GaCLLM: Graph-aware Convolutional Large Language Model for Recommendation

Anonymous ACL submission

Abstract

Leveraging auxiliary textual data can help with user profiling and item characterization in recommender systems (RSs). However, incomplete item descriptions and the subjectivity of 004 user-uploaded content limit the potential of textual information in RSs. Although large lan-007 guage models (LLMs) emerge as promising tools for description enhancement, LLMs may suffer from hallucinations without fully exploring user-item collaborative information. To this end, we propose a Graph-aware Convolutional 011 LLM method, which captures fine-grained collaborative information behind high-order relations in the user-item graph. To bridge the gap 015 between graph structures and LLMs, we employ the LLM as an aggregator for graph convolution process, eliciting it to infer the graph-017 based knowledge iteratively. To mitigate the in-019 formation overload associated with large-scale graphs, we segment the graph processing into manageable steps, progressively incorporating multi-hop information in a least-to-most manner. Experiments on three real-world datasets demonstrate that our method consistently outperforms state-of-the-art approaches.¹

1 Introduction

034

Recommender systems (RSs) are pivotal in delivering personalized services to users for their satisfaction and platform profitability. Traditionally, RSs heavily rely on user-item interaction records (Koren et al., 2009) but face challenges with data sparsity (Sun et al., 2019). Recently, there has been a trend towards utilizing auxiliary textual information for recommendation (Torbati et al., 2023). However, texts with users and items often suffer from incompleteness and bias, with users offering vague self-descriptions and providers giving sparse or strategically biased item descriptions. Such texts negatively impact user profiling and item characterization, hindering accurate recommendations. 039

041

043

044

045

047

051

053

054

059

060

061

062

063

064

065

066

067

068

069

070

071

073

074

075

076

078

079

To enhance the reliability and completeness of textual descriptions, recent approaches have employed large language models (LLMs) to generate LLM-driven descriptions based on raw contents and task-specific prompt instructions (Zheng et al., 2023; Wu et al.; Liu et al., 2023; Wang et al., 2024b), such as incorporating users' behaviors as supplemental knowledge for retrieval-augmented generation (Du et al., 2024; Liu et al., 2024b). Nevertheless, these methods still suffer from unreliable and inaccurate textual generation due to the limited scope of information observed by LLMs and the lack of collaborative user-item insights.

To this end, inspired by the success of graph convolutional networks (GCNs) (Kipf and Welling, 2016), we propose Graph-aware Convolutional LLM (GaCLLM) to integrate collaborative useritem information into LLMs to enhance reasoning and mitigate hallucinations. We focus on two main challenges: the constraints of the context length, and the incompatibility between graph structures and LLMs. First, large-scale user-item graphs pose context length limitations for LLM inputs by simply describing them in a textual format. Specifically, LLMs often struggle to robustly access and utilize information from lengthy contextual inputs, particularly when the critical information (e.g., the key entity in the graph) is located in the middle (Liu et al., 2024a). Second, text-based LLMs are inherently ill-suited for processing structured graph data. Existing methods convert graph data into textual form using templates and sampling strategies (Wang et al., 2023; Wu et al., 2024a). However, these methods limit the LLMs' ability to maintain a global perspective on graphs, thereby hindering their full potential in utilizing reasoning skills for graph-based knowledge.

To tackle these challenges, we develop a convolutional inference strategy to integrate high-order

¹Our code is available at https://anonymous.4open. science/r/GaCLLM_code-C326.

relations from the user-item interaction graph into LLMs. Specifically, we segment the graph process-081 ing into manageable steps in a least-to-most (Zhou et al., 2022) manner, iteratively incorporating multihop neighbor information to refine each node's (i.e., user or item) description. Therefore, the overload of describing the graph can be segmented into several steps, drastically reducing the input's context length for LLMs. It can alleviate the limitations of lengthy inputs to capture critical information for LLM-driven reasoning. To align LLMs with graph structures, we employ the LLM as an aggregator function and maintain a global perspective on graphs. Specifically, the LLM assimilates information from neighboring nodes and ensures 094 layer-by-layer propagation throughout the graph. By leveraging high-order relations in the user-item interaction graph, our method enhances reasoning capabilities and mitigates hallucinations in the LLM-driven descriptions. Finally, we fuse these LLM-driven descriptions into behavioral graph em-100 beddings to bridge the gap between text information and structural data in the user-item graph for recommendation. We conduct extensive experi-103 104 ments on multiple real-world datasets to show that our method consistently outperforms state-of-theart approaches, validating the effectiveness of our 106 proposed strategy through comprehensive ablation studies and in-depth analysis. 108

2 RELATED WORK

110

111

112

113

114

115

116

117

118

119

121

122

123

125

126

127

128

2.1 Graph-based Recommendation

Graph-based recommender systems (Kipf and Welling, 2016; Huang et al., 2024; Yan et al., 2024) employ deep neural networks to model the complex user-item interactions within graph structures. LightGCN (He et al., 2020) streamlines GCNs for collaborative filtering with simplicity and effectiveness. Many studies build upon LightGCN using techniques like contrastive learning (Yu et al., 2022; Chen et al., 2023), transformer (Wei et al., 2023), neighborhood-structure (Lin et al., 2022), and selfsupervised learning (Wu et al., 2021). However, they mainly focus on aggregating node embeddings and fail to extract insights from textual descriptions for recommendation.

2.2 LLM for Recommendation

There is increasing interest in leveraging LLMs in recommender systems (Wu et al., 2024b; Lyu et al., 2024). Non-tuning methods (Kuo and Chen, 2023; Senel et al., 2024) assume that LLMs already possess recommendation capabilities and use them to produce results directly through specific prompts (Kang et al., 2023; Zhang et al., 2023) and in-context learning (Hou et al., 2024; Wang and Lim, 2024). The tuning paradigm (Lu et al., 2024) employs LLM as feature extractors for downstream tasks, aiming to capture contextual information for a precise understanding of user profiles (Zheng et al., 2023; Du et al., 2024), user attributes (Wang et al., 2024a), and item descriptions (Liu et al., 2024b). However, relying only on raw text and ignoring graph knowledge leads to hallucinations. To alleviate this, we aggregate additional information from graphs into LLMs for more reliable textual data to enhance recommendation results.

129

130

131

132

133

134

135

136

137

138

139

140

141

142

143

144

145

146

147

148

149

150

151

152

153

154

155

156

157

158

159

160

161

162

163

165

167

169

170

171

172

173

174

175

176

177

2.3 LLM with Graph Data

Integrating LLMs with graph data (Li et al., 2024; Ye et al., 2023) effectively leverages the rich structure and relationships. Supervised methods use LLMs for graph-aware tasks via encoding text into node embeddings (Chen et al., 2024; Zhang et al., 2021) and incorporating graph elements into training (Sun et al., 2021; Yasunaga et al., 2022; Xie et al., 2023; Zhang et al., 2024b). However, they mainly compress graph knowledge into model parameters, overlooking the LLMs' reasoning mechanism. Unsupervised methods (Wang et al., 2023; Andrus et al., 2022; Wu et al., 2024a; Zhang et al., 2024a) convert graph information into text via templates or sampling strategies for LLMs to process. However, they lack a global view of the graph and still fail to fully exploit LLMs' reasoning potential. Thus, we propose to incorporate high-order graph information into LLMs by iteratively distilling information from neighbors, enhancing its reasoning while reducing token overhead.

3 Methodology

3.1 Problem Definition

We denote $\mathcal{U} = \{u_1, \dots, u_N\}$ and $\mathcal{I} = \{i_1, \dots, i_M\}$ as the sets of users and items, where N and M are sizes. The interaction records between users and items can be denoted as an interaction matrix $\mathcal{R} \in \mathbb{R}^{N \times M}$ where $\mathcal{R}_{u,i} = 1$ if user u interacted with item i, and 0 otherwise. We also possess the textual information (e.g., user resumes and job descriptions in online recruitment scenarios) of both users, denoted as $\mathcal{T}_u = [w_1, \dots, w_{l_u}]$ with length l_u for user u, and items, denoted as



Figure 1: The overall architecture of the proposed method GaCLLM.

1

 $\mathcal{T}_i = [w_1, \cdots, w_{l_i}]$ with length l_i for item *i*, and w_k represents the *k*-th word.

In this paper, our goal is to learn a matching function g(u, i) using the interaction records \mathcal{R} and the textual descriptions. Our task is to recommend K items that a user is most likely to prefer, as known as top-K recommendation.

3.2 Overview

178

179

181

183

185

187

188

189

191

192

194

195

197

199

201

206

210

The overall architecture of GaCLLM is shown in Figure 1. First, we perform supervised fine-tuning (SFT) for LLM to strengthen its effectiveness in the task-related domain. Second, we propose an LLM-based graph-aware convolutional inference strategy to enhance user and item descriptions progressively. Third, we align and integrate the generated text with behavioral information captured through graph-based embeddings. Last, we present the objective function and model learning process.

3.3 Supervised Fine-tuning

To fully exploit the potential of the LLM in understanding the task-related domain, we begin with fine-tuning it on domain-specific data. This involves the training of the LLM using descriptions from matched user-item pairs, enabling it to learn the alignment between user and item descriptions. Specifically, we employ the prompt template: "Query: Given an item's description, generate a user's description that fits it. The item's description is [*Item Desc*]. Answer: ", where [*Item Desc*] represents the actual description of the item. The prompt for inferring item descriptions with the provided user description is designed symmetrically. The optimization process involves minimizing the negative log-likelihood loss for these templates:

$$\mathcal{L}_{\text{sft}} = -\sum_{k=1}^{|T_{\text{Answer}}|} \log \Pr(w_k \mid w_{< k}, T_{\text{Query}}), \quad (1)$$

where w_k denotes the k-th word in Answer sentence T_{Answer} , and $\Pr(T_{\text{Answer}}|T_{\text{Query}})$ denotes the generation probability for the produced answer with a given query. This process uses parameterefficient fine-tuning techniques.

3.4 Convolutional Inference Strategy

Graph Construction. To explore the structured graph with high-order descriptive texts for LLMs, we organize the descriptions of users and items into a unified graph $\mathcal{G} = (\mathcal{V}, \mathcal{E})$ using the collaborative information among users and items. Specifically, the nodes \mathcal{V} in the graph represent users and items, i.e., $\mathcal{V} = \{u | u \in \mathcal{U}\} \cup \{i | i \in \mathcal{I}\}$. The edges \mathcal{E} are constructed by the interactions between users and items $\mathcal{R} \in \mathbb{R}^{N \times M}$, i.e., $\mathcal{E} = \{(u, i) | \mathcal{R}_{u,i} = 1\}$. Each node in the graph has a textual description, such as a user profile in a social network or a resume of a job seeker.

Least-to-Most Text Enhancement. Recognizing the extensive knowledge, advanced text comprehension, and reasoning capabilities of LLMs, we propose an LLM-based convolutional inference strategy to explore high-order relations among textual contents in the user-item interaction graph. To make user descriptions more representative, we leverage the LLM to rewrite a user's raw description T_u by the descriptions of items that the user has interacted with:

$$\mathcal{T}'_{u} = \text{LLM}(\mathsf{P}_{user}(\mathcal{T}_{u}, \{\mathcal{T}_{i}: (u, i) \in \mathcal{E}\})), \quad (2)$$

where P_{user} denotes the prompt template for generating user descriptions. Similarly, to enhance 211

212

213



Figure 2: The prompt design for job recommendation (top) and social recommendation (bottom).

item description \mathcal{T}_i , we use the LLM to produce the enhanced version considering the descriptions of users by interaction:

244

245

246

247

248

250

254

255

258

264

265

268

269

270

271

273

$$\mathcal{T}'_i = \text{LLM}(P_{\texttt{item}}(\mathcal{T}_i, \{\mathcal{T}_u : (u, i) \in \mathcal{E}\})), \quad (3)$$

where P_{item} denotes the prompt template for generating item descriptions. The design of a prompt template varies with the tasks. In this paper, we focus on job and social recommendation tasks. The details of the prompts are shown in Figure 2.

To enable LLMs to effectively explore the structured graph, we iteratively use them to refine the descriptions of nodes (users and items) step by step. Specifically, we set the first-layer descriptions of users $\{\mathcal{L}_u^{(1)} | u \in \mathcal{U}\}$ by raw texts provided by users, i.e., $\mathcal{L}_u^{(1)} = \mathcal{T}_u$, and we set the first-layer descriptions of items $\{\mathcal{L}_i^{(1)} | i \in \mathcal{I}\}$ by raw texts given by item providers, i.e., $\mathcal{L}_i^{(1)} = \mathcal{T}_i$. We employ the LLM as an "aggregator" in the graph convolutional process, enhancing its ability to infer graph-based knowledge through iterative steps. The updated user and item descriptions after each iteration are generated as follows:

$$\begin{aligned} \mathcal{L}_{u}^{(l+1)} &= \text{LLM}(\text{P}_{\texttt{user}}(\mathcal{L}_{u}^{(l)}, \{\mathcal{L}_{i}^{(l)} : (u, i) \in \mathcal{E}\})), \end{aligned} \tag{4} \\ \mathcal{L}_{i}^{(l+1)} &= \text{LLM}(\text{P}_{\texttt{item}}(\mathcal{L}_{i}^{(l)}, \{\mathcal{L}_{u}^{(l)} : (u, i) \in \mathcal{E}\})), \end{aligned}$$

where $\mathcal{L}_{u}^{(l+1)}$ and $\mathcal{L}_{i}^{(l+1)}$ denote the descriptions of users and items at (l+1)-th layer after l iterations of generation, capturing l-hop descriptive information within the graph. After L iterations of this LLMbased convolutional inference strategy, we obtain



Figure 3: The comparison of token usage of convolutional inference strategy (left) and plain description strategy (right) in text enhancement.

progressively refined descriptions across multiple layers for both users and items. 274

275

276

277

278

279

281

282

284

285

287

288

290

291

292

293

294

295

296

298

299

300

301

302

303

304

305

306

307

308

309

310

311

312

313

314

Token Effectiveness and Efficiency. Compared with organizing all hierarchical node descriptions in the graph structure into a single, *plain* paragraph of prompt (e.g., listing each node and its multi-hop neighbors along with their descriptions), the proposed convolutional inference strategy improves both effectiveness and efficiency in token usage.

First, it optimizes the capture of graph-related information within the limited context length of LLMs. Specifically, the proposed strategy decomposes the ultimate task of description enhancement into multiple steps, where each step (layer) only integrates the descriptions of direct (one-hop) neighbors for the target node. This step-by-step approach effectively alleviates the issues of hallucination and distraction with long inputs, significantly reducing the number of tokens required for each inference.

Second, our convolutional inference strategy efficiently reduces the redundancy in describing the graph for target nodes. Specifically, when comparing the number of nodes required to capture Lhop graph-based information for each node, the proposed method incorporates $O(|\mathcal{G}| \cdot |\mathcal{N}| \cdot L)$ nodes into LLMs, where $|\mathcal{G}|$ denotes the number of nodes in the graph and $|\mathcal{N}|$ denotes the average number of neighbors of each node. In contrast, the *plain* description strategy needs to incorporate $O(|\mathcal{G}| \cdot (1 + \dots + |\mathcal{N}|^L))$ nodes into LLMs, leading to a significant increase in token usage. Therefore, by minimizing the overlap in node descriptions (such as the redundant description of common neighbors shown in Figure 3), our method enhances token efficiency.

3.5 Text-graph Alignment

To bridge the gap between LLM-driven text information and behavioral-based structural data in the user-item graph for recommendation, we propose to align the user and item descriptions with their corresponding graph embeddings in a unified man315

317

ner. Specifically, the GCN-based embeddings for users and items at the *l*-th layer, denoted as $e_u^{(l)}$ and $e_i^{(l)}$. They can be iteratively updated as follows:

318
$$\boldsymbol{e}_{u}^{(l+1)} = W_{l} \cdot \left[\sum_{(u,i)\in\mathcal{E}} \frac{\boldsymbol{e}_{i}^{(l)}}{\sqrt{|\mathcal{N}_{u}||\mathcal{N}_{i}|}} \oplus f(\mathcal{L}_{u}^{(l)})\right],$$
(6)

319

321

322

324

328

332

334

336

337

338

341

345

346

347

348

 $oldsymbol{e}_i^{(l+1)} = W_l \cdot [\sum_{(u,i) \in \mathcal{E}} rac{oldsymbol{e}_u^{(l)}}{\sqrt{|\mathcal{N}_u||\mathcal{N}_i|}} \oplus f(\mathcal{L}_i^{(l)})].$ (7)

Here, \mathcal{N}_u denotes the set of items that are interacted by user u, and \mathcal{N}_i denotes the set of users that interact with item i. $|\cdot|$ indicates their sizes. We use d to represent the dimension of latent embedding space and \oplus for the fusing function such as concatenation. The matrix $W_l \in \mathbb{R}^{2d \times d}$ denotes the transformation mapping matrix for the *l*-th layer. In the first layer, each user and item is initialized with a graph embedding based on its ID, represented as $e_u^{(1)} \in \mathbb{R}^d$ and $e_i^{(1)} \in \mathbb{R}^d$. To incorporate the textual descriptions associated with users and items, we encode these descriptions into constant text-based embeddings by $f(\cdot)$. In practice, we add a unique token [CLS] before the original text and feed the combined sequence into the simbert-basechinese model. The output of the [CLS] token is used as the semantic embedding for alignment.

To leverage the descriptions of users and items across all layers, we further combine their embeddings from each layer to produce the final embeddings of users and items through mean-pooling:

$$\tilde{e}_{u} = \frac{1}{L} \sum_{l=1}^{L} e_{u}^{(l)}; \quad \tilde{e}_{i} = \frac{1}{L} \sum_{l=1}^{L} e_{i}^{(l)}.$$
 (8)

Objective Function 3.6

To measure the matching scores between users and items for final predictions, we propose to compute the inner product of their representations for recommendation prediction scores by $R_{u,i} = \langle \tilde{e}_u, \tilde{e}_i \rangle$, where $\langle \cdot, \cdot \rangle$ denotes the inner product operation for similarity. It produces a score or probability of item i that user u will engage. For the model training process, we use the pairwise loss to define the recommendation objective function as follows:

$$\max_{\Theta} \sum_{(u,i,j)\in\mathcal{D}} \log \sigma(\hat{R}_{u,i} - \hat{R}_{u,j}) - \lambda ||\Theta||^2, \quad (9)$$

where the train set $\mathcal{D} = \{(u, i, j)\}$ consists of triplets with a user u, an item i with positive feed-354

back from user u, and an item j with negative feedback from user u. Θ denotes all trainable parameters, and λ is the regularization coefficient of L2 norm $|| \cdot ||^2$.

355

357

358

360

362

363

364

365

367

369

370

371

372

373

374

375

376

377

378

379

380

381

382

383

385

386

388

389

390

391

392

393

394

395

Complexity and Applicability 3.7

The model parameter of GaCLLM is approximately $\mathcal{O}((M+N) \cdot d + 2 \cdot L \cdot d^2) = \mathcal{O}((M+N \cdot d))$ as $(M+N) \gg 2 \cdot L \cdot d$. The complexity is similar to the efficient LightGCN (He et al., 2020). As for model training, the time cost (3.76s per update) is slightly higher than LightGCN (2.27s per update) due to the additional text embeddings. The text enhancement in Equation 4, 5 can be done offline in parallel to speed up the generation phase, thus, GaCLLM is scalable in real-world applications.

4 Experiment

Experimental Setup 4.1

Datasets. We investigate two scenarios: job recommendation and social recommendation. For job recommendation, we use two real-world datasets sourced from an online recruiting platform within the **Design** and **Sales** professions with extensive user-job interactions. The user resumes and job descriptions are available as textual document information. For social recommendation, we use a public dataset Pokec Slovakian Social Network (*Pokec*) collected from an online social platform. It contains the friendship relations among users and their self-descriptions. We aim to suggest connections between users based on diverse preferences and attributes. The dataset is divided into subsets *Pokec-A* and *Pokec-B* by different user groups. The statistics of datasets are in Table 1.

Evaluation. We randomly split the dataset equally into training, validation, and test sets. We utilize two well-recognized top-K recommendation metrics, mean average precision (MAP@K) and normalized discounted cumulative gain (NDCG@K), where K is set to 5 empirically. We run five times and take the average performance as experimental results with different random initializations.

Job	# User Resumes	# Job Descriptions	# Interactions
Designs	12,290	9,143	166,270
Sales	15,854	12,772	145,066
Social	# Group A	# Group B	# Connections
Pokec	6,240	6,213	104,152

Table 1: Statistics of datasets.

	Job Recommendation				Social Recommendation			
Models	Models Design		Sales		Pokec-A		Pokec-B	
	MAP@5	NDCG@5	MAP@5	NDCG@5	MAP@5	NDCG@5	MAP@5	NDCG@5
SGPT-BE MF NCF	0.0651 0.2081 0.2100	0.1042 0.3182 0.3258	0.0491 0.0957 0.1468	0.0861 0.1751 0.2678	0.0724 0.2639 0.2969	0.1013 0.3838 0.4270	0.0710 0.2616 0.2930	0.0980 0.3876 0.4273
LightGCN SimGCL UltraGCN SGL	$\begin{array}{c c} 0.2940\\ \hline 0.1471\\ 0.2639\\ \hline 0.2769\end{array}$	0.4697 0.2277 0.4258 0.4418	0.1658 0.0921 0.1469 0.1431	0.3001 0.1658 0.2725 0.2567	$\begin{array}{c} \underline{0.3293}\\ 0.2940\\ 0.3263\\ 0.3047 \end{array}$	$\begin{array}{r} 0.4664 \\ 0.4235 \\ \underline{0.4691} \\ 0.4385 \end{array}$	$\begin{array}{c c} \underline{0.3294} \\ 0.3093 \\ 0.3204 \\ 0.3012 \end{array}$	$\begin{array}{r} \underline{0.4676}\\ 0.4459\\ 0.4623\\ 0.4394\end{array}$
LLM-CS LLM-TES LGIR	0.2669 0.2208 0.2898	0.2190 0.3478 0.4616	0.1530 0.1520 <u>0.1694</u>	0.2803 0.2797 <u>0.3103</u>	0.2569 0.2593 0.3245	0.3478 0.3517 0.4390	0.2527 0.2571 0.3081	0.3468 0.3512 0.4183
GaCLLM Improvement	0.3060* 4.06%	0.4925* 4.85%	0.1750* 3.32%	0.3234* 4.21%	0.3461* 5.10%	0.4798* 2.28%	0.3446* 4.60%	0.4797* 2.60%

Table 2: Performance of GaCLLM and baseline methods. The best results are in **bold** and the runner-up results are <u>underscored</u>. * indicates significant improvements at the level of 0.05 with a paired t-test.

Baselines. We compare our GaCLLM with the following baselines using various approaches. **Content-based and collaborative filtering RS: SGPT-BE** (Muennighoff, 2022) applies GPT models as Bi-Encoders for asymmetric search. MF (Koren et al., 2009) learns low-dimensional representations of users and items by reconstructing their interaction matrix based on the point loss. NCF (He et al., 2017) enhances collaborative filtering with deep neural networks to explore the non-linear interaction between user and item.

397

400

401

402

403

404

405

406

Graph-based RS: For a fair comparison, we en-407 hance all graph-based methods with text informa-408 tion to incorporate an equal amount of utilized in-409 formation. LightGCN (He et al., 2020) simplifies 410 the vanilla GCN's implementation to improve ef-411 ficiency for recommendation. SimGCL (Yu et al., 412 413 2022) adds uniform noises to graph embeddings and conducts contrastive learning for recommenda-414 tion. UltraGCN (Mao et al., 2021) skips infinite 415 layers of message passing of GCN for efficient rec-416 ommendation. SGL (Wu et al., 2021) conducts the 417 self-supervised learning on the user-item graph to 418 improve accuracy and robustness. 419

LLM-based RS: LLM-CS (Chen et al., 2024) 420 directly encodes text attributes into initial node 421 features by LLMs for graph models in a cascad-422 ing structure. LLM-TES (Chen et al., 2024), as 423 another variant, conducts text-level enhancement 424 structure using LLMs and then encodes them as 425 initial node embeddings. LGIR (Du et al., 2024) 426 designs a GAN-based model and infers users' im-427 plicit characteristics from their behaviors for re-428 sume completion. 429

430 **Implementation Details.** For the LLM backbone,

we use ChatGLM2-6B (Du et al., 2022) for its proficiency in handling multilingual tasks including Chinese, as datasets **Design** and **Sales** are in Chinese. For the SFT stage, we use LoRA (Hu et al., 2022) with a learning rate of 10^{-5} , LoRA dimension of 128, batch size of 2, 10^4 training steps, and gradient accumulation of 1. To ensure a fair comparison, we fix the embedding size of all methods to 768, batch size to 1024, and regularization coefficient to 10^{-4} with AdamW (Loshchilov and Hutter, 2019) optimizer. Following (Yang et al., 2022; Du et al., 2024), we use 20 negative instances for every target item during evaluation. 431

432

433

434

435

436

437

438

439

440

441

442

443

444

4.2 Comparison with Baselines

Table 2 shows the overall comparison between 445 GaCLLM and baselines. From the experimental 446 results, we demonstrate that GaCLLM consistently 447 outperforms all baseline methods across all job rec-448 ommendation and social recommendation scenar-449 ios, with average improvements of 4.46%, 3.77%, 450 3.69%, and 3.60%. Besides, interaction-only (i.e., 451 MF and NCF) and text-only (SGPT-BE) methods 452 show inferior performance compared to the other 453 hybrid approaches, indicating the necessity of uti-454 lizing both text and interaction information. In ad-455 dition, the improvements in GCN-based methods 456 prove the value of extracting both graph and text 457 information for better recommendation outcomes. 458 This supports our motivation to combine LLMs 459 with graph structural information to improve the 460 quality of textual descriptions in recommendation 461 systems. SimGCL shows underwhelming results, 462 likely due to the graphical framework's incompati-463 bility with incorporating text-aware information ef-464

	De	sign	Sales		
Models	MAP@5	NDCG@5	MAP@5	NDCG@5	
RAW	0.2951	0.4717	0.1692	0.3082	
w/o-ALIGN GaCLLM	0.2908	0.4654 0.4925	0.1753 0.1750	0.3212 0.3234	
- ouellin	0.0000	01.20	Pokec-B		
	Pol	xec-A	Pol	kec-B	
Models	Pok MAP@5	xec-A NDCG@5	Pol MAP@5	xec-B NDCG@5	

Table 3: Performance of ablation variants.

fectively. Finally, simply adopting the LLM as the encoder (LLM-CS) or zero-shot reasoner (LLM-TES) produces suboptimal performance. LGIR shows stronger performance by inferring from direct neighbors but still overlooks the more complex, high-order relationships within the graph. As a result, by aligning LLM and high-order graph relations, GaCLLM achieves the best performance, validating its effectiveness.

4.3 Ablation Study

465

466

467

468

469

470

471

472

473

474

481

491

To verify the efficacy of the key components of 475 476 GaCLLM, we test the following variants. **RAW** adopts raw descriptions instead of LLM-driven de-477 478 scriptions by the user-item graph. PLAIN removes the convolutional inference strategy, adopting a 479 template to describe all node descriptions related 480 to the target node in a plain way as the inputs of LLMs. w/o-ALIGN excludes the alignment with 482 graph embeddings and simply adopts the enhanced 483 descriptions by L-hop neighbors for node embed-484 dings of the L-th layer. Table 3 shows the per-485 486 formance of variants and original GaCLLM. First, the proposed GaCLLM consistently outperforms 487 RAW across all scenarios, indicating that utiliz-488 ing high-order relations in the interaction graph 489 can improve the textual content and thus lead to 490 more accurate recommendation predictions. Second, GaCLLM significantly outperforms PLAIN. 492 While **PLAIN** struggles to effectively capture the 493 structured graph by describing high-order relations 494 in a single prompt, GaCLLM elicits the reason-495 ing capacity of LLMs more effectively through a 496 step-by-step, graph-based convolutional inference 497 process. This allows GaCLLM to better utilize 498 499 the graph structure for improved recommendations and avoids context length limits. Third, GaCLLM outperforms w/o-ALIGN as the alignment of textual and graphical representations bridges the gap between LLM-driven information and behavioral 503





Figure 4: GaCLLM with varying numbers of layers.

Figure 5: Performance of the proposed method with varying LLM backbones in Designs dataset.

patterns. Thus, we can fully leverage the layered descriptions generated by the LLM for recommendation. As such, the ablation study supports the efficacy of GaCLLM and the underlying motivations presented in this paper.

504

505

506

507

508

509

510

511

512

513

514

515

516

517

518

519

520

521

523

524

525

526

527

528

529

530

531

4.4 In-depth Analysis

In this subsection, we further conduct experiments to analyze the impact of hyper-parameters, the supervised fine-tuning step, and the LLM model selection. We also illustrate the effectiveness of our GaCLLM by both quantitative subgroup analysis and qualitative case study.

Number of Layers. As shown in Figure 4, we observe that the best performance is produced by (4, 3, 2, 2) layers for *Design*, *Sales*, and *Pokec* datasets, respectively. The layered structure need not be deep. For real-world applications, we suggest using the grid search on optimal layer numbers for GaCLLM implementation empirically. Supervised Fine-tuning Study. In Figure 5 (left), we evaluate the variant without supervised fine-tuning in Section 3.3. Using Designs dataset as an example, we notice a limited improvement, which indicates that the overall performance boost by GaCLLM is **not** obtained directly from the SFT, but from the LLM-based convolutional inference strategy and embedding alignment. Though the impact of SFT is not significant, some recommen-

Encoder	MAP@5	Ţ	NDCG@5
simbert-base-chinese	0.3060	1	0.4925
ChatGLM2-transformer	0.2722	1	0.4291

Table 4: Performance of the proposed method withvarying text encoders in *Designs* dataset.

dation scenarios may contain extra domain-specific information beyond LLMs' pre-trained knowledge. Therefore, the SFT step contributes to the adaptability of GaCLLM.

532

533

534

535

536

540

541

542

556

558

561

562

564

LLM Backbone. In Figure 5 (right), we assess GaCLLM using Llama-2-7B as the backbone replacement of the original ChatGLM2-6B with a similar scale. The result shows comparable performance, validating the robustness of our method and the stability of our convolutional inference strategy in description enhancement for recommendation.

543 Text Encoder. To bridge the gap between LLMdriven text information and behavioral-based graph embeddings, we employ simbert-base-chinese to encode user and item text information into latent space. In Table 4, we also explore using other 547 LLM's backbone as text encoder. The results show 548 that ChatGLM2 yields suboptimal results as a text 549 encoder, likely due to its decoder-only structure optimized for text generation rather than under-551 standing. For better performance and parameter efficiency, the encoder-only simbert-base-chinese is a more suitable choice. 554

Subgroup Analysis. We also investigate in the recommendation performance across user groups by the description length ascendingly from G1 to G5, the difference between GaCLLM and *RAW* in Figure 6 shows the significance of refining descriptions for all raw text. Notably, GaCLLM achieves more substantial improvements in groups with less comprehensive descriptions, highlighting the effectiveness of LLM-based convolutional inference strategy by leveraging the graph structure.



Figure 6: Performance across user subgroups for description improvement analysis.



Figure 7: Case study in Design dataset.

565

566

567

568

569

570

571

572

573

574

575

576

577

578

579

580

581

582

583

584

585

586

587

588

589

590

591

592

593

594

595

596

598

Case study. We qualitatively show the efficacy of the convolutional inference strategy by the case study in Figure 7, where we highlight contents relevant to the target job from a user's resume across layers. The raw resume contains some relevant information and some irrelevant words. As layers increase, our method progressively refines the resume, removing irrelevant content and focusing on job-specific details. The text similarity between the user's resume and the job description significantly improves by the third layer, showcasing the LLM's success in reasoning over graph structure. By revising vague information and inferring potential requirements for job matching, we achieve better recommendation outcomes.

5 Conclusion

In this paper, we propose GaCLLM to enhance auxiliary textual information through user-item interactions for recommendation. Our approach bridges the gap between text-based LLMs and graph-based multi-hop relations that contain collaborative information. By employing an iterative convolutional inference strategy, GaCLLM enables efficient propagation of textual information across the graph within constrained token limits to achieve quality improvement. We further align the LLM-driven texts and the behavioral graph embeddings to enhance recommendation performance. Extensive experiments show that GaCLLM consistently outperforms various baseline methods, with ablation studies and in-depth analysis further validating our model design. In future work, we aim to explore using LLMs to handle multi-modality information beyond text for more fine-grained RSs.

6 Limitation

599

613

614

615

616

617

618

619

621

622

623

625

626

627

630

632

635

636

637

644

645

647

600The primary constraints of this paper are as fol-601lows: (1) The training phase requires substantial602computational resources for LLM inference. Since603some users and items may share similar collab-604orative information, it may not be necessary to605make exact inferences for all nodes in the graph.606(2) In real-world scenarios, users often exhibit dy-607namic preferences for items. However, GaCLLM608relies on a static graph, which fails to capture the609dynamic preferences underlying users' sequential610behaviors. To this end, we leave the exploration of611more efficient and dynamic solutions for sequential612recommendation as future work.

References

- Berkeley R Andrus, Yeganeh Nasiri, Shilong Cui, Benjamin Cullen, and Nancy Fulda. 2022. Enhanced story comprehension for large language models through dynamic document-based knowledge graphs. In *Proceedings of the AAAI Conference on Artificial Intelligence (AAAI)*, pages 10436–10444.
- Mengru Chen, Chao Huang, Lianghao Xia, Wei Wei, Yong Xu, and Ronghua Luo. 2023. Heterogeneous graph contrastive learning for recommendation. In *Proceedings of the Sixteenth ACM International Conference on Web Search and Data Mining (WSDM)*, page 544–552.
 - Zhikai Chen, Haitao Mao, Hang Li, Wei Jin, Hongzhi Wen, Xiaochi Wei, Shuaiqiang Wang, Dawei Yin, Wenqi Fan, Hui Liu, and Jiliang Tang. 2024. Exploring the potential of large language models in learning on graphs. *SIGKDD Explor. Newsl.*, page 42–61.
 - Yingpeng Du, Di Luo, Rui Yan, Xiaopei Wang, Hongzhi Liu, Hengshu Zhu, Yang Song, and Jie Zhang. 2024.
 Enhancing job recommendation through llm-based generative adversarial networks. In *Proceedings* of the AAAI Conference on Artificial Intelligence (AAAI), pages 8363–8371.
- Zhengxiao Du, Yujie Qian, Xiao Liu, Ming Ding, Jiezhong Qiu, Zhilin Yang, and Jie Tang. 2022. Glm: General language model pretraining with autoregressive blank infilling. In Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (ACL), pages 320–335.
- Xiangnan He, Kuan Deng, Xiang Wang, Yan Li, Yongdong Zhang, and Meng Wang. 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 (SIGIR)*, pages 639–648.

Xiangnan He, Lizi Liao, Hanwang Zhang, Liqiang Nie, Xia Hu, and Tat-Seng Chua. 2017. Neural collaborative filtering. In *Proceedings of the 26th international conference on world wide web (WWW)*, pages 173–182. 650

651

652

653

654

655

656

657

658

659

660

661

662

663

664

665

666

667

668

669

670

671

672

673

674

675

676

677

678

679

680

681

682

683

684

685

686

687

688

689

690

691

692

693

694

695

696

697

698

699

700

701

702

703

704

- 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.
- Edward J Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, and Weizhu Chen. 2022. LoRA: Low-rank adaptation of large language models. In *International Conference on Learning Representations (ICLR)*.
- Zhen Huang, Zhongchuan Sun, Jiaming Liu, and Yangdong Ye. 2024. Group-aware graph neural networks for sequential recommendation. *Information Sciences*, 670:120623.
- Wang-Cheng Kang, Jianmo Ni, Nikhil Mehta, Maheswaran Sathiamoorthy, Lichan Hong, Ed Chi, and Derek Zhiyuan Cheng. 2023. Do llms understand user preferences? evaluating llms on user rating prediction. *arXiv preprint arXiv:2305.06474*.
- Thomas N Kipf and Max Welling. 2016. Semisupervised classification with graph convolutional networks. In *International Conference on Learning Representations (ICLR)*.
- Yehuda Koren, Robert Bell, and Chris Volinsky. 2009. Matrix factorization techniques for recommender systems. *Computer*, 42(8):30–37.
- Hui-Chi Kuo and Yun-Nung Chen. 2023. Zero-shot prompting for implicit intent prediction and recommendation with commonsense reasoning. In *Findings of the Association for Computational Linguistics* (ACL), pages 249–258.
- Yuhan Li, Zhixun Li, Peisong Wang, Jia Li, Xiangguo Sun, Hong Cheng, and Jeffrey Xu Yu. 2024. A survey of graph meets large language model: Progress and future directions. In *Proceedings of the Thirty-Third International Joint Conference on Artificial Intelligence (IJCAI)*, pages 8123–8131.
- Zihan Lin, Changxin Tian, Yupeng Hou, and Wayne Xin Zhao. 2022. Improving graph collaborative filtering with neighborhood-enriched contrastive learning. In *Proceedings of the ACM Web Conference (TheWeb-Conf)*, pages 2320–2329.
- Junling Liu, Chao Liu, Renjie Lv, Kang Zhou, and Yan Zhang. 2023. Is chatgpt a good recommender? a preliminary study. *arXiv preprint arXiv:2304.10149*.
- Nelson F. Liu, Kevin Lin, John Hewitt, Ashwin Paranjape, Michele Bevilacqua, Fabio Petroni, and Percy Liang. 2024a. Lost in the middle: How language models use long contexts. *Transactions of the Association for Computational Linguistics (TACL)*, 12:157– 173.

799

800

801

802

803

804

805

806

807

808

809

810

811

812

813

814

815

816

817

761

762

709 710 712

713

706

707

- 714 715 716 717 719 721 723
- 725
- 727
- 729 730 731
- 732 733
- 734 735 736 737
- 738 739 740 741
- 743 744

- 745 746
- 747 748
- 750 751
- 752 753 754
- 755
- 756

- 757

- Qijiong Liu, Nuo Chen, Tetsuya Sakai, and Xiao-Ming Wu. 2024b. Once: Boosting content-based recommendation with both open- and closed-source large language models. In Proceedings of the 17th ACM International Conference on Web Search and Data Mining (WSDM), page 452–461.
- Ilya Loshchilov and Frank Hutter. 2019. Decoupled weight decay regularization. In International Conference on Learning Representations (ICLR).
- Wensheng Lu, Jianxun Lian, Wei Zhang, Guanghua Li, Mingyang Zhou, Hao Liao, and Xing Xie. 2024. Aligning large language models for controllable recommendations. In Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (ACL), pages 8159–8172.
- Hanjia Lyu, Song Jiang, Hanqing Zeng, Yinglong Xia, Qifan Wang, Si Zhang, Ren Chen, Chris Leung, Jiajie Tang, and Jiebo Luo. 2024. LLM-rec: Personalized recommendation via prompting large language models. In Findings of the Association for Computational Linguistics: NAACL, pages 583–612.
- Kelong Mao, Jieming Zhu, Xi Xiao, Biao Lu, Zhaowei Wang, and Xiuqiang He. 2021. Ultragen: ultra simplification of graph convolutional networks for recommendation. In Proceedings of the 30th ACM International Conference on Information & Knowledge Management (CIKM), pages 1253-1262.
- Niklas Muennighoff. 2022. Sgpt: Gpt sentence embeddings for semantic search. arXiv preprint arXiv:2202.08904.
- Lütfi Kerem Senel, Besnik Fetahu, Davis Yoshida, Zhiyu Chen, Giuseppe Castellucci, Nikhita Vedula, Jason Ingyu Choi, and Shervin Malmasi. 2024. Generative explore-exploit: Training-free optimization of generative recommender systems using LLM optimizers. In Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (ACL), pages 5396-5420.
- Yu Sun, Shuohuan Wang, Shikun Feng, Siyu Ding, Chao Pang, Junyuan Shang, Jiaxiang Liu, Xuyi Chen, Yanbin Zhao, Yuxiang Lu, et al. 2021. Ernie 3.0: Large-scale knowledge enhanced pre-training for language understanding and generation. arXiv preprint arXiv:2107.02137.
- Zhu Sun, Qing Guo, Jie Yang, Hui Fang, Guibing Guo, Jie Zhang, and Robin Burke. 2019. Research commentary on recommendations with side information: A survey and research directions. Electronic Commerce Research and Applications (ECRA), 37:100879.
- Ghazaleh Haratinezhad Torbati, Anna Tigunova, and Gerhard Weikum. 2023. Unveiling challenging cases in text-based recommender systems. In Proceedings of the 3rd Workshop Perspectives on the Evaluation of Recommender Systems (Perspectives@RecSys).

- Lei Wang and Ee-Peng Lim. 2024. The whole is better than the sum: Using aggregated demonstrations in in-context learning for sequential recommendation. In Findings of the Association for Computational Linguistics (NAACL), pages 876-895.
- Yan Wang, Zhixuan Chu, Xin Ouyang, Simeng Wang, Hongyan Hao, Yue Shen, Jinjie Gu, Siqiao Xue, James Y Zhang, Qing Cui, et al. 2023. Enhancing recommender systems with large language model reasoning graphs. arXiv preprint arXiv:2308.10835.
- Yan Wang, Zhixuan Chu, Xin Ouyang, Simeng Wang, Hongvan Hao, Yue Shen, Jinije Gu, et al. 2024a. Llmrg: Improving recommendations through large language model reasoning graphs. In Proceedings of the AAAI Conference on Artificial Intelligence (AAAI), volume 38, pages 19189–19196.
- Yancheng Wang, Ziyan Jiang, Zheng Chen, Fan Yang, Yingxue Zhou, Eunah Cho, Xing Fan, Yanbin Lu, Xiaojiang Huang, and Yingzhen Yang. 2024b. Rec-Mind: Large language model powered agent for recommendation. In Findings of the Association for Computational Linguistics (NAACL), pages 4351-4364.
- Yinwei Wei, Wengi Liu, Fan Liu, Xiang Wang, Ligiang Nie, and Tat-Seng Chua. 2023. Lightgt: A light graph transformer for multimedia recommendation. In Proceedings of the 46th International ACM SI-GIR Conference on Research and Development in Information Retrieval (SIGIR), page 1508–1517.
- Jiancan Wu, Xiang Wang, Fuli Feng, Xiangnan He, Liang Chen, Jianxun Lian, and Xing Xie. 2021. Selfsupervised graph learning for recommendation. In Proceedings of the 44th international ACM SIGIR conference on research and development in information retrieval (SIGIR), pages 726–735.
- Likang Wu, Zhaopeng Qiu, Zhi Zheng, Hengshu Zhu, and Enhong Chen. 2024a. Exploring large language model for graph data understanding in online job recommendations. In Proceedings of the AAAI Conference on Artificial Intelligence (AAAQI), pages 9178-9186.
- Likang Wu, Zhi Zheng, Zhaopeng Qiu, Hao Wang, Hongchao Gu, Tingjia Shen, Chuan Qin, Chen Zhu, Hengshu Zhu, Qi Liu, et al. 2024b. A survey on large language models for recommendation. World Wide Web, 27(5):60.
- "Xuansheng Wu, Huachi Zhou, Yucheng Shi, Wenlin Yao, Xiao Huang, and year = "2024" Ninghao Liu". "could small language models serve as recommenders? towards data-centric cold-start recommendation". In "Proceedings of the ACM Web Conference (TheWebConf), pages "3566-3575".
- Han Xie, Da Zheng, Jun Ma, Houyu Zhang, Vassilis N. Ioannidis, Xiang Song, Qing Ping, et al. 2023. Graphaware language model pre-training on a large graph corpus can help multiple graph applications. In Proceedings of the 29th ACM SIGKDD Conference

818 on Knowledge Discovery and Data Mining (KDD),
819 pages 5270–5281.

821

822

831

833

836

841

843

845

847

849

852

853

854

855

856

857 858

863

870

871

- Surong Yan, Chongyang Li, Haosen Wang, Bin Lin, and Yixian Yuan. 2024. Feature interactive graph neural network for kg-based recommendation. *Expert Systems with Applications*, 237:121411.
- Chen Yang, Yupeng Hou, Yang Song, Tao Zhang, Ji-Rong Wen, and Wayne Xin Zhao. 2022. Modeling two-way selection preference for person-job fit. In *Proceedings of the 16th ACM Conference on Recommender Systems (RecSys)*, pages 102–112.
- Michihiro Yasunaga, Antoine Bosselut, Hongyu Ren, Xikun Zhang, Christopher D Manning, Percy S Liang, and Jure Leskovec. 2022. Deep bidirectional language-knowledge graph pretraining. *Advances in Neural Information Processing Systems (NeurIPS)*, pages 37309–37323.
- Ruosong Ye, Caiqi Zhang, Runhui Wang, Shuyuan Xu, and Yongfeng Zhang. 2023. Language is all a graph needs. In *Findings of the Association for Computational Linguistics (EACL)*, pages 1955–1973.
- Junliang Yu, Hongzhi Yin, Xin Xia, Tong Chen, Lizhen Cui, and Quoc Viet Hung Nguyen. 2022. Are graph augmentations necessary? simple graph contrastive learning for recommendation. In *Proceedings of the 45th international ACM SIGIR conference on research and development in information retrieval* (*SIGIR*), pages 1294–1303.
- Junjie Zhang, Ruobing Xie, Yupeng Hou, Wayne Xin Zhao, Leyu Lin, and Ji-Rong Wen. 2023. Recommendation as instruction following: A large language model empowered recommendation approach. *arXiv preprint arXiv:2305.07001*.
- Mengmei Zhang, Mingwei Sun, Peng Wang, Shen Fan, Yanhu Mo, Xiaoxiao Xu, Hong Liu, Cheng Yang, and Chuan Shi. 2024a. Graphtranslator: Aligning graph model to large language model for open-ended tasks. In *Proceedings of the ACM Web Conference* (*TheWebConf*), page 1003–1014.
- Xikun Zhang, Antoine Bosselut, Michihiro Yasunaga, Hongyu Ren, Percy Liang, Christopher D Manning, and Jure Leskovec. 2021. Greaselm: Graph reasoning enhanced language models. In *International conference on learning representations (ICLR)*.
- Yichi Zhang, Zhuo Chen, Lingbing Guo, Yajing Xu, Wen Zhang, and Huajun Chen. 2024b. Making large language models perform better in knowledge graph completion. In *Proceedings of the 32nd ACM International Conference on Multimedia (MM)*, page 233–242.
- Zhi Zheng, Zhaopeng Qiu, Xiao Hu, Likang Wu, Hengshu Zhu, and Hui Xiong. 2023. Generative job recommendations with large language model. *arXiv preprint arXiv:2307.02157*.

Denny Zhou, Nathanael Schärli, Le Hou, Jason Wei, Nathan Scales, Xuezhi Wang, Dale Schuurmans, Claire Cui, Olivier Bousquet, Quoc V Le, et al. 2022. Least-to-most prompting enables complex reasoning in large language models. In *The Eleventh International Conference on Learning Representations* (ICLR).

872

873

874

875

876

877