P_U and items T_V , and derive 128-dimensional features for both users F_U and items M_V as well. We employ OpenAI’s “GPT-4o-mini” for preference and intent reasoning and the responses are encoded using OpenAI’s “text-embedding-ada-002,” also with 128-dimensional outputs. The distance encoding $\{e\}_{j=1}^3$ is set to 2 dimensions. The linear layer Linear_θ use a 130×128 -dimensional linear layer. For the main experiment with the GCN method, we use a two-layer GCN with hidden dimension set at 128. GraphTransformer is used as provided in the original paper [58]. We set the size of the negative set to 10 for Eq 6, and employ the Adam optimizer with a learning rate of $1e-6$ for parameter training. Early stopping with a patience setting of 10

steps is also used during training. All experiments are conducted on an A800 GPU with 80GB of memory.

4.2 Main Results (RQ1)

To answer **RQ1**, we conduct recommendation experiments with results shown in Table 1. From the results, we can see that: (1) Our method CORONA with LLM-empowered reasoning consistently outperforms the state-of-the-art baselines on all three datasets. On average, CORONA has 17.66% relative improvement in recall and 16.06% in NDCG compared to the best baseline. This improvement showcases the effectiveness of our framework. (2) The LLM-based recommendation methods are the most competitive baselines, as they also leverage LLMs’ capabilities to enhance recommendation. But our approach still show a relative improvement of 30.67% in recall and 31.24% in NDCG on average.

4.3 Cold-start Results (RQ2)

To answer **RQ2**, we focus on items with no more than two interactions, and set up an item cold-start scenario for evaluation. From the results shown in Table 2, we can observe that: (1) Our method outperforms all baselines, with an average relative improvement of 8.24% in recall and 10.68% in NDCG. This demonstrates that our framework remains high utility in cold-start scenarios by leveraging the reasoning capabilities of LLMs to mitigate data sparsity. (2) Note that LLM-Ins [20] is specialized for only cold-start settings, and performs the best among baselines. While our approach can outperform LLM-Ins across all datasets and metrics, showcasing the value of our framework. (3) The performance of traditional cold-start methods, namely DropoutNet and ALDI, is relatively low, highlighting the advantage of using LLMs to process textual data. Our method relatively improves recall by 35.51% over traditional methods and 10.39% over other LLM baselines, and NDCG by 37.48% over traditional methods and 11.62% over other LLM baselines.

4.4 Ablation Study (RQ3)

To answer **RQ3**, we consider three categories of ablated variants to test whether each design of CORONA is effective and necessary.

4.4.1 Combinations of Different Subgraph Retrieval Methods and GNNs. To show the effectiveness of our proposed LLM-based subgraph retrieval, we consider three other methods for extracting relevant subgraphs: full graph, 1-hop neighbors of target user (fixed 1-hop) and 2-hop neighbors of target user (fixed 2-hop). Also, to show that our CORONA is compatible with different GNNs, we test 2-layer GCN [22] and GraphTransformer (GT) [58] for recommendation based on extracted subgraphs. We present the results in Table 3. From the results, we observe that: among the GCN-based and GT-based methods, our LLM-based subgraph retrieval demonstrates a clear advantage over others. While the same GNN is applied, our method relatively yields a higher recall by 110.73% and a higher NDCG by 112.03% on average. This highlights that our designed retrieval process significantly contributes to more accurate recommendations. Besides, GT-based CORONA yields competitive performance with GCN-based one, showing the compatibility of our proposed framework.

¹<https://www.kaggle.com/datasets/netflix-inc/netflix-prize-data>

²<https://files.grouplens.org/datasets/movielens/ml-10m-README.html>

³<https://cseweb.ucsd.edu/jmcauley/datasets/amazon/links.html>

Table 1: Recommendation performance on three datasets in terms of $Recall@10/20/50$, and $NDCG@10/20/50$.

Datasets	Netflix						MovieLens						Amazon-book					
Methods	R@10	N@10	R@20	N@20	R@50	N@50	R@10	N@10	R@20	N@20	R@50	N@50	R@10	N@10	R@20	N@20	R@50	N@50
NGCF	0.0357	0.0163	0.0693	0.0231	0.1090	0.0335	0.1153	0.1007	0.1305	0.1164	0.1636	0.1261	0.0844	0.0637	0.1294	0.0783	0.2151	0.1023
LightGCN	0.0385	0.0152	0.0661	0.0222	0.1252	0.0336	0.0966	0.1013	0.1314	0.1173	0.1647	0.1297	0.0857	0.0646	0.1303	0.0791	0.2154	0.1028
GraphPro	0.0397	0.0174	0.0701	0.0235	0.1263	0.0350	0.1376	0.1121	0.1462	0.1296	0.1757	0.1411	0.0865	0.0656	0.1345	0.0839	0.2203	0.1085
GRCN	0.0392	0.0169	0.0695	0.0260	0.1318	0.0348	0.1175	0.1017	0.1356	0.1192	0.1654	0.1308	0.0862	0.0653	0.1332	0.0826	0.2179	0.1067
TALLRec	0.0498	0.0251	0.0795	0.0338	0.1303	0.0420	0.2521	0.1104	0.3522	0.1748	0.5017	0.1921	0.0857	0.0729	0.1277	0.0838	0.2420	0.1169
BinLLM	0.0502	0.0254	0.0797	0.0340	0.1312	0.0431	0.2584	0.1187	0.3657	0.1771	0.5094	0.1995	<u>0.0972</u>	<u>0.0741</u>	<u>0.1512</u>	<u>0.0939</u>	<u>0.2451</u>	<u>0.1194</u>
RLMRec	0.0504	0.0257	0.0806	0.0341	0.1314	0.0439	<u>0.2595</u>	<u>0.1195</u>	<u>0.3668</u>	<u>0.1779</u>	<u>0.5108</u>	<u>0.2001</u>	0.0968	0.0738	0.1499	0.0913	0.2432	0.1175
LLMRec	<u>0.0526</u>	<u>0.0271</u>	<u>0.0808</u>	<u>0.0342</u>	<u>0.1317</u>	<u>0.0442</u>	0.2582	0.1187	0.3528	0.1750	0.5050	0.1904	0.0862	0.0733	0.1285	0.0840	0.2424	0.1171
G-Retriever	0.0327	0.0134	0.0678	0.0213	0.1106	0.0323	0.1179	0.1052	0.1294	0.1121	0.1569	0.1225	0.0831	0.0616	0.1248	0.0719	0.2123	0.0969
CORONA	0.0616	0.0279	0.0938	0.0416	0.1452	0.0487	0.3033	0.1565	0.4214	0.2017	0.5745	0.2507	0.1206	0.0855	0.1857	0.1089	0.3048	0.1299
Improv.	17.11%	2.95%	16.09%	21.64%	10.25%	10.18%	16.88%	30.96%	14.89%	13.38%	12.47%	25.29%	24.07%	15.38%	22.82%	15.97%	24.36%	8.79%

Table 2: Recommendation performance under item cold-start setting in terms of $Recall@10/20/50$, and $NDCG@10/20/50$.

Datasets	Netflix						MovieLens						Amazon-book					
Methods	R@10	N@10	R@20	N@20	R@50	N@50	R@10	N@10	R@20	N@20	R@50	N@50	R@10	N@10	R@20	N@20	R@50	N@50
DropoutNet	0.0175	0.0083	0.0323	0.0157	0.0414	0.0195	0.1352	0.0538	0.2107	0.1101	0.3622	0.1495	0.0425	0.0264	0.0832	0.0496	0.1482	0.1007
ALDI	0.0247	0.0099	0.0409	0.0161	0.0527	0.0202	0.1436	0.0552	0.2244	0.1102	0.3613	0.1488	0.0525	0.0307	0.0898	0.0517	0.1511	0.1026
TALLRec	0.0124	0.0075	0.0301	0.0152	0.0322	0.0181	0.1224	0.0525	0.2083	0.1099	0.3597	0.1475	0.0459	0.0278	0.0732	0.0472	0.1368	0.1004
BinLLM	0.0218	0.0086	0.0395	0.0158	0.0513	0.0200	0.01373	0.0547	0.2159	0.1095	0.3638	0.1477	0.0514	0.0296	0.0873	0.0508	0.1485	0.1012
RLMRec	0.0243	0.0096	0.0412	0.0164	0.0535	0.0207	0.1576	0.0588	0.2473	0.1105	0.3626	0.1501	0.0571	0.0319	0.0907	0.0526	0.1524	0.1035
LLMRec	0.0119	0.0073	0.0296	0.0142	0.0315	0.0175	0.1163	0.0517	0.2052	0.1095	0.3594	0.1473	0.0454	0.0273	0.0725	0.0464	0.1357	0.0971
LLM-Ins	<u>0.0281</u>	<u>0.0126</u>	<u>0.0490</u>	<u>0.0182</u>	<u>0.0739</u>	<u>0.0195</u>	<u>0.1826</u>	<u>0.0854</u>	<u>0.2689</u>	<u>0.1372</u>	<u>0.4063</u>	<u>0.1665</u>	<u>0.0875</u>	<u>0.0417</u>	<u>0.1167</u>	<u>0.0704</u>	<u>0.1645</u>	<u>0.1206</u>
CORONA	0.0295	0.0130	0.0519	0.0201	0.0791	0.0217	0.2054	0.0933	0.2917	0.1425	0.4582	0.1931	0.0916	0.0453	0.1212	0.0764	0.1788	0.1314
Improv.	4.98%	3.17%	5.91%	10.44%	7.04%	11.28%	12.49%	9.25%	8.48%	3.86%	12.77%	15.98%	4.69%	8.63%	3.86%	8.52%	8.69%	8.96%

4.4.2 Ablation of Different Components. We ablate each component of CORONA individually to test the effectiveness of the design. The results are shown in Table 4. For “w/o Preference-assisted Retrieval” and “w/o Intent-assisted Retrieval”, we remove the preference-assisted and intent-assisted retrieval components respectively, and use the remaining LLM-assisted retrieval combined with GNN-enhanced retrieval for recommendation. For “w/o GNN-enhanced Retrieval”, we remove the GNN module GNN_ϕ that updates the embedding of target user by message passing. For “w/o Preference Reasoning”, we remove the preference reasoner LLM_{PR} and directly encode user profile to form the query embedding E_{Q_1} . Similarly, for “w/o Intent Reasoning”, we encode textual user history to form query embedding E_{Q_2} .

From the results, we observe that: (1) The full model CORONA always yields the best performance, demonstrating the soundness of our design and the necessity of each component. Compared to the average of variants, the full framework achieves a relative improvement of 20.21% in recall and 19.52% in NDCG. (2) The “w/o Preference Reasoning” and “w/o Intent Reasoning” variants have the weakest performance, highlighting the importance of LLM reasoning in our framework. Note that the above two variants have even worse performance than the “w/o Preference-assisted Retrieval” and “w/o Intent-assisted Retrieval” variants, indicating that LLM’s reasoning abilities and general knowledge indeed contribute effectively to the recommendation task. (3) The variants without intent

reasoning/assistance perform worse than those without preference reasoning/assistance. This observation highlights the importance of using LLMs to handle dynamic user intents, whereas previous methods using LLMs for pre-processing are unable to infer such dynamic information effectively.

4.4.3 Influence of Different LLMs. We also investigate the impact of using different LLMs on model performance, with results shown in Figure 3. Besides GPT-4o-mini used in main experiments, we also test CORONA with a small open-source LLM, *i.e.*, Vicuna-7B-v1.5 [6]. Vicuna is an open-source chatbot developed by fine-tuning LLaMA [43] with conversations shared by users. We deploy the Vicuna-7B-v1.5 model locally for preference and intent reasoning, using the same “text-embedding-ada-002” encoder to ensure the only difference lies in the reasoning content. The results show that CORONA with Vicuna-7B-v1.5 is already slightly better than previous methods, while CORONA with GPT-4o-mini further yields an average improvement of 15.87%. This demonstrates that higher-quality reasoning leads to more significant performance gains.

4.5 Hyperparameter Analysis (RQ4)

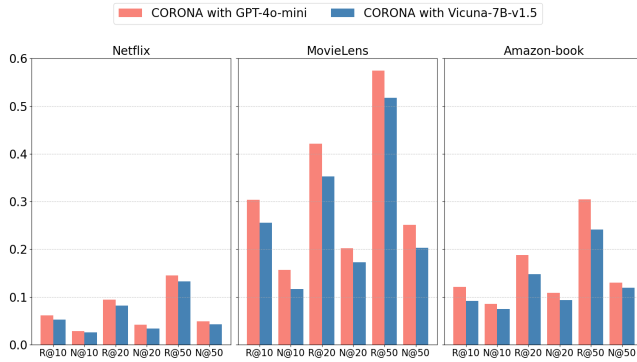
We investigate the impact of key hyperparameters in this subsection, including the size of retrieved subgraph k , the temperature parameter of LLM τ , the dimension of distance encoding $\dim(e)$ and the hidden dimension of GNN. For each dataset, we use $R@20$ as the evaluation metric, and the results are shown in Figure 4.

Table 3: Combinations of different subgraph retrieval methods and GNNs on three datasets in terms of $Recall@10/20/50$, and $NDCG@10/20/50$.

Datasets	Netflix						MovieLens						Amazon-book					
Variants	R@10	N@10	R@20	N@20	R@50	N@50	R@10	N@10	R@20	N@20	R@50	N@50	R@10	N@10	R@20	N@20	R@50	N@50
GCN (full graph)	0.0307	0.0179	0.0662	0.0217	0.1012	0.0247	0.1056	0.0808	0.1371	0.0993	0.1600	0.1208	0.0786	0.0405	0.1121	0.0556	0.1845	0.0797
GCN (fixed 1-hop)	0.0319	0.0107	0.0617	0.0169	0.0981	0.0217	0.1123	0.0804	0.1320	0.1037	0.1599	0.1181	0.0773	0.0427	0.1133	0.0516	0.1925	0.0790
GCN (fixed 2-hop)	0.0305	0.0097	0.0632	0.0184	0.0987	0.0274	0.1092	0.0912	0.1244	0.1086	0.1495	0.1178	0.0726	0.0443	0.1135	0.0582	0.1882	0.0796
GCN (CORONA)	0.0616	0.0279	0.0938	0.0416	0.1452	0.0487	0.3033	0.1565	0.4214	0.2017	0.5745	0.2507	0.1206	0.0855	0.1857	0.1089	0.3048	0.1299
GT (full graph)	0.0302	0.0104	0.0641	0.0192	0.1008	0.0285	0.1055	0.0907	0.1235	0.1074	0.1482	0.1170	0.0719	0.0435	0.1127	0.0571	0.1854	0.0787
GT (fixed 1-hop)	0.0294	0.0081	0.0655	0.0152	0.0988	0.0342	0.1035	0.1039	0.1170	0.0828	0.1501	0.1137	0.0672	0.0330	0.1109	0.0451	0.1835	0.0787
GT (fixed 2-hop)	0.0345	0.0117	0.0665	0.0160	0.1119	0.0305	0.1044	0.0864	0.1370	0.1043	0.1463	0.1193	0.0696	0.0453	0.1187	0.0601	0.1855	0.0790
GT (CORONA)	0.0628	0.0318	0.0931	0.0294	0.1468	0.0471	0.3072	0.1528	0.4165	0.1988	0.5504	0.2619	0.1203	0.1028	0.1928	0.1202	0.3011	0.1221

Table 4: Ablation study of different components on three datasets in terms of $Recall@10/20/50$, and $NDCG@10/20/50$.

Datasets	Netflix						MovieLens						Amazon-book					
Methods	R@10	N@10	R@20	N@20	R@50	N@50	R@10	N@10	R@20	N@20	R@50	N@50	R@10	N@10	R@20	N@20	R@50	N@50
w/o Preference-assisted Retrieval	0.0576	0.0227	0.0848	0.0343	0.1295	0.0431	0.2886	0.1375	0.3873	0.1970	0.5306	0.1829	0.1112	0.0747	0.1535	0.0909	0.2228	0.1036
w/o Intent-assisted Retrieval	0.0529	0.0208	0.0785	0.0329	0.1006	0.0418	0.2694	0.1092	0.3590	0.1913	0.5129	0.1768	0.1037	0.0764	0.1643	0.0937	0.2028	0.0961
w/o GNN-enhanced Retrieval	0.0591	0.0264	0.0917	0.0359	0.1304	0.0457	0.2958	0.1434	0.4088	0.1943	0.5531	0.1905	0.1189	0.0814	0.1729	0.0968	0.2311	0.1106
w/o Preference Reasoning	0.0506	0.0187	0.0759	0.0297	0.0954	0.0386	0.2532	0.1158	0.3316	0.1741	0.4819	0.1703	0.0933	0.0653	0.1490	0.0875	0.1986	0.0887
w/o Intent Reasoning	0.0451	0.0142	0.0735	0.0264	0.0933	0.0371	0.2317	0.1173	0.3221	0.1476	0.4609	0.1659	0.0789	0.0621	0.1215	0.0839	0.1982	0.0835
Full CORONA	0.0616	0.0279	0.0938	0.0416	0.1452	0.0487	0.3033	0.1565	0.4214	0.2017	0.5745	0.2507	0.1206	0.0855	0.1857	0.1089	0.3048	0.1299

**Figure 3: Applying different LLMs in CORONA on three datasets in terms of $Recall@10/20/50$ and $NDCG@10/20/50$.**

4.5.1 Size of Retrieved Subgraph k . We vary the value of k from 1,000 to 4,000, and observe that the performance improves as k increases up to a certain point, after which it begins to decline. For all datasets and metrics the best performance is achieved when k is around 3,000.

4.5.2 Temperature τ of LLM. Follow existing work [51], we conduct experiments on the temperature parameter τ , which controls the randomness of text generation. Higher values lead to greater diversity and creativity, while lower values produce more deterministic outputs. We experiment with τ values of {0, 0.2, 0.4, 0.6, 0.8, 1}. As illustrated in Figure 4, increasing τ slightly enhances most metrics at the beginning, but further increases result in a decline.

4.5.3 Dimension of Distance Encoding $\dim(e)$. We evaluate the impact of different dimensions of distance encoding on model performance, which plays a crucial role in the subgraph retrieval process. A smaller distance encoding dimension weakens the model's ability to differentiate neighbors at varying distances, while a larger dimension may lead to overfitting. We test distance encoding dimensions ranging from 0 to 5, and the results align with expectations, with the best performance achieved at a dimension of 2.

4.5.4 Hidden Dimension of GNN. We test hidden dimensions of {8, 16, 32, 64, 128, 256} on the message-passing component in GNN-enhanced retrieval. The results show improved performance as the hidden dimension increases, peaking at 128 dimensions, after which a slight decline is observed.

4.6 Efficiency Analysis (RQ5)

To answer RQ5 and validate the efficiency of CORONA, we measure the total running time on three datasets, and compare it with three LLM-based methods. Since TALLRec needs to fine-tune LLMs, we only compare it under open-source LLMs. The results are represented in Figure 5. From the results, we can observe that: (1) Although CORONA involves two LLM inference steps, its inference focuses only on users, avoiding the large-scale inference over items or interactions, which ensures better efficiency compared to other methods. (2) CORONA is more efficient with open-source LLMs than GPT-4o-mini that suffers from additional overhead from network transmission. (3) CORONA with GPT-4o-mini consumes approximately 1 US cent and less than 1.5 seconds per user for inference. With locally deployed LLMs, the inference time can be

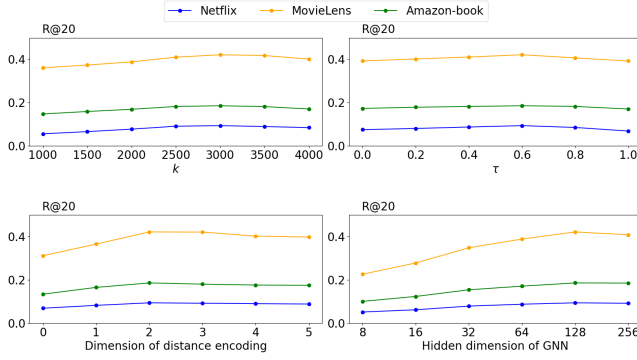


Figure 4: Hyperparameter experiments on three datasets in terms of Recall@20.

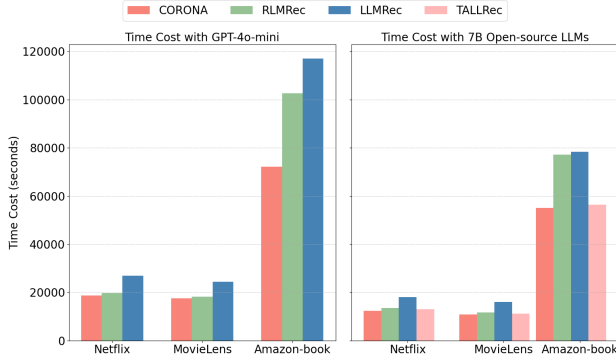


Figure 5: Total time cost (in seconds) of different LLM-based recommendation methods.

further reduced to less than one second. While lightweight open-source LLMs show a slight performance drop, we interpret this as a balance between performance and efficiency, allowing users to choose the version that best fits their needs. (3) LLMRec needs to augment over all items, users and interactions in the dataset, and thus takes longer processing time, especially for the Amazon-book dataset with denser interactions.

4.7 Case Study (RQ6)

To illustrate the inference process of our proposed CORONA framework in an understandable natural language form and further verify whether the LLM-based reasoner can generate plausible user preferences and intents, we present two examples from the test set of Netflix dataset in Figure 6. In the examples, we show the user profile, the preference reasoning, item summary of preference-assisted retrieval, interaction history, intent reasoning, and the final recommendation. The results show that the LLM-based reasoners can effectively infer relevant content leveraging their commonsense knowledge and enhance recommendation performance.

5 Conclusion

In this paper, we present a novel paradigm to incorporate LLMs into recommendation systems, where LLMs are employed for coarse-grained reasoning to assist in retrieval across the entire item set. The proposed CORONA framework has a carefully designed three-stage

User Profile: Age: 40-55; Gender: Male; Country: United States; Language: English; Liked directors: Bernardo Bertolucci

Preference Reasoning: The user likely enjoys **historical dramas**, **epic films** that reflect deep social and personal themes. Prefers visually rich and thematic films (inspired by Bernardo Bertolucci). Likely favors **English-language movies** from the 1970s-1990s.

Summary: Candidate genres include: **historical dramas**, **epic films**, crime dramas, war films, **classic cinema**...

User History: Title: The Last Emperor; Year: 1987; Genre: Historical Drama; Title: The Cook and Other Treasures; Year: 1970; Genre: Classic cinema; Title: Journey: 2001; Year: 2000; Genre: classic cinema; Title: Reds; Year: 1981; Genre: Historical Drama;

Intent Reasoning: The user's viewing preferences lean towards serious, thematic, and genre-specific films, particularly **historical dramas**, **epic films**, music-related films and **international films**.

Final Recommendation @5: Schindler's List (Historical drama, 1993); Crouching Tiger, Hidden Dragon (International film, 2000); Ran (International film, 1985); Lawrence of Arabia (Epic film, 1962); The Godfather (Classic cinema, 1982).

(a) Case of a male user.

User Profile: Age: 30-50; Gender: Female; Country: United States; Language: English; Liked directors: Ingmar Bergman

Preference Reasoning: The user likely enjoys **psychological dramas** and **horror**, artistic, moody, and **emotionally complex films**. The interest in European directors indicates a broader openness to international and foreign films, especially **European cinema**.

Summary: Candidate genres include: **Horror**, Psychological Thriller, **Animation**, **Fantasy**, **Drama**, **Mystery**...

User History: Title: Hour of wolf; Year: 1986; Genre: Horror; Title: Deep Red; Year: 1975; Genre: Horror; Title: Bambi: Platinum Edition; Year: 1972; Genre: animation; Title: Valerie and Her Week of Wonders; Year: 1970; Genre: Fantasy;

Intent Reasoning: The user's viewing preferences lean to something in the **Horror** or **Thriller** genres. They may also be drawn to **Mystery** due to its connection with psychological thrillers and horror.

Final Recommendation @5: The Shining (Horror film, 1980); Psycho (Psychological Thriller, 1960); Picnic as Hanging Rock (Mystery, 1975); The Night Porter (Psychological Thriller, 1977); The Sixth Sence (Fantasy/Horror, 1999).

(b) Case of a female user.

Figure 6: Case study of recommendation with our proposed CORONA. Here we mark the correct items in the final recommendation list in green, and the key information related to the correct items are highlighted in red.

retrieval process that progressively refines the selection at different levels of granularity. CORONA harnesses the reasoning power of LLMs for preference and intent inference, combined with GNNs for efficient recommendations. Extensive experiments confirm the effectiveness of the CORONA framework and validate its design. Future work may extend our framework to larger-scale industrial scenarios with more stages of retrieval. Additionally, exploring different LLMs to find the most cost-effective implementation could be valuable. It is also possible to finetune an open-source LLM for more accurate preference or intent reasoning.

Acknowledgments

This work was supported by the National Key Research and Development Program of China (No.2023YFC3303800).

References

- [1] 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. In *In Proceedings of the 17th ACM Conference on Recommender Systems*. 1007–1014.
- [2] James Bennett, Stan Lanning, et al. 2007. The netflix prize. In *In Proceedings of KDD Cup and Workshop*, Vol. 2007. New York, 35.
- [3] Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. 2020. Language models are few-shot learners. *Advances in Neural Information Processing Systems* 33 (2020), 1877–1901.
- [4] Lei Chen, Le Wu, Richang Hong, Kun Zhang, and Meng Wang. 2020. Revisiting graph based collaborative filtering: A linear residual graph convolutional network approach. In *In Proceedings of the AAAI Conference on Artificial Intelligence*, Vol. 34. 27–34.
- [5] Yuxin Chen, Junfei Tan, An Zhang, Zhengyi Yang, Leheng Sheng, Enzhi Zhang, Xiang Wang, and Tat-Seng Chua. 2024. On Softmax Direct Preference Optimization for Recommendation. *arXiv preprint arXiv:2406.09215* (2024).
- [6] Wei-Lin Chiang, Zhuohan Li, Zi Lin, Ying Sheng, Zhanghao Wu, Hao Zhang, Lianmin Zheng, Siyuan Zhuang, Yonghao Zhuang, Joseph E Gonzalez, et al. 2023. Vicuna: An open-source chatbot impressing gpt-4 with 90%* chatgpt quality. See <https://vicuna.lmsys.org> (accessed 14 April 2023) 2, 3 (2023), 6.
- [7] Sunhao Dai, Ninglu Shao, Haiyuan Zhao, Weijie Yu, Zihua Si, Chen Xu, Zhongxiang Sun, Xiao Zhang, and Jun Xu. 2023. Uncovering chatgpt’s capabilities in recommender systems. In *In Proceedings of the 17th ACM Conference on Recommender Systems*. 1126–1132.
- [8] Abhinandan S Das, Mayur Datar, Ashutosh Garg, and Shyam Rajaram. 2007. Google news personalization: scalable online collaborative filtering. In *In Proceedings of the ACM Web Conference 2007*. 271–280.
- [9] James Davidson, Benjamin Liebald, Junnig Liu, Palash Nandy, Taylor Van Vleet, Ullas Gargi, Sujoy Gupta, Yu He, Mike Lambert, Blake Livingston, et al. 2010. The YouTube video recommendation system. In *In Proceedings of the 4th ACM conference on Recommender systems*. 293–296.
- [10] Chen Gao, Yu Zheng, Nian Li, Yinfeng Li, Yingrong Qin, Jinghua Piao, Yuhuan Quan, Jianxin Chang, Depeng Jin, Xiangnan He, et al. 2023. A survey of graph neural networks for recommender systems: Challenges, methods, and directions. *ACM Transactions on Recommender Systems* 1, 1 (2023), 1–51.
- [11] Yunfan Gao, Tao Sheng, Youlin Xiang, Yun Xiong, Haoften Wang, and Jiawei Zhang. 2023. Chat-rec: Towards interactive and explainable llms-augmented recommender system. *arXiv preprint arXiv:2303.14524* (2023).
- [12] Binzong Geng, Zhaoxin Huan, Xiaolu Zhang, Yong He, Liang Zhang, Fajie Yuan, Jun Zhou, and Linjian Mo. 2024. Breaking the length barrier: Llm-enhanced CTR prediction in long textual user behaviors. In *In Proceedings of the 47th International ACM SIGIR Conference on Research and Development in Information Retrieval*. 2311–2315.
- [13] Shijie Geng, Shuchang Liu, Zuohui Fu, Yingqiang Ge, and Yongfeng Zhang. 2022. Recommendation as language processing (rlp): A unified pretrain, personalized prompt & predict paradigm (p5). In *In Proceedings of the 16th ACM Conference on Recommender Systems*. 299–315.
- [14] F Maxwell Harper and Joseph A Konstan. 2015. The movielens datasets: History and context. *Acm Transactions on Interactive Intelligent Systems* 5, 4 (2015), 1–19.
- [15] Xiangnan He, Kuan Deng, Xiang Wang, Yan Li, Yongdong Zhang, and Meng Wang. 2020. Lightgcn: Simplifying and powering graph convolution network for recommendation. In *In Proceedings of the 43rd International ACM SIGIR conference on Research and development in Information Retrieval*. 639–648.
- [16] Xiaoxin He, Yijun Tian, Yifei Sun, Nitesh Chawla, Thomas Laurent, Yann LeCun, Xavier Bresson, and Bryan Hooi. 2024. G-retriever: Retrieval-augmented generation for textual graph understanding and question answering. *Advances in Neural Information Processing Systems* 37 (2024), 132876–132907.
- [17] Yupeng Hou, Junjie Zhang, Zihan Lin, Hongyu Lu, Ruobing Xie, Julian McAuley, and Wayne Xin Zhao. 2024. Large language models are zero-shot rankers for recommender systems. In *European Conference on Information Retrieval*. Springer, 364–381.
- [18] Wenyue Hua, Lei Li, Shuyuan Xu, Li Chen, and Yongfeng Zhang. 2023. Tutorial on large language models for recommendation. In *In Proceedings of the 17th ACM Conference on Recommender Systems*. 1281–1283.
- [19] Feiran Huang, Zefan Wang, Xiao Huang, Yufeng Qian, Zhetao Li, and Hao Chen. 2023. Aligning distillation for cold-start item recommendation. In *In Proceedings of the 46th International ACM SIGIR Conference on Research and Development in Information Retrieval*. 1147–1157.
- [20] Feiran Huang, Zhenghang Yang, Junyi Jiang, Yuanchen Bei, Yijie Zhang, and Hao Chen. 2024. Large Language Model Interaction Simulator for Cold-Start Item Recommendation. *arXiv preprint arXiv:2402.09176* (2024).
- [21] Sein Kim, Hongseok Kang, Seungyoon Choi, Donghyun Kim, Minchul Yang, and Chanyoung Park. 2024. Large language models meet collaborative filtering: An efficient all-round llm-based recommender system. In *In Proceedings of the 30th ACM SIGKDD Conference on Knowledge Discovery and Data Mining*. 1395–1406.
- [22] Thomas N Kipf and Max Welling. 2016. Semi-supervised classification with graph convolutional networks. *arXiv preprint arXiv:1609.02907* (2016).
- [23] Xiaoyu Kong, Jiancan Wu, An Zhang, Leheng Sheng, Hui Lin, Xiang Wang, and Xiangnan He. 2024. Customizing Language Models with Instance-wise LoRA for Sequential Recommendation. *arXiv preprint arXiv:2408.10159* (2024).
- [24] Yehuda Koren, Steffen Rendle, and Robert Bell. 2021. Advances in collaborative filtering. *Recommender Systems Handbook* (2021), 91–142.
- [25] Haoyang Li, Xin Wang, Ziwei Zhang, Jianxin Ma, Peng Cui, and Wenwu Zhu. 2021. Intention-aware sequential recommendation with structured intent transition. *IEEE Transactions on Knowledge and Data Engineering* 34, 11 (2021), 5403–5414.
- [26] Jiayi Liao, Sihang Li, Zhengyi Yang, Jiancan Wu, Yancheng Yuan, Xiang Wang, and Xiangnan He. 2024. Llara: Large language-recommendation assistant. In *In Proceedings of the 47th International ACM SIGIR Conference on Research and Development in Information Retrieval*. 1785–1795.
- [27] Jianghao Lin, Xinyi Dai, Yunjia Xi, Weiwen Liu, Bo Chen, Hao Zhang, Yong Liu, Chuhan Wu, Xiangyang Li, Chenxu Zhu, et al. 2025. How can recommender systems benefit from large language models: A survey. *ACM Transactions on Information Systems* 43, 2 (2025), 1–47.
- [28] Xinyu Lin, Wenjie Wang, Yongqi Li, Shuo Yang, Fuli Feng, Yinwei Wei, and Tat-Seng Chua. 2024. Data-efficient Fine-tuning for LLM-based Recommendation. In *In Proceedings of the 47th International ACM SIGIR Conference on Research and Development in Information Retrieval*. 365–374.
- [29] Peng Liu, Lemei Zhang, and Jon Atle Gulla. 2023. Pre-train, Prompt, and Recommendation: A Comprehensive Survey of Language Modeling Paradigm Adaptations in Recommender Systems. *Transactions of the Association for Computational Linguistics* 11 (2023), 1553–1571.
- [30] Qidong Liu, Xian Wu, Xiangyu Zhao, Yejing Wang, Zijian Zhang, Feng Tian, and Yefeng Zheng. 2024. Large Language Models Enhanced Sequential Recommendation for Long-tail User and Item. *arXiv preprint arXiv:2405.20646* (2024).
- [31] Weiwen Liu, Qing Liu, Ruiming Tang, Junyang Chen, Xiuqiang He, and Pheng Ann Heng. 2020. Personalized Re-ranking with Item Relationships for E-commerce. In *In Proceedings of the 29th ACM International Conference on Information & Knowledge Management*. 925–934.
- [32] Qiyao Ma, Xubin Ren, and Chao Huang. 2024. XRec: Large Language Models for Explainable Recommendation. *arXiv preprint arXiv:2406.02377* (2024).
- [33] Sheshera Mysore, Andrew McCallum, and Hamed Zamani. 2023. Large language model augmented narrative driven recommendations. In *In Proceedings of the 17th ACM Conference on Recommender Systems*. 777–783.
- [34] Jianmo Ni, Jiacheng Li, and Julian McAuley. 2019. Justifying recommendations using distantly-labeled reviews and fine-grained aspects. In *In Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing*. 188–197.
- [35] Fabio Petroni, Tim Rocktäschel, Patrick Lewis, Anton Bakhtin, Yuxiang Wu, Alexander H Miller, and Sebastian Riedel. 2019. Language models as knowledge bases? *arXiv preprint arXiv:1909.01066* (2019).
- [36] N Reimers. 2019. Sentence-BERT: Sentence Embeddings using Siamese BERT- Networks. *arXiv preprint arXiv:1908.10084* (2019).
- [37] Xubin Ren, Wei Wei, Lianghao Xia, Lixin Su, Suqi Cheng, Junfeng Wang, Dawei Yin, and Chao Huang. 2024. Representation learning with large language models for recommendation. In *In Proceedings of the ACM Web Conference 2024*. 3464–3475.
- [38] Yankun Ren, Zhongde Chen, Xinxing Yang, Longfei Li, Cong Jiang, Lei Cheng, Bo Zhang, Linjian Mo, and Jun Zhou. 2024. Enhancing sequential recommenders with augmented knowledge from aligned large language models. In *In Proceedings of the 47th International ACM SIGIR Conference on Research and Development in Information Retrieval*. 345–354.
- [39] 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 *In Proceedings of the 17th ACM Conference on Recommender Systems*. 890–896.
- [40] Wentao Shi, Xiangnan He, Yang Zhang, Chongming Gao, Xinyue Li, Jizhi Zhang, Qifan Wang, and Fuli Feng. 2024. Enhancing Long-Term Recommendation with Bi-level Learnable Large Language Model Planning. *arXiv preprint arXiv:2403.00843* (2024).
- [41] Zhu Sun, Hongyang Liu, Xinghua Qu, Kaidong Feng, Yan Wang, and Yew Soon Ong. 2024. Large language models for intent-driven session recommendations. In *In Proceedings of the 47th International ACM SIGIR Conference on Research and Development in Information Retrieval*. 324–334.
- [42] Juntao Tan, Shuyuan Xu, Wenyue Hua, Yingqiang Ge, Zelong Li, and Yongfeng Zhang. 2024. Towards llm-recsys alignment with textual id learning. *arXiv preprint arXiv:2403.19021* (2024).
- [43] Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, et al. 2023. Llama: Open and efficient foundation language models. *arXiv preprint arXiv:2302.13971* (2023).
- [44] Maksims Volkovs, Guangwei Yu, and Tomi Poutanen. 2017. Dropoutnet: Addressing cold start in recommender systems. *Advances in Neural Information Processing Systems* 30 (2017).

- [45] Jie Wang, Alexandros Karatzoglou, Ioannis Arapakis, and Joemon M Jose. 2024. Reinforcement Learning-based Recommender Systems with Large Language Models for State Reward and Action Modeling. In *In Proceedings of the 47th International ACM SIGIR Conference on Research and Development in Information Retrieval*. 375–385.
- [46] Jiayin Wang, Fengran Mo, Weizhi Ma, Peijie Sun, Min Zhang, and Jian-Yun Nie. 2024. A User-Centric Multi-Intent Benchmark for Evaluating Large Language Models. In *In Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing*. 3588–3612.
- [47] Xiang Wang, Xiangnan He, Meng Wang, Fuli Feng, and Tat-Seng Chua. 2019. Neural graph collaborative filtering. In *In Proceedings of the 42nd international ACM SIGIR conference on Research and development in Information Retrieval*. 165–174.
- [48] Xiaolei Wang, Xinyu Tang, Wayne Xin Zhao, Jingyuan Wang, and Ji-Rong Wen. 2023. Rethinking the evaluation for conversational recommendation in the era of large language models. *arXiv preprint arXiv:2305.13112* (2023).
- [49] Y. Wang, Z. Chu, X. Ouyang, S. Wang, H. Hao, Y. Shen, J. Gu, S. Xue, J. Zhang, Q. Cui, L. Li, J. Zhou, and S. Li. 2024. LLMRG: Improving Recommendations through Large Language Model Reasoning Graphs. In *In Proceedings of the AAAI Conference on Artificial Intelligence*, Vol. 38. AAAI Press, 19189–19196. doi:10.1609/aaai.v38i17.29887
- [50] Yifan Wang, Suyao Tang, Yuntong Lei, Weiping Song, Sheng Wang, and Ming Zhang. 2020. Disenhan: Disentangled heterogeneous graph attention network for recommendation. In *In Proceedings of the 29th ACM International Conference on Information & Knowledge Management*. 1605–1614.
- [51] Wei Wei, Xubin Ren, Jiabin Tang, Qinyong Wang, Lixin Su, Suqi Cheng, Junfeng Wang, Dawei Yin, and Chao Huang. 2024. Llmrec: Large language models with graph augmentation for recommendation. In *In Proceedings of the 17th ACM International Conference on Web Search and Data Mining*. 806–815.
- [52] Yinwei Wei, Xiang Wang, Liqiang Nie, Xiangnan He, and Tat-Seng Chua. 2020. Graph-refined convolutional network for multimedia recommendation with implicit feedback. In *In Proceedings of the 28th ACM International Conference on Multimedia*. 3541–3549.
- [53] Likang Wu, Zhi Zheng, Zhaopeng Qiu, Hao Wang, Hongchao Gu, Tingjia Shen, Chuan Qin, Chen Zhu, Hengshu Zhu, Qi Liu, et al. 2024. A survey on large language models for recommendation. *World Wide Web* 27, 5 (2024), 60.
- [54] Jiajing Xu, Andrew Zhai, and Charles Rosenberg. 2022. Rethinking personalized ranking at Pinterest: An end-to-end approach. In *In Proceedings of the 16th ACM Conference on Recommender Systems*. 502–505.
- [55] Shenghao Yang, Weizhi Ma, Peijie Sun, Qingyao Ai, Yiqun Liu, Mingchen Cai, and Min Zhang. 2024. Sequential recommendation with latent relations based on large language model. In *In Proceedings of the 47th International ACM SIGIR Conference on Research and Development in Information Retrieval*. 335–344.
- [56] Yuhao Yang, Lianghao Xia, Da Luo, Kangyi Lin, and Chao Huang. 2024. GraphPro: Graph Pre-training and Prompt Learning for Recommendation. In *In Proceedings of the ACM Web Conference 2024*. 3690–3699.
- [57] Rex Ying, Ruining He, Kaifeng Chen, Pong Eksombatchai, William L Hamilton, and Jure Leskovec. 2018. Graph convolutional neural networks for web-scale recommender systems. In *In Proceedings of the 24th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining*. 974–983.
- [58] Seongjun Yun, Minbyul Jeong, Raehyun Kim, Jaewoo Kang, and Hyunwoo J Kim. 2019. Graph transformer networks. *Advances in Neural Information Processing Systems* 32 (2019).
- [59] Jizhi Zhang, Keqin Bao, Yang Zhang, Wenjie Wang, Fuli Feng, and Xiangnan He. 2023. Is chatgpt fair for recommendation? evaluating fairness in large language model recommendation. In *In Proceedings of the 17th ACM Conference on Recommender Systems*. 993–999.
- [60] Junjie Zhang, Ruobing Xie, Yupeng Hou, Xin Zhao, Leyu Lin, and Ji-Rong Wen. 2023. Recommendation as instruction following: A large language model empowered recommendation approach. *ACM Transactions on Information Systems* (2023).
- [61] Kaike Zhang, Qi Cao, Yunfan Wu, Fei Sun, Huawei Shen, and Xueqi Cheng. 2024. Lorec: Large language model for robust sequential recommendation against poisoning attacks. *arXiv preprint arXiv:2401.17723* (2024).
- [62] Xiaoyu Zhang, Yishan Li, Jiayin Wang, Bowen Sun, Weizhi Ma, Peijie Sun, and Min Zhang. 2024. Large language models as evaluators for recommendation explanations. In *In Proceedings of the 18th ACM Conference on Recommender Systems*. 33–42.
- [63] Yang Zhang, Keqin Bao, Ming Yan, Wenjie Wang, Fuli Feng, and Xiangnan He. 2024. Text-like Encoding of Collaborative Information in Large Language Models for Recommendation. *arXiv preprint arXiv:2406.03210* (2024).
- [64] Yang Zhang, Fuli Feng, Jizhi Zhang, Keqin Bao, Qifan Wang, and Xiangnan He. 2025. Collm: Integrating collaborative embeddings into large language models for recommendation. *IEEE Transactions on Knowledge and Data Engineering* (2025).
- [65] Wayne Xin Zhao, Kun Zhou, Junyi Li, Tianyi Tang, Xiaolei Wang, Yupeng Hou, Yingqian Min, Beichen Zhang, Junjie Zhang, Zican Dong, et al. 2023. A survey of large language models. *arXiv preprint arXiv:2303.18223* (2023).
- [66] Yurou Zhao, Yiding Sun, Ruidong Han, Fei Jiang, Lu Guan, Xiang Li, Wei Lin, Weizhi Ma, and Jiaxin Mao. 2024. Aligning Explanations for Recommendation with Rating and Feature via Maximizing Mutual Information. In *In Proceedings of the 33rd ACM International Conference on Information and Knowledge Management*. 3374–3383.
- [67] Zihuai Zhao, Wenqi Fan, Jiatong Li, Yunqing Liu, Xiaowei Mei, Yiqi Wang, Zhen Wen, Fei Wang, Xiangyu Zhao, Jiliang Tang, et al. 2024. Recommender systems in the era of large language models (llms). *IEEE Transactions on Knowledge and Data Engineering* (2024).
- [68] Yu Zheng, Chen Gao, Xiangnan He, Yong Li, and Depeng Jin. 2020. Price-aware recommendation with graph convolutional networks. In *2020 IEEE 36th International Conference on Data Engineering*. IEEE, 133–144.
- [69] Yaochen Zhu, Liang Wu, Qi Guo, Liangjie Hong, and Jundong Li. 2024. Collaborative large language model for recommender systems. In *In Proceedings of the ACM Web Conference 2024*. 3162–3172.