

Retrieval-Augmented Recommendation Explanation Generation with Hierarchical Aggregation

Bangcheng Sun, Yazhe Chen, Jilin Yang, Xiaodong Li, Hui Li

Key Laboratory of Multimedia Trusted Perception and Efficient Computing,
Ministry of Education of China, Xiamen University, China
{36920231153231, 23020241154376, 22920212204262}@stu.xmu.edu.cn
{xdli, hui}@xmu.edu.cn

Abstract

Explainable Recommender System (ExRec) provides transparency to the recommendation process, increasing users' trust and boosting the operation of online services. With the rise of large language models (LLMs), whose extensive world knowledge and nuanced language understanding enable the generation of human-like, contextually grounded explanations, LLM-powered ExRec has gained great momentum. However, existing LLM-based ExRec models suffer from profile deviation and high retrieval overhead, hindering their deployment. To address these issues, we propose Retrieval-Augmented Recommendation Explanation Generation with Hierarchical Aggregation (REXHA). Specifically, we design a hierarchical aggregation based profiling module that comprehensively considers user and item review information, hierarchically summarizing and constructing holistic profiles. Furthermore, we introduce an efficient retrieval module using two types of pseudo-document queries to retrieve relevant reviews to enhance the generation of recommendation explanations, effectively reducing retrieval latency and improving the recall of relevant reviews. Extensive experiments demonstrate that our method outperforms existing approaches by up to 12.6% w.r.t. the explanation quality while achieving high retrieval efficiency.

1 Introduction

The exponential growth of digital content and products across online platforms (e.g., Yahoo News, Amazon and TikTok) has intensified the information overload problem [1], where users face overwhelming challenges in identifying relevant items amid vast information. To mitigate cognitive burdens in decision-making, Recommender System (RecSys) has emerged as an essential tool and is prevalently incorporated into many online services. RecSys effectively filters irrelevant information and delivers tailored recommendations. This way, it not only helps users better find their desired targets but also facilitates the operation of platforms. Hence, RecSys has attracted much attention and progressed actively over the past few decades [2, 3].

As users increasingly rely on recommendations for high-stakes decisions (e.g., financial investments and healthcare choices), mere predictive accuracy proves insufficient without explainable justification for recommendations. This challenge has catalyzed the emergence of Explainable Recommender System (ExRec) [4]. ExRec explains the recommendation results to users and enhances the transparency of the decision-making process. Therefore, it increases users' trust in RecSys, further boosting the operation of online services.

Prior works on ExRec mainly study how to generate interpretable explanations for user-item interactions [5]. For instance, some works adopt deep learning techniques like RNNs [6], GNNs [7] and Transformer [8] for capturing intricate patterns in user behaviors and item attributes, enabling the

generation of persuasive explanations. Although these methods are effective in some cases, they naturally suffer from limited language competence as they are trained over the restricted recommendation data. Hence, a branch of work on ExRec opts to use language models (e.g., BERT [9]) pre-trained on vast, diverse textual corpora to generate human-like explanations. This paradigm shift has gained momentum with the rise of large language models (LLMs), whose extensive world knowledge and nuanced language understanding enable the generation of human-like, contextually grounded explanations [10–12].

XRec is one representative LLM-based ExRec [11]. It enables LLMs to better understand complex user-item interaction patterns and offer explanations via integrating collaborative filtering (CF) signals. Building on this foundation, G-Refer [12] introduces a hybrid graph retrieval method to enhance ExRec by extracting both structural and semantic CF information from the user-item interaction graph. Despite these notable advancements, two critical limitations persist in LLM-powered ExRec: **(C1) Profile Deviation**. To assist with explanation generation, XRec and G-Refer construct user/item profiles by randomly sampling a small subset of user/item reviews, neglecting the broader contextual information embedded in the remaining data. This selective sampling introduces information bias, causing LLMs to generate explanations misaligned with holistic user preferences or item characteristics. **(C2) High Retrieval Overhead**. G-Refer employs the Dijkstra algorithm in path-level retrieval and can only be calculated on the CPU. It is not friendly to parallelism and has a high time complexity of $\mathcal{O}(N^2)$, resulting in a prohibitive retrieval cost.

To alleviate the limitations of existing ExRec works, we present Retrieval-Augmented Recommendation Explanation Generation with Hierarchical Aggregation (REXHA), a framework designed to improve the credibility and efficiency of recommendation explanation generation. The contributions of this work are summarized as follows:

- **Holistic User/Item Profiling:** Unlike prior methods (e.g., XRec and G-Refer) that construct incomplete profiles via random review sampling, we propose a hierarchical review aggregation module to systematically encode all user/item review data. This approach mitigates profile deviation **(C1)** by capturing nuanced user preferences and item characteristics through multi-layered summarization, ensuring LLMs generate explanations grounded in comprehensive contextual signals.
- **Efficient Retrieval via Pseudo-Document Queries:** To overcome the computational bottlenecks of retrieval **(C2)**, we introduce aggregated pseudo-document query generation. By synthesizing latent queries from user/item reviews, REXHA efficiently retrieves semantically rich opinions from vectorized historical reviews, drastically reducing latency compared to prior LLM-powered ExRec.
- **Extensive and Rigorous Evaluation:** We have conducted extensive experiments on public datasets. The results demonstrate that REXHA achieves significant performance gains, outperforming a range of state-of-the-art baselines by up to 12.6% on recommendation explanation generation.

2 Our Method REXHA

Fig. 1 provides an overview of REXHA. The primary idea of REXHA is to efficiently retrieve and harness reviews relevant to the target user-item interaction to enhance the credibility and personalization of the generated explanation. To achieve this, REXHA uses three main modules, collaborative signal extraction (Sec. 2.1), hierarchical aggregation based profiling (Sec. 2.2) and review retrieval for data augmentation (Sec. 2.3), to prepare rich and beneficial background data for the target user-item interaction that is further fed into LLM for retrieval-augmented explanation generation (Sec. 2.4).

2.1 Collaborative Signal Extraction

Collaborative information is crucial to model recommendation process [2]. Constructing a user-item interaction graph and then capturing interaction patterns from the perspective of graph semantics and structure is a prevalent method, providing insights into understanding user preferences and item characteristics [13]. Hence, in REXHA, we leverage the capability of Graph Convolutional Networks (GCNs) to encode collaborative information in the user-item interaction graph into latent embeddings,

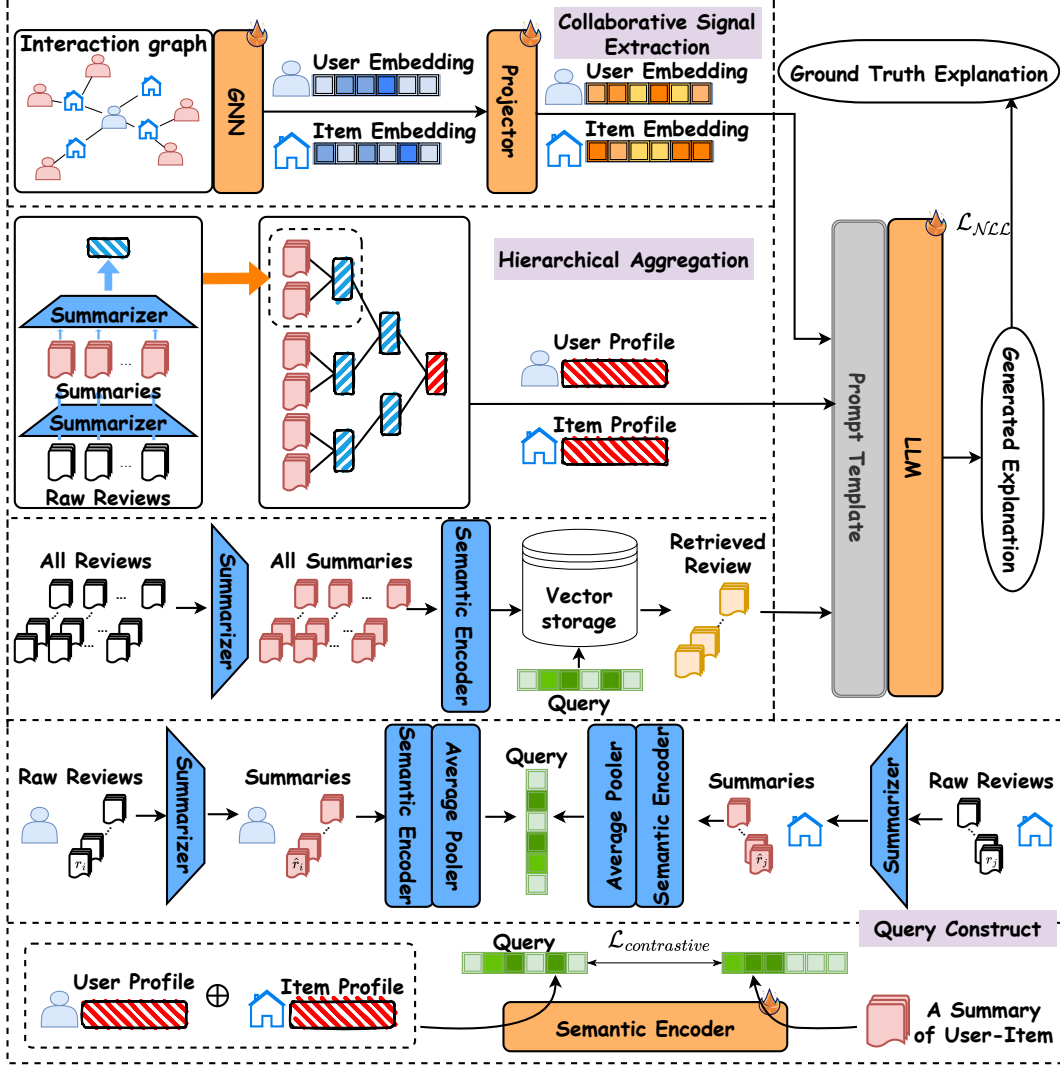


Figure 1: Overview of REXHA. It contains three key components: (1) **Collaborative Signal Extraction** provides collaborative filtering information to LLMs. (2) **Hierarchical Aggregation** compresses and summarizes reviews layer by layer, finally construct textual profiles for user/item. (3) **Review Retrieval** module retrieves relative reviews to enhance LLM generating explanations.

complementing later LLM-based explanation generation by incorporating collaborative signals in recommendation scenarios.

To be specific, we adopt LightGCN [14] to encode the collaborative information:

$$\mathbf{e}_u^{(l+1)} = \sum_{i \in \mathcal{N}_u} \frac{1}{\sqrt{|\mathcal{N}_u|} \sqrt{|\mathcal{N}_i|}} \mathbf{e}_i^{(l)}, \quad \mathbf{e}_i^{(l+1)} = \sum_{u \in \mathcal{N}_i} \frac{1}{\sqrt{|\mathcal{N}_i|} \sqrt{|\mathcal{N}_u|}} \mathbf{e}_u^{(l)} \quad (1)$$

where $\mathbf{e}_u^{(l)}$ and $\mathbf{e}_i^{(l)}$ denote the embedding of user u and item i after l -layer propagation, respectively. \mathcal{N}_u denotes the set of users who have interacted with item i , and \mathcal{N}_i indicates the set of items interacted by user u .

The user and item embeddings are extracted by averaging their embeddings from all GCN layers:

$$\mathbf{e}_u = \sum_{l=0}^L \frac{1}{L+1} \mathbf{e}_u^{(l)}, \quad \mathbf{e}_i = \sum_{l=0}^K \frac{1}{L+1} \mathbf{e}_i^{(l)} \quad (2)$$

where L denotes the number of layers. To align with the dimensionality of LLM’s representation space, we further employ a three-layer feedforward neural network to project user and item embeddings. We utilize the negative log-likelihood (NLL) as our training loss, while the parameters of LLM are frozen.

$$\mathcal{L} = -\frac{1}{N} \sum_{i=1}^N \sum_{c=1}^{C_i} y_{i,c} \log(\hat{y}_{i,c}). \quad (3)$$

Here N denotes the total number of user-item pairs to be explained, and C denotes the token count in each explanations. $y_{i,c}$ and $\hat{y}_{i,c}$ correspond to the ground-truth and predicted tokens at position c of i -th sample, respectively.

2.2 Hierarchical Aggregation Based Profiling

Item textual attributes (e.g., item title, item category, etc) contain important features for capturing item characteristics. Nevertheless, as many items share similar textual attributes, there is insufficient information for constructing distinguishing profiles, making later explanation generation lack useful supporting inputs.

Instead, reviews written by humans reflect user preferences and item characteristics. Therefore, incorporating and summarizing fragmented human-written reviews using natural language as user and item profiles can enhance the richness and usefulness of the constructed profiles, providing better supporting inputs for LLM to generate persuasive and comprehensive recommendation explanations.

Existing LLM-based ExRec methods XRec and G-Refer utilize LLM as the summarizer. Since processing all the reviews of a popular item may exceed the context length of LLM, they randomly sample a few reviews for constructing profiles, losing quite a lot of information beneficial for generating explanations. Furthermore, randomly sampled reviews lack connections to each other, making it more difficult for LLM to summarize and produce comprehensive and correct profiles. Despite that LLMs with a longer context window can be applied for profile summarization, the lost-in-middle issue [15] still affects profile construction.

To conquer the above problem, we opt to construct user and item profiles via hierarchical aggregating their reviews instead of feeding reviews together to LLM. REXHA compresses and summarizes reviews layer by layer in parallel, and finally constructs a refined and informative profile.

2.2.1 Raw Review Summarization

Given a user-item interaction, we conduct hierarchical aggregation twice: one for user profile and the other for item profile. At the bottom of the hierarchical aggregation tree, each leaf contains one raw item review and each review is randomly placed in a leaf node. When constructing a user profile, the raw review comes from one item interacted by that user. When constructing an item profile, we use this item’s reviews as raw reviews. As shown in Fig. 1, the first step is to aggregate reviews of two adjacent sibling nodes and use LLM to summarize the two raw reviews.

The primary motivation for employing LLM to summarize raw item reviews lies in addressing two critical challenges: excessive information redundancy and semantic ambiguity. By condensing user feedback into summaries, this approach establishes a robust semantic foundation that enhances the later hierarchical aggregation process.

Raw review data often suffers from verbosity, noise, and conflicting signals, such as redundant feature mentions (e.g., “camera quality”), irrelevant details, or contradictory sentiments (“loves the design but hates the price”). Directly aggregating unprocessed text risks diluting actionable insights or introducing logical inconsistencies, which undermines both analytical precision and contextual coherence. LLM excels at distilling these noisy inputs into concise, semantically dense representations. This process not only overcomes the limitation of context length but also extracts nuanced user preferences (e.g., “prioritizes camera performance over battery life”) and item characteristics (e.g., “high-resolution sensor with limited battery capacity”). The resulting summaries act as high-fidelity intermediaries, filtering out extraneous details while preserving critical semantic relationships.

Based on the above motivation, we design the instruction prompt to order the LLM to preliminarily summarize each raw reviews. The prompt can be found in Appendix F. By combining the instruction prompt with a raw review r_i as the input to LLM, we can leverage LLM to generate review summarization and similar processes are conducted in parallel for raw reviews. This is similar to a k -ary-tree merging operation, where one raw review takes part in only one summarization. Assume there are m raw reviews at the bottom level. Then, after raw review summarization, there are m/k brief but more information-dense summaries and they will be engaged in the next step of hierarchical aggregation.

2.2.2 Hierarchical Aggregation

When constructing an item profile for an item, after generating summaries from raw reviews, we concatenate every p summaries into a group and conduct summarization similar to the summarization operation at the bottom level. This process is conducted recursively until reaching the root level where we only have one summarization (i.e., the constructed item profile), forming a hierarchical aggregation tree as shown in Fig. 1. When integrated into the hierarchical aggregation process, the refined summaries from each level enable progressive fusion of granular features into the holistic profile. Similarly, a user profile can be constructed in a hierarchical manner.

2.3 Review Retrieval for Data Augmentation

2.3.1 Review Retrieval Augmentation

Hierarchical aggregation only utilizes reviews directly related to the target user/item to construct the user/item profile. In practice, users often refer to the reviews of other similar items or consider reviews made by people with similar preferences before making decisions. Such reviews may contain auxiliary and use information that may support their decisions.

Based on the above observation, REXHA leverages an embedding model f as the semantic encoder to transfer raw review summaries obtained in Sec. 2.2.1 into representation vectors. Then, for a target user-item interaction that requires explanation generation, REXHA uses a retrieval query to find the most relevant raw review summaries based on cosine similarity. The query construction method will be introduced in the next section. Finally, top- q relevant raw review summaries to the target user-item interaction are listed as data augmentation for generating explanations.

2.3.2 Retrieval Query Construction

Unlike general information retrieval scenarios, the recommendation system domain inherently lacks readily available natural language queries that can systematically retrieve review-based evidence for specific user-item interactions to enhance later explanation generation. In conventional recommendation frameworks, there exists no straightforward mechanism to translate implicit user preferences and item characteristics into search-compatible queries that directly enable contextual raw review summaries retrieval.

To address this problem, we propose two types of pseudo-document query construction strategies: the latent representation query and the profile query. The former leverages the global information of the user-item pair, aiming to retrieve diverse and informative reviews. The latter uses user and item profiles—constructed to emphasize high-frequency, long-term preference attributes—as queries. By retrieving based on these specialized profiles, the resulting reviews exhibit a more concentrated semantic distribution and are more closely aligned with the underlying preferences of the user or item.

Latent Representation Query A latent representation query encodes the information of the target user-item interaction. Concretely, REXHA encodes all raw reviews of the target user and item using the embedding model f (semantic encoder) and then aggregates all the encoded representations as

the latent query for the target user-item interaction:

$$\hat{r}_i^u = f(r_i^u), \hat{r}_j^v = f(r_j^v), \quad (4)$$

$$\hat{R}^u = \{\hat{r}_1^u, \hat{r}_2^u, \dots, \hat{r}_n^u\}, \hat{R}^v = \{\hat{r}_1^v, \hat{r}_2^v, \dots, \hat{r}_n^v\}, \quad (5)$$

$$\mathbf{Q}^u = f(\hat{R}^u), \mathbf{Q}^v = f(\hat{R}^v), \quad (6)$$

$$\mathbf{q}^u = \frac{1}{N} \sum_{q_i^u \in \mathbf{Q}^u} \mathbf{q}_i^u, \mathbf{q}^v = \frac{1}{N} \sum_{q_i^v \in \mathbf{Q}^v} \mathbf{q}_i^v, \quad (7)$$

$$\mathbf{q}_{latent} = \frac{\mathbf{q}^u + \mathbf{q}^v}{2}. \quad (8)$$

where r_i^u is the raw review text written by the user, and \hat{r}_i^u denotes its corresponding summarized opinion. Let \hat{R}^u be the set of all user opinions \hat{r}^u . After mapping \hat{R}^u to higher-dimension space by embedding model f , we obtain \mathbf{Q}^u , the set of embedding \mathbf{q}^u . N denotes the number of opinions. The same applies to the item side.

Profile Query Since user and item profiles constructed in Sec. 2.2 contain precise and comprehensive information of user preferences and item characteristics, we design the profile query that directly uses the textual descriptions in profiles of the target user/item as the retrieval query. Because user/item profiles misalign with raw review summaries due to filtering information during the hierarchical aggregation step. In this scenarios, fine-tuning the embedding model is necessary. We designed a contrastive learning approach that minimizes the distance between the user u_i and item v_j profiles p_{u_i, v_j} and their corresponding opinions \hat{r}_{u_i, v_j} , while increasing the distance from irrelevant opinions.

$$\mathcal{L}_{contrastive} = \frac{e^{sim(\mathbf{p}_{u_i, v_j}, \hat{\mathbf{r}}_{u_i, v_j})/\tau}}{e^{sim(\mathbf{p}_{u_i, v_j}, \hat{\mathbf{r}}_{u_i, v_j})/\tau} + \sum_{(i', j') \neq (i, j)} e^{sim(\mathbf{p}_{u_i', v_{j'}}, \hat{\mathbf{r}}_{u_i', v_{j'}})/\tau}}. \quad (9)$$

where $sim(\cdot, \cdot)$ indicates cosine similarity, $\mathbf{p}_{u, v}$ and $\mathbf{r}_{u, v}$ denote embeddings of profiles and summarized reviews of the user-item pair, respectively. The embedding $\mathbf{p}_{u, v}$ and $\mathbf{r}_{u, v}$ are computed as follows:

$$\mathbf{p}_{u, v} = f'(p_{u, v}), \mathbf{r}_{u, v} = f'(\hat{r}_{u, v}). \quad (10)$$

where f' denotes fine-tuned embedding model. $p_{u, v}$ and $\hat{r}_{u, v}$ denote the profile and summarized reviews of user-item pair (u, v) .

2.4 Recommendation Explanation Generation

For a target user-item interaction, we concatenate the corresponding constructed user and item profiles (Sec. 2.2), the extracted user and item embeddings (Sec. 2.1) and the relevant reviews for data augmentation (Sec. 2.3) and use the prompt template shown in Fig. 8 to instruct LLM to generate the recommendation explanation.

3 Experiments

3.1 Experimental Settings

3.1.1 Datasets

We conduct experiments on three public datasets from different domains: **Yelp**¹, **Amazon-books**², and **Google-reviews** [16, 17]. More details of datasets can be found in Tab. 3 in Appendix A.

3.1.2 Evaluation Metrics

For evaluating the semantic explainability and stability of the generated explanation, we follow XRec [11] and G-Refer [12] to employ a suite of metrics that assessing the semantic explainability and stability of generated explanations. Concretely, we use GPT_{score} [18], BERT_{score}^{Precision}, BERT_{score}^{Recall} and BERT_{score}^{F1} [19] to measure the explainability. To evaluate the consistency, we also report the standard deviations of these metrics. The details of metric can be found in Appendix B.

¹<https://www.yelp.com/dataset/challenge>

²<https://jmcauley.ucsd.edu/data/amazon/>

Table 1: Overall performance. Superscripts “P”, “R” and “F1” respectively denote precision, recall, F1-score. Subscripts “std” denotes the standard deviation of each metric. Numbers in bold indicate the best results, while underlined numbers mean the second-best results. “REXHA-P” and “REXHA-L” correspond to REXHA with latent representation queries and profile queries, respectively.

Models	Explainability \uparrow				Stability \downarrow			
	GPT _{score}	BERT _{score} ^P	BERT _{score} ^R	BERT _{score} ^{F1}	GPT _{std}	BERT _{std} ^P	BERT _{std} ^R	BERT _{std} ^{F1}
Amazon-books								
NRT	75.63	0.3444	0.3440	0.3443	12.82	0.1804	0.1035	0.1321
Att2Seq	76.08	0.3746	0.3624	0.3687	12.56	0.1691	0.1051	0.1275
PETER	77.65	<u>0.4279</u>	0.3799	0.4043	11.21	0.1334	0.1035	0.1098
PEPLER	78.77	0.3506	0.3569	0.3543	11.38	0.1105	0.0935	0.0893
XRec	82.57	0.4193	0.4038	0.4122	9.60	0.0836	0.0920	<u>0.0800</u>
G-Refer (7B)	<u>82.70</u>	0.4076	0.4476	<u>0.4282</u>	9.04	0.0937	0.0845	0.0820
REXHA-P	83.28	0.3995	0.3752	0.3881	9.69	<u>0.0844</u>	<u>0.0896</u>	0.0788
REXHA-L	81.44	0.4722	<u>0.4070</u>	0.4400	8.95	0.0908	0.0989	0.0862
Yelp								
NRT	61.94	0.0795	0.2225	0.1495	16.81	0.2293	0.1134	0.1581
Att2Seq	63.91	0.2099	0.2658	0.2379	15.62	0.1583	0.1074	0.1147
PETER	67.00	0.2102	0.2983	0.2513	15.57	0.3315	0.1298	0.2230
PEPLER	67.54	0.2920	0.3183	0.3052	14.18	0.1476	<u>0.1044</u>	0.1050
XRec	74.53	0.3946	0.3506	0.3730	11.45	0.0969	0.1048	0.0852
G-Refer (7B)	<u>74.91</u>	0.3573	0.4264	0.3922	10.88	0.1050	0.0952	<u>0.0862</u>
REXHA-P	76.25	0.4879	<u>0.3604</u>	<u>0.4237</u>	<u>10.64</u>	0.1033	0.1153	0.0933
REXHA-L	74.32	0.5005	0.3603	0.4298	10.40	<u>0.0994</u>	0.1160	0.0938
Google-reviews								
NRT	58.27	0.3509	0.3495	0.3496	19.16	0.2176	0.1267	0.1571
Att2Seq	61.31	0.3619	0.3653	0.3636	17.47	0.1855	0.1247	0.1403
PETER	65.16	0.3892	0.3905	0.3881	17.00	0.2819	0.1356	0.2005
PEPLER	61.58	0.3373	0.3711	0.3546	17.17	0.1134	0.1161	0.0999
XRec	69.12	0.4546	0.4069	0.4311	14.24	0.0972	0.1163	<u>0.0938</u>
G-Refer (7B)	71.47	0.4253	0.4873	0.4566	13.46	0.1184	0.0872	0.0921
REXHA-P	<u>70.35</u>	<u>0.4565</u>	0.4200	0.4385	14.52	0.1060	<u>0.1130</u>	0.0940
REXHA-L	69.91	0.4884	<u>0.4259</u>	0.4573	<u>14.23</u>	<u>0.1001</u>	0.1179	0.0957

3.1.3 Baselines

We compare our method with the following state-of-the-art methods, including NRT [6], Att2Seq [20], PETER [21], PEPLER [22], XRec [11] and G-Refer [12]. The details of baseline can be found in Appendix C.

3.1.4 Implementation Details

For the collaborative signal extractor module, we train LightGCN with a learning rate of 1e-3 and a batch size of 1024. For hierarchical aggregation (HA) module, we set 4 reviews as a set to be summarized. For the retrieval module, we use llm-embedder model [23] to generate the embeddings of the reviews. For the generation module, we use the LLaMA-2-7B model as the base model. We set the learning rate, epochs, and batch size as 8e-4, 2 and 12, respectively. For inference, we set the temperature to 0 and the max output tokens as 128. We run the experiments on a machine with 8 NVIDIA A800 GPUs.

3.2 Overall Performance

Tab. 1 provides the overall results. It can be observed that REXHA achieves outstanding performance in explanation quality across evaluation metrics based on GPT and BERT. XRec and G-Refer take different approaches: the former leverages graph neural networks to capture collaborative filtering information and integrates it into large models, while the latter further exploits graph-structured data using a hybrid retrieval method that combines node-level and path-level strategies to more precisely utilize user-item interaction graphs. In contrast, REXHA-L, which uses latent representation query, focuses on mining the semantic information in reviews, and achieves notable improvements on the BERT_{score}^{Precision} metrics across three datasets with increases of 12.6%, 26.8% and 7.46%. Meanwhile, comparing with G-Refer on BERT_{score}^{F1}, REXHA achieves increases of 2.76%, 9.59% and 0.15%. We

Table 2: Ablation study. The best results are highlighted in bold and the worst in **RED**. “RR” represents the retrieval module, and “HA” is the hierarchical aggregation module. “Random” indicates that user/item profiles are generated by randomly sampling the user/item reviews.

Datasets		Amazon-books \uparrow		Yelp \uparrow	
RR	HA	BERT $^{\text{Precision}}$ _{score}	BERT $^{\text{F1}}$ _{score}	BERT $^{\text{Precision}}$ _{score}	BERT $^{\text{F1}}$ _{score}
w/o RR	Random	0.4570	0.4319	0.4860	0.4192
w/o RR	w/ HA	0.4563	0.4347	0.4930	0.4250
w/ RR	Random	0.4417	0.4085	0.4791	0.4191
REXHA		0.4726	0.4403	0.5031	0.4310

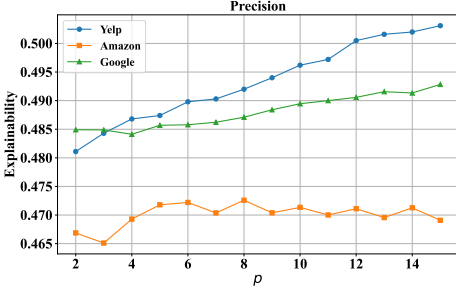


Figure 2: Performance when using different p .

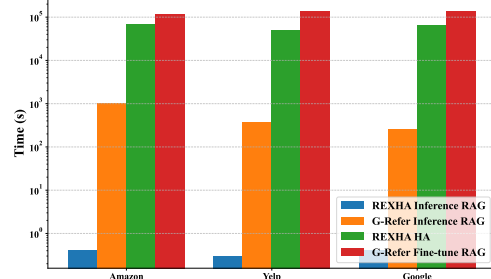


Figure 3: Comparisons of efficiency between REXHA and G-Refer.

observe that REXHA-P with profile query achieves higher GPT scores compared to REXHA-L across all three datasets, with improvements of 2.26%, 2.60%, and 0.64% on Amazon-books, Yelp, and Google-reviews, respectively. However, it performs slightly worse in terms of BERT scores.

3.3 Ablation Study

We provide the ablation study in Tab. 2. It is observed that the hierarchical aggregation module improves the performance, which proves that the profile generated by this module exactly extracts the important information from user/item reviews, benefiting the generation step. Rarely using review retrieval module without the hierarchical aggregation module get the worst results. However, when two modules are combined, the best results are obtained. This observation indicates that the deviation of profiles introduces noise and makes LLM confused when faced with the contradiction between the profile and the retrieved reviews. In contrast, with the combination of the two modules, LLM can better utilize the retrieved reviews according to correct profiles. By integrating these two modules, we achieve a compounded benefit that exceeds what each could offer alone.

3.4 Analysis of Hyperparameter Sensitivity

We further investigate how the performance of REXHA varies with the number of retrieved reviews p . Fig. 2 illustrates the impact of increasing p from 2 to 15 on BERT $^{\text{precision}}$ _{score} and BERT $^{\text{recall}}$ _{score}. We can observe that on both Yelp, BERT $^{\text{precision}}$ _{score} consistently increases as p grows. When $p > 4$, BERT $^{\text{Precision}}$ _{Score} keeps increases on Google-reviews. When $p = 8$ REXHA achieves best performance on Amazon-books. We notice that BERT $^{\text{recall}}$ _{score} remains relatively insensitive to changes of p , indicating that the profile already contains the essential information required for explanation generation.

3.5 Analysis of Retrieval Efficiency

We evaluate the preprocessing time required by REXHA and G-Refer on three datasets. Across all three datasets, the retrieval time during the inference stage for REXHA is consistently under 1 second, significantly faster than G-Refer, which requires over 4 minutes. This highlights the superior inference efficiency of REXHA. Unlike G-Refer, REXHA does not perform any retrieval during the training phase. However, the profile preprocessing step in REXHA involves the Hierarchical Aggregation

(HA) module, which incurs a relatively high computational cost—taking up to 20 hours. Despite this, it is still considerably more efficient compared to G-Refer, whose total retrieval time during training exceeds 40 hours across all three datasets. These results demonstrate that the HA module, while computationally intensive, remains more efficient than G-Refer’s training-time retrieval operations.

4 Related Work

4.1 Explainable Recommender Systems (ExRec)

With the rise of LLMs, many recent works leverage their world knowledge to generate more fluent and informative explanations. Explainable recommendation (ExRec) enhances user trust by revealing the rationale behind recommendations [4, 5], and has received increasing attention. Existing methods fall into three categories: generation-based, extraction-based, and hybrid approaches [5]. Generation-based methods (e.g., NRT [6], PETER [8], SEQUER [24]) produce explanations word-by-word from user/item representations. Extraction-based methods (e.g., ESCOFILT [9], GREENer [7]) select sentences directly from reviews. Hybrid methods (e.g., ERRA [25], ExBERT [26]) integrate both paradigms using retrieval-augmented generation. With the advent of LLMs, recent works utilize their world knowledge for more fluent and informative explanations. PEPLER applies prompt learning [22]; POD adopts prompt distillation [27]; XRec integrates collaborative signals via a lightweight adaptor [11]; and G-Refer enhances explanation quality by retrieving structured and semantic CF signals [12].

4.2 Retrieval-Augmented Generation (RAG)

Retrieval-Augmented Generation (RAG) enhances generative models by incorporating external information through retrieval. Recent research focuses on improving retrievers, generators, and their interaction. Recent RAG methods enhance LLMs by retrieving relevant knowledge and incorporating it into the input. A common approach is to concatenate retrieved documents with the original prompt. For example, In-Context RALM [28] appends retrieved content without modifying LLM parameters. SKR [29] lets the LLM choose between internal knowledge and retrieval. HyDE [30] generates a hypothetical answer to improve retrieval relevance. RichRAG [31] retrieves from multiple aspects and fuses the results for richer input.

5 Limitations

Although REXHA achieves significant improvements in explanation quality and retrieval efficiency, it still has some limitations. For HA module, the time overhead is relatively high, requiring further optimizations such as dynamic pruning of the hierarchical structure based on different user-item pairs to reduce noise. For the retrieval module, we have not yet implemented deeper post-processing of retrieved reviews, including re-ranking the review list and capturing key information.

6 Conclusion

We analyze the limitations of existing explainable recommendation methods in terms of profile generation and retrieval efficiency, and propose REXHA to address these challenges. The key components of REXHA include a hierarchical aggregation module that generates comprehensive yet concise user/item profiles, and two different query construction strategies that capture both global and fine-grained features of users and items to retrieve relevant reviews. Extensive experiments demonstrate the efficiency and effectiveness of REXHA, highlighting its promising potential for real-world applications. Future work will focus on deeper utilization of the retrieved reviews and accelerating the profile generation process.

References

- [1] Al Borchers, Jonathan L. Herlocker, and John Riedl. Ganging up on information overload. *Computer*, 31(4):106–108, 1998.
- [2] Le Wu, Xiangnan He, Xiang Wang, Kun Zhang, and Meng Wang. A survey on accuracy-oriented neural recommendation: From collaborative filtering to information-rich recommendation. *IEEE Trans. Knowl. Data Eng.*, 35(5):4425–4445, 2023.
- [3] Francesco Ricci, Lior Rokach, and Bracha Shapira, editors. *Recommender Systems Handbook*. Springer US, 2022.
- [4] Yongfeng Zhang and Xu Chen. Explainable recommendation: A survey and new perspectives. *Found. Trends Inf. Retr.*, 14(1):1–101, 2020.
- [5] Alejandro Ariza-Casabona, Ludovico Boratto, and Maria Salamó. A comparative analysis of text-based explainable recommender systems. In *RecSys*, pages 105–115, 2024.
- [6] Piji Li, Zihao Wang, Zhaochun Ren, Lidong Bing, and Wai Lam. Neural rating regression with abstractive tips generation for recommendation. In *SIGIR*, pages 345–354, 2017.
- [7] Peng Wang, Renqin Cai, and Hongning Wang. Graph-based extractive explainer for recommendations. In *WWW*, pages 2163–2171, 2022.
- [8] Lei Li, Yongfeng Zhang, and Li Chen. Personalized transformer for explainable recommendation. In *ACL/IJCNLP*, pages 4947–4957, 2021.
- [9] Reinald Adrian Pugoy and Hung-Yu Kao. Unsupervised extractive summarization-based representations for accurate and explainable collaborative filtering. In *ACL/IJCNLP*, pages 2981–2990, 2021.
- [10] Yuxuan Lei, Jianxun Lian, Jing Yao, Xu Huang, Defu Lian, and Xing Xie. Recexplainer: Aligning large language models for explaining recommendation models. In *KDD*, pages 1530–1541, 2024.
- [11] Qiyao Ma, Xubin Ren, and Chao Huang. Xrec: Large language models for explainable recommendation. In *EMNLP (Findings)*, pages 391–402, 2024.
- [12] Yuhan Li, Xinni Zhang, Linhao Luo, Heng Chang, Yuxiang Ren, Irwin King, and Jia Li. G-refer: Graph retrieval-augmented large language model for explainable recommendation. In *WWW*, pages 240–251, 2025.
- [13] Shiwen Wu, Fei Sun, Wentao Zhang, Xu Xie, and Bin Cui. Graph neural networks in recommender systems: A survey. *ACM Comput. Surv.*, 55(5):97:1–97:37, 2023.
- [14] Xiangnan He, Kuan Deng, Xiang Wang, Yan Li, Yong-Dong Zhang, and Meng Wang. Lightgcn: Simplifying and powering graph convolution network for recommendation. In *SIGIR*, pages 639–648, 2020.
- [15] Nelson F. Liu, Kevin Lin, John Hewitt, Ashwin Paranjape, Michele Bevilacqua, Fabio Petroni, and Percy Liang. Lost in the middle: How language models use long contexts. *Trans. Assoc. Comput. Linguistics*, 12:157–173, 2024.
- [16] Jiacheng Li, Jingbo Shang, and Julian J. McAuley. Uctopic: Unsupervised contrastive learning for phrase representations and topic mining. In *ACL*, pages 6159–6169, 2022.
- [17] An Yan, Zhankui He, Jiacheng Li, Tianyang Zhang, and Julian John McAuley. Personalized showcases: Generating multi-modal explanations for recommendations. In *SIGIR*, pages 2251–2255, 2023.
- [18] Jiaan Wang, Yunlong Liang, Fandong Meng, Haoxiang Shi, Zhixu Li, Jinan Xu, Jianfeng Qu, and Jie Zhou. Is chatgpt a good NLG evaluator? A preliminary study. *arXiv Preprint*, 2023. URL <https://arxiv.org/abs/2303.04048>.

- [19] Tianyi Zhang, Varsha Kishore, Felix Wu, Kilian Q. Weinberger, and Yoav Artzi. Bertscore: Evaluating text generation with BERT. In *ICLR*, 2020.
- [20] Li Dong, Shaohan Huang, Furu Wei, Mirella Lapata, Ming Zhou, and Ke Xu. Learning to generate product reviews from attributes. In *EACL*, pages 623–632, 2017.
- [21] Lei Li, Yongfeng Zhang, and Li Chen. Personalized transformer for explainable recommendation. In *ACL/IJCNLP*, pages 4947–4957, 2021.
- [22] Lei Li, Yongfeng Zhang, and Li Chen. Personalized prompt learning for explainable recommendation. *ACM Trans. Inf. Syst.*, 41(4):103:1–103:26, 2023.
- [23] Tianhua Zhang, Hongyin Luo, Yung-Sung Chuang, Wei Fang, Luc Gaitskell, Thomas Hartvigsen, Xixin Wu, Danny Fox, Helen Meng, and James R. Glass. Interpretable unified language checking. *arXiv Preprint*, 2023. URL <https://arxiv.org/abs/2304.03728>.
- [24] Alejandro Ariza-Casabona, Maria Salamó, Ludovico Boratto, and Gianni Fenu. Towards self-explaining sequence-aware recommendation. In *RecSys*, pages 904–911, 2023.
- [25] Hao Cheng, Shuo Wang, Wensheng Lu, Wei Zhang, Mingyang Zhou, Kezhong Lu, and Hao Liao. Explainable recommendation with personalized review retrieval and aspect learning. In *ACL*, pages 51–64, 2023.
- [26] Huijing Zhan, Ling Li, Shaohua Li, Weide Liu, Manas Gupta, and Alex C. Kot. Towards explainable recommendation via bert-guided explanation generator. In *ICASSP*, pages 1–5, 2023.
- [27] Lei Li, Yongfeng Zhang, and Li Chen. Prompt distillation for efficient llm-based recommendation. In *CIKM*, pages 1348–1357, 2023.
- [28] Ori Ram, Yoav Levine, Itay Dalmedigos, Dor Muhlgay, Amnon Shashua, Kevin Leyton-Brown, and Yoav Shoham. In-context retrieval-augmented language models. *Trans. Assoc. Comput. Linguistics*, 11:1316–1331, 2023.
- [29] Yile Wang, Peng Li, Maosong Sun, and Yang Liu. Self-knowledge guided retrieval augmentation for large language models. In *EMNLP (Findings)*, pages 10303–10315, 2023.
- [30] Luyu Gao, Xueguang Ma, Jimmy Lin, and Jamie Callan. Precise zero-shot dense retrieval without relevance labels. In *ACL*, pages 1762–1777, 2023.
- [31] Shuting Wang, Xin Yu, Mang Wang, Weipeng Chen, Yutao Zhu, and Zhicheng Dou. Richrag: Crafting rich responses for multi-faceted queries in retrieval-augmented generation. In *COLING*, pages 11317–11333, 2025.

A Datasets

Table 3: Statistics of the experimental datasets

Datasets	#Users	#Items	#Interactions
Amazon	15,349	15,247	360,839
Yelp	15,942	14,085	393,680
Google	22,582	16,557	411,840

We use three public datasets in the experiments and Tab. 3 provides the data statistics.

- **Amazon-books** comes from the Amazon book platform and contains user reviews, ratings, metadata, and other information on books. It can reflect user reading preferences and product popularity.
- **Yelp** comes from the American review website Yelp, where local businesses (such as restaurants and bars) are considered as items. The dataset provides rich information about local businesses, including user review text and ratings under multiple categories.
- **Google-reviews** comes from user reviews of businesses (such as restaurants, stores, attractions, etc.) on Google, including ratings, review text and timestamps, as well as business metadata.

B Metrics

- **GPT_{score}** [18] utilizes LLMs to evaluate text quality. We regard GPT-3.5-Turbo ChatGPT as a human evaluator and give task-specific (e.g., summarization) and aspect-specific (e.g., relevance) instruction to prompt ChatGPT to evaluate the generated results of NLG models.
- **BART_{score}** [19] assesses the similarity between each token in the text by using the contextual embeddings generated by the pre-trained BERT model and calculating the sum of the cosine similarities of the token embeddings between two sentences.

Given a reference sentence $x = \langle x_1, x_2, \dots, x_n \rangle$ and a candidate sentence $\hat{x} = \langle \hat{x}_1, \hat{x}_2, \dots, \hat{x}_m \rangle$. Bert computes word embeddings sequence for the reference sentence and the candidate sentence as follows:

$$\mathbf{x} = \langle \mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n \rangle = \text{BERT}(x = \langle x_1, x_2, \dots, x_n \rangle) \quad (11)$$

$$\hat{\mathbf{x}} = \langle \hat{\mathbf{x}}_1, \hat{\mathbf{x}}_2, \dots, \hat{\mathbf{x}}_m \rangle = \text{BERT}(x = \langle \hat{x}_1, \hat{x}_2, \dots, \hat{x}_m \rangle) \quad (12)$$

The cosine similarity of a reference token x_i and a candidate token \hat{x}_j is $\frac{\mathbf{x}_i^\top \hat{\mathbf{x}}_j}{\|\mathbf{x}_i\| \|\hat{\mathbf{x}}_j\|}$. As both embeddings are pre-normalized, the formula is reduced to the inner product $\mathbf{x}_i^\top \hat{\mathbf{x}}_j$. Then the $\text{BERT}_{\text{score}}^{\text{Precision}}$, $\text{BERT}_{\text{score}}^{\text{Recall}}$ and $\text{BERT}_{\text{score}}^{\text{F1}}$ are computed as follows:

$$\text{BERT}_{\text{score}}^{\text{Precision}} = \frac{1}{|x|} \sum_{x_i \in x} \max_{\hat{x}_j \in \hat{x}} \mathbf{x}_i^\top \hat{\mathbf{x}}_j, \quad (13)$$

$$\text{BERT}_{\text{score}}^{\text{Recall}} = \frac{1}{|\hat{x}|} \sum_{\hat{x}_j \in \hat{x}} \max_{x_i \in x} \mathbf{x}_i^\top \hat{\mathbf{x}}_j, \quad (14)$$

$$\text{BERT}_{\text{score}}^{\text{F1}} = 2 \frac{P_{\text{BERT}} \cdot R_{\text{BERT}}}{P_{\text{BERT}} + R_{\text{BERT}}}. \quad (15)$$

C Baselines

- **NRT** [6] tackles both rating prediction and tip generation by learning shared representations from user and item IDs through a joint optimization approach. It utilizes a GRU to produce concise, abstractive tips.

- **Att2Seq** [20] builds upon an attribute-to-sequence framework where an attention mechanism helps the model focus on relevant input features. A stacked LSTM is employed for decoding the review text.
- **PETER** [21] introduces a Transformer-based approach that personalizes review generation by linking ID-based representations of users and items with natural language output. Due to dataset limitations, the basic version without auxiliary word features is applied.
- **PEPLER** [22] leverages a pretrained language model to generate textual explanations. It refines this process with techniques like sequential adaptation and regularization to narrow the gap between the prompt structure and the language model’s expectations.
- **XRec** [11] enhances text generation by injecting collaborative filtering signals from GNN-based user and item encoders into every layer of a language model, enabling more contextually relevant and personalized content.
- **G-Refer** [12] strengthens explainable recommendation by retrieving and translating collaborative filtering signals from user-item graphs into human-readable text, enabling large language models to generate personalized explanations through retrieval-augmented fine-tuning.

D Analysis of Hierarchical Aggregation

To validate the effectiveness of the Hierarchical Aggregation (HA) method, we compare it with the “Direct” method, where the LLM is invoked only once to generate the profile based on all reviews at once. We also compare the performance of Qwen2.5-7B with Qwen2.5-7B-1M, which supports a 1M-token context length.

Tab. 4 reports the results. We observe that Qwen2.5-7B fails to generate profiles using the Direct method since the input exceeds its context length. In contrast, the HA method outperforms the Direct method, achieving improvements of 1.57%, 0.17%, and 1.02% on $BERT_{score}^P$, $BERT_{score}^R$, and $BERT_{score}^{F1}$, respectively. Additionally, we evaluate an alternative approach in which second-layer summaries are directly concatenated to form the profiles. This method achieves slightly higher BERT precision than HA (an increase of 0.08%), but results in a slight drop in recall (a decrease of 0.19%).

Table 4: Contrast study for Hierarchical Aggregation module. We evaluate the performance without Review Retrieval module. In this table, “Directly” denotes that for all reviews of user/item, LLMs are invoked once to generate profiles, and “Second Layer” denotes that sub-nodes of the final node are used as profiles.

Models	Type	$BERT_{score}^P$	$BERT_{score}^R$	$BERT_{score}^{F1}$
Qwen2.5-7B-1M	HA	0.4712	0.3576	0.4142
Qwen2.5-7B-1M	Directly	0.4639	0.3570	0.4100
Qwen2.5-7B-1M	Second Layer	0.4716	0.3557	0.4134
Qwen2.5-7B	Directly	—	—	—
Qwen2.5-7B	HA	0.4930	0.3580	0.4250

E Comparison of Retrieval Results between Latent Representation Query and Profile Query

To examine the differences between the latent representation query and the profile query, we sampled a set of user-item pairs and calculated the pairwise similarity among the reviews retrieved by each query type. The results are shown in Fig. 4 and Fig. 5. As observed, the reviews retrieved by the latent representation query exhibit lower similarity compared to those retrieved by the profile query, indicating that the former leads to more diverse retrievals. In contrast, the profile query yields more semantically concentrated results, better aligning with long-term user and item preferences.

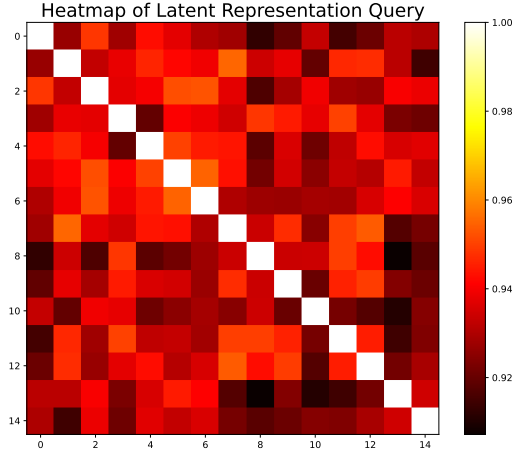


Figure 4: Similarity of each reviews retrieved by latent representation query.

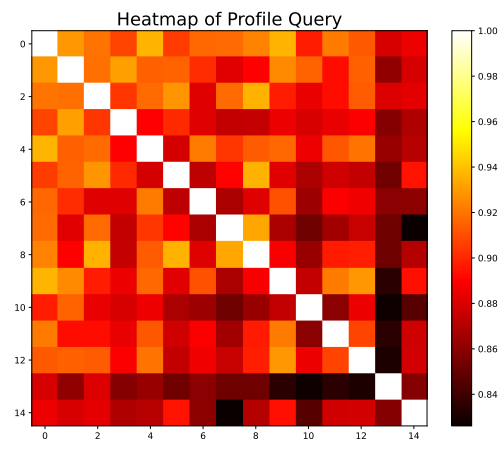


Figure 5: Similarity of each reviews retrieved by profile query.

F Prompt Template

Fig. 6 demonstrates how LLMs generate item profiles on Yelp by combining item metadata and user reviews. This enables the model to identify key user characteristics associated with the item, offering deeper insights into user preferences and improving the understanding of user-item interactions.

System prompt:

You will serve as an assistant to help me summarize which types of users would enjoy a specific business. I will provide you with the basic information (name) of that business and also some reviews from users for it.

Here are the instructions:

1. The basic information will be described in JSON format, with the following attributes:

{ "name": "the name of the business" }

2. Reviews from users will be managed in the following List format:

```
[
  "the first review",
  "the second review",
  ....
]
```

3. The information I will give you:

BASIC INFORMATION: a JSON string describing the basic information about the business.

USER REVIEWS: a List object containing some reviews from users about the business.

Requirements:

1. Please provide your answer in JSON format, following this structure:

{ "summarization": "A summarization of what types of users would enjoy this business." (if you are unable to summarize it, please set this value to "None") }

2. Please ensure that the "summarization" is no longer than 50 words.

3. Do not provide any other text outside the JSON string.

Output:

{ "summarization": "Users who appreciate ... would enjoy this business." }

Figure 6: Prompt for Item Profile

Fig. 7 presents a method for generating user profiles on the Yelp dataset using LLMs. By combining various item metadata, description (i.e., generated item profiles) and user reviews, the system builds detailed user profiles that capture preferences and support personalized recommendations.

The generation process employs a structured prompt, as illustrated in Fig. 8, combining multiple data components such as retrieved reviews, user/item embeddings, and user/item metadata. The prompt is tokenized and transformed into an embedding space representation. To distinguish special tokens within the prompt as distinct entities, we integrate them into the LLM's tokenizer. These tokens are subsequently replaced with their corresponding modified embeddings in the final token embedding representation.

System prompt:

You will serve as an assistant to help me determine which types of business a specific user is likely to enjoy. I will provide you with information about businesses that the user has interacted with, including their reviews of those businesses.

Here are the instructions:

1. Each interacted business will be described in DICTIONARY format, with the following attributes:

```
{
  "title": "the name of the business", (if there is no business, I will set this value to "None")
  "description": "a description of what types of users will like this business",
  "reviews": "the user's review on the business" (if there is no review, I will set this value to "None")
}
```

2. The information I will give you:

PURCHASED BUSINESSES: a list of dictionaries describing the businesses that the user has interacted with.

Requirements:

1. Please provide your answer in JSON format, following this structure:

```
{ "summarization": "A summarization of what types of business this user is likely to enjoy." (if you are unable to summarize it, please set this value to "None") }
```
2. Please ensure that the "summarization" is no longer than 50 words.
3. Do not provide any other text outside the JSON string.

Output:

```
{ "summarization": "This user enjoys ... experiences, ... service, ... atmospheres." }
```

Figure 7: Prompt for User Profile

System prompt:

Explain why the user would buy with the book within 50 words.

User prompt:

Here are some comments from similar users and items, you can use them to help you write the review.

1. The user would enjoy the business because of the ...
2. This business would be enjoyed by the user ...

...

user record: <USER_EMBED> item record: <ITEM_EMBED>

item name: ... user profile: ... item profile: ...

Output:

The user would enjoy the business because of...

Figure 8: Prompt for Reviews Summarization.

G Case Study

We present two cases in Tab. 5 and Tab. 6 to demonstrate the effectiveness of our method and illustrate how hierarchical aggregation based profile and retrieved reviews benefit the generated explanations. The table only shows some of the reviews on the retrieved content. We also provide the ground truth explanations and explanations generated by XRec [11] and G-Refer [12] for comparison.

From Tab. 5, we can observe that our method hierarchically aggregates the item profile to highlight “high quality seafood and steak” and “good ambiance”. This allows our model to integrate the item profile and the retrieved related reviews to generate an explanation that covers both explicit and inferred preferences. In contrast, XRec incorrectly explains that the restaurant atmosphere is not as good as expected, failing to capture true interests, while G-Refer’s explanation is too general, focusing only on common aspects such as “food service” without touching on users’ preferences for specific food or atmosphere. This illustrates that the profile generated by aggregated reviews captures richer and more accurate personalized information than that generated by randomly sampled reviews.

Tab. 6 provides another case that leverages retrieved similar reviews to enrich explanations. While the user profile and item profile simply state their preference for emotionally rich "romance genres", the retrieved reviews reveal deeper elements that appeal to users with similar interests, such as "intense moments", "satisfying endings", and "captivating love stories". Our model successfully incorporates these signals, capturing both emotional depth and structural highlights that users may value. In contrast, the explanations generated by XRec are more vague, and G-Refer only captures "intense" information and introduces irrelevant information about the author, which is not relevant to the result. Together, these two cases demonstrate how retrieved similar reviews and hierarchical aggregate profiles complement each other to produce accurate, diverse, and user-consistent explanations.

Table 5: A case from Yelp, which shows how REXHA leverages HA profile to generate better explanations.

Inputs	
User Profile	Enjoys craft beers, lively atmospheres, diverse food and drink options, and family-friendly attractions with interactive experiences.
Item Profile	Users who enjoy high-quality seafood and steaks , seek good service and value, and appreciate unique dining experiences will enjoy McCormick & Schmick's. Those looking for affordable yet tasty options and a good ambiance during special event days will also find it enjoyable.
Retrieved Reviews	["The user would enjoy the business because it offers a wide variety of options from cheap produce, meat, and seafood to diverse dining choices with the added convenience of a bar, making it a great place for a fun day out with friends.", "The user would enjoy the business for its diverse beer selection, great customer service, and delicious food options, making it a welcoming and enjoyable spot for sampling different beers and relaxing outdoors.", "The user would enjoy the business for its interesting beer selection, tasty food like fish and chips, and a cozy atmosphere, despite its modest setting, making it a great spot to explore different dishes and enjoy a casual dining experience.", ...]
Explanations	
Ground Truth	The user would enjoy the business for its upscale ambiance, delicious steak and seafood options , excellent cocktails, and the chance to spot a celebrity, making it a versatile spot for meeting friends, dates, or work colleagues.
XRec	The user would enjoy the business for its great food and service, even though the atmosphere is not as good as expected. The menu has a variety of options, and the crab cakes are a highlight. The server was attentive and helpful, making the dining experience enjoyable overall.
G-Refer	The user would enjoy the business for its delicious food, especially the crab cakes, and the attentive service provided by the server, making the dining experience enjoyable and memorable.
REXHA (ours)	The user would enjoy the business for its delicious food options, including seafood and steak, and the great atmosphere , making it a perfect spot for a special occasion or a casual dining experience.

Table 6: A case from Amazon, which shows how REXHA generates better explanation with HA profile and retrieved reviews.

Inputs	
User Profile	This user prefers emotionally rich, character-driven stories with strong emotional impacts, often in romance genres .
Item Profile	Fans of romantic suspense , emotional stories, and captivating duets interested in thrilling, suspenseful tales with deep characters and unexpected twists would buy this book.
Retrieved Reviews	["The user would buy the book because they enjoy a captivating love story with complex characters and passionate dynamics, as well as the emotional roller coaster ride provided by the book's plot and the strong emotional investment in the characters.", "The user would buy the book because it features captivating characters that evoke deep emotional connections, a compelling storyline that elicits intense reactions , and a satisfying ending that leaves readers feeling ecstatic and emotionally fulfilled.", "The user would buy the book because it provides a captivating and emotional love story , filled with intense moments and deep character development, which keeps the reader completely engaged and eager for more.", ...]
Explanations	
<u>Ground Truth</u>	The user would buy the book because it promises a thrilling and unexpected roller-coaster ride of a love story with intense twists and turns, a satisfying ending , and relatable themes that make readers question their own love lives.
<u>XRec</u>	The user would buy the book because it offers a compelling and emotionally charged storyline, with well-developed characters and a captivating plot that keeps the reader engaged and invested in the characters' journey, making it a must-read.
<u>G-Refer</u>	The user would buy the book because they are a fan of the author, Jay Crownover, and have enjoyed her previous works. They are drawn to the intense and emotional storyline, the complex characters, and the captivating writing style, making them eager to experience the full impact of the story.
<u>REXHA (ours)</u>	The user would buy the book because it offers a captivating and emotional love story with complex characters, intense moments , and a satisfying ending , making it a must-read for fans of romance novels.