

LLM Collaborative Filtering: User-Item Graph as New Language

Huachi Zhou¹, Yujing Zhang¹, Hao Chen^{2*},
Qinggang Zhang¹, Qijie Shen³, Feiran Huang⁴, Xiao Huang¹

¹The Hong Kong Polytechnic University, Hong Kong

²City University of Macau, Macao

³Alibaba Group, China

⁴Jinan University, China

{huachi.zhou, yu-jing.zhang, qinggang.zhang}@connect.polyu.hk
sundaychenhao@gmail.com, qjshenxd@gmail.com
huangfr@jnu.edu.cn, xiaohuang@comp.polyu.edu.hk

Abstract

In collaborative filtering, learning effective embeddings for users and items from interaction data remains a central challenge. While recent efforts leverage large language models (LLMs) to enhance collaborative filtering, two critical limitations persist: (1) Efficiency: LLM-based inference is significantly slower than traditional embedding-based search; and (2) Topological Modeling: LLMs struggle to capture graph structures, which are essential for modeling multi-order user-item interactions. To address these limitations, we propose New Language Collaborative Filtering (NLCF), a framework that aligns LLMs with collaborative filtering by conceptualizing user-item graphs as new languages. This approach is based on two key insights: (1) LLMs excel at mastering new languages when trained on suitable corpora, and (2) the empirical conditional probability between tokens in corpora converges to the transition probabilities between nodes in graphs. NLCF translates user-item graphs into corpora, where users and items are treated as tokens. These corpora are used to fine-tune LLMs, and the learned representations are aggregated to construct user and item embeddings that encode multi-order interactions. Unlike methods that deploy LLMs for inference, NLCF distills LLM knowledge learned from corpora into compact embeddings, enabling both efficient training and real-time inference. The framework has been deployed on a billion-scale e-commerce platform for several months. Extensive experiments demonstrate that NLCF outperforms traditional graph CF models and LLM-based baselines while achieving significant training and inference efficiency improvement over LLM-based baselines.

Introduction

Collaborative Filtering (CF) plays a pivotal role in recommender systems by suggesting new items to users based on the collaborative information that users with similar interactions share similar preferences (Yuan et al. 2023). Multi-order interactions which represent the multi-order relationships connecting users and items through various paths in the user-item graph enrich this information by uncovering

*Corresponding author.

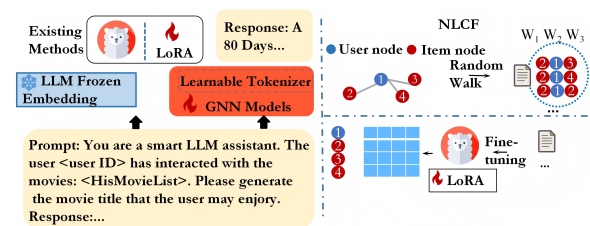


Figure 1: Comparison of pipelines between existing methods and NLCF. NLCF achieves efficiency through using short graph sentences for training and embeddings for inference.

preferences beyond immediate interactions. To model multi-order interactions, traditional models like LightGCN (He et al. 2020) achieve this by iteratively aggregating embeddings from immediate interactions in the user-item graphs to learn user and item embeddings. Recently, the capabilities of large language models (LLMs) (Chen et al. 2024; Hong et al. 2024) have inspired their application in CF. However, aligning LLMs with CF tasks remains challenging, as LLMs are optimized for natural language tasks rather than learning user and item embeddings based on user-item graphs.

Existing methods attempt to adapt LLMs for CF by transforming a user’s immediate interactions into natural language graph descriptions. These approaches can be broadly categorized into two classes: (i) **Description Reasoning-based** methods: These methods leverage LLMs’ natural language understanding to reason about collaborative information. For instance, TransRec (Lin et al. 2024b) constructs descriptive sequences from item identifiers, while LLM-CF (Sun et al. 2024) employs Chain-of-Thought prompting to distill collaborative information explicitly. (ii) **Embedding Injection-based** methods: Methods like LC-Rec (Zheng et al. 2024) and LETTER (Wang et al. 2024) inject user and item embeddings from traditional models (e.g., LightGCN) into prompts. By tokenizing collaborative embeddings and combining them with natural language descriptions, these methods aim to capture both semantic and collaborative information.

Despite their potential, natural language graph description-based methods face one or both of the following limitations in efficiency and topological modeling for CF tasks.

- (1) **Multi-order Collaborative Information Loss:** LLMs can process only a limited number of interactions at a time, making it extremely expensive to model multi-order interactions and the underlying collaborative information during training. While some methods inject collaborative embeddings into prompts (Zheng et al. 2024; Wang et al. 2024), fine-tuning LLMs on precomputed embeddings, rather than the original user-item graph, inevitably leads to information loss (Yang et al. 2024).
- (2) **Slow Inference:** LLMs generate items token by token, with each token depending on the entire prompt (Zhao et al. 2023). The computational cost of generation scales with the prompt length. This sequential generation process, especially when using lengthy graph description prompts, significantly slows down inference compared to embedding-based search.

Although recent methods, e.g., LLMEmb (Liu et al. 2025) and LLM-CF (Sun et al. 2024), have proposed efficient inference mechanisms, they still focus on generating accurate user and item profile representations rather than modeling the user-item graph with LLMs. Given that user-item graphs have been proven highly effective for capturing multi-order collaborative information (Wei et al. 2024), modeling the graph with LLMs powerful learning ability is a promising direction. However, LLMs excel at learning a new language but find it challenging to efficiently learn from the graph with two key reasons:

- (1) **Structural Mismatch between Graph and LLM Inputs:** User-item graphs are irregular, two-dimensional structures, while LLMs are designed to process one-dimensional sequential data. Unlike graph CF models, which aggregate multi-order interactions through graph structure (Wu et al. 2020), LLMs lack a built-in mechanism to handle such graph structures.
- (2) **Computational Constraints:** As the order increases, the number of multi-order interactions grows geometrically. Modeling these interactions while maintaining a concise representation is necessary since LLMs, containing billions of parameters, are computationally expensive during both fine-tuning and inference. And deploying LLMs in real-world recommender systems is particularly challenging due to latency constraints.

To address these challenges, we propose New Language Collaborative Filtering (NLCF), a novel framework that enables LLMs to efficiently learn collaborative embeddings by treating the graph as a new language. As shown in Figure 1, NLCF transforms the graph into concise corpus, where empirical conditional probabilities between tokens in the corpus converge to the transition probabilities between nodes in the graph. This transformation enables NLCF to model multi-order interactions, and then NLCF balances common and rare interactions in the corpus through similarity-based

sampling. The resulting graph corpus is then used to fine-tune the LLM efficiently, allowing it to construct collaborative embeddings. NLCF achieves efficient training by using compact graph corpus rather than lengthy natural language graph descriptions, and efficient inference by leveraging collaborative embeddings instead of token-by-token generation. Our contributions are as follows:

- We propose NLCF, a novel framework that efficiently integrates LLMs with CF tasks to learn user and item collaborative embeddings by treating the user-item graph as a new language.
- We design two core modules: (i) a graph corpus collection module that transforms the graph into concise corpus, modeling multi-order interactions; and (ii) a collaborative embedding construction module that fine-tunes LLMs on this corpus to construct collaborative embeddings for efficient item search.
- Extensive experiments on three datasets show performance gains over both traditional graph CF and LLM-based baselines, with significant efficiency improvement over LLM-based baselines. Online A/B tests further validate NLCF’s effectiveness in industrial applications.

Preliminary

Notation. We represent the user-item graph as a tuple $\mathcal{G} = (\mathcal{V}, \mathcal{E})$, where $\mathcal{V} = \mathcal{U} \cup \mathcal{I}$ denotes the union of user nodes \mathcal{U} and item nodes \mathcal{I} . The node set is indexed as $\{v_1, v_2, \dots, v_{|\mathcal{V}|}\}$, where $|\mathcal{V}|$ is the total number of nodes. The edge set $\mathcal{E} \subseteq \mathcal{U} \times \mathcal{I}$ represents user-item interactions, with cardinality $|\mathcal{E}|$. These interactions are encoded in the adjacency matrix $\mathbf{A} \in \mathbb{R}^{|\mathcal{V}| \times |\mathcal{V}|}$:

$$\mathbf{A}_{ui} = \begin{cases} 1, & \text{if } (v_u, v_i) \in \mathcal{E}, \\ 0, & \text{otherwise.} \end{cases} \quad (1)$$

The bipartite structure ensures $\mathbf{A}_{ui} = 0$ for all user-user and item-item pairs. We fine-tune LLMs on user-item graph \mathcal{G} with LoRA (Hu et al. 2021) $\hat{\mathbf{W}}$ to learn user embeddings \mathbf{h}_{v_u} and item embeddings \mathbf{h}_{v_i} . More preliminary details are put in Appendix D.

New Language Collaborative Filtering

In this section, we present NLCF, an efficient framework that applies LLMs to learn collaborative embeddings from the user-item graph through new language learning. As shown in Figure 2, NLCF consists of two primary modules: (i) **Graph Corpus Collection:** The user-item graph is transformed into a new language corpus that encodes multi-order interactions. We retrieve this corpus by employing similarity-based sampling to reduce the computational burden. (ii) **Collaborative Embedding Construction:** The retrieved corpus is used to fine-tune LLMs, enabling the model to capture multi-order collaborative information. Hidden representations from the fine-tuned model are then aggregated to construct collaborative user and item embeddings for inference.

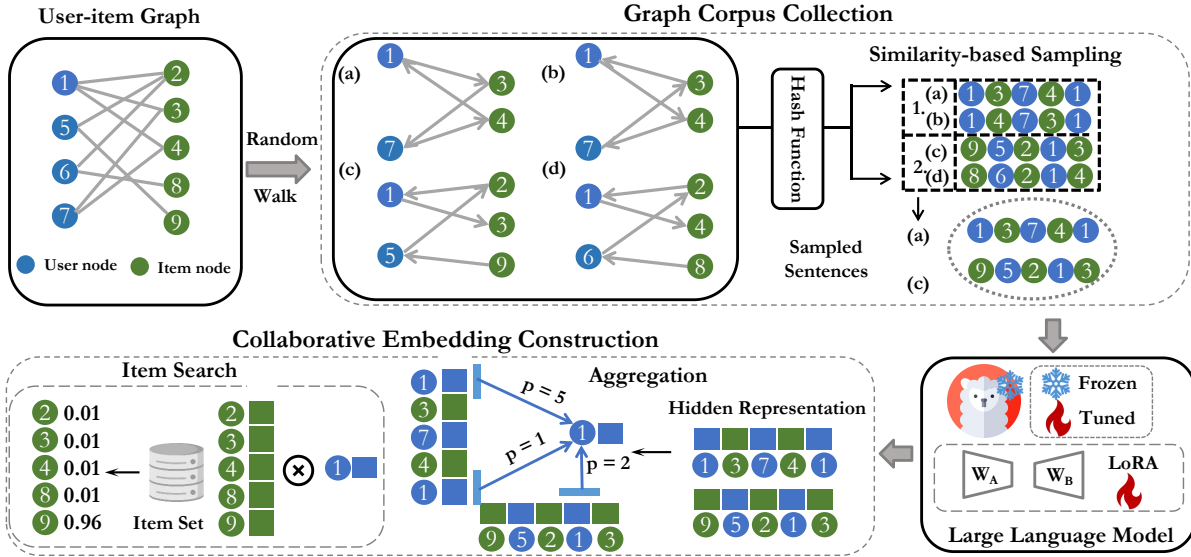


Figure 2: The overall pipeline of the proposed NLCF framework. NLCF treats the user-item graph as a new language. The framework consists of two key modules: in the top part, NLCF first collects graph corpus through random walks and then applies similarity-based sampling; in the bottom part, NLCF constructs collaborative user-item embeddings using fine-tuned LLMs for efficient item search.

Graph Corpus Collection

We begin by mapping basic concepts in user-item graph to their counterparts in graph language and then describe our similarity-based sampling strategy for reducing corpus size.

Definition of Graph Language-related Concepts. User and Item Nodes as Graph Tokens. We define the union of user and item nodes \mathcal{V} as a set of unique graph tokens in this new language and extend the LLM tokenizer’s vocabulary accordingly. These newly added tokens have randomly initialized embeddings.

Graph Path as Graph Sentence. A connected sequence of nodes in the graph, or a path, forms a graph sentence in the corpus. For example, as shown in Figure 2, the path $s_1 = \{v_9, v_5, v_2, v_1, v_3\}$ represents a graph sentence. These paths model multi-order interactions and capture collaborative information by revealing behavior similarity. For instance, v_1 is likely to interact with v_9 because the second-order neighbor user v_5 has previously interacted with v_9 . The initial corpus $\mathcal{S} = \{s_1, s_2, s_3, \dots\}$ is composed of such sentences extracted from the user-item graph.

To extract graph sentences from the user-item graph, we employ random walks (Grover and Leskovec 2016), whose transition probability between nodes is defined as:

$$P_g(v_i | v_u) = \begin{cases} \frac{\mathbf{A}_{ui}}{\sum_{j \in \mathcal{N}(v_u)} \mathbf{A}_{uj}}, & v_i \in \mathcal{N}(v_u), \\ 0, & (v_u, v_i) \notin \mathcal{E}, \end{cases} \quad (2)$$

where $\mathcal{N}(v_u)$ denotes the neighbor set of node v_u , and \mathbf{A}_{ui} represents the edge weight between nodes v_u and v_i . Each random walk continues walking until reaching a predefined length l . This process is designed to ensure that empirical conditional probability between tokens in resulting corpus

converges to the transition probabilities between nodes in graphs. The guarantee is proved by the following theorem:

Theorem 1. Let $P_s(v_i|v_u)$ be the empirical conditional probability computed from the corpus generated by a sufficiently large number of random walks. As the number of walks approaches infinity, $P_s(v_i|v_u)$ converges in probability to $P_g(v_i|v_u)$.

The proof is provided in Appendix A. To sufficiently capture collaborative information, we generate d_u sequences starting at each node v_u , where $d_u = \sum \mathbf{A}_u$: is the degree of v_u . This strategy yields the initial corpus \mathcal{S} , which serves as the foundation for the following processing.

Similarity-based Sampling for Graph Corpus Reduction. The scale of possible graph corpus is proportional to the number of nodes and the average node degree. To limit the corpus size, traditional method, e.g., random sampling or node degree-based sampling does not account for subgraph density around each node. In dense subgraphs, random walks tend to generate highly overlapping sentences due to frequent visits to common nodes. Some existing LLM-based recommendation sampling methods (Lin et al. 2024a,c) focus on selecting important interactions for a small set of candidate items, which are incompatible with our goal of efficient item search from the entire item set (Zhou et al. 2025d). Other approaches rely on LLM fine-tuning (Zhou et al. 2025c) or LLM data integration to identify important samples (Wu et al. 2023) which would be computationally expensive to use.

We propose a similarity-based sampling strategy that addresses these limitations. Our approach groups similar graph sentences into clusters based on their sentence overlapping and assigns lower sampling probabilities to densely popu-

lated clusters. By sampling representative sentences within a pre-defined budget, this strategy reduces redundancy while preserving essential collaborative information.

We measure the overlapping between any two graph sentences using the Jaccard similarity:

$$S_{\text{similarity}}(\mathbf{s}_p, \mathbf{s}_q) = \frac{|\mathbf{s}_p \cap \mathbf{s}_q|}{|\mathbf{s}_p \cup \mathbf{s}_q|}, \quad (3)$$

where \mathbf{s}_p and \mathbf{s}_q represent two graph sentences. We then group sentences into clusters based on their pairwise similarities. A cluster \mathcal{C}_i is defined as a set of sentences where each sentence shares a similarity above threshold t with all other sentences in the cluster:

$$\mathcal{C}_i = \mathbf{s}_p \in \mathcal{S} : S_{\text{similarity}}(\mathbf{s}_p, \mathbf{s}_q) \geq t, \forall \mathbf{s}_q \in \mathcal{C}_i. \quad (4)$$

The size of each cluster $|\mathcal{C}_i|$ reflects the density of similar sentences within it. To achieve balanced representation, we assign higher sampling probabilities to sentences from larger clusters. The sampling probability for each sentence is normalized across all clusters:

$$P(\mathbf{s}_q) = \frac{|\mathcal{C}_i|}{\sum_{j=1}^m |\mathcal{C}_j|} : \mathbf{s}_q \in \mathcal{C}_i, \quad (5)$$

where m is the cluster number automatically determined by the algorithm. Then given a pre-defined sampling ratio $\alpha \in (0, 1]$, we sample sentences according to these normalized probabilities to create a representative subset:

$$\mathcal{S}' = \mathbf{s}_q \sim P(\mathbf{s}_q) : \mathbf{s}_q \in \mathcal{S}, |\mathcal{S}'| = \alpha |\mathcal{S}|. \quad (6)$$

This stratified sampling approach ensures balanced representation across clusters while maintaining computational efficiency within the specified budget constraints.

While computing pairwise similarities using Eq. (3) provides fine-grained clustering, it introduces a significant computational complexity of $O(|\mathcal{S}|^2)$, which is impractical for large corpus. To address this challenge, we employ MinHash to efficiently estimate Jaccard similarities, whose efficiency is guaranteed by the following theorem:

Theorem 2. Let $J(s_p, s_q)$ be the Jaccard similarity between two graph sentences s_p and s_q . The MinHash estimator, $\hat{J}(s_p, s_q)$, constructed from k independent hash functions, is an unbiased estimator of $J(s_p, s_q)$ with a variance of $\frac{J(s_p, s_q)(1-J(s_p, s_q))}{k}$.

The proof is detailed in the Appendix B. This theorem demonstrates that MinHash provides a statistically sound approximation of the true Jaccard similarity. By generating compact signatures for each sentence and storing them in hash tables, we can identify similar sentences with an expected complexity of $O(|\mathcal{S}|)$, thus making the similarity estimation feasible for large-scale graph corpora.

Theoretical Analyses about Connection between Graph Corpus and Collaborative Information. The graph corpus collection module models multi-order interactions by transforming the two-dimensional user-item graph into a one-dimensional graph corpus. It is important to formally analyze whether this transformation quantitatively guarantees the preservation of multi-order collaborative information in the graph corpus. To address this, we present the following theorem:

Theorem 3. Let w denote the importance of node v_i to node v_u , as measured by the gradient norm of a GCN, and let $w' = \frac{n'}{\hat{n}}$ represent the empirical co-occurrence ratio, where n' is the number of graph sentences containing both nodes v_u and v_i , and \hat{n} is the number of graph sentences containing node v_u . The standard error of $|w - w'|$, is bounded by $O\left(\frac{1}{\hat{n}}\right)$.

Proof. Let $h_{v_u}^{(l)}$ denote the hidden feature learned by GCN (Hamilton, Ying, and Leskovec 2017), defined as $\text{ReLU}\left(\frac{1}{d_u} \cdot \sum_{v_j \in \mathcal{N}(v_u)} \mathbf{W}_l h_{v_j}^{(l-1)}\right)$. For simplicity, we remove the non-linear activation function and GCN weights from the GCN aggregation in this proof. However, the theorem still holds when these components are included. The gradient of $h_{v_u}^{(l)}$ with respect to $h_{v_i}^{(0)}$ is given by:

$$\frac{\partial h_{v_u}^{(l)}}{\partial h_{v_i}^{(0)}} = \frac{1}{d_u} \cdot \sum_{v_j \in \mathcal{N}(v_u)} \frac{\partial h_{v_j}^{(l-1)}}{\partial h_{v_i}^{(0)}}.$$

We iteratively expand this formula using the chain rule to get:

$$\frac{\partial h_{v_u}^{(l)}}{\partial h_{v_i}^{(0)}} = \sum_{p=1}^n \left[\frac{\partial h_{v_u}^{(l)}}{\partial h_{v_i}^{(0)}} \right]_p = \sum_{p=1}^n \prod_{q=1}^{l'} \frac{1}{d_{v_{q|p}}} = w,$$

where n is the number of paths containing both nodes v_i and v_u , l' is the length of the current path p , $d_{v_{q|p}}$ represents the degree of the q -th node in the current path p . Since we use Eq. (2) and perform random walk, the probability of node v_u visiting v_i is exactly the same as the sum of probabilities for all paths shown above.

However, it is computationally impractical to retrieve all such paths for every pair of nodes, especially when similarity-based sampling is used. Let X_p be an indicator Bernoulli variable for the p -th walk in practice, which shows whether the walk successfully reaches node v_i starting from v_u :

$$X_p = \begin{cases} 1 & \text{if the walk reaches } v_i, \\ 0 & \text{otherwise.} \end{cases}$$

From the previous theorem, the following expectation and variance properties hold: $E[X_p] = w$, and $\text{Var}(X_p) = w(1-w)$. Now, if there are \hat{n} paths from node v_u to node v_i , the empirical probability satisfies the following properties:

$E\left[\frac{\sum_{p=1}^{\hat{n}} X_p}{\hat{n}}\right] = w$, and $\text{Var}\left[\frac{\sum_{p=1}^{\hat{n}} X_p}{\hat{n}}\right] = \frac{w(1-w)}{\hat{n}}$. This variance is bounded by 1. Therefore, the sampling standard error of the empirical probability w' , i.e., $\frac{n'}{\hat{n}}$ encoded in the sampled corpus, satisfies:

$$|w - w'| = O\left(\frac{1}{\hat{n}}\right).$$

In the proof, w represents collaborative information captured by GCNs in prediction by modeling multi-order interactions, while w' denotes the importance derived from

the sampled graph corpus as reflected by the empirical occurrence ratio. The theorem demonstrates that the sampled graph corpus encodes collaborative information in a manner consistent with GCNs, with the difference diminishing as the number of sampled graph sentences increases. This result provides a theoretical foundation for using graph sentences as a proxy for mining collaborative information underlying multi-order interactions.

Collaborative Embedding Construction

Having transformed the user-item graph into a new language corpus rich in collaborative information, we proceed to fine-tune LLMs to incorporate this information. After fine-tuning, to enable efficient inference across the entire item set, we construct collaborative user and item embeddings by aggregating the hidden representations derived from the curated corpus.

Fine-tuning LLMs on the Corpus. Within these graph sentences, multi-order collaborative information enables distant tokens to influence the prediction of the next token in the sentence. To capture this collaborative information, we fine-tune LLMs by maximizing the likelihood of predicting the next token within the graph corpus. Formally, the fine-tuning objective is defined as:

$$\mathcal{L}_{\text{pre}} = - \sum_{q=1}^{|\mathcal{S}'|} \sum_{j=1}^{|\mathbf{s}_q|} \log P(s_{q,j} | \mathbf{s}_{q,<j}, \mathbf{W}_p, \hat{\mathbf{W}}), \quad (7)$$

where $\mathbf{W}_p \in \mathbb{R}^{|\mathcal{V}| \times d}$ is the learnable head layer parameter for predicting the next token. By optimizing this objective, the fine-tuned LLM learns to capture multi-order collaborative information embedded in the graph sentences. To control the memory usage, we incorporate the dynamic memory bank mechanism and the details are put in the Appendix E.

Graph Sentence Representation Aggregation. After fine-tuning, we obtain a well-trained LLM. However, directly deploying LLMs in an online recommender system would incur prohibitive latency. To address this issue, we precompute collaborative user and item embeddings offline, enabling efficient recommendation at inference time.

LLM is used to compute hidden representations for each graph token in the vocabulary:

$$\mathbf{h}_{q,j} = \text{LLM}(\{s_{q,1}, s_{q,2}, \dots, s_{q,j-1}\}), \quad (8)$$

where $\mathbf{h}_{q,j}$ represents the hidden representation of the j -th token in sentence \mathbf{s}_q , computed based on its preceding tokens $s_{q,1}, s_{q,2}, \dots, s_{q,j-1}$.

To comprehensively encode multi-order collaborative information, we aggregate the hidden representations of the target user v_u across multiple graph sentences. The aggregation is defined as:

$$\mathbf{h}_{v_u} = \sum_{p \in \mathcal{K}} \mathbf{h}_{q,p},$$

$$\mathcal{K} = \{p : \mathbf{s}_q \in \mathcal{S}', s_{q,p} = v_u, p = |\mathbf{s}_q| - k, 0 < k < |\mathbf{s}_q|\}$$
(9)

where \mathcal{K} specifies the valid positions of the token v_u across the sampled corpus \mathcal{S}' , and k is a hyperparameter controlling

the allowed positions of user v_u within a sentence. The item embeddings are computed analogously through the same aggregation process.

By precomputing user and item embeddings offline, NLCF enables efficient real-time inference through simple inner product search with these embeddings, eliminating the need for costly LLM computations during serving.

Experiments

We conduct extensive experiments on three real-world datasets to evaluate the effectiveness and efficiency of the NLCF framework. Our experimental study aims to address the following research questions: **RQ1:** How does NLCF perform compared to state-of-the-art graph CF and LLM-based baselines? **RQ2:** How do different design choices affect NLCF’s performance, particularly regarding sampling strategies and LLM backbone selections? **RQ3:** How sensitive is NLCF to key hyperparameters, such as sampling ratio and sentence length? **RQ4:** How does NLCF perform in real-world recommendation applications?

Experimental Settings

Datasets. We evaluate NLCF on three real-world datasets: Steam (Kang and McAuley 2018), ML-1M and ML-10M (Harper and Konstan 2016). Details about these datasets are shown in Appendix F.1. Specifically, Steam contains 918,951 interactions, 41,008 users, and 2,438 items. ML-10M contains 2,340,369 interactions, 69,428 users, and 5,180 items. ML-1M contains 370,647 interactions, 4,869 users, and 1,818 items.

Baseline Methods. We compare NLCF with the following groups of baselines: Traditional Graph CF Baselines: (i) LightGCN (He et al. 2020), (ii) LightGCL (Cai et al. 2023), (iii) HMLET (Kong et al. 2022), and (iv) AFDGCF (Wu et al. 2024b); LLM-based Baselines: TransRec (Lin et al. 2024b), LLM-CF (Sun et al. 2024), LETTER (Wang et al. 2024), LC-Rec (Zheng et al. 2024) and LLMEmb (Liu et al. 2025). Details about these baselines are shown in Appendix F.2. To evaluate the effectiveness of the similarity-based sampling approach, we compare it with random sampling.

Implementation Details. All experiments are conducted using publicly released codes, with each baseline running on **one dedicated** NVIDIA A100-SXM4-40GB GPU. The software environment is based on 20.04.6. The Python version is 3.9.22. We use the Hugging Face Transformers library 4.45.2. And we use metric Precision@ N (Zhuang et al. 2025; Zhang et al. 2025) and NDCG@ N (Zhou et al. 2023). More implementation details are shown in Appendix F.3.

Main Comparison (RQ1)

Tables 1 and 4 present the overall comparison across three datasets, revealing two key observations:

First, LLM-based approaches do not consistently achieve performance improvement over traditional graph CF approaches across all metrics. LLM-based approaches face

challenges in effectively incorporating collaborative information from user-item graph. While methods like LC-Rec and LETTER attempt to integrate collaborative embeddings from traditional graph CF models, these methods that rely on intermediate embeddings instead of directly modeling multi-order interactions suffer from collaborative information loss. This limitation may explain their performance degradation as metric k increases and their inability to consistently outperform state-of-the-art graph CF models, particularly in Precision metrics.

Second, NLCF achieves the most superior performance across all three datasets, demonstrating a successful paradigm shift. Rather than relying on traditional graph CF methods to model multi-order interactions, NLCF enables LLMs to directly capture the collaborative information from the user-item graph, marking a paradigm shift.

Efficiency Comparison (RQ1)

Table 2 presents the training and testing efficiency results across three datasets, revealing that most LLM-based baselines exhibit significantly longer training and inference times compared to NLCF. This performance gap is expected, as LLM-specific strategies—such as user sequence augmentation and diverse prompt template designs—add considerable computational overhead. During inference, many LLM-based approaches still rely on active LLM computations, further increasing inference costs. While LLMEmb fine-tunes LLMs to generate side information from user-item interactions, it struggles with the inefficiency of lengthy natural language graph descriptions during training. As a result, it only improves inference efficiency compared to earlier methods. Notably, LLMEmb does not model user-item interactions directly; instead, it fine-tunes LLMs based on user-item attributes, which fails to encode multi-order user-item interactions and limits potential performance gains over other LLM-based approaches. In contrast, NLCF achieves remarkable efficiency through two key design choices: (i) during training, it fine-tunes LLMs on concise graph sentences derived from the user-item graph; and (ii) during inference, it constructs collaborative embeddings, enabling efficient item retrieval across the entire item set. These designs significantly reduce both training and inference computational costs, making NLCF more efficient than most LLM-based baselines.

Ablation Study (RQ2)

We examine different variants of NLCF through experiments shown in Figure 4 and 6. Since NLCF employs a straightforward architecture without composite components or multiple training objectives, we focus our analysis on two key factors beyond hyper-parameters: the LLM backbone and sampling strategy. Our experiments reveal two significant observations:

First, NLCF’s performance aligns with empirical neural scaling laws. As the LLM parameter size increases, the model’s performance improves substantially. For instance, the Llama 1B model performs well, highlighting the power of LLMs, while the 7B model achieves better results. However, the performance improvement from 7B to 8B models

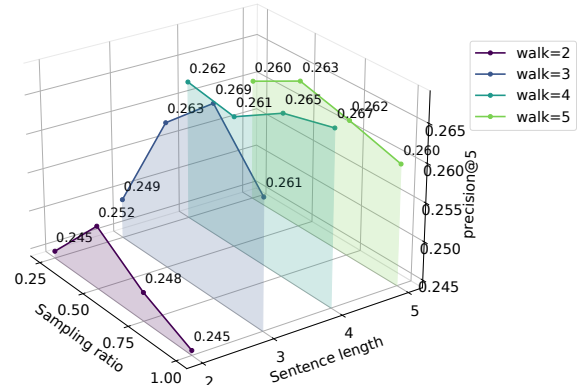


Figure 3: The impact of graph sentence length l and sampling ratio α on Precision@5 on ML-1M dataset.

is relatively modest, suggesting a saturation point in model scaling benefits.

Second, our comparison with alternative sampling strategies demonstrates the superiority of our similarity-based approach across almost all sampling ratios. This advantage likely stems from two factors: random sampling fails to account for subgraph density and struggles to preserve representative samples from each subgraph, while our similarity-based sampling controls the granularity needed for effective graph sentence grouping, particularly at lower sampling ratios where performance gap becomes significant.

Hyper-parameter Sensitivity (RQ3)

We extensively evaluate the effect of hyper-parameters on the performance of NLCF, including sampling ratio α , graph sentence length l , and similarity threshold t . The results, presented in Figure 3, Figure 5, and Table 5, reveal two observations:

First, NLCF achieves optimal performance with a graph sentence length $l = 4$. This length indicates that third-order collaborative information suffices for high-quality recommendations, while higher-order information may introduce noise without contributing positively to performance. Increasing the sampling ratio improves model performance. Notably, performance declines only slightly when the ratio is lowered from 1 to 0.75, indicating the effectiveness of our sampling method.

Second, the similarity threshold demonstrates a critical role in sampling effectiveness. A moderate threshold value enables NLCF to maintain a more representative subset of samples. High threshold values impose overly strict similarity constraints, effectively reducing the sampling to random selection as most graph sentences are deemed dissimilar. Conversely, low threshold values lack discriminative power, marking most samples as similar and diminishing the sampling strategy’s effectiveness.

Model	Steam		ML-10M		ML-1M	
	NDCG@5	NDCG@10	NDCG@5	NDCG@10	NDCG@5	NDCG@10
LightGCN (He et al. 2020)	0.2744	0.2790	0.1988	0.2007	0.2695	0.2436
LightGCL (Cai et al. 2023)	0.2623	0.2807	0.1944	<u>0.2065</u>	0.2112	0.2045
HMLET (Kong et al. 2022)	0.2730	0.2853	0.1919	0.1950	0.2723	<u>0.2503</u>
AFDGCF (Wu et al. 2024b)	<u>0.2809</u>	<u>0.2897</u>	0.2002	0.2004	0.2694	0.2484
TransRec (Lin et al. 2024b)	0.2796	0.2843	0.1961	0.1975	0.2706	0.2435
LLM-CF (Sun et al. 2024)	0.2772	0.2829	0.1976	0.1993	<u>0.2735</u>	0.2478
LETTER (Wang et al. 2024)	0.2717	0.2776	0.1928	0.1883	0.2683	0.2350
LC-Rec (Zheng et al. 2024)	0.2615	0.2682	0.1865	0.1822	0.2626	0.2321
LLMEmb (Liu et al. 2025)	0.2785	0.2858	<u>0.2042</u>	0.2034	0.2711	0.2446
NLCF	0.2864	0.2948	0.2171	0.2147	0.2796	0.2558

Table 1: NDCG performance with $N = 5$ and 10 across Steam, ML-10M, and ML-1M datasets.

Model	Steam		ML-10M		ML-1M	
	Train	Test	Train	Test	Train	Test
TransRec	16h5m	3h25m	18h44m	4h18m	4h32m	52m42s
LLM-CF	13h24m	0.029s	16h58m	0.101s	3h47m	0.003s
LETTER	5h6m	3m31s	7h27m	5m35s	1h43m	1m16s
LC-Rec	11h49m	18m59s	13h22m	26m07s	3h18m	7m21s
LLMEmb	41m36s	0.029s	48m15s	0.101s	10m23s	0.003s
NLCF	42m43s	0.029s	44m19s	0.101s	8m12s	0.003s

Table 2: Training and test time comparison among LLM-based models (in seconds [s], minutes [m], and hours [h]).

A/B Test	PCTR	UCTR	GMV	ResTime
v.s. LightGCN	+4.15%	+3.11%	+5.78%	+ 1.46%
v.s. LLM-CF	+2.51%	+2.47%	+3.14%	- 20.17%

Table 3: Online A/B tests on the industrial platform.

Online Evaluation (RQ4)

We deploy NLCF on a billion-scale online shopping platform and conduct A/B testing to evaluate its performance. The platform serves hundreds of millions of users and billions of items, supported by two main components: an offline computing center and an online service center. The offline computing center processes user logs and generates a graph corpus from processed user interactions through distributed jobs. It trains NLCF and generates collaborative embeddings for each user and item. Notably, we do not require a tokenizer to handle billions of tokens for mapping. Another offline computing center task is the routine update of collaborative embeddings. This process involves generating a new graph corpus from recent user interactions on a daily basis. Because the final embeddings are constructed via summation, they can be updated incrementally by simply adding the hidden representations derived from this new corpus. The generated collaborative embeddings are then transmitted to the online service center, which uses these pre-computed user embeddings to efficiently retrieve items from a massive item pool, thereby eliminating the need for costly real-time LLM reasoning. NLCF is deployed as a recall model, replacing two baselines: the traditional graph CF model LightGCN and the LLM-based approach LLM-CF.

The performance metrics shown in Table 3 are averaged over eight consecutive weeks, with each model allocated 5% of online traffic. Compared to LightGCN, NLCF achieves significant improvements: +4.15% in PCTR, +3.11% in UCTR, and +5.78% in GMV, demonstrating its superiority in capturing user preferences and driving item consumption. Despite using higher-dimensional embeddings, NLCF incurs only a 1.46% increase in latency, ensuring scalability. Against LLM-CF, NLCF shows moderate gains: +2.51% in PCTR, +2.47% in UCTR, and +3.14% in GMV, confirming the benefits of learning collaborative embeddings over data augmentation. Additionally, NLCF reduces latency by 20.17% compared to LLM-CF, highlighting embedding-only inference efficiency.

Conclusion

This study introduces a novel paradigm for fine-tuning LLMs for CF task, enabling efficient modeling of multi-order interactions in the user-item graph to learn effective user and item embeddings. Existing methods that prompt LLMs with graph descriptions face two major limitations: (i) Efficiency – slower inference due to token-by-token item generation compared to embedding-based search, and (ii) Topological Modeling – difficulty encoding multi-order interactions and collaborative information from user-item graphs. To address these challenges, we propose NLCF, which treats the user-item graph as a new language. This approach is built on two insights: (1) LLMs excel at learning new languages with suitable corpora, and (2) token transition probabilities in language align with node transition probabilities in graphs. NLCF operates in two stages: (i) transforming the user-item graph into a language corpus that encodes multi-order interactions, and (ii) fine-tuning LLMs on this corpus to capture underlying multi-order collaborative information and construct collaborative user and item embeddings for efficient search. Extensive offline experiments demonstrate NLCF’s superior performance over both LLM-based and traditional graph CF baselines while significantly improve computational efficiency over LLM-based baselines. And online A/B tests conducted on a world-leading shopping platform validate NLCF’s effectiveness in real-world applications.

Acknowledgements

The work described in this paper was fully supported by a grant from the Innovation and Technology Commission of the Hong Kong Special Administrative Region, China (Project No. GHP/391/22).

References

- Bao, K.; Zhang, J.; Zhang, Y.; Wang, W.; Feng, F.; and He, X. 2023. Tallrec: An effective and efficient tuning framework to align large language model with recommendation. In *Proceedings of the 17th ACM Conference on Recommender Systems*, 1007–1014.
- Cai, X.; Huang, C.; Xia, L.; and Ren, X. 2023. LightGCL: Simple Yet Effective Graph Contrastive Learning for Recommendation. In *The Eleventh International Conference on Learning Representations*.
- Chen, S.; Zhang, Q.; Dong, J.; Hua, W.; Li, Q.; and Huang, X. 2024. Entity Alignment with Noisy Annotations from Large Language Models. In *The Thirty-eighth Annual Conference on Neural Information Processing Systems*.
- Chen, S.; Zhou, C.; Yuan, Z.; Zhang, Q.; Cui, Z.; Chen, H.; Xiao, Y.; Cao, J.; and Huang, X. 2025. You Don't Need Pre-built Graphs for RAG: Retrieval Augmented Generation with Adaptive Reasoning Structures. In *The Fortieth AAAI Conference on Artificial Intelligence*.
- Grover, A.; and Leskovec, J. 2016. node2vec: Scalable feature learning for networks. In *Proceedings of the 22nd ACM SIGKDD international conference on Knowledge discovery and data mining*, 855–864.
- Hamilton, W.; Ying, Z.; and Leskovec, J. 2017. Inductive representation learning on large graphs. *Advances in neural information processing systems*, 30.
- Harper, F. M.; and Konstan, J. A. 2016. The MovieLens Datasets: History and Context. *ACM Trans. Interact. Intell. Syst.*, 5(4): 19:1–19:19.
- He, X.; Deng, K.; Wang, X.; Li, Y.; Zhang, Y.; and Wang, M. 2020. Lightgcn: Simplifying and powering graph convolution network for recommendation. In *Proceedings of the 43rd International ACM SIGIR conference on research and development in Information Retrieval*, 639–648.
- Hong, Z.; Yuan, Z.; Chen, H.; Zhang, Q.; Huang, F.; and Huang, X. 2024. Knowledge-to-SQL: Enhancing SQL Generation with Data Expert LLM. In *Findings of the Association for Computational Linguistics: ACL 2024*.
- Hong, Z.; Yuan, Z.; Zhang, Q.; Chen, H.; Dong, J.; Huang, F.; and Huang, X. 2025. Next-generation database interfaces: A survey of llm-based text-to-sql. *IEEE Transactions on Knowledge and Data Engineering*.
- Hu, E. J.; Shen, Y.; Wallis, P.; Allen-Zhu, Z.; Li, Y.; Wang, S.; Wang, L.; and Chen, W. 2021. Lora: Low-rank adaptation of large language models. *arXiv preprint arXiv:2106.09685*.
- Huang, F.; Yang, Z.; Jiang, J.; Bei, Y.; Zhang, Y.; and Chen, H. 2024. Large Language Model Interaction Simulator for Cold-Start Item Recommendation. *arXiv preprint arXiv:2402.09176*.
- Jiang, Y.; Yang, Y.; Xia, L.; Luo, D.; Lin, K.; and Huang, C. 2024. RecLM: Recommendation Instruction Tuning. *arXiv preprint arXiv:2412.19302*.
- Kang, W.; and McAuley, J. J. 2018. Self-Attentive Sequential Recommendation. In *IEEE International Conference on Data Mining, ICDM 2018, Singapore, November 17-20, 2018*, 197–206. IEEE Computer Society.
- Kim, S.; Kang, H.; Choi, S.; Kim, D.; Yang, M.; and Park, C. 2024. Large language models meet collaborative filtering: An efficient all-round llm-based recommender system. In *Proceedings of the 30th ACM SIGKDD Conference on Knowledge Discovery and Data Mining*, 1395–1406.
- Kong, T.; Kim, T.; Jeon, J.; Choi, J.; Lee, Y.-C.; Park, N.; and Kim, S.-W. 2022. Linear, or non-linear, that is the question! In *Proceedings of the fifteenth ACM international conference on web search and data mining*, 517–525.
- Liao, J.; Li, S.; Yang, Z.; Wu, J.; Yuan, Y.; Wang, X.; and He, X. 2024. Llara: Large language-recommendation assistant. In *Proceedings of the 47th International ACM SIGIR Conference on Research and Development in Information Retrieval*, 1785–1795.
- Lin, J.; Shan, R.; Zhu, C.; Du, K.; Chen, B.; Quan, S.; Tang, R.; Yu, Y.; and Zhang, W. 2024a. Rella: Retrieval-enhanced large language models for lifelong sequential behavior comprehension in recommendation. In *Proceedings of the ACM on Web Conference 2024*, 3497–3508.
- Lin, X.; Wang, W.; Li, Y.; Feng, F.; Ng, S.-K.; and Chua, T.-S. 2024b. Bridging items and language: A transition paradigm for large language model-based recommendation. In *Proceedings of the 30th ACM SIGKDD Conference on Knowledge Discovery and Data Mining*, 1816–1826.
- Lin, X.; Wang, W.; Li, Y.; Yang, S.; Feng, F.; Wei, Y.; and Chua, T.-S. 2024c. Data-efficient Fine-tuning for LLM-based Recommendation. In *Proceedings of the 47th International ACM SIGIR Conference on Research and Development in Information Retrieval*, 365–374.
- Lin, Z.; Tian, C.; Hou, Y.; and Zhao, W. X. 2022. Improving graph collaborative filtering with neighborhood-enriched contrastive learning. In *Proceedings of the ACM web conference 2022*, 2320–2329.
- Liu, Q.; Wu, X.; Wang, W.; Wang, Y.; Zhu, Y.; Zhao, X.; Tian, F.; and Zheng, Y. 2025. LLMEmb: Large Language Model Can Be a Good Embedding Generator for Sequential Recommendation. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 39, 12183–12191.
- Qu, H.; Fan, W.; Zhao, Z.; and Li, Q. 2024. TokenRec: Learning to Tokenize ID for LLM-based Generative Recommendation. *arXiv preprint arXiv:2406.10450*.
- Sun, Z.; Si, Z.; Zang, X.; Zheng, K.; Song, Y.; Zhang, X.; and Xu, J. 2024. Large language models enhanced collaborative filtering. In *Proceedings of the 33rd ACM International Conference on Information and Knowledge Management*, 2178–2188.
- Wang, W.; Bao, H.; Lin, X.; Zhang, J.; Li, Y.; Feng, F.; Ng, S.-K.; and Chua, T.-S. 2024. Learnable item tokenization for generative recommendation. In *Proceedings of the 33rd*

- ACM International Conference on Information and Knowledge Management, 2400–2409.
- Wang, X.; He, X.; Wang, M.; Feng, F.; and Chua, T.-S. 2019. Neural graph collaborative filtering. In *Proceedings of the 42nd international ACM SIGIR conference on Research and development in Information Retrieval*, 165–174.
- Wei, W.; Ren, X.; Tang, J.; Wang, Q.; Su, L.; Cheng, S.; Wang, J.; Yin, D.; and Huang, C. 2024. Llmrec: Large language models with graph augmentation for recommendation. In *Proceedings of the 17th ACM International Conference on Web Search and Data Mining*, 806–815.
- Wu, J.; Liu, Q.; Hu, H.; Fan, W.; Liu, S.; Li, Q.; Wu, X.-M.; and Tang, K. 2023. Leveraging Large Language Models (LLMs) to Empower Training-Free Dataset Condensation for Content-Based Recommendation. *arXiv preprint arXiv:2310.09874*.
- Wu, L.; Qiu, Z.; Zheng, Z.; Zhu, H.; and Chen, E. 2024a. Exploring large language model for graph data understanding in online job recommendations. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 38, 9178–9186.
- Wu, W.; Wang, C.; Shen, D.; Qin, C.; Chen, L.; and Xiong, H. 2024b. Afdgcf: Adaptive feature de-correlation graph collaborative filtering for recommendations. In *Proceedings of the 47th International ACM SIGIR Conference on Research and Development in Information Retrieval*, 1242–1252.
- Wu, Z.; Pan, S.; Chen, F.; Long, G.; Zhang, C.; and Philip, S. Y. 2020. A comprehensive survey on graph neural networks. *IEEE transactions on neural networks and learning systems*, 32(1): 4–24.
- Yang, H.; Wang, X.; Tao, Q.; Hu, S.; Lin, Z.; and Zhang, M. 2024. GL-Fusion: Rethinking the Combination of Graph Neural Network and Large Language model. *arXiv preprint arXiv:2412.06849*.
- Yang, Z.; Wu, J.; Luo, Y.; Zhang, J.; Yuan, Y.; Zhang, A.; Wang, X.; and He, X. 2023. Large language model can interpret latent space of sequential recommender. *arXiv preprint arXiv:2310.20487*.
- Yuan, Z.; Chen, H.; Hong, Z.; Zhang, Q.; Huang, F.; Li, Q.; and Huang, X. 2025. Knapsack optimization-based schema linking for llm-based Text-to-SQL generation. *arXiv preprint arXiv:2502.12911*.
- Yuan, Z.; Yuan, F.; Song, Y.; Li, Y.; Fu, J.; Yang, F.; Pan, Y.; and Ni, Y. 2023. Where to go next for recommender systems? id-vs. modality-based recommender models revisited. In *Proceedings of the 46th International ACM SIGIR Conference on Research and Development in Information Retrieval*, 2639–2649.
- Zhang, A.; Deng, Y.; Lin, Y.; Chen, X.; Wen, J.-R.; and Chua, T.-S. 2024a. Large Language Model Powered Agents for Information Retrieval. In *Proceedings of the 47th International ACM SIGIR Conference on Research and Development in Information Retrieval*, 2989–2992.
- Zhang, Q.; Chen, S.; Bei, Y.; Yuan, Z.; Zhou, H.; Hong, Z.; Chen, H.; Xiao, Y.; Zhou, C.; Dong, J.; et al. 2025. A survey of graph retrieval-augmented generation for customized large language models. *arXiv preprint arXiv:2501.13958*.
- Zhang, Y.; Bao, K.; Yan, M.; Wang, W.; Feng, F.; and He, X. 2024b. Text-like Encoding of Collaborative Information in Large Language Models for Recommendation. *arXiv preprint arXiv:2406.03210*.
- Zhang, Y.; Feng, F.; Zhang, J.; Bao, K.; Wang, Q.; and He, X. 2023. Collm: Integrating collaborative embeddings into large language models for recommendation. *arXiv preprint arXiv:2310.19488*.
- Zhao, W. X.; Zhou, K.; Li, J.; Tang, T.; Wang, X.; Hou, Y.; Min, Y.; Zhang, B.; Zhang, J.; Dong, Z.; et al. 2023. A survey of large language models. *arXiv preprint arXiv:2303.18223*.
- Zheng, B.; Hou, Y.; Lu, H.; Chen, Y.; Zhao, W. X.; Chen, M.; and Wen, J.-R. 2024. Adapting large language modelA survey of large language models by integrating collaborative semantics for recommendation. In *2024 IEEE 40th International Conference on Data Engineering (ICDE)*, 1435–1448. IEEE.
- Zhou, C.; Du, J.; Zhou, H.; Chen, H.; Huang, F.; and Huang, X. 2025a. Text-Attributed Graph Learning with Coupled Augmentations. In *Proceedings of the 31st International Conference on Computational Linguistics*, 10865–10876.
- Zhou, C.; Wang, Z.; Chen, S.; Du, J.; Zheng, Q.; Xu, Z.; and Huang, X. 2025b. Taming language models for text-attributed graph learning with decoupled aggregation. In *Proceedings of the 63rd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, 3463–3474.
- Zhou, H.; Chen, H.; Dong, J.; Zha, D.; Zhou, C.; and Huang, X. 2023. Adaptive popularity debiasing aggregator for graph collaborative filtering. In *Proceedings of the 46th international ACM SIGIR conference on research and development in information retrieval*, 7–17.
- Zhou, H.; Du, J.; Zhou, C.; Yang, C.; Xiao, Y.; Xie, Y.; and Huang, X. 2025c. Each Graph is a New Language: Graph Learning with LLMs. *arXiv preprint arXiv:2501.11478*.
- Zhou, H.; Yu, K.; Zhang, Q.; Chen, H.; Zha, D.; Pei, W.; Kong, A.; and Huang, X. 2025d. Self-Monitoring Large Language Models for Click-Through Rate Prediction. *ACM Transactions on Information Systems*, 44(1): 1–25.
- Zhou, H.; Zhou, S.; Chen, H.; Liu, N.; Yang, F.; and Huang, X. 2024. Enhancing explainable rating prediction through annotated macro concepts. In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, 11736–11748.
- Zhuang, L.; Chen, S.; Xiao, Y.; Zhou, H.; Zhang, Y.; Chen, H.; Zhang, Q.; and Huang, X. 2025. LinearRAG: Linear Graph Retrieval Augmented Generation on Large-scale Corpora. *arXiv preprint arXiv:2510.10114*.