

Adaptive and Representative Multi-Interest Modeling for Recommendation with Large Language Model

Anonymous ACL submission

Abstract

Large language models (LLMs) show potential for multi-interest analysis of users in recommender systems, going beyond heuristic assumptions in existing methods, e.g., co-occurring items indicate the same interest. Despite the effectiveness, two key challenges remain. First, the granularity of raw generation of LLMs for multi-interests is agnostic, possibly leading to overly fine or coarse interest grouping. Second, adopting LLM to analyze individual user behaviors lacks a global perspective on how items relate across users. In this paper, we propose an LLM-driven adaptive and representative multi-interest modeling framework to address the challenges. At the user-individual level, we exploit LLM analysis and alleviate the agnostic granularity by adaptively aggregating semantic clusters to collaborative multi-interests. At the user-crowd level, to mitigate the limited insights in individual behaviors, we formulate a max covering problem to expand the scope of LLM analysis with compactness and representativeness, disentangling interest representations from global perspectives. Experiments on real-world datasets show that our approach outperforms various baselines.

1 Introduction

Recommender systems (RSs) play a crucial role in personalized user experience, while modeling users' multi-faceted and dynamic interests remains challenging. To bridge this gap, multi-interest methods (Li et al., 2019; Zhang et al., 2022) use multiple representations for each user to capture users' interest facets behind their behaviors.

Despite their effectiveness, existing methods often rely on heuristic assumptions, such as similar items indicate the same interest for users, where similarity can be measured through item embeddings (Ma et al., 2020), co-occurrence statistics (Du et al., 2024b), or auxiliary information (Chai et al., 2022). However, due to the sparsity in user-

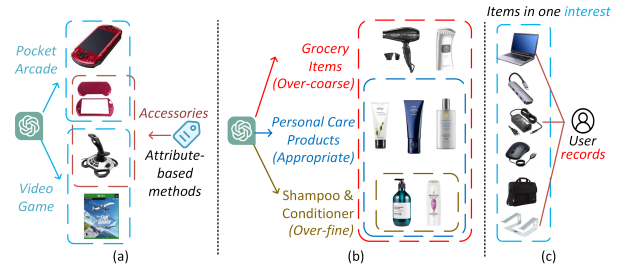


Figure 1: (a) Difference between attribute- and LLM-driven multi-interest analysis. (b) LLM-driven multi-interest analysis leads to varying granularity. (c) Individual user's behaviors lack global item relationships.

item interactions and incompleteness in auxiliary data, accurately measuring item similarity and user interests remains a significant challenge. As in Figure 1(a), though the items within the red box belong to the 'Accessories' attribute, they significantly differ in functionalities and appeal to different user interests. To address these problems, we leverage large language models (LLMs) with their rich knowledge and powerful reasoning capabilities, analyzing users' semantic multi-interests beyond their interaction records and auxiliary information. LLMs can offer semantic guidance for multi-interest extraction, providing more accurate guidance (blue boxes) on multi-interest analysis.

Although leveraging LLMs offers a promising way for multi-interest modeling, directly using them as a black-box is not a one-size-fits-all solution, as there remain two significant challenges. **First**, the granularity of LLM-driven multi-interest analysis is agnostic, i.e., overly-fine or overly-coarse item division among users' behaviors, making it hard to model their multi-interests accurately. As in Figure 1(b), LLMs might categorize interests with overly-fine distinctions among items that should belong to the same interest ('Shampoo & Conditioner' in brown), or with overly-coarse groupings that fail to capture interest discrimina-

tion among items (‘Grocery Items’ in red), while only the ‘Personal Care Products’ (in blue) is the desired or ideal interest. **Second**, multi-interest analysis for individual users lacks a global perspective on item relations across entire population. Specifically, relying solely on individual user interactions analyzed by LLMs, which are inherently sparse, may lead to incomplete modeling of multi-interests. As in Figure 1 (c), besides the engaged items in users’ behaviors connected by red lines, there remain multiple non-co-occurring items that may also reflect the same user interests.

To address the first challenge, we guide the LLM-driven multi-interest analysis at user-individual level using a tailored prompt to cluster each user’s interaction sequence into distinct semantic interest groups. We then introduce an interpretable alignment module that dynamically aggregates LLM-driven semantic clusters and maps them into collaborative interests learned by capsule network. By doing so, we adaptively adjust the granularity to fit user patterns. For the second challenge, we generate synthetic users for LLM-driven analysis at the user-crowd level, ensuring the compactness (limited number of interests with coherent behaviors) and representativeness (covering multiple items). We achieve this by clustering real users with similar preference and solving a max covering problem (MCP) to select synthetic users that span the item space. Finally, we introduce contrastive learning to encourage item concentration within interests and dispersion across interests, enhancing multi-interest representations.

Our key contributions are three-fold. Firstly, we propose a novel LLM-driven Adaptive and Representative Multi-Interest (LARMI) modeling framework to explore semantic information from both the user-individual level and user-crowd level, for more effective multi-interest recommendation. Secondly, We address the issues of LLM’s agnostic granularity by designing an adaptive alignment module and MCP optimization with contrastive learning to ensure the compactness and representativeness from a global perspective. Thirdly, We evaluate the proposed method across three real-world datasets to demonstrate its effectiveness in multi-interest modeling.

2 Literature Review

Single- and Multi-Interest Modeling for Recommendation. Single-interest modeling in recom-

mendation include methods built upon recurrent neural networks (Hidasi et al., 2016; Guo et al., 2020), self-attention (Kang and McAuley, 2018; Zhang et al., 2019), transformers (Sun et al., 2019; Xia et al., 2021), and graph convolutional networks (He et al., 2020; Wang et al., 2025), but these methods overlook the diversity in user interests. Current multi-interest modeling methods that adopt capsule networks (Sabour et al., 2017; Li et al., 2019; Xie et al., 2023; Tian et al., 2022) solely rely on users’ engaged items. Attention mechanisms are also incorporated (Cen et al., 2020; Xiao et al., 2020). Regularization strategies (Zhang et al., 2022; Lee et al., 2024) stabilize the learning of multiple embeddings. Diffusion model (Le et al., 2025) and item partition objective (Du et al., 2024b) are applied for interest-aware denoising and enhancement. Other methods utilize auxiliary sources such as users’ profiles (Chai et al., 2022), timestamps (Chen et al., 2021), items’ categories (Liu et al., 2024), and knowledge graphs (Liu et al., 2022).

Large Language Models for Recommendation.

LLMs’ success has inspired their incorporation in recommendation pipelines (Wu et al., 2024). LLM-as-recommender methods employ LLMs as scoring functions or rankers. Methods fully fine-tune the LLM parameters (Geng et al., 2022; Qu et al., 2024) or conduct parameter-efficient fine-tuning (PEFT) (Bao et al., 2023; Jiang et al., 2025) bridge the gap between LLMs and recommendation tasks. Non-tuning methods align the recommendation objectives for LLMs through zero-shot prompting (Dai et al., 2023; Hou et al., 2024) and in-context learning (Sanner et al., 2023; Bao et al., 2025) strategies. LLM-as-extractor methods apply LLMs for data augmentation. Studies focus on encoding historical behaviors and item attributes produce expressive embeddings (Wang et al., 2024; Harte et al., 2023) to capture complex semantic information. Besides, several methods adopt LLMs to extract additional knowledge such as user profiles (Zheng et al., 2023; Du et al., 2024a), item descriptions (Ren et al., 2024; Wei et al., 2024), and other textual data (Mohbat and Zaki, 2025) through semantic mining.

3 Methodology

This section presents **LARMI**, the proposed LLM-based Adaptive and Representative Multi-Interest modeling framework, which leverages LLM-driven semantics for multi-interest modeling in RSs. We denote the set of M users as $\mathcal{U} = \{u_1, \dots, u_M\}$

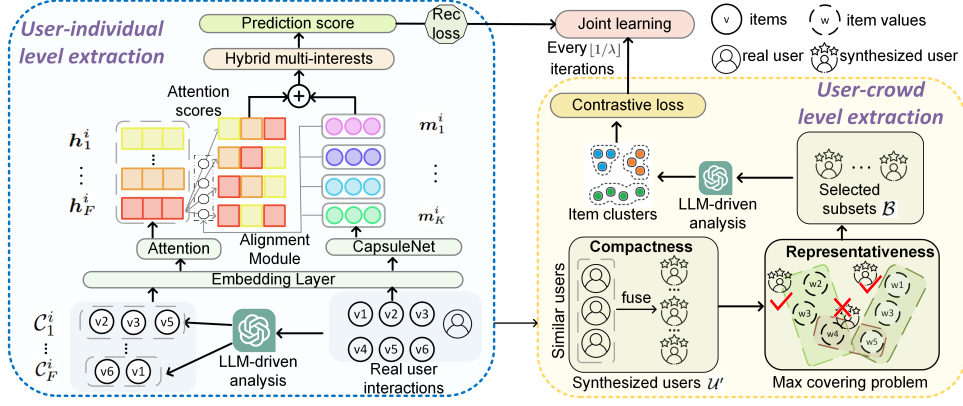


Figure 2: The overall architecture of our proposed LARMI.

and the set of N items as $\mathcal{V} = \{v_1, \dots, v_N\}$. Each user $u_i \in \mathcal{U}$ has a sequential behavior series sorted by timestamps denoted as $s(u_i) = \{v_1^i, v_2^i, \dots, v_L^i\}$, where v_j^i denotes the j -th item engaged by the user u_i and L is the length of the sequence. Besides, we suppose to know item titles in $s(u_i)$ denoted as $t(u_i) = \{t_1^i, t_2^i, \dots, t_L^i\}$. Given the user's historical behavior, we aim to generate a top- n ranking list containing items that the user is likely to engage in the recent future.

3.1 Model Overview

For the user-individual level, we employ the LLM to analyze sequential behaviors for each user, infer distinctive and meaningful semantic multi-interests, and adaptively align with collaborative interests for proper granularity. For the user-crowd level, we synthesize compact and representative users with the MCP optimization, and then bridge the gap between real and synthesized users to expand the LLM analysis scope beyond individual users. Figure 2 shows the overall framework of LARMI.

3.2 User-individual Multi-interest Extraction

We propose to leverage the semantic knowledge of the LLM to guide multi-interest extraction, overcoming the limitation of heuristic assumptions such as co-occurring items implying the same interest of users. Specifically, we prompt the LLM to conduct multi-interest analysis as follows, generating distinctive semantic clusters, each representing a cohesive set of items with shared characteristics.

$$C_1^i, \dots, C_F^i = LLM(prompt, t(u_i)), \quad (1)$$

where $prompt$ denotes the multi-interest analysis prompt. C_f^i denotes the f -th cluster that contains items belonging to the same semantic group by

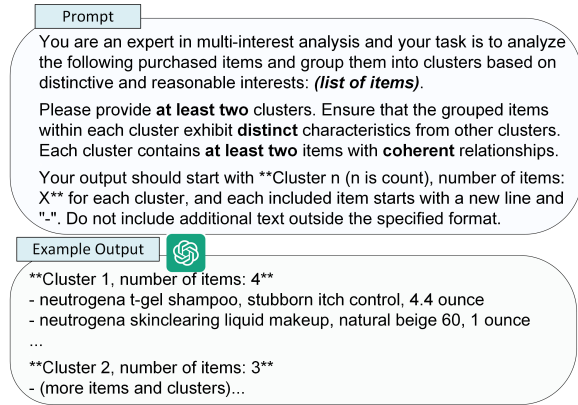


Figure 3: Example output for LLM analysis.

LLM for user u_i . F is an unknown varying number relying on LLMs' analysis and reasoning. We show the prompt and example output in Figure 3.

Intuitively, different clusters reflect distinct aspects (e.g., functionality, preference) for the user. However, the granularity of LLM-driven clusters is agnostic, making it uninterpretable and hard to effectively model multi-interests with overly-fine or overly-coarse clusters. To address this, we propose an adaptive alignment module consisting of an attention mechanism and a projection layer. For over-coarse clusters in LLMs' analysis, we use an attention mechanism to dynamically aggregate the items' representations in a cluster C_f^i , allowing more specific signals to dominate the cluster and sharpen the encoding. The LLM-driven multi-interest representation is defined as:

$$h_f^i = \sum_{v_j \in C_f^i} \alpha_j \cdot v_j, \quad \alpha_j = \frac{\exp(\mathbf{w}^T v_j + b)}{\sum_{v_k \in C_f^i} \exp(\mathbf{w}^T v_k + b)}, \quad (2)$$

where v_j denotes item v_j 's learned ID embedding, and $\mathbf{w} \in \mathbb{R}^d$, $b \in \mathbb{R}$ are learnable weights. Second, the over-fine clusters can be averaged in the pre-

vious step, and further merged with collaborative interests based on their similarity. We use capsule network to model the collaborative multi-interests:

$$\mathbf{m}_1^i, \dots, \mathbf{m}_K^i = \text{CapsuleNet}([\mathbf{v}_1^i, \dots, \mathbf{v}_L^i]), \quad (3)$$

where \mathbf{v}_j^i denotes the ID embedding of the j -th item engaged by the user u_i , and K is the predefined number of users' multi-interests. Then, we compute attention scores between pairs of semantic and collaborative interest facets for alignment:

$$\mathbf{z}_k^i = \sum_{f=1}^F \alpha_{kf} \cdot \mathbf{h}_f^i, \quad \alpha_{kf} = \frac{\exp(\mathbf{m}_k^i \cdot \tanh(\mathbf{W}_1 \mathbf{h}_f^i))}{\sum_{f'} \exp(\mathbf{m}_k^i \cdot \tanh(\mathbf{W}_1 \mathbf{h}_{f'}^i))},$$

where \mathbf{z}_k^i is the aggregation of the semantic clusters related to interests \mathbf{m}_k^i of user u_i , α_{kf} is the attention score between LLM-derived semantic clusters \mathbf{h}_f^i and collaborative interests \mathbf{m}_k^i , $\mathbf{W}_1 \in \mathbb{R}^{d \times d}$ is a learnable projection matrix, and $\tanh(\cdot)$ introduces non-linearity to better capture complex cross-representation relationships.

Therefore, our alignment module allows the cluster granularity to be adaptively adjusted, where specific items can be emphasized through attention weights for sharper interest representation for over-coarse clusters, while similar attention weights can be learned to enable merging over-fine clusters. To take advantage of LLM-driven semantics and collaborative signals, we aggregate them as follows:

$$\mathbf{o}_k^i = \mathbf{m}_k^i + \mathbf{z}_k^i. \quad (4)$$

Thus, we obtain the final hybrid multi-interest representations $\{\mathbf{o}_1^i, \dots, \mathbf{o}_K^i\}$ for each user u_i .

3.3 User-crowd Multi-interest Extraction

While user-individual multi-interest modeling leverages LLMs for semantic analysis, it lacks a representative global perspective. Individual users typically interact with only a small subset of items, inevitably preventing all related items from being grouped to the same interest. To address this, we synthesize users with richer behaviors for a more comprehensive LLM analysis. However, simply synthesizing users through random item selection produces dispersed interests and largely increases the scale and cost of LLM inference. To this end, we propose to synthesize users with the consideration of compactness and representativeness principles. **Compactness** ensures that synthesized users have focused interests, with each containing a cohesive set of semantically related items. Otherwise,

aggregating unrelated behaviors can lead to fragmented interests and generate sparse, ineffective interest clusters. **Representativeness** maximizes the coverage of unique items across all interests. This avoids redundant LLM analysis and enhances the generalization capability to better represent real-world interest diversity by the synthesized users.

To achieve **compactness** for user synthesis, we formulate a max covering problem (MCP) by grouping cliques of users with overlapping preferences and combine their behaviors to synthesize a user. Specifically, a clique $c(u'_i)$ can be generated by clustering similar users w.r.t. a real user u_i , i.e.,

$$c(u'_i) = \bigcup_{u_g \in \mathcal{N}(u_i)} s(u_g), \quad (5)$$

where $\mathcal{N}(u_i)$ denotes users who share the most overlapped behaviors to the user u_i . Therefore, each clique can be regarded as the union set of items of a compact synthesized user u'_i . We allow items to belong to multiple interest clusters, enabling LLM-driven analysis to successfully detect distinct intents. However, prompting LLM for all synthesized users leads to redundancy and inefficiency. To this end, we propose selecting representative users for LLM-driven multi-interest analysis. To further achieve **representativeness**, we select a small portion of synthesized users covering as many valuable items as possible across all interests, i.e.,

$$\max_{|\mathcal{B}| \leq Z, \mathcal{B} \subset \mathcal{U}'} \sum_{j=1}^N \mathbb{I}(v_j \in \bigcup_{u'_i \in \mathcal{B}} c(u'_i)) \cdot w_{v_j}, \quad (6)$$

where \mathcal{U}' is the set of synthesized users, \mathcal{B} is the selected representative synthesized users with maximal size Z , and w_{v_j} is the value for covering item v_j . Assuming popular items have higher values because they have more impacts, we formulate values and construct a synthesized user-item interaction matrix $\mathbf{A} \in \mathbb{R}^{M \times N}$ based on compactness:

$$w_{v_j} = 1 + \frac{\sum_{i=1}^M \mathbb{I}(v_j \in s(u_i))}{\sum_{i=1}^M |s(u_i)|}, \quad \mathbf{A}_{ij} = \begin{cases} 1 & \text{if } v_j \in c(u'_i), \\ 0 & \text{otherwise,} \end{cases} \quad (7)$$

where each row indicates the behaviors of a synthesized user based on the corresponding real user. We then formally formulate Equation (6) into a standard MCP. Specifically, we propose to represent the solution \mathcal{B} with an indicator vector $\mathbf{x} \in \{0, 1\}^M$, where x_i denotes the i -th element of \mathbf{x} , showing whether the synthesized user u'_i is included in \mathcal{B} .

Then, Equation (6) can be reformulated as:

$$\begin{aligned}
& \max_{|\mathcal{B}| \leq Z, \mathcal{B} \subset \mathcal{U}'} \sum_{j=1}^N \mathbb{I}(v_j \in \bigcup_{u'_i \in \mathcal{B}} c(u'_i)) \cdot w_{v_j} \\
& \Leftrightarrow \max_{|\mathcal{B}| \leq Z, \mathcal{B} \subset \mathcal{U}'} \sum_{j=1}^N \mathbb{I} \left(\sum_{u'_i \in \mathcal{B}} \mathbb{I}(v_j \in c(u'_i)) \right) \cdot w_{v_j} \\
& \Leftrightarrow \max_{|\mathcal{B}| \leq Z, \mathcal{B} \subset \mathcal{U}'} \sum_{j=1}^N \mathbb{I} \left(\sum_{u'_i \in \mathcal{B}} \mathbf{A}_{ij} \right) \cdot w_{v_j} \\
& \Leftrightarrow \max_{\mathbf{x} \in \{0,1\}^M, \|\mathbf{x}\|_0 \leq Z} \sum_{j=1}^N \mathbb{I} \left(\sum_{i=1}^M \mathbf{x}_i \cdot \mathbf{A}_{ij} \right) \cdot w_{v_j}.
\end{aligned} \tag{8}$$

To solve the MCP, we adopt a differentiable optimal transport model to collect the compact and representative synthesized users \mathcal{B} following (Wang et al., 2022). Given these synthesized users with rich behaviors, we propose to trigger the LLM to generate distinctive and comprehensive interest clusters for each synthetic user $u'_i \in \mathcal{B}$, i.e.,

$$\mathcal{C}_1^i, \dots, \mathcal{C}_F^i = LLM(pmt, t(u'_i)) \tag{9}$$

where $t(u'_i)$ represents the titles of items within the synthesized user u'_i 's behaviors $c(u'_i)$, and \mathcal{C}_f^i is the f -th cluster generated by the LLM.

As simply producing multi-interest representations of synthesized users has limited impact on real users, we leverage their LLM-driven multi-interests through contrastive learning to refine item distributions, eventually contributing to real users bridged by item representations in a global view. Specifically, we encourage item similarity within the same clusters and dispersion among different clusters in the representation space:

$$\mathcal{L}_{u'_i}^{cst} = - \sum_{v_j \in c(u'_i)} \sum_{v_j^* \in \mathcal{C}^i(v_j)} \log \frac{e^{(v_j^\top \cdot v_j^* / \tau)}}{\sum_{v_j' \in c(u'_i) / \mathcal{C}^i(v_j)} e^{(v_j^\top \cdot v_j' / \tau)}},$$

where for each synthesized user u'_i , $\mathcal{C}^i(v_j)$ denotes the cluster which contains the item v_j , v_j^* denotes a positive sample that belongs to one same cluster as v_j , and v_j' denotes negative samples that do not occur with v_j . τ controls the sharpness of the similarity distribution. Therefore, for each item in an interest cluster, intra-cluster items are positive instances, and inter-cluster items are hard negative instances. We aggregate the contrastive learning for all synthesized users as the overall loss,

$$\mathcal{L}^{cst} = - \frac{1}{|\mathcal{B}|} \sum_{u'_i \in \mathcal{B}} \mathcal{L}_{u'_i}^{cst}. \tag{10}$$

In summary, we synthesize a compact and representative subset of users with rich behaviors and formulate an MCP optimization for LLM-driven analysis, enabling global perspective for improved multi-interest modeling.

3.4 Multi-task Objective Function

To effectively bridge the adaptive user-individual and representative user-crowd multi-interest modeling, we propose a multi-task learning framework. For real users, we employ hard readout to predict user-item score for recommendation:

$$f(u_i, v_j) = \max_{1 \leq k \leq K} (\mathbf{o}_k^\top \mathbf{v}_j), \tag{11}$$

where the interest is selected among all interests by the maximum score. For the recommendation task, the objective function can be formulated with the InfoNCE loss as follows:

$$\mathcal{L}^{rec} = - \sum_{(u_i, v_j) \in \mathcal{D}} \log \frac{\exp(f(u_i, v_j))}{\sum_{v_j' \in \mathcal{V}} \exp(f(u_i, v_j'))}, \tag{12}$$

where \mathcal{D} is the training set for user-item interactions. For selected synthesized users, we employ contrastive learning in Equation (10) to bridge them with real users through item representation learning to enhance multi-interest modeling.

To perform multi-task learning for user-individual and user-crowd multi-interest modeling, our overall objective function aggregates these two goals in a weighted way:

$$\mathcal{L} = \mathcal{L}^{rec} + \lambda \cdot \mathcal{L}^{cst}, \tag{13}$$

where λ controls the trade-off between user-individual and user-crowd multi-interest modeling. As synthesized users usually contain rich behaviors leading to high-computation of contrastive learning, we conduct backpropagation on \mathcal{L}^{cst} every $\lfloor 1/\lambda \rfloor$ iterations ($\lambda > 0$) for efficiency consideration.

4 Experiments

4.1 Experimental Setup

Datasets: We use three subcategories of Amazon Review Data (Ni et al., 2019) with varying scales: *Beauty*, *Books*, and *Video Games*, abbreviated as Beauty, Book, and Game, respectively. All three datasets contain users' ratings on items with timestamps and item title information. Following prior studies (Xie et al., 2023; Du et al., 2024b), we filter out users and items with less than 5 records, then we convert ratings as implicit feedback for these three datasets. The statistical details of datasets are summarized in Appendix Table 4.

Evaluation Settings: To ensure a fair comparison, we follow prior studies (Xie et al., 2023), we chronologically split the user interactions with maximum length of 20 into training, validation,

	Metrics	Pop	GRU4Rec	LLMBRec	MIND	ComiRec	Re4	REMI	DisMIR	EIMF	LARMI
Beauty	<i>R@20</i>	0.0228	0.0349	0.0289	0.0477	0.0367	0.0550	0.0616	0.0702	<u>0.0765</u>	0.0872
	<i>N@20</i>	0.0161	0.0180	0.0205	0.0248	0.0176	0.0271	0.0320	0.0364	<u>0.0390</u>	0.0443
	<i>H@20</i>	0.0351	0.0454	0.0580	0.0669	0.0500	0.0715	0.0838	0.1057	<u>0.1126</u>	0.1380
	<i>R@50</i>	0.0391	0.0452	0.0473	0.0646	0.0519	0.0751	0.0817	0.0955	<u>0.0982</u>	0.1092
	<i>N@50</i>	0.0209	0.0186	0.0272	0.0257	0.0194	0.0274	0.0325	0.0433	<u>0.0427</u>	0.0505
	<i>H@50</i>	0.0593	0.0613	0.0943	0.0897	0.0702	0.0987	0.1099	0.1360	<u>0.1429</u>	0.1552
Book	<i>R@20</i>	0.0075	0.0215	0.0187	0.0236	0.0275	0.0298	0.0441	0.0639	<u>0.0722</u>	0.0804
	<i>N@20</i>	0.0052	0.0112	0.0132	0.0154	0.0166	0.0187	0.0293	0.0401	<u>0.0413</u>	0.0485
	<i>H@20</i>	0.0121	0.0314	0.0311	0.0334	0.0382	0.0420	0.0639	0.1034	<u>0.1060</u>	0.1228
	<i>R@50</i>	0.0133	0.0296	0.0256	0.0313	0.0381	0.0425	0.0592	0.0898	<u>0.0951</u>	0.1036
	<i>N@50</i>	0.0070	0.0113	0.0148	0.0158	0.0169	0.0194	0.0305	0.0404	<u>0.0443</u>	0.0519
	<i>H@50</i>	0.0217	0.0431	0.0538	0.0482	0.0570	0.0630	0.0915	0.1367	<u>0.1436</u>	0.1594
Game	<i>R@20</i>	0.0226	0.0751	0.0661	0.0950	0.0751	0.0967	0.1082	0.1221	0.1172	0.1305
	<i>N@20</i>	0.0139	0.0393	0.0420	0.0514	0.0384	0.0533	0.0543	0.0618	0.0604	0.0713
	<i>H@20</i>	0.0372	0.1099	0.1137	0.1401	0.1050	0.1465	0.1571	<u>0.1689</u>	0.1593	0.2133
	<i>R@50</i>	0.0435	0.1073	0.0849	0.1387	0.1145	0.1409	0.1510	<u>0.1597</u>	0.1571	0.1742
	<i>N@50</i>	0.0206	0.0423	0.0445	0.0552	0.0401	0.0558	0.0581	<u>0.0689</u>	0.0627	0.0774
	<i>H@50</i>	0.0710	0.1552	0.1576	0.2048	0.1496	0.2007	0.2175	<u>0.2479</u>	0.2315	0.2662

Table 1: Performance comparison of baseline methods and our proposed LARMI on three datasets. The best results are in **bold** and the runner-up results are underlined. The improvements are significant on the t-test ($p \leq 0.05$).

and test sets by the proportion of 6:2:2 and test last 20% items in each sequence. We adopt three widely used top- n evaluation metrics, i.e., Recall (R), Hit Rate (H), and Normalized Discounted Cumulative Gain (N), to evaluate all methods with $n = \{20, 50\}$.

Baseline Methods: We compare our model LARMI with the following baseline methods. **Pop** takes the most popular items as the recommendation results. **GRU4Rec** (Hidasi et al., 2016) models sequential behaviors through RNN structure. **LLM-BRec** (Harte et al., 2023) leverages an LLM to produce expressive embeddings by the BERT4Rec structure. **MIND** (Li et al., 2019) uses dynamic routing with a capsule network for multi-interest learning. **ComiRec-SA** (Cen et al., 2020) allows diversity control and introduces multi-head attention to model users’ multi-interests. **Re4** (Zhang et al., 2022) leverages the backward flow to re-examine and regulate interest representations. **REMI** (Xie et al., 2023) introduces interest-aware hard negative sampling with routing variation regularization for multi-interest learning. **DisMIR** (Du et al., 2024b) formulates an item partition problem to encourage items in each group to focus on a discriminated interest. **EIMF** (Qiao et al., 2024) uses an LLM to extract similar items for multi-interest modeling.

4.2 Comparison with Baselines

From the results in Table 1, we summarize our key findings to answer RQ1. First, LARMI consistently outperforms baselines across all three datasets, highlighting the effectiveness of our LLM-driven

analysis. In addition, the improvements demonstrate the advantage of integrating LLM analysis with our adaptive and representative multi-interest modeling framework by capturing diverse and meaningful user interests. We mitigate issues of uninterpretable agnostic granularity in LLM generation and the lack of a representative global perspective. Second, we observe that multi-interest baselines generally outperform single-interest baselines, showing that capturing multiple facets of user interests is beneficial for better results. Third, the superior performance of LARMI over the LLM-based baseline demonstrates that our approach achieves effective and coherent integration of the LLM-driven analysis and conventional multi-interest recommendation methods. Specifically, LARMI outperforms the relatively strong EIMF due to its emphasis on personalization and the granularity issue. Last, the relatively strong performance of DisMIR shows the positive impact of a global perspective item partition task. However, relying solely on the sparse co-occurrence of items limits the insights about item relationships. In comparison, LARMI produces synthesized users with compact and representative item subsets by formulating and solving the MCP, thus achieving improvements.

4.3 Ablation Studies

To validate the effects of key components, we conduct ablation studies as follows. **w/o-sem** removes the semantic-based multi-interest modeling with LLM-driven analysis at user-individual level, i.e., $h_f^i = 0$. **w/o-col** removes the collaborative-based

Metrics	w/o-sem	w/o-col	w/o-com	w/o-rep	LARMI
$R@20$	0.0544	0.0604	0.0841	0.0850	0.0872
$N@20$	0.0272	0.0310	0.0403	0.0425	0.0443
$H@20$	0.0745	0.0920	0.1310	0.1344	0.1380
$R@50$	0.0747	0.0889	0.1042	0.1058	0.1092
$N@50$	0.0285	0.0338	0.0471	0.0479	0.0505
$H@50$	0.1013	0.1279	0.1473	0.1525	0.1552

Table 2: Ablation study results with best scores in **bold**.

multi-interest modeling at the user-individual level, i.e., $\mathbf{o}_k^i = \mathbf{h}_k^i$. **w/o-com** removes the compactness rule for user synthesis, e.g., MCP formulation and user synthesize based on users’ similarities. **w/o-rep** removes the representativeness rule for user synthesis. Instead, it randomly selects cliques as synthesized users for user-crowd level analysis.

From Table 2, first, at the user-individual level, the lack of semantic multi-interests (w/o-sem) results in the worst performance, proving the effectiveness of LLM-driven analysis in handling the limitations of existing multi-interest modeling assumptions like co-occurring items simply indicating same interests. Second, w/o-col also shows inferior performance, showing our integration with collaborative interests and the alignment module is successful in alleviating the agnostic granularity issue in LLM-driven multi-interest analysis. Third, w/o-com and w/o-rep shows degraded performance, indicating that analyzing multi-interests for individual users only provides limited insights. On the one hand, generating user cliques from similar preferences ensures a moderate number of interests for synthesized users, thus each containing a rich set of cohesive items. On the other hand, the formulated MCP is crucial for selecting a representative subset that reduces redundancy and enhances representativeness. Thus, LARMI achieves representative multi-interest modeling by our user-crowd level multi-interest extraction.

4.4 Hyper-parameter Analysis

4.4.1 Loss Weight and Update Strategy

Figure 4 investigates the impact of (a) the loss weight λ on model accuracy and (b) the training time, with an update conducted every $\lfloor 1/\lambda \rfloor$ iterations. Higher λ (lower $\lfloor 1/\lambda \rfloor$) usually leads to higher model accuracy but requires a significantly longer time for model training. To balance effectiveness and efficiency, we suggest selecting $\lfloor 1/\lambda \rfloor = 100$. Generally, we observe that their performance improves as λ increases, with only slight differences between the two update strate-

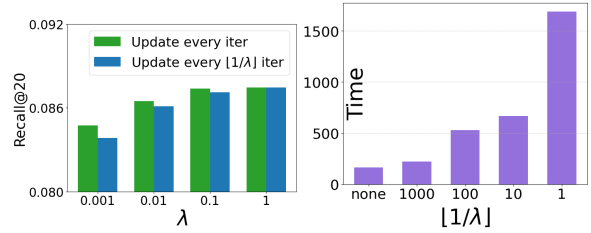


Figure 4: The model performance (a) and training time (b) with varying loss weights.

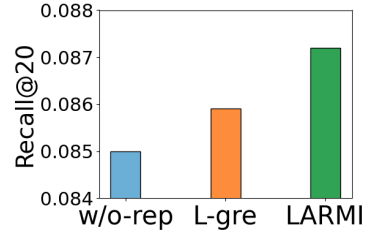


Figure 5: Model performance across MCP solve types.

gies when $\lambda \geq 1 \times 10^{-2}$. Therefore, the efficient update strategy with a proper setting of λ significantly accelerates the training while maintaining comparable performance.

4.4.2 MCP Solver

In Figure 5, we test the performance of different solvers in MCP for the multi-interest extraction at the user-crowd level. Specifically, we replace the MCP solver (Wang et al., 2022) with a greedy search of representative synthesized users, L-gre. First, LARMI shows an advantage over L-gre, reflecting the effectiveness of the MCP solver in selecting representative synthesized users. Second, both of the two approaches (LARMI and L-gre) outperform the w/o-rep variant, indicating the necessity of selecting representative synthetic users at the user-crowd level.

4.4.3 Number of Interests

Table 3 shows the effect of the number of multi-interests K for LARMI, where the optimal number of interests is 4 for *Beauty* dataset, matching the average numbers of LLM-derived clusters. Results for other two datasets are in Appendix Table 6. Specifically, insufficient interest numbers make it hard to capture the diverse facets of user preferences, while too large interest numbers may lead to overly-fine multi-interest modeling. As a result, we suggest a moderate interest number $K = 4$ for real-world applications.

Metrics	2	4	6	8
$R@20$	0.0869	0.0872	0.0847	0.0849
$N@20$	0.0438	0.0443	0.0423	0.0417
$H@20$	0.1374	0.1380	0.1295	0.1334
$R@50$	0.1077	0.1092	0.1009	0.1006
$N@50$	0.0497	0.0505	0.0469	0.0486
$H@50$	0.1530	0.1552	0.1523	0.1537

Table 3: LARMI performance of different interest nums.

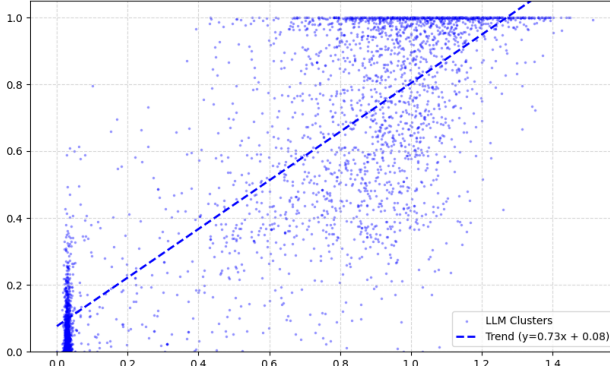


Figure 6: Positive correlation occurs between LLM cluster coarseness (x-axis) and the alignment module’s attention concentration (y-axis).

4.5 Adaptive Granularity Control Analysis

To empirically validate how our model adaptively addresses the agnostic granularity issue with raw LLM outputs, we visualize the alignment module’s behavior. Specifically, we define cluster coarseness (x-axis) as the semantic divergence of a raw LLM-derived semantic cluster C_f^i , computed by the average pairwise distance of the item embeddings within that cluster. Interest Concentration (y-axis) measures the sharpening effect of our alignment module, which is computed according to the negative entropy of the attention distribution a_j .

As shown in Figure 6, there is a strong positive correlation between the two metrics, demonstrating our module’s adaptive capability. Coarse clusters exhibit high attention concentration, indicating that the model selectively emphasizes specific items within a broad semantic group to filter noise and align with the user’s certain interest. Conversely, over-fine clusters are effectively merged into a unified, broader representation for further alignment with collaborative signals. Thus, LARMI successfully mitigates the issue of agnostic LLM semantics, ensuring multi-interest profiles are dynamically adjusted to the optimal level of granularity.

4.6 Case Study

We illustrate a case study to qualitatively investigate the effect of our model for multi-interest ex-

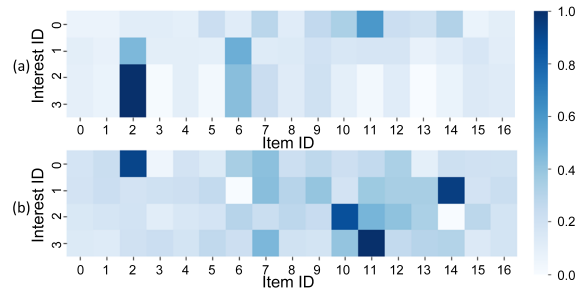


Figure 7: Multi-interest heatmap of baseline method MIND (a) and our proposed LARMI (b).

traction via LLMs. To visualize the cosine similarity between interests and items, we plot the multi-interest heatmap in Figure 7 for a classic baseline MIND and our LARMI model, which corresponds to a user sequence randomly sampled from *Beauty* dataset containing 17 items across 4 interests. In the heatmap, darker colors indicate higher cosine similarity between interests and items. We notice that the heatmaps with interest ID 2 and 3 in MIND exhibit significant overlap, indicating collapsed facets. Rather than being dominated mainly by one item for each interest, each interest in LARMI contains multiple relevant items, indicating discriminated and balanced multi-interest modeling. This improvement validates LARMI’s ability to adaptively capture diverse and fine-grained interests with the LLM and our model design.

5 Conclusion

In this paper, we present LARMI, the LLM-based Adaptive and Representative Multi-Interest approach that leverages the semantic knowledge and reasoning skills of LLMs to address critical challenges in multi-interest modeling for recommendation. At the user-individual level, LARMI provides an adaptive solution to the agnostic granularity in raw LLM generations, which merges and aligns semantic clusters with collaborative interests. At the user-crowd level, LARMI leverages MCP optimization and contrastive learning to mitigate the limitation in individual user behaviors and extend the scope of LLM analysis to a representative global perspective. Extensive experiments validate the superiority of LARMI over single- and multi-interest baselines. Further analysis supports the effectiveness of our model design. Future work includes fine-tuning open-source LLMs to align semantics with collaborative behaviors to further improve the scalability for multi-interest modeling.

598 Limitations

599 The limitations in our current framework include
600 the following points. First, our method relies on the
601 semantic richness of item metadata to prompt the
602 LLM for multi-interest analysis. However, it may
603 be less effective if the data has noise or incomplete-
604 ness, whereas ID-based methods would be unaf-
605 fected. Second, our reliance on external APIs may
606 be unstable due to changes in backbones. Third, the
607 inference cost of querying Large Language Models
608 remains higher than traditional lightweight recom-
609 mendation models, especially with larger datasets.

610 References

611 Keqin Bao, Ming Yan, Yang Zhang, Jizhi Zhang, Wenjie
612 Wang, Fuli Feng, and Xiangnan He. 2025. Customiz-
613 ing in-context learning for dynamic interest adap-
614 tion in LLM-based recommendation. In *Findings of
615 the Association for Computational Linguistics: ACL
616 2025*, pages 14278–14291.

617 Keqin Bao, Jizhi Zhang, Yang Zhang, Wenjie Wang,
618 Fuli Feng, and Xiangnan He. 2023. Tallrec:an ef-
619 fective and efficient tuning framework to align large
620 language model with recommendation. In *Proceed-
621 ings of the 17th ACM Conference on Recommender
622 Systems (RecSys)*, page 1007–1014.

623 Yukuo Cen, Jianwei Zhang, Xu Zou, Chang Zhou,
624 Hongxia Yang, and Jie Tang. 2020. Controllable
625 multi-interest framework for recommendation. In
626 *Proceedings of the 26th ACM SIGKDD International
627 Conference on Knowledge Discovery & Data Mining
628 (KDD)*, page 2942–2951.

629 Zheng Chai, Zhihong Chen, Chenliang Li, Rong Xiao,
630 Houyi Li, Jiawei Wu, Jingxu Chen, and Haihong
631 Tang. 2022. User-aware multi-interest learning for
632 candidate matching in recommenders. In *Proceed-
633 ings of the 45th International ACM SIGIR Confer-
634 ence on Research and Development in Information
635 Retrieval (SIGIR)*, page 1326–1335.

636 Gaode Chen, Xinghua Zhang, Yanyan Zhao, Cong Xue,
637 and Ji Xiang. 2021. Exploring periodicity and in-
638 teractivity in multi-interest framework for sequential
639 recommendation. In *Proceedings of the 30th Inter-
640 national Joint Conference on Artificial Intelligence
641 (IJCAI)*, pages 1426–1433.

642 Sunhao Dai, Ninglu Shao, Haiyuan Zhao, Weijie Yu,
643 Zihua Si, Chen Xu, Zhongxiang Sun, Xiao Zhang,
644 and Jun Xu. 2023. Uncovering chatgpt’s capabilities
645 in recommender systems. In *Proceedings of the 17th
646 ACM Conference on Recommender Systems (RecSys)*,
647 page 1126–1132.

648 Yingpeng Du, Di Luo, Rui Yan, Xiaopei Wang, Hongzhi
649 Liu, Hengshu Zhu, Yang Song, and Jie Zhang. 2024a.
650 Enhancing job recommendation through llm-based

generative adversarial networks. *Proceedings of the
AAAI Conference on Artificial Intelligence (AAAI)*,
pages 8363–8371. 651
652
653

Yingpeng Du, Ziyang Wang, Zhu Sun, Yining Ma,
Hongzhi Liu, and Jie Zhang. 2024b. Disentangled
multi-interest representation learning for sequential
recommendation. In *Proceedings of the 30th ACM
SIGKDD Conference on Knowledge Discovery and
Data Mining (KDD)*, page 677–688. 654
655
656
657
658
659

Shijie Geng, Shuchang Liu, Zuohui Fu, Yingqiang Ge,
and Yongfeng Zhang. 2022. Recommendation as
language processing (rlp): A unified pretrain, person-
alized prompt & predict paradigm (p5). In *Proceed-
ings of the 16th ACM Conference on Recommender
Systems (RecSys)*, page 299–315. 660
661
662
663
664
665

Qing Guo, Zhu Sun, Jie Zhang, and Yin-Leng Theng.
2020. An attentional recurrent neural network for
personalized next location recommendation. In *Pro-
ceedings of the AAAI Conference on artificial intelli-
gence*, volume 34, pages 83–90. 666
667
668
669
670

William L. Hamilton, Rex Ying, and Jure Leskovec.
2017. Inductive representation learning on large
graphs. In *Proceedings of the 31st International Con-
ference on Neural Information Processing Systems
(NeurIPS)*, page 1025–1035. 671
672
673
674
675

Jesse Harte, Wouter Zorgdrager, Panos Louridas, As-
terios Katsifodimos, Dietmar Jannach, and Marios
Fragkoulis. 2023. Leveraging large language models
for sequential recommendation. In *Proceedings of
the 17th ACM Conference on Recommender Systems
(RecSys)*. 676
677
678
679
680
681

Xiangnan He, Kuan Deng, Xiang Wang, Yan Li, Yong-
dong Zhang, and Meng Wang. 2020. Lightgcn: Sim-
plifying and powering graph convolution network for
recommendation. In *Proceedings of the 43rd Inter-
national ACM SIGIR conference on research and de-
velopment in Information Retrieval*, pages 639–648. 682
683
684
685
686
687

Balazs Hidasi, Alexandros Karatzoglou, Linas Bal-
trunas, and Domonkos Tikk. 2016. Session-based
recommendations with recurrent neural networks. In
*International Conference on Learning Representa-
tions (ICLR)*. 688
689
690
691
692

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 (ECIR)*, pages 364–381. 693
694
695
696
697

Yangqin Jiang, Yuhao Yang, Lianghao Xia, Da Luo,
Kangyi Lin, and Chao Huang. 2025. RecLM: Recom-
mendation instruction tuning. In *Proceedings of the
63rd Annual Meeting of the Association for Compu-
tational Linguistics (Volume 1: Long Papers)*, pages
15443–15459. 698
699
700
701
702
703

Wang-Cheng Kang and Julian McAuley. 2018. Self-
attentive sequential recommendation. In *2018 IEEE
International Conference on Data Mining (ICDM)*,
pages 197–206. 704
705
706
707

708	Samir Khuller, Anna Moss, and Joseph Seffi Naor. 1999.	Rethinking the usage of pre-trained language model	764
709	The budgeted maximum coverage problem. <i>Information</i>	in sequential recommendation. In <i>Proceedings of</i>	765
710	<i>processing letters</i> , 70(1):39–45.	<i>the 18th ACM Conference on Recommender Systems</i>	766
711	Yankun Le, Haoran Li, Baoyuan Ou, Yingjie Qin, Zhixuan	(<i>RecSys</i>), page 53–62.	767
712	Yang, Ruilong Su, and Fu Zhang. 2025. Diffusion	Xubin Ren, Wei Wei, Lianghao Xia, Lixin Su, Suqi	768
713	model for interest refinement in multi-interest	Cheng, Junfeng Wang, Dawei Yin, and Chao Huang.	769
714	recommendation. <i>Preprint</i> , arXiv:2502.05561.	2024. Rlmrec: Representation learning with large	770
715	Jaeri Lee, Jeongin Yun, and U Kang. 2024. Towards true	language models for recommendation. In <i>Proceed-</i>	771
716	multi-interest recommendation: Enhanced scheme	<i>ings of the ACM on Web Conference 2024 (TheWeb-</i>	772
717	for balanced interest training. In <i>2024 IEEE Inter-</i>	<i>Conf)</i> , pages 3464–3475.	773
718	<i>national Conference on Big Data (BigData)</i> , pages	Sara Sabour, Nicholas Frosst, and Geoffrey E. Hinton.	774
719	394–402.	2017. Dynamic routing between capsules. In <i>Pro-</i>	775
720	Chao Li, Zhiyuan Liu, Mengmeng Wu, Yuchi Xu, Huan	<i>ceedings of the 31st International Conference on Neu-</i>	776
721	Zhao, Pipei Huang, Guoliang Kang, Qiwei Chen, Wei	<i>ral Information Processing Systems (NeurIPS)</i> , page	777
722	Li, and Dik Lun Lee. 2019. Multi-interest network	3859–3869.	778
723	with dynamic routing for recommendation at tmall.	Scott Sanner, Krisztian Balog, Filip Radlinski, Ben	779
724	In <i>Proceedings of the 28th ACM International Con-</i>	Wedin, and Lucas Dixon. 2023. Large language mod-	780
725	<i>ference on Information and Knowledge Management</i>	els are competitive near cold-start recommenders for	781
726	(<i>CIKM</i>), page 2615–2623.	language- and item-based preferences. In <i>Proced-</i>	782
727	Danyang Liu, Yuji Yang, Mengdi Zhang, Wei Wu, Xing	<i>ings of the 17th ACM Conference on Recommender</i>	783
728	Xie, and Guangzhong Sun. 2022. Knowledge en-	<i>Systems (RecSys)</i> , page 890–896.	784
729	hanced multi-interest network for the generation of	Fei Sun, Jun Liu, Jian Wu, Changhua Pei, Xiao Lin,	785
730	recommendation candidates. In <i>Proceedings of the</i>	Wenwu Ou, and Peng Jiang. 2019. Bert4rec: Se-	786
731	<i>31st ACM International Conference on Information &</i>	quential recommendation with bidirectional encoder	787
732	<i>Knowledge Management (CIKM)</i> , page 3322–3331.	representations from transformer. In <i>Proceedings</i>	788
733	Yaokun Liu, Xiaowang Zhang, Minghui Zou, and Zhiy-	<i>of the 28th ACM International Conference on Infor-</i>	789
734	ong Feng. 2024. Attribute simulation for item embed-	<i>mation and Knowledge Management (CIKM)</i> , page	790
735	ding enhancement in multi-interest recommendation.	1441–1450.	791
736	In <i>Proceedings of the 17th ACM International Con-</i>	Yu Tian, Jianxin Chang, Yanan Niu, Yang Song, and	792
737	<i>ference on Web Search and Data Mining (WSDM)</i> ,	Chenliang Li. 2022. When multi-level meets multi-	793
738	page 482–491.	interest: A multi-grained neural model for sequential	794
739	Jianxin Ma, Chang Zhou, Hongxia Yang, Peng Cui,	recommendation. In <i>Proceedings of the 45th Inter-</i>	795
740	Xin Wang, and Wenwu Zhu. 2020. Disentangled	<i>national ACM SIGIR Conference on Research and</i>	796
741	self-supervision in sequential recommenders. In <i>Pro-</i>	<i>Development in Information Retrieval (SIGIR)</i> , page	797
742	<i>ceedings of the 26th ACM SIGKDD International</i>	1632–1641.	798
743	<i>Conference on Knowledge Discovery & Data Mining</i>	Runzhong Wang, Li Shen, Yiting Chen, Xiaokang Yang,	799
744	(<i>KDD</i>), page 483–491.	Dacheng Tao, and Junchi Yan. 2022. Towards one-	800
745	Fnu Mohbat and Mohammed J Zaki. 2025. KERL:	shot neural combinatorial solvers: Theoretical and	801
746	Knowledge-enhanced personalized recipe recommen-	empirical notes on the cardinality-constrained case.	802
747	dation using large language models. In <i>Proceedings</i>	In <i>The 11th International Conference on Learning</i>	803
748	<i>of the 63rd Annual Meeting of the Association for</i>	<i>Representations (ICLR)</i> .	804
749	<i>Computational Linguistics (Volume 1: Long Papers)</i> ,	Xinfeng Wang, Jin Cui, Fumiyo Fukumoto, and Yoshimi	805
750	pages 19125–19141.	Suzuki. 2025. AGRec: Adapting autoregressive de-	806
751	Jianmo Ni, Jiacheng Li, and Julian McAuley. 2019. Jus-	coders with graph reasoning for LLM-based sequen-	807
752	tifying recommendations using distantly-labeled re-	tial recommendation. In <i>Findings of the Association</i>	808
753	views and fine-grained aspects. In <i>Proceedings of</i>	<i>for Computational Linguistics: ACL 2025</i> , pages	809
754	<i>the 2019 Conference on Empirical Methods in Natu-</i>	7076–7090.	810
755	<i>ral Language Processing and the 9th International</i>	Yan Wang, Zhixuan Chu, Xin Ouyang, Simeng Wang,	811
756	<i>Joint Conference on Natural Language Processing</i>	Hongyan Hao, Yue Shen, Jinjie Gu, Siqiao Xue,	812
757	(<i>EMNLP-IJCNLP</i>), pages 188–197.	James Zhang, Qing Cui, and 1 others. 2024. Llmrg:	813
758	Shutong Qiao, Chen Gao, Yong Li, and Hongzhi Yin.	Improving recommendations through large language	814
759	2024. Llm-assisted explicit and implicit multi-	model reasoning graphs. In <i>Proceedings of the AAAI</i>	815
760	interest learning framework for sequential recom-	<i>Conference on Artificial Intelligence (AAAI)</i> , vol-	816
761	mendation. <i>arXiv preprint arXiv:2411.09410</i> .	ume 38, pages 19189–19196.	817
762	Zekai Qu, Ruobing Xie, Chaojun Xiao, Zhanhui Kang,	Wei Wei, Xubin Ren, Jiabin Tang, Qinyong Wang, Lixin	818
763	and Xingwu Sun. 2024. The elephant in the room:	Su, Suqi Cheng, Junfeng Wang, Dawei Yin, and Chao	819
		Huang. 2024. Llmrec: Large language models with	820

821 graph augmentation for recommendation. In *Pro-*
822 *ceedings of the 17th ACM International Conference*
823 *on Web Search and Data Mining (WSDM)*, pages
824 806–815.

825 Likang Wu, Zhi Zheng, Zhaopeng Qiu, Hao Wang,
826 Hongchao Gu, Tingjia Shen, Chuan Qin, Chen Zhu,
827 Hengshu Zhu, Qi Liu, Hui Xiong, and Enhong Chen.
828 2024. A survey on large language models for recom-
829 mendation. *World Wide Web*, 27:60.

830 Lianghao Xia, Chao Huang, Yong Xu, Peng Dai, Xiyue
831 Zhang, Hongsheng Yang, Jian Pei, and Liefeng Bo.
832 2021. Knowledge-enhanced hierarchical graph trans-
833 former network for multi-behavior recommendation.
834 In *Proceedings of the AAAI conference on artificial*
835 *intelligence (AAAI)*, volume 35, pages 4486–4493.

836 Zhibo Xiao, Luwei Yang, Wen Jiang, Yi Wei, Yi Hu,
837 and Hao Wang. 2020. Deep multi-interest network
838 for click-through rate prediction. In *Proceedings*
839 *of the 29th ACM International Conference on Infor-*
840 *mation & Knowledge Management (CIKM)*, page
841 2265–2268.

842 Yueqi Xie, Jingqi Gao, Peilin Zhou, Qichen Ye, Yining
843 Hua, Jae Boum Kim, Fangzhao Wu, and Sunghun
844 Kim. 2023. Rethinking multi-interest learning for
845 candidate matching in recommender systems. In
846 *Proceedings of the 17th ACM Conference on Recom-*
847 *mender Systems (RecSys)*, page 283–293.

848 Shengyu Zhang, Lingxiao Yang, Dong Yao, Yujie Lu,
849 Fuli Feng, Zhou Zhao, Tat-Seng Chua, and Fei Wu.
850 2022. Re4: Learning to re-contrast, re-attend, re-
851 construct for multi-interest recommendation. In *The*
852 *ACM Web Conference (TheWebConf)*, pages 2216–
853 2226.

854 Shuai Zhang, Yi Tay, Lina Yao, Aixin Sun, and Jake An.
855 2019. Next item recommendation with self-attentive
856 metric learning. In *Thirty-Third AAAI conference on*
857 *artificial intelligence (AAAI)*, volume 9.

858 Zhi Zheng, Zhaopeng Qiu, Xiao Hu, Likang Wu, Heng-
859 shu Zhu, and Hui Xiong. 2023. Generative job rec-
860 ommendations with large language model. *arXiv*
861 *preprint arXiv:2307.02157*.

862 A Backgrounds

863 A.1 Multi-Interest Modeling for 864 Recommendation

865 Capturing different preference representations of
866 users is essential for multi-interest methods. Cap-
867 sule networks (Sabour et al., 2017) have gained
868 popularity for multi-interest modeling recently (Li
869 et al., 2019; Xie et al., 2023). Specifically, the cap-
870 sule network *CapsuleNet*(\cdot) can generate users’
871 K -interest representations $\{\mathbf{m}_1, \dots, \mathbf{m}_K\}$ de-
872 rived from their sequential behaviors $s_{(u_i)}$, where
873 K denotes the pre-defined number of users’ multi-
874 interests. The k -th collaborative interest capsule
875 $\mathbf{m}_k \in \mathbb{R}^d$ is calculated by:

$$876 \mathbf{m}_k = \sum_{j=1}^L b_{jk} \cdot \mathbf{W} \cdot \mathbf{v}_j^i, \quad k = 1, \dots, K, \quad (14)$$

877 where $\mathbf{v}_j^i \in \mathbb{R}^d$ denotes the embeddings of the j -th
878 item in the user u_i ’s behaviors $s_{(u_i)}$, and d denotes
879 the dimension of embedding space. $\mathbf{W} \in \mathbb{R}^{d \times d}$
880 is the transformation matrix, and b_{jk} denotes the
881 routing weight of item v_j^i to the k -th interest cap-
882 sule. The routing weight b_{jk} is calculated by the
883 softmax operation of the routing logits g_{jk} that mea-
884 sures the similarity between the item embedding
885 and squashed vector \mathbf{e}_k , i.e.,

$$886 \begin{aligned} \mathbf{e}_k &= \frac{\|\mathbf{m}_k\|^2}{\|\mathbf{m}_k\|^2 + 1} \cdot \frac{\mathbf{m}_k}{\|\mathbf{m}_k\|}, \\ g_{jk} &\leftarrow g_{jk} + \mathbf{v}_j^{i\top} \cdot \mathbf{W} \cdot \mathbf{e}_k, \\ b_{jk} &= \frac{\exp(g_{jk})}{\sum_{k=1}^K \exp(g_{jk})}, \end{aligned} \quad (15)$$

887 where the iterative process updates the routing log-
888 its g_{jk} and recalculates the routing weights b_{jk}
889 based on updated \mathbf{m}_k .

890 A.2 Max Covering Problem

891 The max covering problem (MCP) (Khuller et al.,
892 1999) is a combinatorial optimization problem,
893 widely applied in decision-making under uncer-
894 tainty. Given P sets, each set containing an indefi-
895 nite number of objects, and Q objects, each object
896 associated with a specific value, the MCP aims to
897 select a subset of Z sets ($Z \ll P$) such that the
898 union of Z sets maximizes the sum of the associ-
899 ated values of the covered objects. The problem

can be formulated as:

$$\begin{aligned} &\max_{\mathbf{x}} \sum_{j=1}^Q \left(\mathbb{I} \left(\sum_{i=1}^P \mathbf{x}_i \mathbf{A}_{ij} \geq 1 \right) \cdot w_j \right), \\ &\text{s.t. } \mathbf{x} \in \{0, 1\}^P, \|\mathbf{x}\|_0 \leq Z, \end{aligned} \quad (16)$$

900 where $\mathbf{A} \in \{0, 1\}^{P \times Q}$ is the adjacency matrix of a
901 bipartite graph linking the sets and objects, $w_j \in \mathbb{R}$
902 denotes the j -th object value, $\mathbb{I}(\cdot)$ is a condition
903 indicator function, \mathbf{x} is the selection outcome, and
904 each element \mathbf{x}_i is a scalar indicating whether the
905 i -th set is selected in solution. To tackle the MCP,
906 an advanced neural solver encodes the bipartite
907 graph with a three-layer GraphSage model (Hamil-
908 ton et al., 2017), integrates the MCP constraints
909 into a differentiable layer, and then predicts the
910 probabilities of selecting each set. 911 912

913 B Model Complexity and Scalability

914 The computational complexity of LARMI comes
915 from (a) LLM-driven inference and (b) multi-
916 interest model training. For (a), assuming the com-
917 plexity for analyzing each user’s behaviors is \mathcal{T} ,
918 the total cost for user-individual multi-interest ex-
919 traction is $\mathcal{O}(M \cdot \mathcal{T})$. For user-crowd multi-interest
920 extraction, it takes $\mathcal{O}(M^2)$ for MCP and $\mathcal{O}(Z \cdot \mathcal{T})$
921 for LLM-driven analysis, where $Z \ll M$. There-
922 fore, the overall complexity for LLM-driven multi-
923 interest analysis is $\mathcal{O}(M^2 + M \cdot \mathcal{T})$, which is similar
924 to existing LLM-based recommendation methods
925 (Dai et al., 2023; Hou et al., 2024). Regarding
926 (b), computing collaborative multi-interests takes
927 $\mathcal{O}(L \cdot K \cdot d)$, where L, K, d are sequence length,
928 number of routing facets, embedding dimension.
929 Aligning semantic and collaborative multi-interests
930 also takes $\mathcal{O}(L \cdot K \cdot d)$. Therefore, the overall com-
931 plexity of user-individual multi-interest learning is
932 $\mathcal{O}(M \cdot L \cdot K \cdot d)$, which is equivalent to existing
933 multi-interest recommendation models (Li et al.,
934 2019; Xie et al., 2023). For contrastive learning,
935 it takes $\mathcal{O}(Z \cdot \overline{c_{(u')}}^2 \cdot d)$ for each update, where
936 $\overline{c_{(u')}}$ is the average number of behaviors for syn-
937 thesized users. To tackle the high computational
938 complexity in contrastive learning, we only update
939 the loss every $\lceil 1/\lambda \rceil$ iteration. In summary, our
940 method matches existing complexity and avoids
941 online latency, thus ensuring scalability.

942 C Supplementary Experimental Results

943 We provide dataset details, implementation details,
944 and experimental results of our ablation study and

Dataset	# Users	# Items	# Interactions	Avg Len	Density
Beauty	15,097	44,261	100,055	6.63	1.5e-4
Book	99,101	361,002	780,018	7.87	2.2e-5
Game	20,551	27,456	153,541	7.47	2.7e-4

Table 4: Dataset statistics.

	Metrics	w/o-sem	w/o-col	w/o-com	w/o-rep	LARMI
Beauty	<i>R@20</i>	0.0544	0.0604	0.0841	0.0850	0.0872
	<i>N@20</i>	0.0272	0.0310	0.0403	0.0425	0.0443
	<i>H@20</i>	0.0745	0.0920	0.1310	0.1344	0.1380
	<i>R@50</i>	0.0747	0.0889	0.1042	0.1058	0.1092
	<i>N@50</i>	0.0285	0.0338	0.0471	0.0479	0.0505
	<i>H@50</i>	0.1013	0.1279	0.1473	0.1525	0.1552
Book	<i>R@20</i>	0.0504	0.0585	0.0745	0.0783	0.0804
	<i>N@20</i>	0.0319	0.0362	0.0442	0.0461	0.0485
	<i>H@20</i>	0.0775	0.0891	0.1153	0.1197	0.1228
	<i>R@50</i>	0.0624	0.0740	0.0919	0.0937	0.1036
	<i>N@50</i>	0.0352	0.0416	0.0498	0.0510	0.0519
	<i>H@50</i>	0.1079	0.1263	0.1507	0.1558	0.1594
Game	<i>R@20</i>	0.0905	0.0934	0.1186	0.1256	0.1305
	<i>N@20</i>	0.0430	0.0468	0.0677	0.0695	0.0713
	<i>H@20</i>	0.1460	0.1523	0.2012	0.2084	0.2133
	<i>R@50</i>	0.1331	0.1405	0.1679	0.1708	0.1742
	<i>N@50</i>	0.0481	0.0525	0.0703	0.0750	0.0774
	<i>H@50</i>	0.1767	0.1953	0.2560	0.2589	0.2662

Table 5: Ablation study results on Beauty, Book and Game datasets. The best scores are in **bold**.

further analysis as follows. Table 4 shows the details of datasets, Table 5 shows the full ablation study results, and Table 6 shows the performance across interest numbers. They follow similar trend and can validate our motivations as well as model design.

C.1 Implementation Details

For a fair comparison, all methods are optimized by the Adam optimizer with a batch size of 128 and we adopt fixed embedding dimension 64 for all methods following (Xie et al., 2023; Du et al., 2024b). For all methods, we select the best performance by varying the number of interests in $\{2, 4, 6, 8\}$, learning rate in $\{1e^{-2}, 1e^{-3}, 1e^{-4}\}$, and weight decay in $\{1e^{-4}, 1e^{-5}, 1e^{-6}\}$. For hyper-parameters, we set $\tau = 0.1$ and $\lambda = 0.01$ for the contrastive loss \mathcal{L}^{cst} . For the MCP solver, we follow the same hyper-parameters as suggested in the original paper (Wang et al., 2022). For the LLM implementation, we use the *gpt-4o* as the backbone model with temperature set to 0 for reproducibility. For other hyper-parameters in baseline methods, we follow the authors’ implementation if they exist, otherwise we tune them to their best according to their performance on the valida-

	Metrics	2	4	6	8
Beauty	<i>R@20</i>	0.0869	0.0872	0.0847	0.0849
	<i>N@20</i>	0.0438	0.0443	0.0423	0.0417
	<i>H@20</i>	0.1374	0.1380	0.1295	0.1334
	<i>R@50</i>	0.1077	0.1092	0.1009	0.1006
	<i>N@50</i>	0.0497	0.0505	0.0469	0.0486
	<i>H@50</i>	0.1530	0.1552	0.1523	0.1537
Book	<i>R@20</i>	0.0775	0.0804	0.0799	0.0783
	<i>N@20</i>	0.0479	0.0485	0.0497	0.0482
	<i>H@20</i>	0.1195	0.1228	0.1206	0.1190
	<i>R@50</i>	0.1004	0.1036	0.1042	0.0997
	<i>N@50</i>	0.0509	0.0519	0.0517	0.0483
	<i>H@50</i>	0.1588	0.1594	0.1598	0.1586
Game	<i>R@20</i>	0.1251	0.1305	0.1274	0.1240
	<i>N@20</i>	0.0704	0.0713	0.0717	0.0711
	<i>H@20</i>	0.2119	0.2133	0.2115	0.2084
	<i>R@50</i>	0.1689	0.1742	0.1754	0.1712
	<i>N@50</i>	0.0732	0.0774	0.0741	0.0748
	<i>H@50</i>	0.2601	0.2662	0.2620	0.2635

Table 6: LARMI performance of different interest nums.

tion set. Our source code is available at <https://anonymous.4open.science/r/LARMI-26A8>.

970

971