

---

# Language Models Encode Collaborative Signals in Recommendation

---

Anonymous Author(s)

Affiliation

Address

email

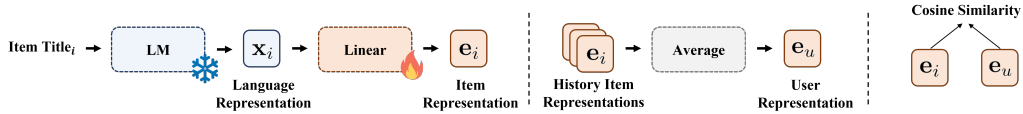
## Abstract

1       Recent studies empirically indicate that language models (LMs) encode rich world  
2       knowledge beyond mere semantics, attracting significant attention across various  
3       fields. However, in the recommendation domain, it remains uncertain whether  
4       LMs implicitly encode user preference information. Contrary to the prevailing  
5       understanding that LMs and traditional recommender models learn two distinct rep-  
6       resentation spaces due to a huge gap in language and behavior modeling objectives,  
7       this work rethinks such understanding and explores extracting a recommendation  
8       space directly from the language representation space. Surprisingly, our findings  
9       demonstrate that item representations, when linearly mapped from advanced LM  
10      representations, yield superior recommendation performance. This outcome sug-  
11      gests a homomorphic relationship between the language representation space and  
12      an effective recommendation space, implying that collaborative signals may indeed  
13      be encoded within advanced LMs. Motivated by these findings, we propose a  
14      simple yet effective collaborative filtering (CF) model named **AlphaRec**, which  
15      utilizes language representations of item textual metadata (*e.g.*, titles) instead of tra-  
16      ditional ID-based embeddings. Specifically, AlphaRec is comprised of three main  
17      components: a multilayer perceptron (MLP), graph convolution, and contrastive  
18      learning (CL) loss function, making it extremely easy to implement and train. Our  
19      empirical results show that AlphaRec outperforms leading ID-based CF models  
20      on multiple datasets, marking the first instance of such a recommender with text  
21      embeddings achieving this level of performance. Moreover, AlphaRec introduces  
22      a new text-based CF paradigm with several desirable advantages: being easy to  
23      implement, lightweight, rapid convergence, superior zero-shot recommendation  
24      abilities in new domains, and being aware of user intention.

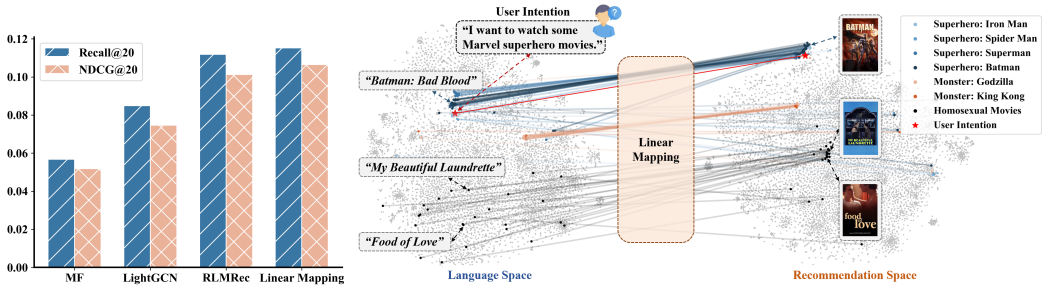
## 25 1 Introduction

26    Language models (LMs) have achieved great success across various domains [3–7], prompting a  
27    critical question about the knowledge encoded within their representation spaces. Recent studies  
28    empirically find that LMs extend beyond semantic understanding to encode comprehensive world  
29    knowledge about various domains, including game states [8], lexical attributes [9], and even concepts  
30    of space and time [10] through language modeling. However, in the domain of recommendation  
31    where the integration of LMs is attracting widespread interest [11–15], it remains unclear whether  
32    LMs inherently encode relevant information on user preferences and behaviors. One possible reason  
33    is the significant difference between the objectives of language modeling for LMs and user behavior  
34    modeling for recommenders [16–19].

35    Currently, one prevailing understanding holds that general LMs and traditional recommenders  
36    encode two distinct representation spaces: the language space and the recommendation space  
37    (*i.e.*, user and item representation space), each offering potential enhancements to the other for



(a) Linearly mapping language representations into the recommendation space



(b) Performance comparison (c) The t-SNE representations of movies and user intention in two spaces.

Figure 1: Linearly mapping item titles in language representation space into recommendation space yields superior recommendation performance on Movies & TV [1] dataset. (1a) The framework of linear mapping. (1b) The recommendation performance comparison between leading CF recommenders and linear mapping. (1c) The t-SNE [2] visualizations of movie representations, with colored lines linking identical movies or user intention across language space (left) and linearly projected recommendation space (right).

38 recommendation tasks [17, 20]. On the one hand, when using LMs as recommenders, aligning the  
 39 language space with the recommendation space could significantly improve the performance of  
 40 LM-based recommendation [14, 21–23]. Various alignment strategies are proposed, including fine-  
 41 tuning LMs with recommendation data [15, 16, 24–26], incorporating embeddings from traditional  
 42 recommenders as a new modality of LMs [17, 20, 27], and extending the vocabulary of LMs with item  
 43 tokens [18, 19, 28–31]. On the other hand, when using LMs as the enhancer, traditional recommenders  
 44 greatly benefit from leveraging text representations [32–45], semantic and reasoning information  
 45 [46–49], and generated user behaviors [50, 51]. Despite these efforts, explicit explorations of the  
 46 relationship between language and recommendation spaces remain largely unexplored.

47 In this work, we rethink the prevailing understanding and explore whether LMs inherently encode  
 48 user preferences through language modeling. Specifically, we test the possibility of directly deriving a  
 49 recommendation space from the language representation space, assessing whether the representations  
 50 of item textual metadata (e.g., titles) obtained from LMs can independently achieve satisfactory  
 51 recommendation performance. Positive results would imply that user behavioral patterns, such as  
 52 **collaborative signals** (i.e., user preference similarities between items) [52, 53], may be implicitly  
 53 encoded by LMs. To test this hypothesis, we employ linear mapping to project the language  
 54 representations of item titles into a recommendation space (see Figure 1a). Our observations include:

- 55 • Surprisingly, this simple linear mapping yields high-quality item representations, which achieve  
 56 exceptional recommendation performance (see Figure 1b and experimental results in Section 2).
- 57 • The clustering of items is generally preserved from the language space to the recommendation  
 58 space (see Figure 1c). For example, movies with the theme of superheroes and monsters are  
 59 gathering in both language and recommendation spaces.
- 60 • Interestingly, the linear mapping effectively reveals preference similarities that may be implicit  
 61 or even obscure in the language space. For instance, while certain movies, such as those of  
 62 homosexual movies (illustrated in Figure 1c), show dispersed representations in the language space,  
 63 their projections through linear mapping tend to cluster together, reflecting their genres affiliation.

64 These findings indicate a homomorphic relationship between the language representation space of  
 65 LMs and an effective item representation space for recommendation. Motivated by this insight, we  
 66 propose a new text-based recommendation paradigm for general collaborative filtering (CF), which  
 67 utilizes the pre-trained language representations of item titles as the item input and the average  
 68 historical interactions’ representations as the user input. Different from traditional ID-based CF  
 69 models [54, 55, 52] that heavily rely on trainable user and item IDs, this paradigm solely uses

70 pre-trained LM embeddings and completely abandons ID-based embeddings. In this paper, to fully  
 71 explore the potential of advanced language representations, we adopt a simple model architecture  
 72 consisting of a two-layer MLP with graph convolution, and the popular contrastive loss, InfoNCE  
 73 [56–58], as the objective function. This model is named **AlphaRec** for its originality and a series of  
 74 good properties.

75 Benefiting from paradigm shifts from ID-based embeddings to language representations, AlphaRec  
 76 presents three desirable advantages. First, AlphaRec is notable for its simplicity, lightweight, rapid  
 77 convergence, and exceptional recommendation performance (see Section 4.1). We empirically  
 78 demonstrate that, for the first time, such a simple model with embeddings from pre-trained LMs can  
 79 outperform leading CF models on multiple datasets. This finding strongly supports the possibility  
 80 for developing language-representation-based recommender systems. Second, AlphaRec exhibits  
 81 a strong zero-shot recommendation capability across untrained domains (see Section 4.2). By  
 82 co-training on three Amazon datasets (Books, Movies & TV, and Video Games) [1], AlphaRec  
 83 can achieve performance comparable to the fully-trained LightGCN on entirely different platforms  
 84 (MovieLens-1M [59] and BookCrossing [60]), and even exceed LightGCN in a completely new  
 85 domain (Amazon Industrial), without additional training on these target datasets. This capability  
 86 underscores AlphaRec’s potential to develop more general recommenders. Third, AlphaRec is user-  
 87 friendly, offering a new research paradigm that enhances recommendation by leveraging language-  
 88 based user feedback (see Section 4.3). Endowed with its inherent semantic comprehension of  
 89 language representations, AlphaRec can refine recommendations based on user intentions expressed  
 90 in natural language, enabling traditional CF recommenders to evolve into intention-aware systems  
 91 through a straightforward paradigm shift.

## 92 2 Uncovering Collaborative Signals in LMs via Linear Mapping

93 In this section, we aim to explore whether LMs implicitly encode collaborative signals in their  
 94 representation spaces. We first formulate the personalized item recommendation task, then detail the  
 95 linear mapping and its empirical findings. Empirical evidence indicates a homomorphic relationship  
 96 between the representation spaces of advanced LMs and effective recommendation spaces.

97 **Task formulation.** Personalized item recommendation with implicit feedback aims to select items  
 98  $i \in \mathcal{I}$  that best match user  $u$ ’s preferences based on binary interaction data  $\mathbf{Y} = [y_{ui}]$ , where  $y_{ui} = 1$   
 99 ( $y_{ui} = 0$ ) indicates user  $u \in \mathcal{U}$  has (has not) interacted with item  $i$  [58]. The primary objective of  
 100 recommendation is to model the user-item interaction matrix  $\mathbf{Y}$  using a scoring function  $\hat{y} : \mathcal{U} \times \mathcal{I} \rightarrow$   
 101  $\mathbb{R}$ , where  $\hat{y}_{ui}$  measures  $u$ ’s preference for  $i$ . The scoring function  $\hat{y}_{ui} = s \circ \phi_\theta(\mathbf{x}_u, \mathbf{x}_i)$  comprises  
 102 three key components: pre-existing features  $\mathbf{x}_u$  and  $\mathbf{x}_i$  for user  $u$  and item  $i$ , a representation learning  
 103 module  $\phi_\theta(\cdot, \cdot)$  parametrized by  $\theta$ , and a similarity function  $s(\cdot, \cdot)$ . The representation learning  
 104 module  $\phi_\theta$  transfers  $u$  and  $i$  into representations  $\mathbf{e}_u$  and  $\mathbf{e}_i$  for similarity matching  $s(\mathbf{e}_u, \mathbf{e}_i)$ , and the  
 105 Top- $K$  highest scoring items are recommended to  $u$ .

106 Different recommenders employ various pre-existing features  $\mathbf{x}_u, \mathbf{x}_i$  and representation learning  
 107 architecture  $\phi_\theta(\cdot, \cdot)$ . Traditional ID-based recommenders use one-hot vectors as pre-existing features  
 108  $\mathbf{x}_u, \mathbf{x}_i$ . The choice of ID-based representation learning architecture  $\phi_\theta$  can vary widely, including  
 109 ID-based embedding matrix [54], multilayer perception [61], graph neural network [52, 62], and  
 110 variational autoencoder [63]. The commonly used similarity function is cosine similarity [64, 57]  
 111  $s(\mathbf{e}_u, \mathbf{e}_i) = \frac{\mathbf{e}_u^\top \mathbf{e}_i}{\|\mathbf{e}_u\| \cdot \|\mathbf{e}_i\|}$ , which we adopt in this paper.

112 **Linear mapping.** Building on the extensive knowledge encoded by LMs, we explore utilizing LMs  
 113 as feature extractors, leveraging the language representations of item titles as initial item feature  $\mathbf{x}_i$ .  
 114 For initial user feature  $\mathbf{x}_u$ , we use the average of the title representations of historically interacted  
 115 items, defined as  $\mathbf{x}_u = \frac{1}{|\mathcal{N}_u|} \sum_{i \in \mathcal{N}_u} \mathbf{x}_i$ , where  $\mathcal{N}_u$  is the set of items user  $u$  has interacted with.  
 116 Detailed procedures for obtaining these language-based features are provided in Appendix B.2.  
 117 We select a trainable linear mapping matrix  $\mathbf{W}$  as the representation learning module  $\phi_\theta$ , setting  
 118  $\mathbf{e}_u = \mathbf{W}\mathbf{x}_u$  and  $\mathbf{e}_i = \mathbf{W}\mathbf{x}_i$ . To learn the linear mapping  $\mathbf{W}$ , we adopt the InfoNCE loss [56] as the  
 119 objective function, which has demonstrated state-of-the-art performance in both ID-based [65, 66]  
 120 and LM-enhanced collaborative filtering (CF) recommendations [47] (refer to Equation (4) for the  
 121 formula). The overall framework of the linear mapping process is illustrated in Figure 1a. We directly  
 122 use linearly mapped representations  $\mathbf{e}_u$  and  $\mathbf{e}_i$  to calculate the user-item similarity  $s(\mathbf{e}_u, \mathbf{e}_i)$  for  
 123 recommendation. High performance on the test set would suggest that collaborative signals (*i.e.*, user

Table 1: The recommendation performance of linear mapping comparing with classical CF baselines.

	Books			Movies & TV			Video Games		
	Recall	NDCG	HR	Recall	NDCG	HR	Recall	NDCG	HR
MF (Rendle et al., 2012)	0.0437	0.0391	0.2476	0.0568	0.0519	0.3377	0.0323	0.0195	0.0864
MultVAE (Liang et al., 2018)	0.0722	0.0597	<u>0.3418</u>	0.0853	0.0776	0.4434	0.0908	0.0531	0.2211
LightGCN (He et al., 2021)	0.0723	<u>0.0608</u>	<b>0.3489</b>	0.0849	0.0747	0.4397	0.1007	0.0590	0.2281
Linear Mapping									
<b>BERT</b>	0.0226	0.0194	0.1240	0.0415	0.0399	0.2362	0.0524	0.0309	0.1245
<b>RoBERTa</b>	0.0247	0.0209	0.1262	0.0406	0.0387	0.2277	0.0578	0.0338	0.1339
<b>Llama2-7B</b>	0.0662	0.0559	0.3176	0.1027	0.0955	0.4952	0.1249	0.0729	0.2746
<b>Mistral-7B</b>	0.0650	0.0544	0.3124	0.1039	0.0963	0.4994	0.1270	0.0687	0.2428
<b>text-embedding-ada-v2</b>	0.0515	0.0436	0.2570	0.0926	0.0874	0.4563	0.1176	0.0683	0.2579
<b>text-embeddings-3-large</b>	<u>0.0735</u>	<u>0.0608</u>	0.3355	<u>0.1109</u>	<u>0.1023</u>	<u>0.5200</u>	<u>0.1367</u>	<b>0.0793</b>	<b>0.2928</b>
<b>SFR-Embedding-Mistral</b>	<b>0.0738</b>	<b>0.0610</b>	0.3371	<b>0.1152</b>	<b>0.1065</b>	<b>0.5327</b>	<b>0.1370</b>	<u>0.0787</u>	<u>0.2927</u>

124 preference similarities between items) have been implicitly encoded in the language representation  
 125 space [67, 10].

126 **Empirical findings.** We compare the recommendation performance of the linear mapping method  
 127 with three classical CF baselines, matrix factorization (MF) [54, 68], MultVAE [63], and LightGCN  
 128 [55] (see more details about baselines in Appendix C.2.1). We report three widely used metrics Hit  
 129 Ratio (HR@K), Recall@K, Normalized Discounted Cumulative Gain (NDCG@K) to evaluate  
 130 the effectiveness of linear mapping, with K set by default at 20. We evaluate a wide range of LMs,  
 131 including BERT-style models [4, 5], decoder-only language models [6, 69], and LM-based text  
 132 embedding models [70, 71] (see Appendix B.1 for details about used LMs).

133 Table 1 reports the recommendation performance yielded by the linear mapping on three Amazon  
 134 datasets [1], comparing with classic CF baselines. We observe that the performance of most advanced  
 135 text embedding models (e.g., text-embeddings-3-large [70] and SFR-Embedding-Mistral [71]) exceed  
 136 LightGCN on all datasets. We further empirically prove that these improvements do not merely  
 137 come from the better feature encoding ability (refer to Appendix B.3). These findings indicate  
 138 the homomorphic relationship between the language representation space of advanced LMs and an  
 139 effective item representation space for recommendation. Moreover, with the advances in LMs, the  
 140 performance of item representation linearly mapped from LMs exhibits a rising trend, gradually  
 141 surpassing traditional ID-based CF models. Representations from early BERT-style models (e.g.,  
 142 BERT [4] and RoBERTa [5]) only show weaker or equal capabilities compared with MF, while the  
 143 performance of decoder-only LMs (e.g., Llama-7B [6]) start to match MultVAE and LightGCN.

### 144 3 AlphaRec

145 This finding of space homomorphic relationship sheds light on building advanced CF models purely  
 146 based on LM representations without introducing ID-based embeddings. To be specific, we try to  
 147 incorporate only three simple components (i.e., nonlinear projection [61], graph convolution [55]  
 148 and contrastive learning (CL) objectives [56]), to develop a simple yet effective CF model called  
 149 AlphaRec. It is important to highlight that our approach is centered on exploring the potential of  
 150 LM representations for CF by integrating essential components from leading CF models, rather than  
 151 deliberately inventing new CF mechanisms. We present the model structure of AlphaRec in Section  
 152 3.1, and compare AlphaRec with two popular recommendation paradigms in Section 3.2.

#### 153 3.1 Method

154 We present how AlphaRec is designed and trained. Generally, the representation learning architecture  
 155  $\phi_\theta(\cdot, \cdot)$  of AlphaRec is simple, which only contains a two-layer MLP and the basic graph convolution  
 156 operation, with language representations as the input features  $\mathbf{x}_u, \mathbf{x}_i$ . The cosine similarity is used as  
 157 the similarity function  $s(\cdot, \cdot)$ , and the contrastive loss InfoNCE [56, 57] is adopted for optimization.  
 158 For simplicity, we consistently adopt text-embeddings-3-large [70] as the language representation  
 159 model, for its excellent language understanding and representation capabilities.

160 **Nonlinear projection.** In AlphaRec, we substitute the linear mapping matrix delineated in Section 2  
 161 with a nonlinear MLP. This conversion from linear to nonlinear is non-trivial, for the paradigm shift  
 162 from ID-based embeddings to LM representations, since nonlinear transformation helps in excavating  
 163 more comprehensive collaborative signals from the LM representation space with rich semantics (see

164 discussions about this in Appendix C.2.3) [61]. Specifically, we project the language representation  
 165  $\mathbf{x}_i$  of the item title to an item space for recommendation with the two-layer MLP, and obtain user  
 166 representations as the average of historical items:

$$\mathbf{e}_i^{(0)} = \mathbf{W}_2 \text{LeakyReLU}(\mathbf{W}_1 \mathbf{x}_i + \mathbf{b}_1) + \mathbf{b}_2, \quad \mathbf{e}_u^{(0)} = \frac{1}{|\mathcal{N}_u|} \sum_{i \in \mathcal{N}_u} \mathbf{e}_i^{(0)}. \quad (1)$$

167 **Graph convolution.** Graph neural networks (GNNs) have shown superior effectiveness for recom-  
 168 mendation [52, 55], owing to the natural user-item graph structure in recommender systems [72].  
 169 In AlphaRec, we employ a minimal graph convolution operation [55] to capture more complicated  
 170 collaborative signals from high-order connectivity [55, 73, 74, 72] as follows:

$$\mathbf{e}_u^{(k+1)} = \sum_{i \in \mathcal{N}_u} \frac{1}{\sqrt{|\mathcal{N}_u|} \sqrt{|\mathcal{N}_i|}} \mathbf{e}_i^{(k)}, \quad \mathbf{e}_i^{(k+1)} = \sum_{u \in \mathcal{N}_i} \frac{1}{\sqrt{|\mathcal{N}_i|} \sqrt{|\mathcal{N}_u|}} \mathbf{e}_u^{(k)}. \quad (2)$$

171 The information of connected neighbors is aggregated with a symmetric normalization term  
 172  $\frac{1}{\sqrt{|\mathcal{N}_u|} \sqrt{|\mathcal{N}_i|}}$ . Here  $\mathcal{N}_u$  ( $\mathcal{N}_i$ ) denotes the historical item (user) set that user  $u$  (item  $i$ ) has inter-  
 173 acted with. The features  $\mathbf{e}_u^{(0)}$  and  $\mathbf{e}_i^{(0)}$  projected from the MLP are used as the input of the first layer.  
 174 After propagating for  $K$  layers, the final representation of a user (item) is obtained as the average of  
 175 features from each layer:

$$\mathbf{e}_u = \frac{1}{K+1} \sum_{k=0}^K \mathbf{e}_u^{(k)}, \quad \mathbf{e}_i = \frac{1}{K+1} \sum_{k=0}^K \mathbf{e}_i^{(k)}. \quad (3)$$

176 **Contrastive learning objective.** The introduction of contrasting learning is another key element for  
 177 the success of leading CF models. Recent research suggests that the contrast learning objective, rather  
 178 than data augmentation, plays a more significant role in improving recommendation performance  
 179 [66, 75, 65]. Therefore, we simply use the contrast learning object InfoNCE [56] as the loss function  
 180 without any additional data augmentation on the graph [76, 57]. With cosine similarity as the  
 181 similarity function  $s(\mathbf{e}_u, \mathbf{e}_i) = \frac{\mathbf{e}_u^\top \mathbf{e}_i}{\|\mathbf{e}_u\| \cdot \|\mathbf{e}_i\|}$ , the InfoNCE loss [56, 76, 77] is written as:

$$\mathcal{L}_{\text{InfoNCE}} = - \sum_{(u,i) \in \mathcal{O}^+} \log \frac{\exp(s(u,i)/\tau)}{\exp(s(u,i)/\tau) + \sum_{j \in \mathcal{S}_u} \exp(s(u,j)/\tau)}. \quad (4)$$

182 Here,  $\tau$  is a hyperparameter called temperature [78],  $\mathcal{O}^+ = \{(u,i) | y_{ui} = 1\}$  denoting the observed  
 183 interactions between users  $\mathcal{U}$  and items  $\mathcal{I}$ . And  $\mathcal{S}_u$  is a randomly sampled subset of negative items  
 184 that user  $u$  does not adopt.

### 185 3.2 Discussion of Recommendation Paradigms

186 We compare the language-representation-based AlphaRec with two popular recommendation  
 187 paradigms in Table 2 (see more discussion about related works in Appendix A).

188 **ID-based recommendation (ID-Rec)** [52, 54]. In the traditional ID-based recommendation paradigm,  
 189 users and items are represented by ID-based learnable embeddings derived from a large number of  
 190 user interactions. While ID-Rec exhibits excellent recommendation capabilities with low training and  
 191 inference costs [62, 76], it also has two significant drawbacks. Firstly, these ID-based embeddings  
 192 learned in specific domains are difficult to transfer to new domains without overlapping users  
 193 and items [37], thereby hindering zero-shot recommendation capabilities. Additionally, there is a  
 194 substantial gap between ID-Rec and natural languages [34], which makes ID-based recommenders  
 195 hard to incorporate language-based user intentions and further refine recommendations accordingly.

196 **LM-based recommendation (LM-Rec)** [15, 16, 24]. Benefiting from the extensive world knowledge  
 197 and powerful reasoning capabilities of LMs [7, 79], the LM-based recommendation paradigm has  
 198 gained widespread attention [11, 13]. LM-Rec tends to convert user interaction history into text  
 199 prompts as input for LMs, utilizing pre-trained or fine-tuned LMs in a text generation pattern to  
 200 recommend items. LM-Rec demonstrates zero-shot and few-shot abilities and can easily understand  
 201 language-based user intentions. However, LM-Rec faces significant challenges. Firstly, the LM-based  
 202 model architecture leads to huge training and inference costs, with real-world deployment difficulties.

Table 2: Comparison of recommendation paradigms

Recommendation Paradigms	Training Cost	Zero-shot Ability	Intention-aware Ability
ID-based	Low	✗	✗
LLM-based	High	✓	✓
Language-representation-based	Low	✓	✓

203 Additionally, limited by the text generation paradigm, LM-based models tend to perform candidate  
 204 selection [17] or generate a single next item [24]. It remains difficult for LM-Rec to comprehensively  
 205 rank the entire item corpus or recommend multiple items that align with user interests.

206 **Language-representation-based recommendation.** We argue that AlphaRec follows a new CF  
 207 paradigm, which we term the language-representation-based paradigm. This paradigm replaces  
 208 the ID-based embeddings in ID-Rec with representations from pre-trained LMs, employing feature  
 209 encoders to map LM representations directly into the recommendation space. Few early studies lie in  
 210 this paradigm, including using BERT-style LMs to learn universal sequence representations [37, 44],  
 211 or adopting the same model architecture as ID-Rec with simple input features replacement [34, 35].  
 212 These early explorations, which are mostly based on BERT-style LMs, are usually only applicable in  
 213 certain specific scenarios, such as the transductive setting with the help of ID-based embeddings [37].  
 214 This phenomenon is consistent with our previous findings in Section 2, indicating that BERT-style  
 215 LMs may fail to effectively encode collaborative signals. We point out that AlphaRec is the first  
 216 recommender in the language-representation-based paradigm to surpass the traditional ID-based  
 217 paradigm on multiple tasks, faithfully demonstrating the effectiveness and potential of this paradigm.

## 218 4 Experiments

219 In this section, we aim to explore the effectiveness of AlphaRec. Specifically, we are trying to answer  
 220 the following research questions:

- 221 • **RQ1:** How does AlphaRec perform compared with leading ID-based CF methods?
- 222 • **RQ2:** Can AlphaRec learn general item representations, and achieve good zero-shot recommenda-  
 223 tion performance on entirely new datasets?
- 224 • **RQ3:** Can AlphaRec capture user intention described in natural language and adjust the recom-  
 225 mendation results accordingly?

### 226 4.1 General Recommendation Performance (RQ1)

227 **Motivation.** We aim to explore whether the language-representation-based recommendation paradigm  
 228 can outperform the ID-Rec paradigm. An excellent performance of AlphaRec would shed light on  
 229 the research line of building representation-based recommenders in the future.

230 **Baselines.** We only consider ID-based baselines in this section. We ignore LM-based methods due to  
 231 two practical difficulties: the huge inference cost on datasets with millions of interactions and the  
 232 task limitation of candidate selection or next item prediction. In addition to classic baselines (*i.e.*, MF,  
 233 MultVAE, and LightGCN) introduced in section 2, we consider two categories of leading ID-based  
 234 CF baselines: CL-based CF methods: SGL [80], BC Loss [76], XSimGCL [66] and LM-enhanced  
 235 methods: KAR [48], RLMRec [47]. See more details about baselines in Appendix C.2.1.

236 **Results.** Table 3 presents the performance of AlphaRec compared with leading CF baselines. The  
 237 best-performing methods are bold, while the second-best methods are underlined. Figure 2a and  
 238 Figure 2b report the training efficiency and ablation results. We observe that:

- 239 • **AlphaRec consistently outperforms leading CF baselines by a large margin across all metrics**  
 240 **on all datasets.** AlphaRec shows an improvement ranging from 6.79% to 9.75% on Recall@20  
 241 compared to the best baseline RLMRec [47]. We further conduct the ablation study to explore the  
 242 reason for its success (see more ablation results in Appendix C.2.2). As shown in Figure 2b, each  
 243 component in AlphaRec contributes positively. Specifically, the performance degradation caused by  
 244 replacing the MLP with a linear weight matrix (w/o MLP) indicates that nonlinear transformations  
 245 can further extract the implicit collaborative signals encoded in the LM representation space.

Table 3: The performance comparison with ID-based CF baselines. The improvement achieved by AlphaRec is significant ( $p$ -value  $\ll 0.05$ ).

	Books			Movies & TV			Video Games		
	Recall	NDCG	HR	Recall	NDCG	HR	Recall	NDCG	HR
MF (Rendle et al., 2012)	0.0437	0.0391	0.2476	0.0568	0.0519	0.3377	0.0323	0.0195	0.0864
MultVAE (Liang et al., 2018)	0.0722	0.0597	0.3418	0.0853	0.0776	0.4434	0.0908	0.0531	0.2211
LightGCN (He et al., 2021)	0.0723	0.0608	0.3489	0.0849	0.0747	0.4397	0.1007	0.0590	0.2281
SGL (Wu et al., 2021)	0.0789	0.0657	0.3734	0.0916	0.0838	0.4680	0.1089	0.0634	0.2449
BC Loss (Zhang et al., 2022)	0.0915	0.0779	0.4045	0.1039	0.0943	0.5037	0.1145	0.0668	0.2561
XSimGCL (Yu et al., 2024)	0.0879	0.0745	0.3918	0.1057	0.0984	0.5128	0.1138	0.0662	0.2550
KAR (Xi et al., 2023)	0.0852	0.0734	0.3834	0.1084	0.1001	0.5134	0.1181	0.0693	0.2571
RLMRec (Ren et al., 2024)	<u>0.0928</u>	<u>0.0774</u>	<u>0.4092</u>	<u>0.1119</u>	<u>0.1013</u>	<u>0.5301</u>	<u>0.1384</u>	<u>0.0809</u>	<u>0.2997</u>
<b>AlphaRec</b>	<b>0.0991*</b>	<b>0.0828*</b>	<b>0.4185*</b>	<b>0.1221*</b>	<b>0.1144*</b>	<b>0.5587*</b>	<b>0.1519*</b>	<b>0.0894*</b>	<b>0.3207*</b>
Imp.% over the best baseline	6.79%	5.34%	2.27%	9.12%	10.75%	5.40%	9.75%	10.51%	7.01%

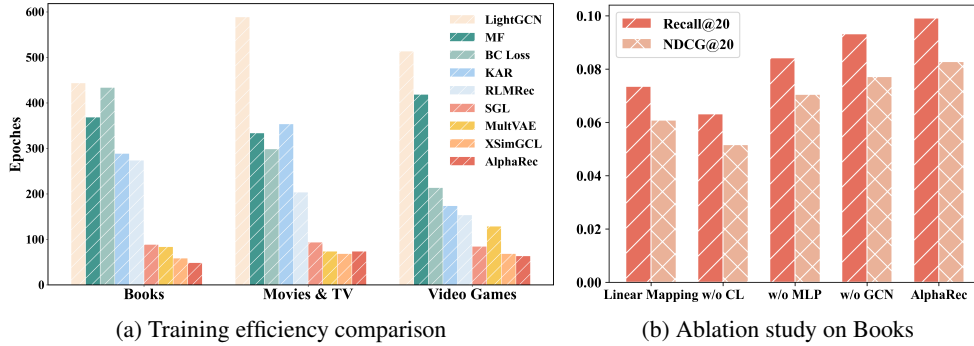


Figure 2: (2a) The bar charts show the number of epochs needed for each model to converge. AlphaRec tends to exhibit an extremely fast convergence speed. (2b) The effect of each component in AlphaRec on Books dataset.

246 Moreover, the performance drop from replacing InfoNCE loss [57] with BPR loss [68] (w/o CL)  
 247 and removing the graph convolution (w/o GCN) suggests that explicitly modeling the collaborative  
 248 relationships through the loss function and model architecture can further enhance recommendation  
 249 performance. These findings suggest that, by carefully designing the model to extract collaborative  
 250 signals, the language-representation-based paradigm can surpass the ID-Rec paradigm.

251 • **The incorporation of semantic LM representations into traditional ID-based CF methods can**  
 252 **lead to significant performance improvements.** We note that two LM-enhanced CF methods,  
 253 KAR and RLMRec, both show improvements over CL-based CF methods. Nevertheless, the combination  
 254 of ID-based embeddings and LM representations in these methods does not yield higher  
 255 results than purely language-representation-based AlphaRec. We attribute this phenomenon to the  
 256 fact that the performance contribution of these methods mainly comes from the LM representations,  
 257 which is consistent with the previous findings [34, 44].

258 • **AlphaRec exhibits fast convergence speed.** We find that the convergence speed of AlphaRec is  
 259 comparable with, or even surpasses, CL-based methods with data augmentation (*e.g.*, SGL [80]  
 260 and XSimGCL [66]). Meanwhile, methods based solely on graph convolution (LightGCN [55]) or  
 261 CL objective (BC Loss [76]) show relatively slow convergence speed, indicating that introducing  
 262 these modules may not lead to convergence speed improvement. Therefore, we attribute the fast  
 263 convergence speed of AlphaRec to the homomorphic relationship between the LM representation  
 264 space and a good recommendation space, so only minor adjustments to the LM representations are  
 265 needed for recommendation.

## 266 4.2 Zero-shot Recommendation Performance on Entirely New Datasets (RQ2)

267 **Motivation.** We aim to explore whether AlphaRec has learned general item representations [37],  
 268 which enables it to perform well on entirely new datasets without any user and item overlap.

269 **Task and datasets.** In zero-shot recommendation [38], there is not any item or user overlap between  
 270 the training set and test set [38, 33], which is different from the research line of cross-domain  
 271 recommendation in ID-Rec [81]. We jointly train AlphaRec on three source datasets (*i.e.*, Books,  
 272 Movies & TV, and Video Games), while testing it on three completely new target datasets (*i.e.*,

Table 4: The zero-shot recommendation performance comparison on entirely new datasets. The improvement achieved by AlphaRec is significant ( $p$ -value  $\ll 0.05$ ).

	Industrial			MovieLens-1M			Book Crossing			
	Recall	NDCG	HR	Recall	NDCG	HR	Recall	NDCG	HR	
full	MF (Rendle et al., 2012)	0.0344	0.0225	0.0521	0.1855	0.3765	0.9634	0.0316	0.0317	0.2382
	MultVAE (Liang et al., 2018)	0.0751	0.0459	0.1125	0.2039	0.3741	0.9740	0.0736	0.0634	0.3716
	LightGCN (He et al., 2021)	0.0785	0.0533	0.1078	0.2019	0.4017	0.9715	0.0630	0.0588	0.3475
zero-shot	Random	0.0148	0.0061	0.0248	0.0068	0.0185	0.2611	0.0039	0.0036	0.0443
	Pop	0.0216	0.0087	0.0396	0.0253	0.0679	0.5439	0.0119	0.0101	0.1157
	ZESRec (Ding et al., 2021)	0.0326	0.0272	0.0628	0.0274	0.0787	0.5786	0.0155	0.0143	0.1347
	UniSRec (Hou et al., 2022)	0.0453	0.0350	0.0863	0.0578	0.1412	0.7135	0.0396	0.0332	0.2454
	<b>AlphaRec</b>	<b>0.0913*</b>	<b>0.0573</b>	<b>0.1277*</b>	<b>0.1486*</b>	<b>0.3215*</b>	<b>0.9296*</b>	<b>0.0660*</b>	<b>0.0545*</b>	<b>0.3381*</b>
	Imp.% over the best zero-shot baseline	157.09%	127.69%	30.29%	66.67%	64.16%	37.78%	101.55%	63.71%	47.97%

273 MovieLens-1M [59], Book Crossing [60], and Industrial [1]) without further training on these new  
 274 datasets. (see more details about how we train AlphaRec on multiple datasets in Appendix C.3.1).

275 **Baselines.** Due to the lack of zero-shot recommenders in the field of general recommendation, we  
 276 slightly modify two zero-shot methods in the sequential recommendation [82], ZESRec [37] and  
 277 UniSRec [37], as baselines. We also incorporate two strategy-based CF methods, Random and Pop  
 278 (see more details about these baselines in Appendix C.3.2).

279 **Results.** Table 4 presents the zero-shot recommendation performance comparison on entirely new  
 280 datasets. The best-performing methods are bold and starred, while the second-best methods are  
 281 underlined. We observe that:

- 282 • **AlphaRec demonstrates strong zero-shot recommendation capabilities, comparable to or even**  
 283 **surpassing the fully trained LightGCN.** On datasets from completely different platforms (*e.g.*,  
 284 MovieLens-1M and Book Crossing), AlphaRec is comparable with the fully trained LightGCN.  
 285 On the same Amazon platform dataset, Industrial, AlphaRec even surpasses LightGCN, which we  
 286 attribute to the possibility that AlphaRec implicitly learns unique user behavioral patterns on the  
 287 Amazon platform [1]. Conversely, ZESRec and UniSRec exhibit a marked performance decrement  
 288 compared with AlphaRec. We attribute this phenomenon to two aspects. On the one hand, BERT-  
 289 style LMs [4, 5] used in these works may not have effectively encoded collaborative signals, which  
 290 is consistent with our findings in Section 2. On the other hand, components designed for the  
 291 next item prediction task in sequential recommendation [83] may not be suitable for capturing the  
 292 general preferences of users in CF scenarios.
- 293 • **The zero-shot recommendation capability of AlphaRec generally benefits from an increased**  
 294 **amount of training data, without harming the performance on source datasets.** As illustrated  
 295 in Figure 8, the zero-shot performance of AlphaRec, when trained on a mixed dataset, is generally  
 296 superior to training on one single dataset [37]. Additionally, we also note that training data with  
 297 themes similar to the target domain contributes more to the zero-shot performance. For instance,  
 298 the zero-shot capability on MovieLens-1M may primarily stem from Movies & TV. Furthermore, we  
 299 discover that AlphaRec, when trained jointly on multiple datasets, hardly experiences a performance  
 300 decline on each source dataset. These findings further point to the general recommendation  
 301 capability of a single pre-trained AlphaRec across multiple datasets. The above findings also offer  
 302 a potential research path to achieve general recommendation capabilities, by incorporating more  
 303 training data with more themes. See more details about these results in Appendix C.3.3.

### 304 4.3 User Intention Capture Performance (RQ3)

305 **Motivation.** We aim to investigate whether a straightforward paradigm shift enables pre-trained  
 306 AlphaRec to perceive text-based user intentions and refine recommendations.

307 **Task and datasets.** We test the user intention capture ability of AlphaRec on MovieLens-1M and  
 308 Video Games. In the test set, only one target item remains for each user [84], with one intention  
 309 query generated by ChatGPT [85, 40] (see the details about how to generate and check these intention  
 310 queries in Appendix C.4.1). In the training stage, we follow the same procedure as illustrated in  
 311 Section 2 to train AlphaRec. In the inference stage, we obtain the LM representation  $e_u^{Intention}$   
 312 for each user intention query and combine it with the original user representation to get a new user  
 313 representation as  $\tilde{e}_u^{(0)} = (1 - \alpha)e_u^{(0)} + \alpha e_u^{Intention}$  [84]. This new user representation is sent into the  
 314 frozen AlphaRec for recommendation. We report a relatively small  $K = 5$  for all metrics to better  
 315 reflect the intention capture accuracy.



Table 5: The performance comparison in user intention capture.

	MovieLens-1M		Video Games	
	HR@5	NDCG@5	HR@5	NDCG@5
TEM (Bi et al., 2020)	0.2738	0.1973	0.2212	0.1425
AlphaRec (w/o Intention)	0.0793	0.0498	0.0663	0.0438
AlphaRec (w Intention)	<b>0.4704*</b>	<b>0.3738*</b>	<b>0.2569*</b>	<b>0.1862*</b>

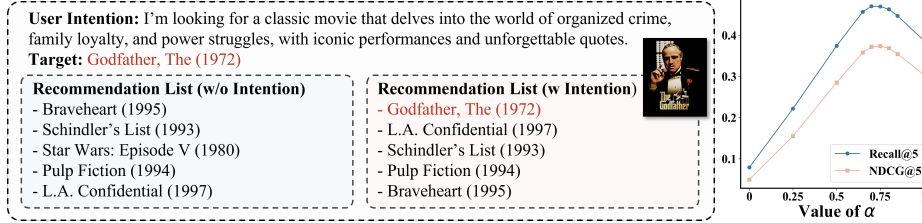


Figure 3: User intention capture experiments on MovieLens-1M. (3a) AlphaRec refines the recommendations according to language-based user intention. (3b) The effect of user intention strength  $\alpha$ .

316 **User intention capture results.** Table 5 represents the user intention capture experiment results, compared with the baseline TEM [86]. Clearly, the introduction of user intention (w Intention) significantly refines the recommendations of the pre-trained AlphaRec (w/o Intention). Moreover, AlphaRec outperforms the baseline model TEM by a large margin, even without additional training on search tasks. We further conduct a case study on MovieLens-1M to demonstrate how AlphaRec captures the user (see more case study results in Appendix C.4.3). As shown in Figure 3a, AlphaRec accurately captures the hidden user intention for “Godfather”, while keeping most of the recommendation results unchanged. This indicates that AlphaRec captures the user intention and historical interests simultaneously.

325 **Effect of the intention strength  $\alpha$ .** By controlling the value of  $\alpha$ , AlphaRec can provide better recommendation results, with a balance between user historical interests and user intent capture. Figure 3b depicts the effect of  $\alpha$ . Initially, as  $\alpha$  increases, the recommendation performance rises accordingly, indicating that incorporating user intention enables AlphaRec to provide better recommendation results. However, as the  $\alpha$  approaches 1, the recommendation performance starts to decrease, which suggests that the user historical interests learned by AlphaRec also play a vital role. The similar effect of  $\alpha$  on Video Games is discussed in Appendix C.4.4.

## 332 5 Limitations

333 There are several limitations not addressed in this paper. On the one hand, although we have demonstrated the excellence of AlphaRec for multiple tasks on various offline datasets, the effectiveness of online employment remains unclear. On the other hand, although we have successfully explored the potential of language-representation-based recommenders by incorporating essential components in leading CF models, we do not elaboratively focus on designing new components for CF models.

## 338 6 Conclusion

339 In this paper, we explored what knowledge about recommendations has been encoded in the LM representation space. Specifically, we found that the advanced LMs representation space exhibits a homomorphic relationship with an effective recommendation space. Based on this finding, we developed a simple yet effective CF model called AlphaRec, which exhibits good recommendation performance with zero-shot recommendation and user intent capture ability. We pointed out that AlphaRec follows a new recommendation paradigm, language-representation-based recommendation, which uses language representations from LMs to represent users and items and completely abandons ID-based embeddings. We believed that AlphaRec is an important stepping stone towards building general recommenders in the future.<sup>1</sup>

<sup>1</sup>The broader impact of AlphaRec will be detailed in Appendix E

348 **References**

- 349 [1] Jianmo Ni, Jiacheng Li, and Julian J. McAuley. Justifying recommendations using distantly-  
350 labeled reviews and fine-grained aspects. In *EMNLP*, 2019.
- 351 [2] Laurens Van der Maaten and Geoffrey Hinton. Visualizing data using t-sne. *Journal of machine*  
352 *learning research*, 9(11), 2008.
- 353 [3] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N. Gomez,  
354 Lukasz Kaiser, and Illia Polosukhin. Attention is all you need. In *NeurIPS*, 2017.
- 355 [4] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. BERT: pre-training of  
356 deep bidirectional transformers for language understanding. In *ACL*, 2019.
- 357 [5] Yinhan Liu, Myle Ott, Naman Goyal, Jingfei Du, Mandar Joshi, Danqi Chen, Omer Levy,  
358 Mike Lewis, Luke Zettlemoyer, and Veselin Stoyanov. Roberta: A robustly optimized BERT  
359 pretraining approach. *CoRR*, abs/1907.11692, 2019.
- 360 [6] Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei,  
361 Nikolay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti Bhosale, Dan Bikel, Lukas  
362 Blecher, Cristian Canton-Ferrer, Moya Chen, Guillem Cucurull, David Esiobu, Jude Fernandes,  
363 Jeremy Fu, Wenyin Fu, Brian Fuller, Cynthia Gao, Vedanuj Goswami, Naman Goyal, Anthony  
364 Hartshorn, Saghar Hosseini, Rui Hou, Hakan Inan, Marcin Kardas, Viktor Kerkez, Madian  
365 Khabsa, Isabel Kloumann, Artem Korenev, Punit Singh Koura, Marie-Anne Lachaux, Thibaut  
366 Lavril, Jenya Lee, Diana Liskovich, Yinghai Lu, Yuning Mao, Xavier Martinet, Todor Mi-  
367 haylov, Pushkar Mishra, Igor Molybog, Yixin Nie, Andrew Poulton, Jeremy Reizenstein, Rashi  
368 Rungta, Kalyan Saladi, Alan Schelten, Ruan Silva, Eric Michael Smith, Ranjan Subramanian,  
369 Xiaoqing Ellen Tan, Binh Tang, Ross Taylor, Adina Williams, Jian Xiang Kuan, Puxin Xu,  
370 Zheng Yan, Iliyan Zarov, Yuchen Zhang, Angela Fan, Melanie Kambadur, Sharan Narang,  
371 Aurélien Rodriguez, Robert Stojnic, Sergey Edunov, and Thomas Scialom. Llama 2: Open  
372 foundation and fine-tuned chat models. *CoRR*, abs/2307.09288, 2023.
- 373 [7] Tom B. Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared Kaplan, Prafulla  
374 Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, Sandhini Agar-  
375 wal, Ariel Herbert-Voss, Gretchen Krueger, Tom Henighan, Rewon Child, Aditya Ramesh,  
376 Daniel M. Ziegler, Jeffrey Wu, Clemens Winter, Christopher Hesse, Mark Chen, Eric Sigler,  
377 Mateusz Litwin, Scott Gray, Benjamin Chess, Jack Clark, Christopher Berner, Sam McCand-  
378 lish, Alec Radford, Ilya Sutskever, and Dario Amodei. Language models are few-shot learners.  
379 In *NeurIPS*, 2020.
- 380 [8] Kenneth Li, Aspen K. Hopkins, David Bau, Fernanda B. Viégas, Hanspeter Pfister, and  
381 Martin Wattenberg. Emergent world representations: Exploring a sequence model trained on a  
382 synthetic task. In *ICLR*, 2023.
- 383 [9] Ivan Vulic, Edoardo Maria Ponti, Robert Litschko, Goran Glavas, and Anna Korhonen. Probing  
384 pretrained language models for lexical semantics. In *EMNLP*, 2020.
- 385 [10] Wes Gurnee and Max Tegmark. Language models represent space and time. *CoRR*,  
386 abs/2310.02207, 2023.
- 387 [11] Wenqi Fan, Zihuai Zhao, Jiatong Li, Yunqing Liu, Xiaowei Mei, Yiqi Wang, Jiliang Tang,  
388 and Qing Li. Recommender systems in the era of large language models (llms). *CoRR*,  
389 abs/2307.02046, 2023.
- 390 [12] Lei Li, Yongfeng Zhang, Dugang Liu, and Li Chen. Large language models for generative  
391 recommendation: A survey and visionary discussions. *CoRR*, abs/2309.01157, 2023.
- 392 [13] Likang Wu, Zhi Zheng, Zhaopeng Qiu, Hao Wang, Hongchao Gu, Tingjia Shen, Chuan Qin,  
393 Chen Zhu, Hengshu Zhu, Qi Liu, Hui Xiong, and Enhong Chen. A survey on large language  
394 models for recommendation. *CoRR*, abs/2305.19860, 2023.
- 395 [14] Jianghao Lin, Xinyi Dai, Yunjia Xi, Weiwen Liu, Bo Chen, Xiangyang Li, Chenxu Zhu,  
396 Hui Feng Guo, Yong Yu, Ruiming Tang, and Weinan Zhang. How can recommender systems  
397 benefit from large language models: A survey. *CoRR*, abs/2306.05817, 2023.

- 398 [15] Junjie Zhang, Ruobing Xie, Yupeng Hou, Wayne Xin Zhao, Leyu Lin, and Ji-Rong Wen. Rec-  
399 ommendation as instruction following: A large language model empowered recommendation  
400 approach. *CoRR*, abs/2305.07001, 2023.
- 401 [16] Keqin Bao, Jizhi Zhang, Yang Zhang, Wenjie Wang, Fuli Feng, and Xiangnan He. Tallrec: An  
402 effective and efficient tuning framework to align large language model with recommendation.  
403 In *RecSys*, 2023.
- 404 [17] Jiayi Liao, Sihang Li, Zhengyi Yang, Jiancan Wu, Yancheng Yuan, and Xiang Wang. Llara:  
405 Aligning large language models with sequential recommenders. *CoRR*, abs/2312.02445, 2023.
- 406 [18] Yaochen Zhu, Liang Wu, Qi Guo, Liangjie Hong, and Jundong Li. Collaborative large language  
407 model for recommender systems. *CoRR*, abs/2311.01343, 2023.
- 408 [19] Bowen Zheng, Yupeng Hou, Hongyu Lu, Yu Chen, Wayne Xin Zhao, Ming Chen, and  
409 Ji-Rong Wen. Adapting large language models by integrating collaborative semantics for  
410 recommendation. *CoRR*, abs/2311.09049, 2023.
- 411 [20] Yang Zhang, Fuli Feng, Jizhi Zhang, Keqin Bao, Qifan Wang, and Xiangnan He. Collm:  
412 Integrating collaborative embeddings into large language models for recommendation. *CoRR*,  
413 abs/2310.19488, 2023.
- 414 [21] Arpita Vats, Vinija Jain, Rahul Raja, and Aman Chadha. Exploring the impact of large  
415 language models on recommender systems: An extensive review. *CoRR*, abs/2402.18590,  
416 2024.
- 417 [22] Chengkai Huang, Tong Yu, Kaige Xie, Shuai Zhang, Lina Yao, and Julian J. McAuley. Founda-  
418 tion models for recommender systems: A survey and new perspectives. *CoRR*, abs/2402.11143,  
419 2024.
- 420 [23] Lanling Xu, Junjie Zhang, Bingqian Li, Jinpeng Wang, Mingchen Cai, Wayne Xin Zhao, and  
421 Ji-Rong Wen. Prompting large language models for recommender systems: A comprehensive  
422 framework and empirical analysis. *CoRR*, abs/2401.04997, 2024.
- 423 [24] Shijie Geng, Shuchang Liu, Zuohui Fu, Yingqiang Ge, and Yongfeng Zhang. Recommendation  
424 as language processing (RLP): A unified pretrain, personalized prompt & predict paradigm  
425 (P5). In *RecSys*, 2022.
- 426 [25] Zeyu Cui, Jianxin Ma, Chang Zhou, Jingren Zhou, and Hongxia Yang. M6-rec: Generative  
427 pretrained language models are open-ended recommender systems. *CoRR*, abs/2205.08084,  
428 2022.
- 429 [26] Jianghao Lin, Rong Shan, Chenxu Zhu, Kounianhua Du, Bo Chen, Shigang Quan, Ruiming  
430 Tang, Yong Yu, and Weinan Zhang. Rella: Retrieval-enhanced large language models for  
431 lifelong sequential behavior comprehension in recommendation. *CoRR*, abs/2308.11131, 2023.
- 432 [27] Zhengyi Yang, Jiancan Wu, Yanchen Luo, Jizhi Zhang, Yancheng Yuan, An Zhang, Xiang  
433 Wang, and Xiangnan He. Large language model can interpret latent space of sequential  
434 recommender. *CoRR*, abs/2310.20487, 2023.
- 435 [28] Shashank Rajput, Nikhil Mehta, Anima Singh, Raghunandan Hulikal Keshavan, Trung Vu,  
436 Lukasz Heldt, Lichan Hong, Yi Tay, Vinh Q. Tran, Jonah Samost, Maciej Kula, Ed H. Chi,  
437 and Mahesh Sathiamoorthy. Recommender systems with generative retrieval. In Alice Oh,  
438 Tristan Naumann, Amir Globerson, Kate Saenko, Moritz Hardt, and Sergey Levine, editors,  
439 *NeurIPS*, 2023.
- 440 [29] Qijiong Liu, Nuo Chen, Tetsuya Sakai, and Xiao-Ming Wu. ONCE: boosting content-based  
441 recommendation with both open- and closed-source large language models. In Luz Angelica  
442 Caudillo-Mata, Silvio Lattanzi, Andrés Muñoz Medina, Leman Akoglu, Aristides Gionis, and  
443 Sergei Vassilvitskii, editors, *WSDM*, 2024.
- 444 [30] Jiaqi Zhai, Lucy Liao, Xing Liu, Yueming Wang, Rui Li, Xuan Cao, Leon Gao, Zhaojie Gong,  
445 Fangda Gu, Michael He, Yinghai Lu, and Yu Shi. Actions speak louder than words: Trillion-  
446 parameter sequential transducers for generative recommendations. *CoRR*, abs/2402.17152,  
447 2024.

- 448 [31] Wenyue Hua, Shuyuan Xu, Yingqiang Ge, and Yongfeng Zhang. How to index item ids for  
449 recommendation foundation models. In Qingyao Ai, Yiqin Liu, Alistair Moffat, Xuanjing  
450 Huang, Tetsuya Sakai, and Justin Zobel, editors, *SIGIR-AP*, 2023.
- 451 [32] Guanghu Yuan, Fajie Yuan, Yudong Li, Beibei Kong, Shujie Li, Lei Chen, Min Yang, Chenyun  
452 Yu, Bo Hu, Zang Li, Yu Xu, and Xiaohu Qie. Tenrec: A large-scale multipurpose benchmark  
453 dataset for recommender systems. In Sanmi Koyejo, S. Mohamed, A. Agarwal, Danielle  
454 Belgrave, K. Cho, and A. Oh, editors, *NeurIPS*, 2022.
- 455 [33] Jiaqi Zhang, Yu Cheng, Yongxin Ni, Yunzhu Pan, Zheng Yuan, Junchen Fu, Youhua Li,  
456 Jie Wang, and Fajie Yuan. Ninerec: A benchmark dataset suite for evaluating transferable  
457 recommendation. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2024.
- 458 [34] Zheng Yuan, Fajie Yuan, Yu Song, Youhua Li, Junchen Fu, Fei Yang, Yunzhu Pan, and Yongxin  
459 Ni. Where to go next for recommender systems? ID- vs. modality-based recommender models  
460 revisited. In *SIGIR*, 2023.
- 461 [35] Ruyi Li, Wenhao Deng, Yu Cheng, Zheng Yuan, Jiaqi Zhang, and Fajie Yuan. Exploring the  
462 upper limits of text-based collaborative filtering using large language models: Discoveries and  
463 insights. *CoRR*, abs/2305.11700, 2023.
- 464 [36] Youhua Li, Hanwen Du, Yongxin Ni, Pengpeng Zhao, Qi Guo, Fajie Yuan, and Xiaofang Zhou.  
465 Multi-modality is all you need for transferable recommender systems. *CoRR*, abs/2312.09602,  
466 2023.
- 467 [37] Yupeng Hou, Shanlei Mu, Wayne Xin Zhao, Yaliang Li, Bolin Ding, and Ji-Rong Wen.  
468 Towards universal sequence representation learning for recommender systems. In *KDD*, pages  
469 585–593. ACM, 2022.
- 470 [38] Hao Ding, Yifei Ma, Anoop Deoras, Yuyang Wang, and Hao Wang. Zero-shot recommender  
471 systems. *CoRR*, abs/2105.08318, 2021.
- 472 [39] Jiacheng Li, Ming Wang, Jin Li, Jinmiao Fu, Xin Shen, Jingbo Shang, and Julian J. McAuley.  
473 Text is all you need: Learning language representations for sequential recommendation. In  
474 *KDD*, 2023.
- 475 [40] Yupeng Hou, Jiacheng Li, Zhankui He, An Yan, Xiusi Chen, and Julian J. McAuley. Bridging  
476 language and items for retrieval and recommendation. *CoRR*, abs/2403.03952, 2024.
- 477 [41] Yupeng Hou, Zhankui He, Julian J. McAuley, and Wayne Xin Zhao. Learning vector-quantized  
478 item representation for transferable sequential recommenders. In Ying Ding, Jie Tang, Juan F.  
479 Sequeda, Lora Aroyo, Carlos Castillo, and Geert-Jan Houben, editors, *WWW*, 2023.
- 480 [42] Zhiming Mao, Huimin Wang, Yiming Du, and Kam-Fai Wong. Unitrec: A unified text-to-text  
481 transformer and joint contrastive learning framework for text-based recommendation. In *ACL*,  
482 2023.
- 483 [43] Zhaopeng Qiu, Xian Wu, Jingyue Gao, and Wei Fan. U-BERT: pre-training user representations  
484 for improved recommendation. In *Thirty-Fifth AAAI Conference on Artificial Intelligence*,  
485 *AAAI 2021*, *Thirty-Third Conference on Innovative Applications of Artificial Intelligence*, *IAAI*  
486 *2021*, *The Eleventh Symposium on Educational Advances in Artificial Intelligence*, *EAAI 2021*,  
487 *Virtual Event, February 2-9, 2021*. AAAI, 2021.
- 488 [44] Lingzi Zhang, Xin Zhou, Zhiwei Zeng, and Zhiqi Shen. Are ID embeddings necessary?  
489 whitening pre-trained text embeddings for effective sequential recommendation. *CoRR*,  
490 abs/2402.10602, 2024.
- 491 [45] Fajie Yuan, Xiangnan He, Alexandros Karatzoglou, and Liguang Zhang. Parameter-efficient  
492 transfer from sequential behaviors for user modeling and recommendation. In *SIGIR*, 2020.
- 493 [46] Wei Wei, Xubin Ren, Jiabin Tang, Qinyong Wang, Lixin Su, Suqi Cheng, Junfeng Wang,  
494 Dawei Yin, and Chao Huang. Llmrec: Large language models with graph augmentation for  
495 recommendation. In Luz Angelica Caudillo-Mata, Silvio Lattanzi, Andrés Muñoz Medina,  
496 Leman Akoglu, Aristides Gionis, and Sergei Vassilvitskii, editors, *WSDM*, 2024.

- 497 [47] Xubin Ren, Wei Wei, Lianghao Xia, Lixin Su, Suqi Cheng, Junfeng Wang, Dawei Yin, and  
498 Chao Huang. Representation learning with large language models for recommendation. *CoRR*,  
499 abs/2310.15950, 2024.
- 500 [48] Yunjia Xi, Weiwen Liu, Jianghao Lin, Jieming Zhu, Bo Chen, Ruiming Tang, Weinan Zhang,  
501 Rui Zhang, and Yong Yu. Towards open-world recommendation with knowledge augmentation  
502 from large language models. *CoRR*, abs/2306.10933, 2023.
- 503 [49] Binzong Geng, Zhaoxin Huan, Xiaolu Zhang, Yong He, Liang Zhang, Fajie Yuan, Jun Zhou,  
504 and Linjian Mo. Breaking the length barrier: Llm-enhanced CTR prediction in long textual  
505 user behaviors. *CoRR*, abs/2403.19347, 2024.
- 506 [50] An Zhang, Leheng Sheng, Yuxin Chen, Hao Li, Yang Deng, Xiang Wang, and Tat-Seng Chua.  
507 On generative agents in recommendation. *CoRR*, abs/2310.10108, 2023.
- 508 [51] Junjie Zhang, Yupeng Hou, Ruobing Xie, Wenqi Sun, Julian J. McAuley, Wayne Xin Zhao,  
509 Leyu Lin, and Ji-Rong Wen. Agentcf: Collaborative learning with autonomous language  
510 agents for recommender systems. *CoRR*, abs/2310.09233, 2023.
- 511 [52] Xiang Wang, Xiangnan He, Meng Wang, Fuli Feng, and Tat-Seng Chua. Neural graph  
512 collaborative filtering. In *SIGIR*, 2019.
- 513 [53] Yang Li, Tong Chen, Yadan Luo, Hongzhi Yin, and Zi Huang. Discovering collaborative  
514 signals for next POI recommendation with iterative seq2graph augmentation. In *IJCAI*, 2021.
- 515 [54] Yehuda Koren, Robert M. Bell, and Chris Volinsky. Matrix factorization techniques for  
516 recommender systems. *Computer*, 42(8):30–37, 2009.
- 517 [55] Xiangnan He, Kuan Deng, Xiang Wang, Yan Li, Yong-Dong Zhang, and Meng Wang. Lightgcn:  
518 Simplifying and powering graph convolution network for recommendation. In *SIGIR*, 2021.
- 519 [56] Aäron van den Oord, Yazhe Li, and Oriol Vinyals. Representation learning with contrastive  
520 predictive coding. *CoRR*, abs/1807.03748, 2018.
- 521 [57] Jiancan Wu, Xiang Wang, Xingyu Gao, Jiawei Chen, Hongcheng Fu, Tianyu Qiu, and Xi-  
522 angnan He. On the effectiveness of sampled softmax loss for item recommendation. *CoRR*,  
523 abs/2201.02327, 2022.
- 524 [58] Steffen Rendle. Item recommendation from implicit feedback. In *Recommender Systems*  
525 *Handbook*. Springer US, 2022.
- 526 [59] F. Maxwell Harper and Joseph A. Konstan. The movielens datasets: History and context. *ACM*  
527 *Trans. Interact. Intell. Syst.*, 5(4):19:1–19:19, 2016.
- 528 [60] Hyeop Lee, Jinbae Im, Seongwon Jang, Hyunsook Cho, and Sehee Chung. Melu: Meta-  
529 learned user preference estimator for cold-start recommendation. In *KDD*, 2019.
- 530 [61] Xiangnan He, Lizi Liao, Hanwang Zhang, Liqiang Nie, Xia Hu, and Tat-Seng Chua. Neural  
531 collaborative filtering. In *WWW*, 2017.
- 532 [62] Xuheng Cai, Chao Huang, Lianghao Xia, and Xubin Ren. Lightgcl: Simple yet effective graph  
533 contrastive learning for recommendation. In *ICLR*, 2023.
- 534 [63] Dawen Liang, Rahul G. Krishnan, Matthew D. Hoffman, and Tony Jebara. Variational  
535 autoencoders for collaborative filtering. In *WWW*, 2018.
- 536 [64] Jiawei Chen, Junkang Wu, Jiancan Wu, Xuezhi Cao, Sheng Zhou, and Xiangnan He. Adap- $\tau$  :  
537 Adaptively modulating embedding magnitude for recommendation. In *WWW*, 2023.
- 538 [65] An Zhang, Leheng Sheng, Zhibo Cai, Xiang Wang, and Tat-Seng Chua. Empowering collabo-  
539 rative filtering with principled adversarial contrastive loss. In *NeurIPS*, 2023.
- 540 [66] Junliang Yu, Xin Xia, Tong Chen, Lizhen Cui, Nguyen Quoc Viet Hung, and Hongzhi Yin.  
541 Xsimgcl: Towards extremely simple graph contrastive learning for recommendation. *IEEE*  
542 *Trans. Knowl. Data Eng.*, 36(2):913–926, 2024.

- 543 [67] Abhilasha Ravichander, Yonatan Belinkov, and Eduard H. Hovy. Probing the probing paradigm:  
544 Does probing accuracy entail task relevance? In *EACL*, 2021.
- 545 [68] Steffen Rendle, Christoph Freudenthaler, Zeno Gantner, and Lars Schmidt-Thieme. BPR:  
546 bayesian personalized ranking from implicit feedback. *CoRR*, abs/1205.2618, 2012.
- 547 [69] Albert Q. Jiang, Alexandre Sablayrolles, Arthur Mensch, Chris Bamford, Devendra Singh  
548 Chaplot, Diego de Las Casas, Florian Bressand, Gianna Lengyel, Guillaume Lample, Lu-  
549 cile Saulnier, L elio Renard Lavaud, Marie-Anne Lachaux, Pierre Stock, Teven Le Scao,  
550 Thibaut Lavril, Thomas Wang, Timoth e Lacroix, and William El Sayed. Mistral 7b. *CoRR*,  
551 abs/2310.06825, 2023.
- 552 [70] Arvind Neelakantan, Tao Xu, Raul Puri, Alec Radford, Jesse Michael Han, Jerry Tworek,  
553 Qiming Yuan, Nikolas Tezak, Jong Wook Kim, Chris Hallacy, Johannes Heidecke, Pranav  
554 Shyam, Boris Power, Tyna Eloundou Nekoul, Girish Sastry, Gretchen Krueger, David Schnurr,  
555 Felipe Petroski Such, Kenny Hsu, Madeleine Thompson, Tabarak Khan, Toki Sherbakov,  
556 Joanne Jang, Peter Welinder, and Lilian Weng. Text and code embeddings by contrastive  
557 pre-training. *CoRR*, abs/2201.10005, 2022.
- 558 [71] Rui Meng, Ye Liu, Shafiq Rayhan Joty, Caiming Xiong, Yingbo Zhou, and Semih Yavuz. Sfr-  
559 embedding-mistral:enhance text retrieval with transfer learning. Salesforce AI Research Blog,  
560 2024. URL <https://blog.salesforceairesearch.com/sfr-embedded-mistral/>.
- 561 [72] Shiwen Wu, Fei Sun, Wentao Zhang, Xu Xie, and Bin Cui. Graph neural networks in  
562 recommender systems: A survey. *ACM Comput. Surv.*, 55(5):97:1–97:37, 2023.
- 563 [73] Felix Wu, Amauri H. Souza Jr., Tianyi Zhang, Christopher Fifty, Tao Yu, and Kilian Q.  
564 Weinberger. Simplifying graph convolutional networks. In *ICML*, 2019.
- 565 [74] Lei Chen, Le Wu, Richang Hong, Kun Zhang, and Meng Wang. Revisiting graph based  
566 collaborative filtering: A linear residual graph convolutional network approach. In *AAAI*, 2020.
- 567 [75] Junliang Yu, Hongzhi Yin, Xin Xia, Tong Chen, Lizhen Cui, and Quoc Viet Hung Nguyen.  
568 Are graph augmentations necessary?: Simple graph contrastive learning for recommendation.  
569 In *SIGIR*, 2022.
- 570 [76] An Zhang, Wenchang Ma, Xiang Wang, and Tat-Seng Chua. Incorporating bias-aware margins  
571 into contrastive loss for collaborative filtering. In *NeurIPS*, 2022.
- 572 [77] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini  
573 Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, Gretchen Krueger, and  
574 Ilya Sutskever. Learning transferable visual models from natural language supervision. In  
575 *ICML*, 2021.
- 576 [78] Feng Wang and Huaping Liu. Understanding the behaviour of contrastive loss. In *CVPR*,  
577 2021.
- 578 [79] Wayne Xin Zhao, Kun Zhou, Junyi Li, Tianyi Tang, Xiaolei Wang, Yupeng Hou, Yingqian  
579 Min, Beichen Zhang, Junjie Zhang, Zican Dong, Yifan Du, Chen Yang, Yushuo Chen, Zhipeng  
580 Chen, Jinhao Jiang, Ruiyang Ren, Yifan Li, Xinyu Tang, Zikang Liu, Peiyu Liu, Jian-Yun Nie,  
581 and Ji-Rong Wen. A survey of large language models. *CoRR*, abs/2303.18223, 2023.
- 582 [80] Jiancan Wu, Xiang Wang, Fuli Feng, Xiangnan He, Liang Chen, Jianxun Lian, and Xing Xie.  
583 Self-supervised graph learning for recommendation. In *SIGIR*, 2021.
- 584 [81] Feng Zhu, Yan Wang, Chaochao Chen, Jun Zhou, Longfei Li, and Guanfeng Liu. Cross-domain  
585 recommendation: Challenges, progress, and prospects. In *IJCAI*, 2021.
- 586 [82] Shoujin Wang, Liang Hu, Yan Wang, Longbing Cao, Quan Z. Sheng, and Mehmet A. Orgun.  
587 Sequential recommender systems: Challenges, progress and prospects. In *IJCAI*, 2019.
- 588 [83] Wang-Cheng Kang and Julian J. McAuley. Self-attentive sequential recommendation. In  
589 *ICDM*, 2018.

- 590 [84] Qingyao Ai, Yongfeng Zhang, Keping Bi, Xu Chen, and W. Bruce Croft. Learning a hierarchi-  
591 cal embedding model for personalized product search. In *SIGIR*, 2017.
- 592 [85] OpenAI. GPT-4 technical report. *CoRR*, 2023.
- 593 [86] Keping Bi, Qingyao Ai, and W. Bruce Croft. A transformer-based embedding model for  
594 personalized product search. In *SIGIR*, 2020.
- 595 [87] Guillaume Alain and Yoshua Bengio. Understanding intermediate layers using linear classifier  
596 probes. In *ICLR (Workshop)*, 2017.
- 597 [88] Roma Patel and Ellie Pavlick. Mapping language models to grounded conceptual spaces. In  
598 *ICLR*, 2022.
- 599 [89] Kiho Park, Yo Joong Choe, and Victor Veitch. The linear representation hypothesis and the  
600 geometry of large language models. *CoRR*, abs/2311.03658, 2023.
- 601 [90] Xubin Ren, Wei Wei, Lianghao Xia, and Chao Huang. A comprehensive survey on self-  
602 supervised learning for recommendation. *arXiv preprint arXiv:2404.03354*, 2024.
- 603 [91] Zheng Chen. PALR: personalization aware llms for recommendation. *CoRR*, abs/2305.07622,  
604 2023.
- 605 [92] Yuling Wang, Changxin Tian, Binbin Hu, Yanhua Yu, Ziqi Liu, Zhiqiang Zhang, Jun Zhou,  
606 Liang Pang, and Xiao Wang. Can small language models be good reasoners for sequential  
607 recommendation? *CoRR*, abs/2403.04260, 2024.
- 608 [93] Qingxiu Dong, Lei Li, Damai Dai, Ce Zheng, Zhiyong Wu, Baobao Chang, Xu Sun, Jingjing  
609 Xu, and Zhifang Sui. A survey on in-context learning. *arXiv preprint arXiv:2301.00234*, 2022.
- 610 [94] Yupeng Hou, Junjie Zhang, Zihan Lin, Hongyu Lu, Ruobing Xie, Julian J. McAuley, and  
611 Wayne Xin Zhao. Large language models are zero-shot rankers for recommender systems. In  
612 *ECIR*, 2024.
- 613 [95] Junling Liu, Chao Liu, Renjie Lv, Kang Zhou, and Yan Zhang. Is chatgpt a good recommender?  
614 A preliminary study. *CoRR*, abs/2304.10149, 2023.
- 615 [96] Sunhao Dai, Ninglu Shao, Haiyuan Zhao, Weijie Yu, Zihua Si, Chen Xu, Zhongxiang Sun,  
616 Xiao Zhang, and Jun Xu. Uncovering chatgpt’s capabilities in recommender systems. In  
617 *RecSys*, 2023.
- 618 [97] Yunfan Gao, Tao Sheng, Youlin Xiang, Yun Xiong, Haofen Wang, and Jiawei Zhang. Chat-  
619 rec: Towards interactive and explainable llms-augmented recommender system. *CoRR*,  
620 abs/2303.14524, 2023.
- 621 [98] Xiangyang Li, Bo Chen, Lu Hou, and Ruiming Tang. Ctrl: Connect tabular and language  
622 model for ctr prediction. *arXiv preprint arXiv:2306.02841*, 2023.
- 623 [99] Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timo-  
624 thée Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, Aurélien Rodriguez,  
625 Armand Joulin, Edouard Grave, and Guillaume Lample. Llama: Open and efficient foundation  
626 language models. *CoRR*, abs/2302.13971, 2023.
- 627 [100] Eric Todd, Millicent L. Li, Arnab Sen Sharma, Aaron Mueller, Byron C. Wallace, and David  
628 Bau. Function vectors in large language models. *CoRR*, abs/2310.15213, 2023.
- 629 [101] Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Brian Ichter, Fei Xia, Ed H. Chi,  
630 Quoc V. Le, and Denny Zhou. Chain-of-thought prompting elicits reasoning in large language  
631 models. In *NeurIPS*, 2022.
- 632 [102] Julian J. McAuley, Rahul Pandey, and Jure Leskovec. Inferring networks of substitutable and  
633 complementary products. In *KDD*, 2015.
- 634 [103] Walid Krichene and Steffen Rendle. On sampled metrics for item recommendation. In *KDD*,  
635 2020.

## 636 A Related Works

637 **Representations in LMs.** The impressive capabilities demonstrated by LMs across various tasks  
638 raise a wide concern about what they have learned in the representation space. An important and  
639 effective approach for interpreting and analyzing representations of LMs is linear probing [67, 87].  
640 The main idea of linear probing is simple: training linear classifiers to predict some specific attributes  
641 or concepts (*e.g.*, lexical structure [9]) from the representations in the hidden layers of LMs. A high  
642 probing result (*e.g.*, classification accuracy on the out-of-sample test set) tends to imply relevant  
643 information has been implicitly encoded in the representation space of LMs, although this does  
644 not imply LMs directly use these representations [67, 10]. Recent studies empirically demonstrate  
645 that concepts such as color [88], game states [8], and geographic position are encoded in LMs.  
646 Furthermore, these concepts may even be linearly encoded in the representation space of LMs [8, 89].

647 **Collaborative filtering.** Collaborative filtering (CF) [90] is an advanced technique in modern  
648 recommender systems. The prevailing CF methods tend to adopt an ID-based paradigm, where users  
649 and items are typically represented as one-hot vectors, with an embedding table used for lookup [54].  
650 Usually, these embedding parameters are learned by optimizing specific loss functions to reconstruct  
651 the history interaction pattern [68]. Recent advances in CF mainly benefit from two aspects, graph  
652 convolution [72] and contrastive learning [90]. These CF models exhibit superior recommendation  
653 performance by conducting the embedding propagation [52, 55] and applying contrastive learning  
654 objectives [80, 62, 66]. However, although effective, these methods are still limited, due to the  
655 ID-based paradigm. Since one-hot vectors contain no feature information beyond being identifiers, it  
656 is challenging to transfer pre-trained ID embeddings to other domains [37] or to leverage leading  
657 techniques from computer vision (CV) and natural language processing (NLP) [34].

658 **LMs for recommendation.** The remarkable language understanding and reasoning ability shown by  
659 LMs has attracted extensive attention in the field of recommendation. The application of LMs in rec-  
660 ommendation can be categorized into three main approaches: LM-enhanced recommendation, LM as  
661 the modality encoder, and LLM-based recommendation. The first research direction, LLM-enhanced  
662 recommendation, focuses on empowering traditional recommenders with the semantic representations  
663 from LMs [48, 47, 46, 49, 91, 92]. Specifically, these methods introduce representations from LMs as  
664 additional features for traditional ID-based recommenders, to capture complicated user preferences.  
665 The second research line lies in adopting the LM as the text modality encoder, which is also known  
666 as a kind of modality-based recommendation (MoRec) [34, 35]. These methods tend to train the  
667 LM as the text modality encoder together with the traditional recommender. In previous studies,  
668 BERT-style LMs are widely used as the text modality encoder. The third research line, LLM-based  
669 recommendation, directly uses LLMs as the recommender and recommends items in a text generation  
670 paradigm. Early attempts focus on adopting in-context learning (ICL) [93] and prompting pre-trained  
671 LLMs [94–97]. However, such naive methods tend to yield poor performance compared to traditional  
672 models. Therefore, recent studies concentrate on fine-tuning LLMs on recommendation-related corpus  
673 [16, 15, 26, 25, 29] and align the LLMs with the representations from traditional recommenders  
674 as the additional modality [17, 20, 27, 98].

## 675 B Linear Mapping

### 676 B.1 Brief of Used LMs

677 We briefly introduce the LMs we use for linear mapping in Section 2.

- 678 • **BERT** [4] is an encoder-only language model based on the transformer architecture [3], pre-trained  
679 on text corpus with unsupervised tasks. BERT adopts bidirectional self-attention heads to learn  
680 bidirectional representations.
- 681 • **RoBERTa** [5] is an enhanced version of BERT. RoBERTa preserves the architecture of BERT but  
682 improves it by training with more data and large batches, adopting dynamic masking, and removing  
683 the next sentence prediction objective.
- 684 • **Llama2-7B** [6] is an open-source decoder-only LLM with 7 billion parameters. Llama2 adopts  
685 grouped-query attention, with longer context length and larger size of the pre-training corpus  
686 compared with Llama-7B [99].



Table 6: Linear mapping performance of randomly shuffled item representations

	Books			Movies & TV			Video Games		
	Recall	NDCG	HR	Recall	NDCG	HR	Recall	NDCG	HR
<b>BERT</b>	0.0226	0.0194	0.1240	0.0415	0.0399	0.2362	0.0524	0.0309	0.1245
<b>text-embeddings-3-large (Random)</b>	0.0200	0.0197	0.1316	0.0559	0.0528	0.3204	0.0562	0.0328	0.1351
<b>text-embeddings-3-large</b>	0.0735	0.0608	0.3355	0.1109	0.1023	0.5200	0.1367	0.0793	0.2928

Table 7: Dataset statistics.

	Books	Movies & TV	Video Games	Industrial	MovieLens-1M	Book Crossing
#Users	7,176	14,382	40,834	15,141	6,040	6,273
#Items	10,728	1,000	14,344	5,163	3,043	5,335
#Interactions	1,304,453	129,748	390,013	82,578	995,492	253,057
Density	0.0169	0.0090	0.0701	0.0010	0.0542	0.0076

- 687 • **Mistral-7B** [69] is an open-source pre-trained decoder-only LLM with 7 billion parameters. Mistral  
688 7B leverages grouped-query attention, coupled with sliding window attention for faster and lower  
689 cost inference.
- 690 • **text-embedding-ada-v2 & text-embeddings-3-large** [70] are leading text embedding models  
691 released by OpenAI. These models are built upon decoder-only GPT models, pre-trained on  
692 unsupervised data at scale with contrastive learning objectives.
- 693 • **SFR-Embedding-Mistral** [71] is a decoder-based text embedding model built upon the open-  
694 source LLM Mistral-7B [69]. SFR-Embedding-Mistral introduces task-homogeneous batching and  
695 computes contrastive loss on “hard negatives”, which brings a better performance than the vanilla  
696 Mistral-7B model.

## 697 B.2 Extracting Representations from LMs

698 We present how to extract representations from LMs. For encoder-based LMs (*e.g.*, BERT [4]  
699 and RoBERTa [5]), we use the output representation corresponding to the [CLS] token [40]. For  
700 decoder-based models (*e.g.*, Llama-7B [6, 69], Mistral-7B, and SFR-Embedding-Mistral [71]),  
701 we use the representation in the last transformer block [3], corresponding to the last input token  
702 [10, 100, 70]. Especially, for the commercial closed-source model (*e.g.*, text-embedding-ada-v2 and  
703 text-embeddings-3-large<sup>2</sup> [70]), we directly call the API interface to obtain representations.

## 704 B.3 Empirical Findings

705 We find more evidence about representations in leading LM encode collaborative signals beyond  
706 better feature encoding ability. We randomly shuffle item representations and conduct the same linear  
707 mapping experiment. As illustrated in Table 6, randomly shuffled representations, text-embeddings-  
708 3-large (Random), yield similar performance with BERT, lagging largely behind the vanilla linear  
709 mapping method. These results indicate that BERT may only serve as a good feature encoder, while  
710 the latest LM may further encode collaborative signals beyond naive feature encoding.

## 711 C Experiments

### 712 C.1 Datasets

713 We incorporate six datasets in this paper, including four datasets from the Amazon platform<sup>3</sup> [1]  
714 (*i.e.*, Books, Movies & TV, Video Games, and Industrial), and two datasets from other platforms (*i.e.*,  
715 MovieLens-1M and Book Crossing). Table 7 reports the data statistics of each dataset.

716 We divide the history interaction of each user into training, validation, and testing sets with a ratio  
717 of 4:3:3, and remove users with less than 20 interactions following previous studies [50]. We also  
718 remove items from the testing and validation sets that do not appear in the training set, to address the  
719 cold start problem.

<sup>2</sup><https://platform.openai.com/docs/guides/embeddings>

<sup>3</sup>[www.amazon.com](http://www.amazon.com)

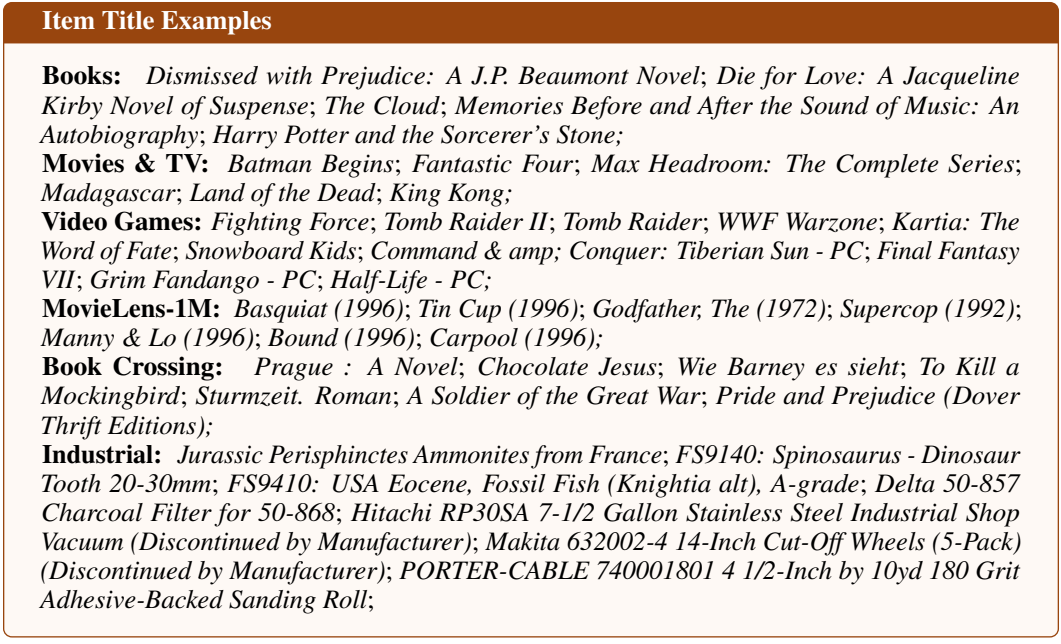


Figure 4: Example of item titles.

720 In this paper, we only use the item titles as the text description. Figure 4 gives some item title  
 721 examples from different datasets.

722 **C.2 General Recommendation**

723 **C.2.1 Baselines**

724 We incorporate a series of CF models as our baselines for general recommendation. These models  
 725 are classified as classical CF methods (MF, MultVAE, and LightGCN), CL-based CF methods (SGL,  
 726 BC Loss, and XSimGCL), and LM-enhanced CF methods (KAR, RLMRec). For these LM-enhanced  
 727 CF methods, we adopt the leading CF method XSimGCL as the backbone.

- 728 • **MF** [54, 68] is the most basic CF model. It denotes users and items with ID-based embeddings and  
 729 conducts matrix factorization with Bayesian personalized ranking (BPR) loss.
- 730 • **MultVAE** [63] is a traditional CF model based on the variational autoencoder (VAE). It regards the  
 731 item recommendation as a generative process from a multinomial distribution and uses variational  
 732 inference to estimate parameters. We adopt the same model structure as suggested in the paper:  
 733  $600 \rightarrow 200 \rightarrow 600$ .
- 734 • **LightGCN** [55] is a light graph convolution network tailored for the recommendation, which  
 735 deletes redundant feature transformation and activation function in NGCF [52].
- 736 • **SGL** [80] introduces graph contrastive learning into recommender models for the first time. By  
 737 employing node or edge dropout to generate augmented graph views and conduct contrastive  
 738 learning between two views, SGL achieves better performance than LightGCN.
- 739 • **BC Loss** [76] introduces a robust and model-agnostic contrastive loss, handling various data biases  
 740 in recommendation, especially for popularity bias.
- 741 • **XSimGCL** [66] directly generates augmented views by adding noise into the inner layer of  
 742 LightGCN without graph augmentation. The simplicity of XSimGCL leads to a faster convergence  
 743 speed and better performance.
- 744 • **KAR** [48] enhances recommender models by integrating knowledge from large language models  
 745 (LLMs). It generates textual descriptions of users and items and combine the LM representations  
 746 with traditional recommenders using a hybrid-expert adaptor.

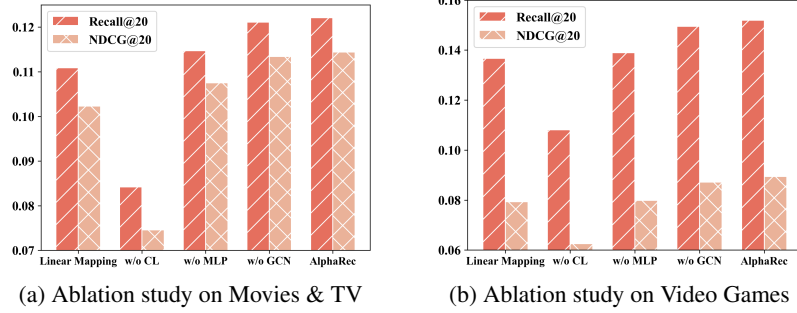


Figure 5: Ablation study

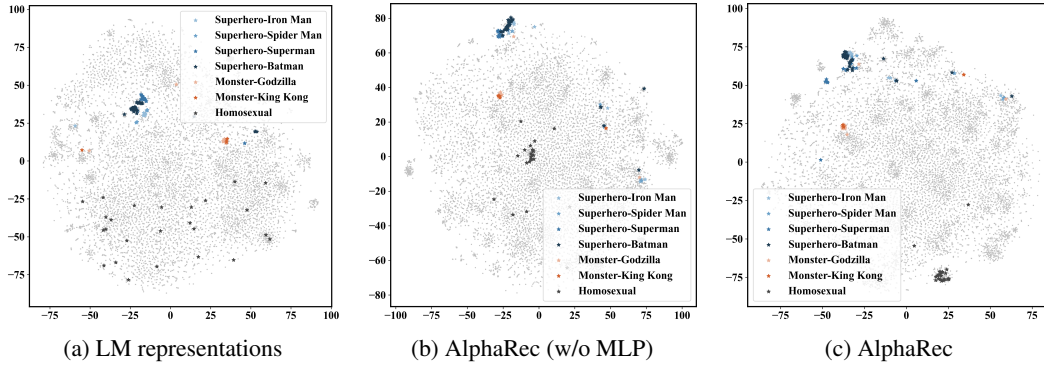


Figure 6: The t-SNE visualization of representations on Movies & TV. (6a) The item representations in the LM space. (6b) The item representations obtained by replacing the MLP with a linear mapping matrix in AlphaRec. (6c) The item representations obtained from AlphaRec.

747 • **RLMRec** [47] aligns semantic representations of users and items with the representations in CF  
 748 models through a contrastive loss, as an additional loss trained together with the CF model. The  
 749 fusion of semantic information and collaborative information brings performance improvement.

### 750 C.2.2 Ablation Study

751 We conduct the same ablation study as introduced in Section 4.1 on Movies & TV and Video Games  
 752 datasets. As illustrated in Figure 5, each component in AlphaRec contributes positively, which is  
 753 consistent with our findings in Section 4.1.

### 754 C.2.3 The t-SNE Visualization Comparison

755 In this section, we aim to intuitively explore how the MLP in AlphaRec further helps in excavating  
 756 collaborative signals in language representations, compared to the linear mapping matrix. We  
 757 visualize the item representations from LMs, AlphaRec (w/o MLP), and AlphaRec in Figure 6, where  
 758 AlphaRec (w/o MLP) denotes replacing the MLP with a linear mapping matrix. We observed that  
 759 movies about superhero and monster cluster in all representation spaces, indicating both AlphaRec  
 760 (w/o MLP) and AlphaRec capture the preference similarities between these items and preserve  
 761 the clustering relationship. The difference between AlphaRec (w/o MLP) and AlphaRec may lie  
 762 in the ability to capture obscure preference similarities among items. As shown in Figure 6a,  
 763 homosexual movies are dispersed in the language space, indicating the possible semantic differences  
 764 between them. AlphaRec successfully captures the preference similarities and gathers these items  
 765 in the representation space, while AlphaRec (w/o MLP) remains some items dispersed. Moreover,  
 766 AlphaRec outperforms AlphaRec (w/o MLP) by a large margin, as indicated in Figure 5a. These  
 767 results indicate that AlphaRec exhibits a more fine-grained preference capture ability with the help of  
 768 nonlinear transformation.

Table 8: The effect of the training dataset on zero-shot recommendation

	Industrial			MovieLens-1M			Book Crossing		
	Recall	NDCG	HR	Recall	NDCG	HR	Recall	NDCG	HR
AlphaRec (trained on Books)	0.0896	0.0562	0.1256	0.1218	0.2619	0.8942	<u>0.0646</u>	<u>0.0532</u>	<u>0.3346</u>
AlphaRec (trained on Movies & TV)	<u>0.0909</u>	<b>0.0581</b>	<u>0.1266</u>	<u>0.1438</u>	<u>0.3122</u>	<u>0.9200</u>	0.0471	0.0406	0.2600
AlphaRec (trained on Video Games)	0.0905	0.0567	0.1225	0.1221	0.2313	0.9034	0.0412	0.0378	0.2585
AlphaRec (trained on mixed dataset)	<b>0.0913</b>	<u>0.0573</u>	<b>0.1277</b>	<b>0.1486</b>	<b>0.3215</b>	<b>0.9296</b>	<b>0.0660</b>	<b>0.0545</b>	<b>0.3381</b>

Table 9: Performance comparison between training on the single dataset and the mixed dataset

	Books			Movies & TV			Video Games		
	Recall	NDCG	HR	Recall	NDCG	HR	Recall	NDCG	HR
AlphaRec (trained on single dataset)	<b>0.0991</b>	<b>0.0828</b>	<b>0.4185</b>	<b>0.1221</b>	<b>0.1144</b>	<b>0.5587</b>	<b>0.1519</b>	<b>0.0894</b>	<b>0.3207</b>
AlphaRec (trained on mixed dataset)	0.0979	0.0818	0.4147	0.1194	0.1107	0.5463	0.1381	0.0827	0.2985

769 **C.3 Zero-shot Recommendation**

770 **C.3.1 Co-training on Multiple Datasets**

771 Co-training on multiple datasets is similar to training on one single dataset, where the only difference  
 772 lies in the negative sampling. When co-training on multiple datasets, the negative items are restricted  
 773 to the same dataset as the positive item rather than the full item pool. The other training procedures  
 774 remain the same with training on one single dataset.

775 **C.3.2 Baselines**

776 Since previous works about zero-shot recommendation mostly focus on sequential recommendation  
 777 [83, 82], we slightly modify two methods in sequential recommendation, ZESRec [38] and UniSRec  
 778 [37] as our baselines. Specifically, we maintain the model structure as provided in the paper, and  
 779 adopt the training paradigm of CF.

- 780 • **Random** denotes randomly recommending items from the entire item pool.
- 781 • **Pop** denotes randomly recommending from the most popular items. Here popularity denotes the  
 782 number of users that have interacted with the item.
- 783 • **ZESRec** [38] is the first work that defines the problem of zero-shot recommendation. To address  
 784 this problem, this work introduces a hierarchical Bayesian model with representations from the  
 785 pre-trained BERT.
- 786 • **UniSRec** [37] aims to learn universal item representations from BERT, with parametric whitening  
 787 and a MoE-enhanced adaptor. By pre-training on multiple source datasets, UniSRec can conduct  
 788 zero-shot recommendation on various datasets in a transductive or inductive paradigm.

789 **C.3.3 The Effect of Training Datasets**

790 **The effect of the training dataset on zero-shot recommendation.** We report the zero-shot  
 791 recommendation performance differences trained on different datasets in Table 8. Here AlphaRec  
 792 (trained on Books) denotes training on a single Books dataset, while AlphaRec (trained on mixed  
 793 dataset) denotes co-training on three Amazon datasets. Generally, training on more datasets lead to a  
 794 better zero-shot performance.

795 **The performance comparison between training on the single dataset and the mixed dataset.** In  
 796 Table 9, AlphaRec (trained on single dataset) denotes training and testing on the same single dataset,  
 797 while AlphaRec (trained on mixed dataset) denotes training on three Amazon datasets and testing  
 798 on one single dataset. Generally, co-training on three Amazon datasets yields similar performance  
 799 compared with training on one single dataset. The only exception lies in Video Games, which shows  
 800 some performance degradation. We attribute this to the difference between the selection of  $\tau$ . We use  
 801  $\tau = 0.15$  when trained on the mixed dataset, while the optimal  $\tau$  for Video Games lies around 0.2.  
 802 These results indicate that a single AlphaRec can capture user preferences among various datasets,  
 803 showcasing a general collaborative signal capture ability.

## 804 C.4 User Intention Capture

### 805 C.4.1 Intention Query Generation

**Intention Query Generation**

**Input**  
You are an expert in generating queries for a target movie. Please help me generate the most suitable query for the target movie within one sentence, following the given example.  
Example:  
TARGET: [BUG-A-SALT 3.0 Black Fly Edition](#).  
QUERY: *I want a gun that I can use while gardening to get rid of stink bugs, ants, flies, and spiders in my house. It needs to be amazing and help me feel less scared.*  
TARGET: [Toy Story \(1995\)](#).

**Output**  
QUERY: *I'm looking for a heartwarming animated movie that follows the adventures of a group of toys who come to life when their owner is not around.*

Figure 7: Example of item query generation.

806 The user intention query is a natural language sentence implying the target item of interest. For  
807 each item in the dataset, we generate a fixed user intention query. Following the previous work  
808 [40], we generate user intention queries with the help of ChatGPT [85]. As shown in Figure 7, we  
809 prompt ChatGPT in a Chain-of-Thought (CoT) [101] paradigm and adopt the output as the user  
810 intention query. We adopt a rule-based strategy to ensure that the output query is in first person, and  
811 regenerate the wrong query. Considering the huge amount of item title text, we use ChatGPT3.5 API  
812 for generating all queries for the budget's sake.

### 813 C.4.2 Baseline

814 AlphaRec exhibits user intention capture abilities, although not specially designed for search tasks.  
815 We compare AlphaRec with TEM [86] which falls in the field of personalized search [84, 102].

- 816 • TEM [86] uses a transformer to encode the intention query together with user history behaviors,  
817 which enables it to achieve better search results by considering the user's historical interest.

### 818 C.4.3 Case Study

819 We conduct two more case studies to verify the user intention capture ability of AlphaRec. As  
820 illustrated in Figure 8 and Figure 9, AlphaRec provides proper recommendation results, including the  
821 target item for the user intention at the top.

### 822 C.4.4 Effect of the Intention Strength Alpha

823 The value of  $\alpha$  controls the balance between the user's historical interests and the user intention  
824 query. A larger  $\alpha$  incorporates more about the user intention while considering less about the user's  
825 historical interests. As shown in Figure 10, the effect of  $\alpha$  on Video Games shows a similar trend  
826 with MovieLens-1M.

## 827 C.5 Training Cost

828 We report the training cost of AlphaRec in this section. Table 10 reports the seconds needed per  
829 epoch and the total training cost until convergence. Here Amazon-Mix denotes the mixed dataset of  
830 Books, Movies & TV, and Video Games. It's worth noting that AlphaRec converges quickly and only  
831 requires a small amount of training time.

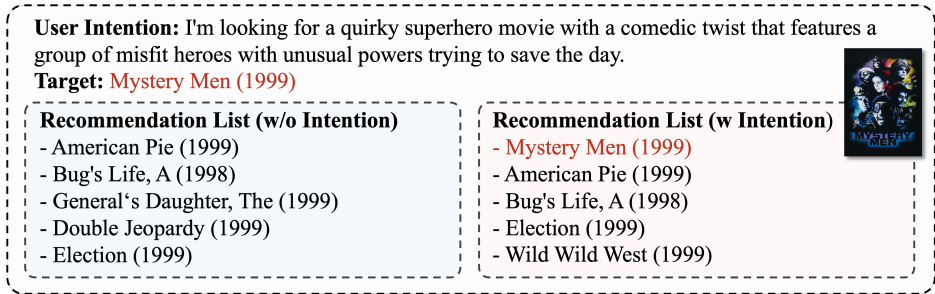


Figure 8: Case study of user intention capture on MovieLens-1M

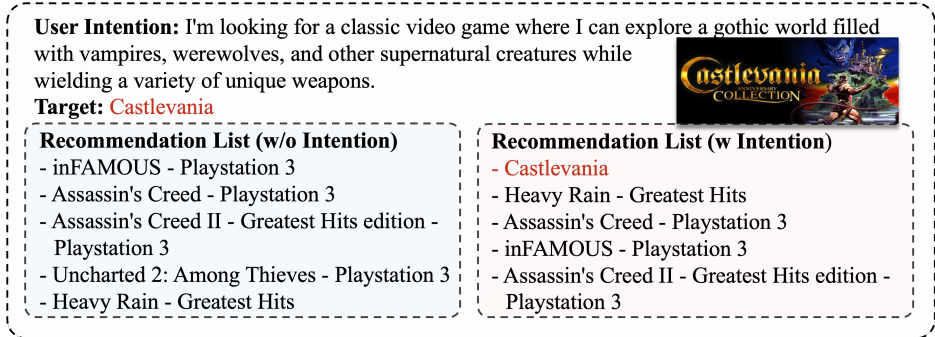


Figure 9: Case study of user intention capture on Video Games

## 832 D Hyperparameter Settings and Implementation Details

833 We conduct all the experiments in PyTorch with a single NVIDIA RTX A5000 (24G) GPU and a  
 834 64 AMD EPYC 7543 32-Core Processor CPU. We optimize all methods with the Adam optimizer.  
 835 For all ID-based CF methods, we set the layer numbers of graph propagation by default at 2, with  
 836 the embedding size as 64 and the size of sampled negative items  $|\mathcal{S}_u|$  as 256. We use the early stop  
 837 strategy to avoid overfitting. We stop the training process if the Recall@20 metric on the validation  
 838 set does not increase for 20 successive evaluations. In AlphaRec, the dimensions of the input and  
 839 output in the two-layer MLP are 3072 and 64 respectively, with the hidden layer dimension as 1536.  
 840 We apply the all-ranking strategy [103] for all experiments, which ranks all items except positive ones  
 841 in the training set for each user. We search hyperparameters for baselines according to the suggestion  
 842 in the literature. The hyperparameter search space is reported in Table 11. For these LM-enhanced  
 843 models, KAR and RLMRec, we also search the hyperparameter of their backbone XSimGCL.

844 For AlphaRec, the only hyperparameter is the temperature  $\tau$  and we search it in [0.05, 2]. We report  
 845 the temperature  $\tau$  we used for each dataset in Table 12. For the mixed dataset Amazon-Mix in  
 846 Section 4.2, we use a universal  $\tau = 0.15$ . We adopt  $\tau = 0.2$  for the MovieLens-1M dataset for the user  
 847 intention capture experiment in Section 4.3.

## 848 E Broader Impact

849 The proposed AlphaRec can significantly improve the performance of zero-shot recommendation  
 850 and the capability of user intent capture, offering a good approach to crafting more personalized  
 851 recommendation results. One concern of AlphaRec is the potential for the representations generated  
 852 by language models can be maliciously attacked, which may result in erroneous or unexpected  
 853 recommendations. Therefore, we kindly advise researchers to cautiously check the quality of the  
 854 language representations before using AlphaRec.

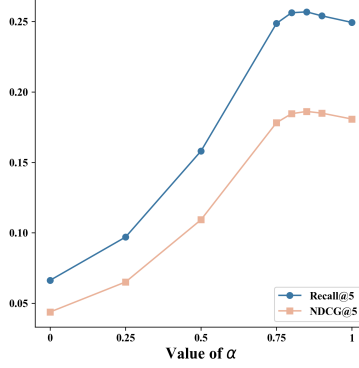


Figure 10: Effect of  $\alpha$  on Video Games

Table 10: Training cost of AlphaRec (seconds per epoch/in total).

	Books	Movies & TV	Video Games	Amazon-Mix
AlphaRec	40.1 / 1363.4	12.3 / 479.7	7.4 / 214.6	107.2 / 5788.8

Table 11: Hyperparameters search spaces for baselines.

	Hyperparameter space
<b>MF &amp; LightGCN</b>	$lr \sim \{1e-5, 3e-5, 5e-5, 1e-4, 3e-4, 5e-4, 1e-3\}$
<b>MultVAE</b>	dropout ratio $\sim \{0, 0.2, 0.5\}$ , $\beta \sim \{0.2, 0.4, 0.6, 0.8\}$
<b>SGL</b>	$\tau \sim [0.05, 2]$ , $\lambda_1 \sim \{0.005, 0.01, 0.05, 0.1, 0.5, 1.0\}$ , $\rho \sim \{0, 0.1, 0.2, 0.3, 0.4, 0.5\}$
<b>BC Loss</b>	$\tau_1 \sim [0.05, 3]$ , $\tau_2 \sim [0.05, 2]$
<b>XSimGCL</b>	$\tau \sim [0.05, 2]$ , $\epsilon \sim \{0.01, 0.05, 0.1, 0.2, 0.5, 1.0\}$ , $\lambda \sim \{0.005, 0.01, 0.05, 0.1, 0.5, 1.0\}$ , $l^* = 1$
<b>KAR</b>	No. shared experts $\sim \{3, 4, 5\}$ , No. preference experts $\sim \{4, 5\}$
<b>RLMRec</b>	kd weight $\sim [0.05, 2]$ , kd temperature $\sim [0.01, 0.05, 0.1, 0.15, 0.2, 0.5, 1]$
<b>ZESRec</b>	$\lambda_u \sim \{0.01, 0.05, 0.1, 0.5, 1.0\}$ , $\lambda_v \sim \{0.01, 0.05, 0.1, 0.5, 1.0\}$
<b>UniSRec</b>	$lr \sim \{3e-4, 1e-3, 3e-3, 1e-2\}$
<b>TEM</b>	$l \sim \{2, 3\}$ , head $h \sim \{4, 8\}$
<b>AlphaRec</b>	$\tau \sim [0.05, 2]$

Table 12: The hyperparameters of AlphaRec

	Books	Movies & TV	Video Games	Amazon-Mix
$\tau$	0.15	0.15	0.2	0.15

## 855 **NeurIPS Paper Checklist**

856 The checklist is designed to encourage best practices for responsible machine learning research,  
857 addressing issues of reproducibility, transparency, research ethics, and societal impact. Do not remove  
858 the checklist: **The papers not including the checklist will be desk rejected.** The checklist should  
859 follow the references and follow the (optional) supplemental material. The checklist does NOT count  
860 towards the page limit.

861 Please read the checklist guidelines carefully for information on how to answer these questions. For  
862 each question in the checklist:

- 863 • You should answer [Yes], [No], or [NA].
- 864 • [NA] means either that the question is Not Applicable for that particular paper or the  
865 relevant information is Not Available.
- 866 • Please provide a short (1–2 sentence) justification right after your answer (even for NA).

867 **The checklist answers are an integral part of your paper submission.** They are visible to the  
868 reviewers, area chairs, senior area chairs, and ethics reviewers. You will be asked to also include it  
869 (after eventual revisions) with the final version of your paper, and its final version will be published  
870 with the paper.

871 The reviewers of your paper will be asked to use the checklist as one of the factors in their evaluation.  
872 While "[Yes]" is generally preferable to "[No]", it is perfectly acceptable to answer "[No]" provided a  
873 proper justification is given (e.g., "error bars are not reported because it would be too computationally  
874 expensive" or "we were unable to find the license for the dataset we used"). In general, answering  
875 "[No]" or "[NA]" is not grounds for rejection. While the questions are phrased in a binary way, we  
876 acknowledge that the true answer is often more nuanced, so please just use your best judgment and  
877 write a justification to elaborate. All supporting evidence can appear either in the main paper or the  
878 supplemental material, provided in appendix. If you answer [Yes] to a question, in the justification  
879 please point to the section(s) where related material for the question can be found.

880 IMPORTANT, please:

- 881 • **Delete this instruction block, but keep the section heading “NeurIPS paper checklist”.**
- 882 • **Keep the checklist subsection headings, questions/answers and guidelines below.**
- 883 • **Do not modify the questions and only use the provided macros for your answers.**

### 884 **1. Claims**

885 Question: Do the main claims made in the abstract and introduction accurately reflect the  
886 paper’s contributions and scope?

887 Answer: [Yes]

888 Justification: We clearly state the claims made in the abstract and introduction.

889 Guidelines:

- 890 • The answer NA means that the abstract and introduction do not include the claims  
891 made in the paper.
- 892 • The abstract and/or introduction should clearly state the claims made, including the  
893 contributions made in the paper and important assumptions and limitations. A No or  
894 NA answer to this question will not be perceived well by the reviewers.
- 895 • The claims made should match theoretical and experimental results, and reflect how  
896 much the results can be expected to generalize to other settings.
- 897 • It is fine to include aspirational goals as motivation as long as it is clear that these goals  
898 are not attained by the paper.

### 899 **2. Limitations**

900 Question: Does the paper discuss the limitations of the work performed by the authors?

901 Answer: [Yes]

902 Justification: We discuss the limitations of this work in the Section 5.



903  
904  
905  
906  
907  
908  
909  
910  
911  
912  
913  
914  
915  
916  
917  
918  
919  
920  
921  
922  
923  
924  
925  
926  
927  
928  
929  
930  
931  
932  
933  
934  
935  
936  
937  
938  
939  
940  
941  
942  
943  
944  
945  
946  
947  
948  
949  
950  
951  
952  
953  
954  
955

Guidelines:

- The answer NA means that the paper has no limitation while the answer No means that the paper has limitations, but those are not discussed in the paper.
- The authors are encouraged to create a separate "Limitations" section in their paper.
- The paper should point out any strong assumptions and how robust the results are to violations of these assumptions (e.g., independence assumptions, noiseless settings, model well-specification, asymptotic approximations only holding locally). The authors should reflect on how these assumptions might be violated in practice and what the implications would be.
- The authors should reflect on the scope of the claims made, e.g., if the approach was only tested on a few datasets or with a few runs. In general, empirical results often depend on implicit assumptions, which should be articulated.
- The authors should reflect on the factors that influence the performance of the approach. For example, a facial recognition algorithm may perform poorly when image resolution is low or images are taken in low lighting. Or a speech-to-text system might not be used reliably to provide closed captions for online lectures because it fails to handle technical jargon.
- The authors should discuss the computational efficiency of the proposed algorithms and how they scale with dataset size.
- If applicable, the authors should discuss possible limitations of their approach to address problems of privacy and fairness.
- While the authors might fear that complete honesty about limitations might be used by reviewers as grounds for rejection, a worse outcome might be that reviewers discover limitations that aren't acknowledged in the paper. The authors should use their best judgment and recognize that individual actions in favor of transparency play an important role in developing norms that preserve the integrity of the community. Reviewers will be specifically instructed to not penalize honesty concerning limitations.

**3. Theory Assumptions and Proofs**

Question: For each theoretical result, does the paper provide the full set of assumptions and a complete (and correct) proof?

Answer: [NA]

Justification: This is an empirical article and contains no theoretical results.

Guidelines:

- The answer NA means that the paper does not include theoretical results.
- All the theorems, formulas, and proofs in the paper should be numbered and cross-referenced.
- All assumptions should be clearly stated or referenced in the statement of any theorems.
- The proofs can either appear in the main paper or the supplemental material, but if they appear in the supplemental material, the authors are encouraged to provide a short proof sketch to provide intuition.
- Inversely, any informal proof provided in the core of the paper should be complemented by formal proofs provided in appendix or supplemental material.
- Theorems and Lemmas that the proof relies upon should be properly referenced.

**4. Experimental Result Reproducibility**

Question: Does the paper fully disclose all the information needed to reproduce the main experimental results of the paper to the extent that it affects the main claims and/or conclusions of the paper (regardless of whether the code and data are provided or not)?

Answer: [Yes]

Justification: We present all the experiment details and datasets in Appendix C, and Hyperparameters settings are reported in Appendix D. Moreover, we have uploaded the code and data we used in the supplementary material.

Guidelines:

- The answer NA means that the paper does not include experiments.

- 956 • If the paper includes experiments, a No answer to this question will not be perceived  
957 well by the reviewers: Making the paper reproducible is important, regardless of  
958 whether the code and data are provided or not.
- 959 • If the contribution is a dataset and/or model, the authors should describe the steps taken  
960 to make their results reproducible or verifiable.
- 961 • Depending on the contribution, reproducibility can be accomplished in various ways.  
962 For example, if the contribution is a novel architecture, describing the architecture fully  
963 might suffice, or if the contribution is a specific model and empirical evaluation, it may  
964 be necessary to either make it possible for others to replicate the model with the same  
965 dataset, or provide access to the model. In general, releasing code and data is often  
966 one good way to accomplish this, but reproducibility can also be provided via detailed  
967 instructions for how to replicate the results, access to a hosted model (e.g., in the case  
968 of a large language model), releasing of a model checkpoint, or other means that are  
969 appropriate to the research performed.
- 970 • While NeurIPS does not require releasing code, the conference does require all submis-  
971 sions to provide some reasonable avenue for reproducibility, which may depend on the  
972 nature of the contribution. For example
  - 973 (a) If the contribution is primarily a new algorithm, the paper should make it clear how  
974 to reproduce that algorithm.
  - 975 (b) If the contribution is primarily a new model architecture, the paper should describe  
976 the architecture clearly and fully.
  - 977 (c) If the contribution is a new model (e.g., a large language model), then there should  
978 either be a way to access this model for reproducing the results or a way to reproduce  
979 the model (e.g., with an open-source dataset or instructions for how to construct  
980 the dataset).
  - 981 (d) We recognize that reproducibility may be tricky in some cases, in which case  
982 authors are welcome to describe the particular way they provide for reproducibility.  
983 In the case of closed-source models, it may be that access to the model is limited in  
984 some way (e.g., to registered users), but it should be possible for other researchers  
985 to have some path to reproducing or verifying the results.

## 986 5. Open access to data and code

987 Question: Does the paper provide open access to the data and code, with sufficient instruc-  
988 tions to faithfully reproduce the main experimental results, as described in supplemental  
989 material?

990 Answer: [Yes]

991 Justification: We provide access to the data and code we used in the supplementary material.

992 Guidelines:

- 993 • The answer NA means that paper does not include experiments requiring code.
- 994 • Please see the NeurIPS code and data submission guidelines ([https://nips.cc/  
995 public/guides/CodeSubmissionPolicy](https://nips.cc/public/guides/CodeSubmissionPolicy)) for more details.
- 996 • While we encourage the release of code and data, we understand that this might not be  
997 possible, so “No” is an acceptable answer. Papers cannot be rejected simply for not  
998 including code, unless this is central to the contribution (e.g., for a new open-source  
999 benchmark).
- 1000 • The instructions should contain the exact command and environment needed to run to  
1001 reproduce the results. See the NeurIPS code and data submission guidelines ([https://  
1002 //nips.cc/public/guides/CodeSubmissionPolicy](https://nips.cc/public/guides/CodeSubmissionPolicy)) for more details.
- 1003 • The authors should provide instructions on data access and preparation, including how  
1004 to access the raw data, preprocessed data, intermediate data, and generated data, etc.
- 1005 • The authors should provide scripts to reproduce all experimental results for the new  
1006 proposed method and baselines. If only a subset of experiments are reproducible, they  
1007 should state which ones are omitted from the script and why.
- 1008 • At submission time, to preserve anonymity, the authors should release anonymized  
1009 versions (if applicable).

- Providing as much information as possible in supplemental material (appended to the paper) is recommended, but including URLs to data and code is permitted.

## 6. Experimental Setting/Details

Question: Does the paper specify all the training and test details (e.g., data splits, hyperparameters, how they were chosen, type of optimizer, etc.) necessary to understand the results?

Answer: [Yes]

Justification: Datasets and data split are presented in Appendix C.1, and hyperparameters are searched according to the suggestion in the literature. See more details in Appendix D.

Guidelines:

- The answer NA means that the paper does not include experiments.
- The experimental setting should be presented in the core of the paper to a level of detail that is necessary to appreciate the results and make sense of them.
- The full details can be provided either with the code, in appendix, or as supplemental material.

## 7. Experiment Statistical Significance

Question: Does the paper report error bars suitably and correctly defined or other appropriate information about the statistical significance of the experiments?

Answer: [Yes]

Justification: We validate the p-value to support the main claims of this paper.

Guidelines:

- The answer NA means that the paper does not include experiments.
- The authors should answer "Yes" if the results are accompanied by error bars, confidence intervals, or statistical significance tests, at least for the experiments that support the main claims of the paper.
- The factors of variability that the error bars are capturing should be clearly stated (for example, train/test split, initialization, random drawing of some parameter, or overall run with given experimental conditions).
- The method for calculating the error bars should be explained (closed form formula, call to a library function, bootstrap, etc.)
- The assumptions made should be given (e.g., Normally distributed errors).
- It should be clear whether the error bar is the standard deviation or the standard error of the mean.
- It is OK to report 1-sigma error bars, but one should state it. The authors should preferably report a 2-sigma error bar than state that they have a 96% CI, if the hypothesis of Normality of errors is not verified.
- For asymmetric distributions, the authors should be careful not to show in tables or figures symmetric error bars that would yield results that are out of range (e.g. negative error rates).
- If error bars are reported in tables or plots, The authors should explain in the text how they were calculated and reference the corresponding figures or tables in the text.

## 8. Experiments Compute Resources

Question: For each experiment, does the paper provide sufficient information on the computer resources (type of compute workers, memory, time of execution) needed to reproduce the experiments?

Answer: [Yes]

Justification: We conduct all the experiments in PyTorch with a single NVIDIA RTX A5000 (24G) GPU and a 64 AMD EPYC 7543 32-Core Processor CPU. And Detailed time costs are shown in Appendix C.5.

Guidelines:

- The answer NA means that the paper does not include experiments.

- 1061
- The paper should indicate the type of compute workers CPU or GPU, internal cluster, or cloud provider, including relevant memory and storage.
  - The paper should provide the amount of compute required for each of the individual experimental runs as well as estimate the total compute.
  - The paper should disclose whether the full research project required more compute than the experiments reported in the paper (e.g., preliminary or failed experiments that didn't make it into the paper).

1068 **9. Code Of Ethics**

1069 Question: Does the research conducted in the paper conform, in every respect, with the  
1070 NeurIPS Code of Ethics <https://neurips.cc/public/EthicsGuidelines>?

1071 Answer: [Yes]

1072 Justification: The research adheres to all ethical guidelines outlined by NeurIPS. Specifically,  
1073 we have ensured that our data collection methods are ethical, our experiments are conducted  
1074 responsibly, and all potential biases are addressed. Additionally, we have considered the  
1075 broader impacts of our work and have taken steps to mitigate any negative consequences.

1076 Guidelines:

- The answer NA means that the authors have not reviewed the NeurIPS Code of Ethics.
- If the authors answer No, they should explain the special circumstances that require a deviation from the Code of Ethics.
- The authors should make sure to preserve anonymity (e.g., if there is a special consideration due to laws or regulations in their jurisdiction).

1082 **10. Broader Impacts**

1083 Question: Does the paper discuss both potential positive societal impacts and negative  
1084 societal impacts of the work performed?

1085 Answer: [Yes]

1086 Justification: We consider both the potential societal impacts and negative societal impacts,  
1087 and also discuss possible mitigation strategies. Details are shown in Appendix E.

1088 Guidelines:

- The answer NA means that there is no societal impact of the work performed.
- If the authors answer NA or No, they should explain why their work has no societal impact or why the paper does not address societal impact.
- Examples of negative societal impacts include potential malicious or unintended uses (e.g., disinformation, generating fake profiles, surveillance), fairness considerations (e.g., deployment of technologies that could make decisions that unfairly impact specific groups), privacy considerations, and security considerations.
- The conference expects that many papers will be foundational research and not tied to particular applications, let alone deployments. However, if there is a direct path to any negative applications, the authors should point it out. For example, it is legitimate to point out that an improvement in the quality of generative models could be used to generate deepfakes for disinformation. On the other hand, it is not needed to point out that a generic algorithm for optimizing neural networks could enable people to train models that generate Deepfakes faster.
- The authors should consider possible harms that could arise when the technology is being used as intended and functioning correctly, harms that could arise when the technology is being used as intended but gives incorrect results, and harms following from (intentional or unintentional) misuse of the technology.
- If there are negative societal impacts, the authors could also discuss possible mitigation strategies (e.g., gated release of models, providing defenses in addition to attacks, mechanisms for monitoring misuse, mechanisms to monitor how a system learns from feedback over time, improving the efficiency and accessibility of ML).

1111 **11. Safeguards**

1112 Question: Does the paper describe safeguards that have been put in place for responsible  
1113 release of data or models that have a high risk for misuse (e.g., pretrained language models,  
1114 image generators, or scraped datasets)?

1115  
1116  
1117  
1118  
1119  
1120  
1121  
1122  
1123  
1124  
1125  
1126  
1127  
1128  
1129  
1130  
1131  
1132  
1133  
1134  
1135  
1136  
1137  
1138  
1139  
1140  
1141  
1142  
1143  
1144  
1145  
1146  
1147  
1148  
1149  
1150  
1151  
1152  
1153  
1154  
1155  
1156  
1157  
1158  
1159  
1160  
1161  
1162  
1163  
1164  
1165

Answer: [NA]

Justification: The paper poses no such risks.

Guidelines:

- The answer NA means that the paper poses no such risks.
- Released models that have a high risk for misuse or dual-use should be released with necessary safeguards to allow for controlled use of the model, for example by requiring that users adhere to usage guidelines or restrictions to access the model or implementing safety filters.
- Datasets that have been scraped from the Internet could pose safety risks. The authors should describe how they avoided releasing unsafe images.
- We recognize that providing effective safeguards is challenging, and many papers do not require this, but we encourage authors to take this into account and make a best faith effort.

## 12. Licenses for existing assets

Question: Are the creators or original owners of assets (e.g., code, data, models), used in the paper, properly credited and are the license and terms of use explicitly mentioned and properly respected?

Answer: [Yes]

Justification: We incorporate six datasets, including four datasets from the Amazon platform[1](Books, Movies & TV, Video Games, and Industrial), Movielens-1M[59], and Book Crossing[60], all of which are open-source. The backend language models used in our research are BERT [4], RoBERTa [5], Llama2-7B [6], Mistral-7B [69], text-embedding-ada-v2 & text-embeddings-3-large [70], and SFR-Embedding-Mistral [71].

Guidelines:

- The answer NA means that the paper does not use existing assets.
- The authors should cite the original paper that produced the code package or dataset.
- The authors should state which version of the asset is used and, if possible, include a URL.
- The name of the license (e.g., CC-BY 4.0) should be included for each asset.
- For scraped data from a particular source (e.g., website), the copyright and terms of service of that source should be provided.
- If assets are released, the license, copyright information, and terms of use in the package should be provided. For popular datasets, [paperswithcode.com/datasets](https://paperswithcode.com/datasets) has curated licenses for some datasets. Their licensing guide can help determine the license of a dataset.
- For existing datasets that are re-packaged, both the original license and the license of the derived asset (if it has changed) should be provided.
- If this information is not available online, the authors are encouraged to reach out to the asset's creators.

## 13. New Assets

Question: Are new assets introduced in the paper well documented and is the documentation provided alongside the assets?

Answer: [NA]

Justification: This paper does not release new assets.

Guidelines:

- The answer NA means that the paper does not release new assets.
- Researchers should communicate the details of the dataset/code/model as part of their submissions via structured templates. This includes details about training, license, limitations, etc.
- The paper should discuss whether and how consent was obtained from people whose asset is used.

- 1166                   • At submission time, remember to anonymize your assets (if applicable). You can either  
1167                   create an anonymized URL or include an anonymized zip file.

1168 **14. Crowdsourcing and Research with Human Subjects**

1169                   Question: For crowdsourcing experiments and research with human subjects, does the paper  
1170                   include the full text of instructions given to participants and screenshots, if applicable, as  
1171                   well as details about compensation (if any)?

1172                   Answer: [NA]

1173                   Justification: The paper does not involve crowdsourcing nor research with human subjects.

1174                   Guidelines:

- 1175                   • The answer NA means that the paper does not involve crowdsourcing nor research with  
1176                   human subjects.
- 1177                   • Including this information in the supplemental material is fine, but if the main contribu-  
1178                   tion of the paper involves human subjects, then as much detail as possible should be  
1179                   included in the main paper.
- 1180                   • According to the NeurIPS Code of Ethics, workers involved in data collection, curation,  
1181                   or other labor should be paid at least the minimum wage in the country of the data  
1182                   collector.

1183 **15. Institutional Review Board (IRB) Approvals or Equivalent for Research with Human  
1184                   Subjects**

1185                   Question: Does the paper describe potential risks incurred by study participants, whether  
1186                   such risks were disclosed to the subjects, and whether Institutional Review Board (IRB)  
1187                   approvals (or an equivalent approval/review based on the requirements of your country or  
1188                   institution) were obtained?

1189                   Answer: [NA]

1190                   Justification: The paper does not involve crowdsourcing nor research with human subjects.

1191                   Guidelines:

- 1192                   • The answer NA means that the paper does not involve crowdsourcing nor research with  
1193                   human subjects.
- 1194                   • Depending on the country in which research is conducted, IRB approval (or equivalent)  
1195                   may be required for any human subjects research. If you obtained IRB approval, you  
1196                   should clearly state this in the paper.
- 1197                   • We recognize that the procedures for this may vary significantly between institutions  
1198                   and locations, and we expect authors to adhere to the NeurIPS Code of Ethics and the  
1199                   guidelines for their institution.
- 1200                   • For initial submissions, do not include any information that would break anonymity (if  
1201                   applicable), such as the institution conducting the review.