

# ClusterRAG: Cluster-Based Collaborative Filtering for Personalized Retrieval-Augmented Generation

Anonymous ACL submission

## Abstract

Personalized Retrieval-Augmented Generation (RAG) relies on accurately selecting user-relevant documents. In practice, existing approaches often suffer from high retrieval costs and overlook that collaborative signals from similar users can enhance personalized generation for the current user. We propose **ClusterRAG**, a **Cluster**-Based Collaborative Filtering for Personalized **R**etrieval-**A**ugmented **G**eneration. ClusterRAG represents users through their profile documents, organizes users into semantically coherent clusters using density-based clustering, and performs retrieval at both the cluster and document levels via cluster-level similarity and fine-grained ranking. Extensive experiments on the LaMP benchmark demonstrate that jointly leveraging the target user’s profile and profiles from top similar users consistently yields the best performance across diverse tasks. Further analysis shows that ClusterRAG integrates seamlessly with different dense retrievers and rankers, and remains effective when paired with both fine-tuned and zero-shot language models.

## 1 Introduction

Retrieval-Augmented Generation (RAG) has emerged as a powerful paradigm for knowledge-intensive language tasks by combining parametric knowledge in large language models (LLMs) with non-parametric retrieval over external documents, significantly reducing hallucinations and improving factuality in text generation (Lewis et al., 2020). RAG systems typically retrieve documents related to the immediate query, then condition a generative model on those documents to produce responses (Fan et al., 2024; Huang and Huang, 2024; Li et al., 2025b). Despite impressive gains, current RAG pipelines often ignore long-term user information and inter-user relationships when constructing retrieval contexts, limiting personalization and the

ability to leverage analogous users’ knowledge for improved generation quality.

Personalization is crucial in many real-world applications (e.g., personal assistants, tutoring systems, and personalized search) because user history and preferences strongly influence what information is relevant and how it should be framed (Li et al., 2025a; Ahmad et al., 2025). Existing personalization strategies in RAG largely fall into two extremes: (1) user-only approaches that condition retrieval and prompts solely on a user’s own profile, which can be sparse or noisy (Salemi et al., 2024b; Zerhoubi and Granitzer, 2024; Dong et al., 2025), and (2) non-personalized approaches that ignore user history altogether (Lewis et al., 2020; Asai et al., 2024; Yang et al., 2025; Zhang, 2025). Both approaches miss a middle ground where signals from similar users can enrich prompts while preserving user-specific nuance. Recent surveys (Xu et al., 2025; Li et al., 2025a) highlight the potential of end-to-end personalization across the RAG pipeline, from pre-retrieval user modeling to retrieval and generation, but also emphasize practical challenges such as balancing personal vs. collaborative signals.

Collaborative filtering (CF), long established in recommender systems (Xue et al., 2017; Wang et al., 2019; Sgardelis et al., 2025; Shi et al., 2025), naturally complements personalization by exploiting similarities between users to infer missing preferences or relevant content. However, direct application of CF to RAG introduces new questions: (1) *how should users be represented for retrieval tasks?* (2) *how to retrieve similar users to capture heterogeneous behavior at scale?* and (3) *how to leverage both collaborative documents and target user’s profile when forming prompts for LLMs?*

In this paper, we propose **Cluster**-Based Collaborative Filtering for Personalized **R**etrieval-**A**ugmented **G**eneration (**ClusterRAG**), a practical pipeline that (1) constructs compact user repre-

083 presentations by aggregating each user’s profile doc- 131  
084 uments into embeddings, (2) groups users into 132  
085 clusters using HDBSCAN (McInnes et al., 2017) 133  
086 to reveal cohorts of similar users and builds a 134  
087 cluster-level ranking matrix by scoring intra-cluster 135  
088 user similarities with effective rankers, e.g., Col- 136  
089 BERT (Khattab and Zaharia, 2020), and (3) clusters 137  
090 and retrieves candidate profile documents from the 138  
091 top  $k$  similar users to form collaborative, user-only, 139  
092 or hybrid prompts for downstream generation. Our 140  
093 method explicitly leverages cluster structure to re- 141  
094 duce search complexity, provide robust neighbor se- 142  
095 lection in variable-density settings, and enable prin- 143  
096 ciple mixing of collaborative and individual sig- 144  
097 nals when constructing prompts. We evaluate multi- 145  
098 ple retrievers (ColBERT, Contriever (Izacard et al., 146  
099 2022), BGE (Xiao et al., 2024), BM25 (Robert- 147  
100 son et al., 1995), Recency and Random). Beyond 148  
101 methodological novelty, we also demonstrate that 149  
102 ClusterRAG is model-agnostic and robust across 150  
103 architectures and retrieval backbones. Cluster- 151  
104 RAG integrates seamlessly with both fine-tuned 152  
105 sequence-to-sequence (seq2seq) encoder-decoder 153  
106 language models and zero-shot LLMs, without re- 154  
107 quiring any model-specific adaptation. 155

108 We validate ClusterRAG on the LaMP bench- 156  
109 mark (Salemi et al., 2024b) and report improve- 157  
110 ments in personalized generation quality compared 158  
111 to non-personalized and naive user-only baselines. 159  
112 We also provide a link to an anonymous project 160  
113 repository for ClusterRAG <sup>1</sup>. 161

114 The remainder of this paper unfolds as follows. 162  
115 The next section presents related work; Section 3 163  
116 provides problem formulation; Section 4 describes 164  
117 the ClusterRAG framework; Section 5 presents the 165  
118 experimental evaluation; and Section 6 presents the 166  
119 conclusion, limitations, and ethical considerations. 167  
120 Additional analysis and results are provided in the 168  
121 Appendix. 169

## 122 2 Related Work

123 **Retrieval-Augmented Generation (RAG).** The 172  
124 adoption of RAG has demonstrated improvements 173  
125 across a range of tasks, including question an- 174  
126 swering, dialogue comprehension, and code gen- 175  
127 eration (Lewis et al., 2020; Xu et al., 2023; Zer- 176  
128 houdi and Granitzer, 2024; Fan et al., 2024). Early 177  
129 RAG pipelines typically index a shared corpus (e.g., 178  
130 Wikipedia or domain corpora) and retrieve pas- 179

180 sages conditioned only on the current query; the 181  
182 retrieved passages are then used to condition an 182  
LLM at generation time (Lewis et al., 2020; Zhang, 183  
2025). Subsequent work has explored numerous 184  
improvements to retriever architectures, reranking 185  
strategies, and retrieval-generation coupling mech- 186  
anisms (Gao et al., 2023; Siriwardhana et al., 2023; 187  
Fan et al., 2024). Even though these works es- 188  
tablish strong task-agnostic baselines and sophisti- 189  
cated retriever-generator interfaces, they typically 190  
operate in a *user-agnostic* manner and do not lever- 191  
age a user’s long-term profile or cross-user signals 192  
when selecting documents for conditioning. 193

194 **Personalized Retrieval-Augmented Genera- 194  
195 tion.** A growing body of work studies how RAG 195  
196 can be adapted for personalized applications (e.g., 196  
197 personal assistants, tutoring, and individualized 197  
198 question answering) by incorporating user history, 198  
199 preferences, or authoring signals into retrieval and 199  
200 prompting (Salemi et al., 2024b; Zerhoudi and 200  
201 Granitzer, 2024; Li et al., 2025a). Some systems 201  
202 personalize LLMs by (1) fine-tuning model pa- 202  
203 rameters (either fully or selectively) for individual 203  
204 users (Li and Liang, 2021; Hu et al., 2022; Zollo 204  
205 et al., 2025), (2) incorporating latent user repre- 205  
206 sentations into the model (Ning et al., 2025; Qiu 206  
207 et al., 2025; Huber et al., 2025), and (3) augmenting 207  
208 model prompts with a user profile or a small set of 208  
209 user documents (Zamani et al., 2022; Salemi et al., 209  
210 2024b,a). The first two strategies require modify- 210  
211 ing the model’s architecture or parameters, which 211  
212 can be expensive, or in some cases infeasible, due 212  
213 to storage, computational, and time constraints. In 213  
214 addition, they cannot perform well for cold-start 214  
215 users. In contrast, the third approach, which is 215  
216 adopted in this work, can be applied to any gener- 216  
217 ative model (Salemi et al., 2024a). User-specific 217  
218 personalization works well when profiles are dense 218  
219 and representative; this notwithstanding, user-only 219  
220 personalization suffers when profiles are sparse, 220  
221 noisy, or unrepresentative of the current intent. 221

222 **Collaborative Personalized Retrieval Aug- 222  
223 mented Generation.** Collaborative filtering has 223  
224 been extensively studied in recommender systems 224  
225 and consistently shown to be effective (Xue et al., 225  
226 2017; Wang et al., 2019; Zhang et al., 2024; Shen 226  
227 et al., 2024; Shi et al., 2024; Tang et al., 2025; Xin 227  
228 et al., 2025; Zhang et al., 2025). The core premise 228  
229 is that users with comparable interaction histories 229  
230 tend to exhibit similar preferences; therefore, lever- 230  
231 aging items favored by similar users can help gener- 231  
232 ate relevant recommendations for a target user. 232

<sup>1</sup>[https://github.com/academicprojects44/  
anonymous](https://github.com/academicprojects44/anonymous)

Recent efforts have started to marry CF and retrieval for generation (Shi et al., 2025; Zhu et al., 2025). For instance, Shi et al. (2025) employs contrastive learning to generate user embeddings that retrieve similar users and incorporate collaborative signals with a user input for prompt creation.

Despite this progress, two practical challenges remain unresolved in current collaborative RAG research. First, naively computing pairwise similarities across millions of users is costly; clustering users into cohorts can reduce search complexity. ClusterRAG is designed to address this issue by combining user-level clustering with document-level collaborative retrieval and introducing cluster-level ranking matrices that summarize intra-cluster user similarity and thereby enabling robust neighbor selection even in variable-density cohorts. Second, once neighbor users are found, selecting which of their documents to include and how to merge collaborative documents with a target user’s own profile remains an open design choice with direct impact on generation quality. ClusterRAG uses flexible prompt fusion modes and evaluates multiple retrievers, showing empirical robustness across both fine-tuned and training-free generative models.

To the best of our knowledge, ClusterRAG is the first framework to integrate user-level clustering with collaborative document retrieval for personalized RAG, explicitly leveraging cross-user similarity to enrich sparse user profiles for personalized generation.

### 3 Problem Formulation

A standard RAG setup consists of two core components, retrieval and generation: given an input query  $x$ , the model predicts the most probable output sequence  $y$  conditioned on  $x$  and a retrieved document  $d$ . Personalized RAG extends this formulation by conditioning generation on a user  $u$ , typically represented through a user profile. Formally, we let  $T = \{(u_1, x_1, y_1), (u_2, x_2, y_2), \dots, (u_N, x_N, y_N)\}$  denote a set of  $N$  training instances, where each tuple consists of a user  $u$ , a user-issued input query  $x$ , and a corresponding personalized ground-truth output  $y$ . For each user  $u$ , a user profile  $U_p$  is available and serves as an auxiliary context for personalized generation. The profile  $U_p = \{d_1, d_2, \dots, d_n\}$  is a collection of personal documents or historical records associated with  $u$ , such as past queries and

generated outputs.

In ClusterRAG, given a target user  $u$ , our objective is to identify the top  $k$  most similar users using a clustering-assisted ranking strategy. We then retrieve and rank the top  $m$  documents from these users using a retriever (ranker)  $R$ .

## 4 ClusterRAG Framework

ClusterRAG is built on the intuition that users with similar behaviors and preferences can provide valuable contextual signals for one another. By combining information from a target user’s own profile with carefully selected profiles from similar users, ClusterRAG enhances an LLM’s ability to generate accurate and personalized responses. As shown in Figure 1, ClusterRAG framework consists of three main stages: (1) user representation and retrieval, (2) profile retrieval, and (3) personalized generation.

### 4.1 User Representation and Retrieval

Since explicit user representations are typically unavailable, we first construct user embeddings from observed interaction data. All data instances associated with a user are aggregated into a user-level profile  $U_p$ . Each document  $d_i \in U_p$  is encoded using a dense embedding model, specifically ColBERTv2 (Santhanam et al., 2022). A compact user representation is then obtained by averaging document embeddings:

$$\mathbf{z}_u = \frac{1}{n_u} \sum_{i=1}^{n_u} f(d_i), \quad (1)$$

where  $f(\cdot)$  denotes the embedding function.

Computing similarities between all user pairs is prohibitively expensive at scale. To address this, ClusterRAG groups users into similarity-based cohorts using a hierarchical density-based clustering method, HDBSCAN (McInnes et al., 2017), which automatically identifies clusters of varying density, making it well-suited for collaborative filtering in our setting, where the number of user groups in the user-document collection is unknown a priori.

Clustering alone does not quantify the relative similarity of users within each cluster. Therefore, we compute intra-cluster similarities using a modern reranker, ColBERTv2 (Santhanam et al., 2022), which provides fine-grained token-level interactions inherited from ColBERT (Khattab and Zaharia, 2020) while incorporating residual compression and denoised supervision for improved effi-

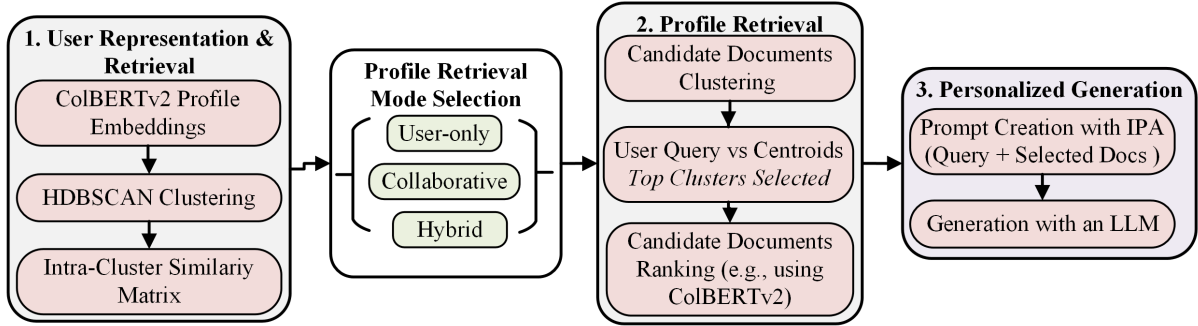


Figure 1: Overview of the ClusterRAG framework.

ciency and generalization. These properties make ColBERTv2 well-suited for robust similarity estimation between user profiles.

**Cluster-Level Similarity Ranking.** This step aims at restricting similarity computation to cluster members to improve robustness and scalability by focusing comparisons on behaviorally consistent cohorts. For each cluster  $C$ , we construct an intra-cluster similarity matrix  $R^C$  defined as:

$$R_{u,v}^C = \text{ColBERTv2}(\mathbf{z}_u, \mathbf{z}_v), \quad (2)$$

where  $u, v \in C$  and  $v \neq u$ . The diagonal entries in  $R^C$  correspond to self-similarity and are discarded. For each user  $u$ , we finally retain an ordered list of the top  $k$  most similar users within the same cluster.

## 4.2 Profile Retrieval

The profile retrieval stage integrates search and ranking to identify documents that are most beneficial for personalized generation (Huang and Huang, 2024). Incorporating collaborative filtering introduces two key challenges: selecting relevant documents from similar users and effectively leveraging both collaborative documents and the target user’s documents for personalized RAG. ClusterRAG addresses these challenges by leveraging cluster structure to provide both topical coherence and retrieval efficiency and organizes profile retrieval as follows.

**Profile Retrieval Modes.** First, given a query  $q$  from user  $u$ , we retrieve candidate profile documents using one of the following three retrieval modes. (1) *User-only retrieval*: the simplest strategy, which considers only the user’s own profile. (2) *Collaborative retrieval*: this mode retrieves documents from the profiles of the top  $k$  most similar users. It is particularly beneficial for sparse or cold-start users, whose own profiles may not adequately

capture the intent of the current query. (3) *Hybrid retrieval*: this mode combines both user-only and collaborative profiles. As a result, the effective profile  $U_p$  used for generation may consist of: (i) documents from the target user only, (ii) documents from similar users only, or (iii) similar users alone, or (iii) a combination of documents from both sources.

**Clustering for Topical Organization.** After selecting a profile retrieval mode, all candidate profile documents are encoded using a dense retriever, such as ColBERTv2, and partitioned into clusters using HDBSCAN, producing a set of clusters  $\mathcal{C} = \{C_1, \dots, C_K\}$  with corresponding computed centroids  $\mathcal{M} = \{\mu_1, \dots, \mu_K\}$ . This clustering captures latent topical structure and enables ClusterRAG to automatically infer the number of topics present in a user’s profile without requiring prior specification.

**Cluster-Level Indexing and Retrieval.** Finally, ClusterRAG employs a two-stage retrieval strategy. First, a cluster index stores centroid embeddings: given a query  $q$ , its embedding  $e_q$  is computed similarly as document embeddings, and compared against all centroids, and the top  $B$  clusters are selected from  $\mathcal{C}$ . Second, within each selected cluster, documents are retrieved and reranked by similarity to  $e_q$ , and the top  $m$  documents are selected for generation.

This hierarchical retrieval reduces complexity from  $\mathcal{O}(N)$  to  $\mathcal{O}(K + B \cdot N/K)$ , where  $N$  is the total number of documents, each cluster  $C_i$  contains  $\approx |N/K|$  documents, and  $B \leq K \leq N$ .

We use ColBERTv2 as the primary model to retrieve and rerank profile documents; even so, our experiments evaluate additional dense, sparse, heuristic, and random retrievers to demonstrate the framework’s retriever-agnostic design.

### 4.3 Personalized Generation

Effective generation by LLMs critically depends on well-engineered prompts that are tailored to the downstream task. Prompts allow seamless integration of pre-trained models into downstream tasks by eliciting desired model behaviors solely based on the given prompt (Sahoo et al., 2025). To effectively leverage selected user documents from  $U_p$ , ClusterRAG adopts *In-Prompt Augmentation (IPA)* (Salemi et al., 2024b), which integrates the user query with all relevant retrieved documents directly within the prompt. IPA is particularly well suited for ClusterRAG, as it can be applied to both training-free (zero-shot) and fine-tuned settings and is compatible with a wide range of model architectures. Accordingly, ClusterRAG supports both fine-tuned LLMs and zero-shot generative models and can be combined with a variety of state-of-the-art document retrievers.

Specifically, to balance user profile context and query specificity, given a maximum prompt length  $L_{\max}$ , the allocated profile length is computed as:

$$|U_p| = \mathcal{G}_t(L_{\max} - \min(|q|, \lfloor \gamma L_{\max} \rfloor)), \quad (3)$$

where  $\gamma \in [0, 1]$  is a tunable mixing parameter,  $|q|$  is query length, and  $\mathcal{G}_t(\cdot)$  is a task-specific prompt generator (we provide detailed task-specific prompt generators in Section 5.2 and Appendix A). This formulation allows ClusterRAG to incorporate strong personalization signals while preserving the relevance of the current query.

## 5 Experiments

In this section, we present the experimental setup employed with ClusterRAG, including the datasets, baseline models, evaluation metrics, implementation details, and experimental results. We also provide an overview of the prompts used in our experiments.

### 5.1 Experimental Setup

**Datasets.** Our experiments employ the LaMP benchmark (Salemi et al., 2024b), a publicly available dataset that covers a broad range of personalized text generation tasks. The benchmark consists of three personalized text classification tasks and four personalized text generation tasks. One of the four text generation tasks, **LaMP-6** (Personalized Email Subject Generation), is excluded in this work because its data is not publicly available. Specifically, the remaining tasks include: **LaMP-1:**

Task	#users	#train	#dev	#test
LaMP-1	6542	6542	1500	1500
LaMP-2	929	5073	1410	1557
LaMP-3	20000	20000	2500	2500
LaMP-4	1643	12500	1500	1800
LaMP-5	14682	14682	1500	1500
LaMP-7	13437	13437	1498	1500

Table 1: Statistics of the LaMP benchmark with time-based data split.

**Personalized Citation Identification**, formulated as a binary classification task; **LaMP-2: Personalized Movie Tagging**, a 15-class categorical classification task; **LaMP-3: Personalized Product Rating**, an ordinal classification task predicting ratings from one to five stars for e-commerce products; **LaMP-4: Personalized News Headline Generation**; **LaMP-5: Personalized Scholarly Title Generation**; and **LaMP-7: Personalized Tweet Paraphrasing**. ClusterRAG experiments follow the time-based LaMP split to partition the data into training, validation, and test sets. We provide statistics of the dataset in Table 1, detailed dataset statistics in Table 6 in Appendix B, and detailed task descriptions in Appendix C. We additionally provide dataset licensing information in Appendix B.

**Evaluation Metrics.** Following previous work (Salemi et al., 2024b,a; Shi et al., 2025), we evaluate LaMP-1 and LaMP-2 using Accuracy and F1-measure, and LaMP-3 using Mean Absolute Error (MAE) and Root Mean Squared Error (RMSE). We evaluate text generation performance on LaMP-4, LaMP-5, and LaMP-7 using ROUGE-1 (R-1) and ROUGE-L (R-L) (Lin, 2004).

**Baseline Models.** We firstly compare ClusterRAG to **no personalization** and call this **vanillaRAG**. In this baseline, the generative model is presented with the original task’s input without any profile documents to assess whether personalization improves the model effectiveness. Then we consider personalized baselines: (1) **User-only** models, which include (a) **ROPG** (Salemi et al., 2024a): it optimizes the dense retrieval model based on the results generated by a LLM and (b) **LaMP-IPA**: it was introduced with the LaMP benchmark, and it uses in-prompt augmentation for prompt creation; (2) **Collaborative** baseline, **CFRAG** (Shi et al., 2025), that uses contrastive learning to find similar users. To the best of our

441 knowledge, CFRAG is the only existing collabora- 490  
442 tive work for personalized RAG; therefore, we com- 491  
443 pare these baselines with three versions of Cluster- 492  
444 RAG: user-only, collaborative, and hybrid profile 493  
445 retrieval modes. 494

446 **Implementation details.** We implement Cluster- 495  
447 RAG using the HuggingFace transformers frame- 496  
448 work (Wolf et al., 2020) and the PyTorch library 497  
449 (Paszke et al., 2019). ClusterRAG adopts a fine-  
450 tuned FlanT5-base (Chung et al., 2024) for gener-  
451 ation; unless explicitly stated otherwise (in ex-  
452 periments with zero-shot LLMs), it uses a causal  
453 Qwen2-7B-Instruct (Yang et al., 2024) or a seq2seq  
454 FlanT5-XXL (Chung et al., 2024); all models  
455 are open-source. Training is performed using the  
456 Trainer or Seq2SeqTrainer APIs<sup>2</sup>, depending on  
457 whether the LLM backbone is a causal or seq2seq  
458 model. FlanT5-base has 250M parameters; Qwen2-  
459 7B-Instruct uses 7.07B parameters; and FlanT5-  
460 XXL has 11B parameters.

461 We use AdamW (Loshchilov and Hutter, 2019)  
462 optimization with a learning rate of  $5 \times 10^{-5}$ ,  
463 weight decay of  $10^{-4}$ , linear learning rate schedul-  
464 ing, and a warm-up ratio of 0.05. Models are  
465 trained for up to 30 epochs, with evaluation and  
466 checkpointing conducted at the end of each epoch.  
467 The maximum prompt and output lengths ( $L_{\max}$   
468 and  $|\bar{y}|$ ) of generative models is set to 512 and 128  
469 tokens, respectively. Maximum sequence length  
470 for embedding models is set to 256 tokens. We set  
471 the number of collaborative (similar) users,  $k$ , to  
472 1 and retrieved profile documents,  $m$ , to 2.  $\gamma$  for  
473 prompt formulation in Equation 3 is set to 0.55,  
474 while batch size is set to 16. Beam search (Freitag  
475 and Al-Onaizan, 2017) with a beam size of 4 is  
476 employed for text generation. All hyperparam-  
477 eters for the baseline models are searched according  
478 to the settings in the original papers. We select  
479 the optimal hyperparameters for ClusterRAG via  
480 grid search and report the corresponding tuning  
481 grid in Table 7 in Appendix D. All experiments are  
482 conducted on a Quadro RTX 8000 GPUs, 48 GB  
483 VRAM for a range of 10-24 hours per experiment  
484 depending on the task.

## 485 5.2 Prompts Used in ClusterRAG

486 This subsection presents the prompt templates em-  
487 ployed during generation. Each prompt contains  
488 an instruction, input, and profile. We provide the  
489 template for LaMP-1 only below; templates for

<sup>2</sup><https://github.com/huggingface/transformers>

other tasks are presented in Appendix A. In the 490  
491 template,  $\{Paper\ abstract\}$  and  $\{Reference\ list\}$ ,  
492  $\{Movie\ description\}$  represent user input for the  
493 corresponding LaMP tasks, while  $\{Paper\ abstract\}$   
494 represent user profile entries. The remaining text  
495 is the instruction guiding an LLM to generate the  
496 intended output.

**LaMP-1 Prompt Template:** Given an author  
who has previously written papers  $\{Paper\ list\}$   
and now has written  $\{Paper\ abstract\}$ . Which  
reference below is related? Just answer with  
[1] or [2] without explanation.  $\{Reference\ list\}$

## 497 5.3 Experimental Results 498

499 When reporting experimental results, we identify 500  
501 statistically significant differences of ClusterRAG 502  
503 performance using a two-tailed paired t-test for 504  
505 generation and ordinal classification evaluation 506  
507 (ROUGE-1, ROUGE-L, MAE, and RMSE) and 508  
509 McNemar test for categorical text classification 509  
510 evaluation (Accuracy and F1). We report results 510  
511 comparing ClusterRAG with baselines, analyzing 511  
512 its retriever-agnostic design, language model versa-  
512 tility, and ablation study. In all tables, the symbol  
512  $\uparrow$  indicates that higher values are better, while the  
512 symbol  $\downarrow$  indicates that lower values are better; all  
512 results presented are obtained from a single experi-  
512 mental run.

### 513 5.3.1 Comparison with Baselines

514 Comparison results are presented in Table 2, which 514  
515 shows that ClusterRAG consistently outperforms 515  
516 all baselines across the LaMP benchmark. Import- 516  
517 tantly, the hybrid variant (*ClusterRAG-H*) achieves 517  
518 the best performance on every task. These im- 518  
519 provements indicate that retrieving and aggregat- 519  
520 ing documents from similar users substantially en- 520  
521 hances RAG, providing more relevant and person- 521  
522 alized evidence than standard RAG pipelines. No- 522  
523 tably, ClusterRAG achieves strong performance 523  
524 using only two profile documents, whereas base- 524  
525 line methods require at least four documents to 525  
526 reach their optimal results, indicating that Cluster- 526  
527 RAG is well suited for low-resource personaliza- 527  
528 tion settings. While *ClusterRAG-C* (collaborative) 528  
529 and *ClusterRAG-U* (user-only) achieve competitive 529  
530 second-best results on several tasks, their combina- 530  
531 tion in the hybrid model yields the most robust and 531  
532 consistent gains across diverse task settings. 532

Models	LaMP-1		LaMP-2		LaMP-3		LaMP-4		LaMP-5		LaMP-7	
	Acc.↑	F1↑	Acc.↑	F1↑	MAE↓	RMSE↓	R-1↑	R-L↑	R-1↑	R-L↑	R-1↑	R-L↑
vanillaRAG	0.630	0.630	0.520	0.440	0.371	0.709	0.171	0.154	0.462	0.413	0.310	0.273
LaMP-IPA	<u>0.674</u>	<u>0.664</u>	<u>0.570</u>	<u>0.522</u>	<u>0.289</u>	<u>0.608</u>	0.175	0.169	0.472	0.423	<u>0.508</u>	<u>0.457</u>
ROPG	0.644	0.322	0.468	0.031	0.346	0.692	<u>0.184</u>	<u>0.163</u>	0.464	0.396	0.353	0.288
CFRAG	0.633	0.327	0.534	0.036	0.354	0.707	0.162	0.141	<u>0.473</u>	<u>0.425</u>	0.375	0.306
ClusterRAG-C	0.674*	0.673*	0.644	0.607	0.284	0.624	0.179	0.157	0.480*	0.430*	0.507	0.454
ClusterRAG-U	0.645	0.645	0.649*	0.612*	0.271*	0.599*	0.184*	0.165*	0.475	0.425	0.514*	0.464*
<b>ClusterRAG-H</b>	<b>0.690</b>	<b>0.690</b>	<b>0.661</b>	<b>0.620</b>	<b>0.270</b>	<b>0.594</b>	<b>0.190</b>	<b>0.176</b>	<b>0.490</b>	<b>0.440</b>	<b>0.521</b>	<b>0.470</b>

Table 2: Comparison of the performance of ClusterRAG with baselines on the LaMP benchmark. Best results are shown in **bold** and second-best in underlined; boldface indicates statistically significant improvements over the second-best ( $p < 0.05$ ). The symbol \* denotes the second-best ClusterRAG variant.

### 5.3.2 ClusterRAG Retriever-Agnostic Design

In addition to *ColBERTv2*, ClusterRAG explores five more retrievers: (1) a dense unsupervised dual-encoder retriever, *Contriever* (Izacard et al., 2022), (2) a fine-tuned multilingual dense retriever optimized for semantic similarity and retrieval tasks, *BGE* (Xiao et al., 2024), (3) a classical sparse lexical retriever based on term frequency-inverse document frequency (TF-IDF), *BM25* (Robertson et al., 1995), (4) a heuristic retriever that ranks documents solely based on temporal proximity to the query time, favoring the most recently published documents, *Recency*, and (5) a non-informative baseline that samples documents uniformly at random, *Random*. We provide retriever-agnostic design results in Table 3 for LaMP-(1,2,7) and Figure 2 for LaMP-5. The table demonstrates that ClusterRAG consistently benefits from stronger retrievers, with dense semantic models outperforming sparse and heuristic baselines across all LaMP tasks. *ColBERTv2* (*ColBERT* or *Col*) achieves the best overall performance on all tasks, highlighting the advantage of late-interaction matching in personalized retrieval. *BGE* and *Contriever* (*Con*) provide competitive performance, confirming the retriever-agnostic nature of ClusterRAG, while *BM25* and *Recency* (*Rec*) offer modest gains over *Random* (*Ran*) but lag behind dense methods. These results indicate that ClusterRAG is robust across retrieval paradigms, yet most effectively leverages high-capacity dense retrievers to maximize personalization and generation quality.

### 5.3.3 ClusterRAG LLM Versatility

Table 4 reports the performance of ClusterRAG on LaMP-(1,2,5) when paired with zero-shot LLMs: FlanT5-XXL and Qwen2-7B-Instruct. For each LLM, we compare a non-personalized variant

Retrievers	LaMP-1		LaMP-2		LaMP-7	
	Acc.↑	F1↑	Acc.↑	F1↑	R-1↑	R-L↑
Random	0.640	0.639	0.609	0.608	0.500	0.449
Recency	0.659	0.650	0.618	0.610	0.507	0.456
BM25	0.662	0.658	0.629	0.621	0.510	0.460
Contriever	0.681	0.681	0.649	<u>0.623</u>	<u>0.511</u>	0.459
BGE	<u>0.684</u>	<u>0.682</u>	<u>0.658</u>	0.613	0.509	<u>0.461</u>
<b>ColBERT</b>	<b>0.690</b>	<b>0.690</b>	<b>0.661</b>	<b>0.620</b>	<b>0.521</b>	<b>0.470</b>

Table 3: Comparison of the performance of ClusterRAG under different retrievers on LaMP-(1,2,7).

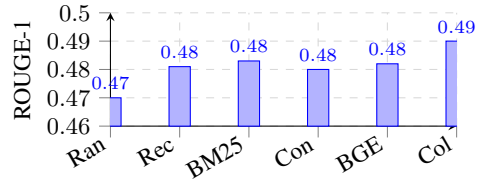


Figure 2: Retrievers' ROUGE-1 scores on LaMP-5.

(*nFlan*, *nQwen2*) against its personalized counterpart (*pFlan*, *pQwen2*). As shown in Table 4, personalized variants consistently outperform their non-personalized counterparts across all tasks, demonstrating the effectiveness of ClusterRAG in injecting user-specific and collaborative signals into generation. Notably, *pFlan* achieves the strongest overall performance on LaMP-1, while *pQwen2* attains the best results on LaMP-2 and LaMP-5, indicating that the benefits of ClusterRAG generalize across model architectures, providing consistent gains without requiring additional model fine-tuning.

### 5.3.4 Ablation Study

The integral components of ClusterRAG are user representation and retrieval and profile retrieval. By systematically removing or replacing these in-

LLMs	LaMP-1		LaMP-2		LaMP-5	
	Acc.↑	F1↑	Acc.↑	F1↑	R-1↑	R-L↑
nFlan	0.546	0.540	0.451	0.448	0.431	0.398
pFlan	<b>0.648</b>	<b>0.647</b>	0.601	0.601	0.484	0.432
nQwen2	0.602	0.600	0.521	0.521	0.457	0.419
pQwen2	0.639	0.635	<b>0.610</b>	<b>0.606</b>	<b>0.488</b>	<b>0.447</b>

Table 4: Performance comparison of ClusterRAG using different LLMs on LaMP-(1,2,5).

Derivatives	LaMP-3		LaMP-7	
	MAE↓	RMSE↓	R-1↑	R-L↑
w/o user clustering	0.320	0.637	0.458	0.371
w/o intra-cluster sim	0.329	0.639	0.501	0.442
w/o doc ranking	0.331	0.642	0.462	0.413
Centroids only	0.400	0.643	0.472	0.438
<i>k</i> -means	0.291	0.610	0.502	0.453
<b>ClusterRAG</b>	<b>0.270</b>	<b>0.594</b>	<b>0.521</b>	<b>0.470</b>

Table 5: Ablation study of ClusterRAG on LaMP-(3,7).

dividual components, we assess their impact on personalized generation performance on selected LaMP tasks in Table 5.

First, we examine the role of collaborative user modeling. Replacing clustering-based neighbor selection with random user sampling (*w/o user clustering*) leads to a substantial degradation across all tasks, highlighting the importance of structured user grouping. Similarly, removing intra-cluster similarity ranking (*w/o intra-cluster sim*) consistently reduces performance, indicating that fine-grained similarity estimation within clusters is critical for identifying truly relevant collaborative signals. Lastly, we analyze the profile retrieval module. Using cluster centroids alone to represent user profiles (*Centroids only*) results in the largest performance drop, demonstrating that document-level evidence is essential. Excluding document ranking (*w/o doc ranking*) further degrades results, confirming that effective reranking is necessary to prioritize high-quality contextual evidence. Replacing HDBSCAN with *k*-means (Na et al., 2010) clustering yields slightly weaker yet competitive performance, suggesting that while ClusterRAG is robust to the choice of clustering algorithm, density-aware clustering provides additional benefits.

## 5.4 Discussion

**Computational Effectiveness.** ClusterRAG consistently improves personalized generation across diverse tasks and evaluation metrics by jointly

leveraging user-specific history and collaborative signals from similar users. The gains observed in both classification and generation settings indicate that clustering-based collaborative filtering provides complementary information beyond individual user profiles, particularly for sparse or ambiguous queries.

**Computational Efficiency.** ClusterRAG is computationally efficient due to its modular and lightweight design. User and document clustering is performed using HDBSCAN, a non-parametric algorithm without learnable parameters, enabling fast and scalable user grouping. The primary retriever, ColBERTv2, maintains a parameter size comparable to BERT by introducing only a small linear projection layer (approximately 0.1M parameters), resulting in a total model size of roughly 110M parameters. Since user and document embeddings can be precomputed offline, inference primarily involves similarity computations, yielding low latency and minimal overhead. This efficiency allows ClusterRAG to scale effectively and generalize rapidly to new datasets.

We further investigate the sensitivity of ClusterRAG to the number of similar users and the size of the retrieved profile in Appendix E. In addition, Appendix F evaluates cluster cohesion for collaborative user retrieval, and Appendix G presents a qualitative case study that illustrates the performance of ClusterRAG.

## 6 Conclusion

This work introduced ClusterRAG, a collaborative framework that organizes users and their documents into semantically coherent clusters and performs retrieval at both the cluster and document levels, effectively reducing search complexity while preserving retrieval quality. Extensive experiments on the LaMP benchmark demonstrate that the hybrid profile retrieval mode, which jointly leverages the target user’s profile and profiles from top similar users, is the most effective configuration, yielding the best overall performance. Additionally, the experiments indicate that ClusterRAG is retriever-agnostic, allowing seamless integration with different dense retrievers and rankers, and remains effective when paired with both fine-tuned and zero-shot language models, highlighting its robustness and generality. Overall, ClusterRAG offers a design approach that improves effectiveness without incurring significant computational overhead.

## 667 Limitations

668 While ClusterRAG’s primary goal is to retrieve  
669 similar-user documents and enhance personalized  
670 RAG, we highlight a few factors that may affect  
671 the model’s performance. First, ClusterRAG relies  
672 on prompt-based generation, and the adopted IPA  
673 strategy may not be optimal; more advanced and  
674 structured prompt formulation techniques could fur-  
675 ther improve personalized generation performance.  
676 Nevertheless, prompt engineering, which is cru-  
677 cial for performance of LLMs, is not the central  
678 objective of this study. Second, the LaMP-1 (Per-  
679 sonalized Citation Identification) and LaMP-5 (Per-  
680 sonalized Scholarly Title Generation) tasks pro-  
681 vide only paper abstracts rather than full-text con-  
682 tent, which may limit the contextual information  
683 available to LLMs when selecting citations or gen-  
684 erating titles, even though abstracts offer useful  
685 sequence constraints. Third, our evaluation is re-  
686 stricted to English, text-only datasets, leaving the  
687 effectiveness of ClusterRAG in multilingual and  
688 multimodal settings unexplored. Finally, Cluster-  
689 RAG’s end performance depends on the underlying  
690 language model, which may introduce additional  
691 limitations inherited from the backend LLM. Fu-  
692 ture work will focus on designing more sophisti-  
693 cated prompt generation strategies and extending  
694 ClusterRAG to multilingual and multimodal per-  
695 sonalized RAG scenarios.

## 696 Ethical Considerations

697 ClusterRAG leverages user interaction histories  
698 and collaborative signals, which raises considera-  
699 tions related to privacy, consent, and potential bias.  
700 Although the framework operates on anonymized  
701 user profiles and does not require access to explicit  
702 personal identifiers, improper handling of user-  
703 generated data could still risk unintended infor-  
704 mation leakage. Moreover, collaborative filtering  
705 may amplify existing biases if certain user groups  
706 or preferences are overrepresented in the data, po-  
707 tentially affecting fairness in personalized outputs.  
708 To mitigate these risks, ClusterRAG can be de-  
709 ployed with standard data governance practices,  
710 including data anonymization, access control, and  
711 bias-aware evaluation. Importantly, ClusterRAG is  
712 a model-agnostic retrieval framework rather than a  
713 user profiling system, and it does not infer sensitive  
714 attributes beyond observed interactions, helping  
715 limit ethical exposure while enabling effective per-  
716 sonalization.

## References

- Aleena Ahmad, Gibson Nkhata, Abdul Rafay Bajwa, Hannah Marsico, Bryan Le, and Susan Gauch. 2025. [Colbert-based user profiles for personalized information retrieval](#). In *Proceedings of the Seventeenth International Conference on Information, Process, and Knowledge Management, eKNOW '25*, pages 51–58, Nice, France. 718–724
- Akari Asai, Zeqiu Wu, Yizhong Wang, Avirup Sil, and Hannaneh Hajishirzi. 2024. [Self-RAG: Learning to retrieve, generate, and critique through self-reflection](#). In *The Twelfth International Conference on Learning Representations*, Vienna, Austria. 725–728
- Hyung Won Chung, Le Hou, Shayne Longpre, Barret Zoph, Yi Tay, William Fedus, Yunxuan Li, Xuezhi Wang, Mostafa Dehghani, Siddhartha Brahma, Albert Webson, Shixiang Shane Gu, Zhuyun Dai, Mirac Suzgun, Xinyun Chen, Aakanksha Chowdhery, Alex Castro-Ros, Marie Pellat, Kevin Robinson, and 16 others. 2024. [Scaling instruction-finetuned language models](#). *Journal of Machine Learning Research*, 25(70):1–53. 730–737
- Guanting Dong, Yutao Zhu, Chenghao Zhang, Zechen Wang, Ji-Rong Wen, and Zhicheng Dou. 2025. [Understand what llm needs: Dual preference alignment for retrieval-augmented generation](#). In *Proceedings of the ACM on Web Conference 2025, WWW '25*, page 4206–4225, New York, NY, USA. Association for Computing Machinery. 739–745
- Wenqi Fan, Yujuan Ding, Liangbo Ning, Shijie Wang, Hengyun Li, Dawei Yin, Tat-Seng Chua, and Qing Li. 2024. [A survey on rag meeting llms: Towards retrieval-augmented large language models](#). In *Proceedings of the 30th ACM SIGKDD Conference on Knowledge Discovery and Data Mining, KDD '24*, page 6491–6501, New York, NY, USA. Association for Computing Machinery. 746–753
- Markus Freitag and Yaser Al-Onaizan. 2017. [Beam search strategies for neural machine translation](#). In *Proceedings of the First Workshop on Neural Machine Translation*, pages 56–60, Vancouver. Association for Computational Linguistics. 754–758
- Yunfan Gao, Yun Xiong, Xinyu Gao, Kangxiang Jia, Jinliu Pan, Yuxi Bi, Yi Dai, Jiawei Sun, Qianyu Guo, Meng Wang, and Haofen Wang. 2023. [Retrieval-augmented generation for large language models: A survey](#). *ArXiv*, abs/2312.10997. 759–763
- Edward J Hu, yelong shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, and Weizhu Chen. 2022. [LoRA: Low-rank adaptation of large language models](#). In *International Conference on Learning Representations*. 764–768
- Yizheng Huang and Jimmy X. Huang. 2024. [A survey on retrieval-augmented text generation for large language models](#). *ArXiv*, abs/2404.10981. 769–774

772	Bernd Huber, Ghazal Fazelnia, Andreas Damianou, Sebastian Peleato, Max Lefarov, Praveen Ravichandran, Marco De Nadai, Mounia Lalmas-Roellke, and Paul N. Bennett. 2025. <a href="#">Embedding-to-prefix: Parameter-efficient personalization for pre-trained large language models</a> . <i>Preprint</i> , arXiv:2505.17051.	829
773		830
774		831
775		832
776		833
777		834
778	Gautier Izacard, Mathilde Caron, Lucas Hosseini, Sebastian Riedel, Piotr Bojanowski, Armand Joulin, and Edouard Grave. 2022. <a href="#">Unsupervised dense information retrieval with contrastive learning</a> . <i>Transactions on Machine Learning Research</i> , 2022.	835
779		
780		
781		
782		
783	Omar Khattab and Matei Zaharia. 2020. <a href="#">Colbert: Efficient and effective passage search via contextualized late interaction over bert</a> . In <i>Proceedings of the 43rd International ACM SIGIR Conference on Research and Development in Information Retrieval</i> , SIGIR '20, page 39–48, New York, NY, USA. Association for Computing Machinery.	836
784		837
785		838
786		839
787		840
788		841
789		842
790	Patrick Lewis, Ethan Perez, Aleksandra Piktus, Vladimir Karpukhin, Naman Goyal, and 1 others. 2020. <a href="#">Retrieval-augmented generation for knowledge-intensive nlp</a> . In <i>Advances in Neural Information Processing Systems 33 (NeurIPS 2020)</i> .	843
791		844
792		845
793		846
794		847
795	Xiang Lisa Li and Percy Liang. 2021. <a href="#">Prefix-tuning: Optimizing continuous prompts for generation</a> . In <i>Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)</i> , pages 4582–4597, Online. Association for Computational Linguistics.	848
796		849
797		
798		
799		
800		
801		
802		
803	Xiaopeng Li, Pengyue Jia, Derong Xu, Yi Wen, Yingyi Zhang, Wenlin Zhang, Wanyu Wang, Yichao Wang, Zhaocheng Du, Xiangyang Li, Yong Liu, Hui Feng Guo, Ruiming Tang, and Xiangyu Zhao. 2025a. <a href="#">A survey of personalization: From rag to agent</a> . <i>Preprint</i> , arXiv:2504.10147.	850
804		851
805		852
806		853
807		854
808		
809	Zongxi Li, Zijian Wang, Weiming Wang, Kevin Hung, Haoran Xie, and Fu Lee Wang. 2025b. <a href="#">Retrieval-augmented generation for educational application: A systematic survey</a> . <i>Computers and Education: Artificial Intelligence</i> , 8:100417.	855
810		856
811		857
812		858
813		
814	Chin-Yew Lin. 2004. <a href="#">ROUGE: A package for automatic evaluation of summaries</a> . In <i>Text Summarization Branches Out</i> , pages 74–81, Barcelona, Spain. Association for Computational Linguistics.	859
815		860
816		861
817		862
818	Ilya Loshchilov and Frank Hutter. 2019. <a href="#">Decoupled weight decay regularization</a> . <i>Preprint</i> , arXiv:1711.05101.	863
819		864
820		865
821	Leland McInnes, John Healy, and Sean Astels. 2017. <a href="#">hdbscan: Hierarchical density based clustering</a> . <i>Journal of Open Source Software (JOSS)</i> , 2(11).	866
822		867
823		868
824	Shi Na, Liu Xumin, and Guan Yong. 2010. <a href="#">Research on k-means clustering algorithm: An improved k-means clustering algorithm</a> . In <i>2010 Third International Symposium on Intelligent Information Technology and Security Informatics</i> , pages 63–67, Jian, China.	869
825		870
826		871
827		872
828		873
		874
		875
		876
		877
		878
		879
		880
		881
		882
		883
		884
		885
		886
		887
		888
		889
		890
		891
		892
		893
		894
		895
		896
		897
		898
		899
		900
		901
		902
		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
		956
		957
		958
		959
		960
		961
		962
		963
		964
		965
		966
		967
		968
		969
		970
		971
		972
		973
		974
		975
		976
		977
		978
		979
		980
		981
		982
		983
		984
		985
		986
		987
		988
		989
		990
		991
		992
		993
		994
		995
		996
		997
		998
		999
		1000

884		<i>the Association for Computational Linguistics: Human Language Technologies</i> , pages 3715–3734, Seattle, United States. Association for Computational Linguistics.		
885				941
886				942
887				
888	Kiriakos Sgardelis, Dionisis Margaritis, Dimitris Spiliotopoulos, and Costas Vassilakis. 2025. <a href="#">An evaluation review of user similarity metrics in sparse collaborative filtering datasets</a> . <i>International Journal of Data Science and Analytics</i> , 20:6665–6693.			943
889				944
890				945
891				946
892				947
893	Chenglei Shen, Xiao Zhang, Teng Shi, Changshuo Zhang, Guofu Xie, and Jun Xu. 2024. <a href="#">A survey of controllable learning: Methods and applications in information retrieval</a> . <i>ArXiv</i> , abs/2407.06083.			948
894				949
895				950
896				
897	Teng Shi, Zihua Si, Jun Xu, Xiao Zhang, Xiaoxue Zang, Kai Zheng, Dewei Leng, Yanan Niu, and Yang Song. 2024. <a href="#">Unisar: Modeling user transition behaviors between search and recommendation</a> . In <i>Proceedings of the 47th International ACM SIGIR Conference on Research and Development in Information Retrieval</i> , SIGIR '24, page 1029–1039, New York, NY, USA. Association for Computing Machinery.			951
898				952
899				953
900				954
901				955
902				956
903				957
904				958
905	Teng Shi, Jun Xu, Xiao Zhang, Xiaoxue Zang, Kai Zheng, Yang Song, and Han Li. 2025. <a href="#">Retrieval augmented generation with collaborative filtering for personalized text generation</a> . In <i>Proceedings of the 48th International ACM SIGIR Conference on Research and Development in Information Retrieval</i> , SIGIR '25, page 1294–1304, New York, NY, USA. Association for Computing Machinery.			959
906				960
907				961
908				962
909				963
910				964
911				965
912				966
913	Shamane Siriwardhana, Rivindu Weerasekera, Elliott Wen, Tharindu Kaluarachchi, Rajib Rana, and Suranga Nanayakkara. 2023. <a href="#">Improving the domain adaptation of retrieval augmented generation (RAG) models for open domain question answering</a> . <i>Transactions of the Association for Computational Linguistics</i> , 11:1–17.			967
914				968
915				969
916				970
917				971
918				972
919				
920	Jiakai Tang, Sunhao Dai, Teng Shi, Jun Xu, Xu Chen, Wen Chen, Wu Jian, and Yuning Jiang. 2025. <a href="#">Think before recommend: Unleashing the latent reasoning power for sequential recommendation</a> . <i>ArXiv</i> , abs/2503.22675.			973
921				974
922				975
923				976
924				977
925	Xiang Wang, Xiangnan He, Meng Wang, Fuli Feng, and Tat-Seng Chua. 2019. <a href="#">Neural graph collaborative filtering</a> . In <i>Proceedings of the 42nd International ACM SIGIR Conference on Research and Development in Information Retrieval</i> , SIGIR'19, page 165–174, New York, NY, USA. Association for Computing Machinery.			978
926				979
927				
928				980
929				981
930				982
931				983
932	Thomas Wolf, Lysandre Debut, Victor Sanh, Julien Chaumond, Clement Delangue, Anthony Moi, Pierric Cistac, Tim Rault, Remi Louf, Morgan Funtowicz, Joe Davison, Sam Shleifer, Patrick von Platen, Clara Ma, Yacine Jernite, Julien Plu, Canwen Xu, Teven Le Scao, Sylvain Gugger, and 3 others. 2020. <a href="#">Transformers: State-of-the-art natural language processing</a> . In <i>Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing: System Demonstrations</i> , pages 38–45, Online. Association for Computational Linguistics.			984
933				985
934				986
935				987
936				988
937				989
938				990
939				991
940				992
				993
				994
				995
				996

997 Changshuo Zhang, Teng Shi, Xiao Zhang, Yanping  
998 Zheng, Ruobing Xie, Qi Liu, Jun Xu, and Jirong  
999 Wen. 2024. Qagcf: Graph collaborative filtering for  
1000 q&a recommendation. *ArXiv*, abs/2406.04828.

1001 Changshuo Zhang, Xiao Zhang, Teng Shi, Jun Xu, and  
1002 Jirong Wen. 2025. Test-time alignment for track-  
1003 ing user interest shifts in sequential recommendation.  
1004 *ArXiv*, abs/2504.01489.

1005 Yangxiao Zhang. 2025. A retrieval-augmented genera-  
1006 tion framework with retriever and generator modules  
1007 for enhancing factual consistency. *Applied and Com-  
1008 putational Engineering*, 166:149–155.

1009 Yaochen Zhu, Chao Wan, Harald Steck, Dawen Liang,  
1010 Yesu Feng, Nathan Kallus, and Jundong Li. 2025.  
1011 Collaborative retrieval for large language model-  
1012 based conversational recommender systems. In  
1013 *Proceedings of the ACM on Web Conference 2025*,  
1014 WWW '25, page 3323–3334, New York, NY, USA.  
1015 Association for Computing Machinery.

1016 Thomas P Zollo, Andrew Wei Tung Siah, Naimeng Ye,  
1017 Ang Li, and Hongseok Namkoong. 2025. Personal-  
1018 LLM: Tailoring LLMs to individual preferences. In  
1019 *The Thirteenth International Conference on Learning  
1020 Representations*.

## 1021 A Prompts Used in ClusterRAG

1022 This subsection presents the prompt templates em-  
1023 ployed during generation for the remaining LaMP  
1024 tasks. As already stated, each prompt contains an  
1025 instruction, input (query), and profile. We provide  
1026 the templates below. In the templates, *{Movie de-  
1027 scription}* and *{Movie tags}*, *{Review}*, *{Article}*,  
1028 *{Paper abstract}*, and *{Tweet}* represent user input  
1029 for the corresponding LaMP tasks, while the rest of  
1030 the italicized text represent user profile entries. The  
1031 remaining text is the instruction guiding an LLM  
1032 to generate the intended output.

**LaMP-2 (Personalized Movie Tagging)**  
**Prompt Template:** Given the user pre-  
vious movie tag pairs: The tag for the  
movie description: *<Movie\_1\_Description>*  
is *<Tag\_1>*, the tag for the movie descrip-  
tion: *<Movie\_2\_Description>* is *<Tag\_2>*,  
..., the tag for the movie description:  
*<Movie\_M\_Description>* is *<Tag\_M>*, which  
tag does the movie description: *{Movie de-  
scription}* relate to among the following tags?  
Just answer with the tag name without further  
explanation. Movie tags: *{Movie tags}*

1033

### LaMP-3 (Personalized Product Rating)

**Prompt Template:** Given the user previous  
review-score pairs: *<Score\_1>* is the score for  
*<Review\_1\_Text>*, *<Score\_2>* is the score for  
*<Review\_2\_Text>*, ..., *<Score\_M>* is the score  
for *<Review\_M\_Text>*. What is the score of  
the following review on a scale of 1 to 5? Just  
answer with 1, 2, 3, 4, or 5 without further  
explanation. Review: *{Review}*

1034

### LaMP-4 (Personalized News Headline Generation)

**Prompt Template:** Given the user’s  
previous article-headline pairs: *<Headline\_1>*  
is the title for *<Article\_1\_Text>*, *<Headline\_2>*  
is the title for *<Article\_2\_Text>*, ..., *<Head-  
line\_M>* is the title for *<Article\_M\_Text>*.  
Generate a headline for the following article.  
Article: *{Article}*

1035

### LaMP-5 (Personalized Scholarly Title Generation)

**Prompt Template:** Given the user’s  
previous abstract-title pairs: *<Title\_1>* is a title  
for *<Abstract\_1\_Text>*, *<Title\_2>* is a title for  
*<Abstract\_2\_Text>*, ..., *<Title\_M>* is a title for  
*<Abstract\_M\_Text>*. Generate a title for the  
following abstract of a paper. Abstract: *{Paper  
abstract}*

1036

### LaMP-7 (Personalized Tweet Paraphrasing)

**Prompt Template:** Given the user’s pre-  
vious tweets: *<Tweet\_1>*, *<Tweet\_2>*, ...,  
*<Tweet\_M>*. Paraphrase the following tweet  
without any explanation before or after it fol-  
lowing the user’s tweeting patterns. Tweet:  
*{Tweet}*

1037

## 1038 B Detailed Dataset Statistics and 1039 Licensing Information

1040 Table 6 below provides detailed statistics of the  
1041 LaMP benchmark based on the time-based split, as  
1042 stated in Section 5.1. We additionally summarize  
1043 the licensing information and terms of use for each  
1044 LaMP task considered in this study as follows:

- 1045 1. Personalized Citation Identification (LaMP-  
1046 1): CC BYNC-SA 4.0.
- 1047 2. Personalized Movie Tagging (LaMP-2): Ed-  
1048 ucational or academic research, NON COM-  
1049 MERCIAL USE.

1050	3. Personalized Product Rating (LaMP-3): CC	<b>LaMP-4: Personalized News Headline Genera-</b>	1094
1051	BY-NC-SA 4.0.	<b>tion.</b> This is a generative task that evaluates the	1095
1052	4. Personalized news Headline Generation	ability of an LLM to capture the stylistic patterns	1096
1053	(LaMP-4): CC BY-NC-SA 4.0.	of an author $u$ by requiring it to generate a headline	1097
1054	5. Personalized Scholarly Title Generation	for an input news article, $x$ , given a user profile of	1098
1055	(LaMP-5): CC BY-NC-SA 4.0.	the authors' historical article-title pairs.	1099
1056	6. Personalized Tweet Paraphrasing (LaMP-7):	<b>LaMP-5: Personalized Scholarly Title Gener-</b>	1100
1057	CC BYNC-SA 4.0.	<b>ation.</b> LaMP-5 is another generative task that	1101
1058	<b>C Task Descriptions</b>	requires a generative model to generate titles for	1102
1059	The LaMP benchmark contains English and text-	an input article $x$ , given a user profile of histor-	1103
1060	only data. The documents used in each LaMP task	ical article-title pairs for an author. It is similar	1104
1061	do not contain personally identifiable information	to LaMP-4, but it uses different corpus domain,	1105
1062	that could otherwise compromise privacy issues.	scientific papers (similar to LaMP-1).	1106
1063	Below, we provide detailed descriptions of each	<b>LaMP-7: Personalized Tweet Paraphrasing.</b>	1107
1064	LaMP task included in our evaluation, outlining	This is framed as a generative personalized tweet	1108
1065	the task objectives, input-output formulations, and	paraphrasing task, which requires an LLM to gener-	1109
1066	the specific aspects of personalization each task is	ate a tweet in the style of a user $u$ given an input	1110
1067	designed to assess.	tweet $x$ , and a user profile of historical tweets by	1111
1068	<b>LaMP-1: Personalized Citation Identification.</b>	the user.	1112
1069	This task frames citation recommendation as a bi-	<b>D Hyperparameter Tuning Grid</b>	1113
1070	nary classification task and assesses the ability of	Table 7 summarizes the hyperparameter search	1114
1071	a language model to identify user preferences for	space explored during model selection, detailing	1115
1072	citations. Specifically, if the user $u$ writes a paper	the ranges and candidate values used to tune Clus-	1116
1073	$x$ , a language model must determine which of two	terRAG across optimization, training, and decod-	1117
1074	given candidate papers ( $a$ or $b$ ) $u$ will cite in $x$ .	ing configurations. As described in Section 5.1, the	1118
1075	The profile of each user encompasses all the papers	optimal hyperparameters were identified via grid	1119
1076	they have authored. Only the title and abstract of	search with early stopping applied on LaMP-1 and	1120
1077	each paper are retained in the user's profile; it uses	LaMP-7, representing classification and generative	1121
1078	scientific papers.	tasks, respectively.	1122
1079	<b>LaMP-2: Personalized Movie Tagging.</b> This	<b>E Impact of Similar User Size and Profile</b>	1123
1080	task recasts movie tagging as a multi-class classi-	<b>Size</b>	1124
1081	fication task. Given a movie description $x$ and a	This section investigates the sensitivity of Clus-	1125
1082	user's historical movie-tag pairs, a language model	terRAG to two key design parameters that govern	1126
1083	must predict one of 15 tags for $x$ . The movie tags	collaborative context construction: the number of	1127
1084	are: sci-fi, based on a book, comedy, action, twist	similar users ( $k$ ) incorporated during retrieval and	1128
1085	ending, dystopia, dark comedy, classic, psychol-	the number of profile documents ( $m$ ) selected per	1129
1086	ogy, fantasy, romance, thought-provoking, social	user. By systematically varying these parameters,	1130
1087	commentary, violence, and true story.	we analyze how the breadth of collaborative sig-	1131
1088	<b>LaMP-3: Personalized Product Rating.</b> LaMP-	nals and the depth of user profile information affect	1132
1089	3 is also framed as a multi-class classification task.	model performance across different LaMP tasks.	1133
1090	In particular, given the user $u$ 's historical review	<b>Impact of the Number of Similar Users (<math>k</math>).</b> Ta-	1134
1091	and rating pairs of products and an input review $x$ ,	ble 8 shows that incorporating a small number of	1135
1092	the model must predict an integer rating (from 1 to	similar users consistently improves ClusterRAG	1136
1093	5) of the review.	performance across all LaMP tasks. Performance	1137
		generally increases as $k$ grows from 1 to 3, where	1138
		most metrics achieve their peak or near-peak val-	1139
		ues, indicating that a limited set of highly simi-	1140

Task	#users	#train	#dev	#test	Input Length	Output Length	#Profile Size	#classes
LaMP-1	6542	6542	1500	1500	$51.43 \pm 5.70$	–	$84.15 \pm 47.54$	2
LaMP-2	929	5073	1410	1557	$92.39 \pm 21.95$	–	$86.76 \pm 189.52$	15
LaMP-3	20000	20000	2500	2500	$128.18 \pm 146.25$	–	$185.40 \pm 129.30$	5
LaMP-4	1643	12500	1500	1800	$29.97 \pm 12.09$	$10.07 \pm 3.10$	$204.59 \pm 250.75$	–
LaMP-5	14682	14682	1500	1500	$162.34 \pm 65.63$	$9.71 \pm 3.21$	$87.88 \pm 53.63$	–
LaMP-7	13437	13437	1498	1500	$29.72 \pm 7.01$	$16.96 \pm 5.67$	$15.71 \pm 14.86$	–

Table 6: Detailed statistics of the LaMP benchmark with time-based data split.

Hyperparameter	Tested values
Learning rate	$5 \times 10^{-5}, 3 \times 10^{-3}, 10^{-3}, 10^{-4}$
Weight decay	$5 \times 10^{-6}, 10^{-4}, 10^{-3}$
Warm-up ratio	0.05 to 0.10
Batch size	8, 16, 32, 64
Epochs	10, 20 30, 50, 70, 100
Max seq length	64, 128, 256, 512
Beam size	1 – 6
$\gamma$	0.1 – 0.9
$k$	1 – 5
$m$	1 – 12
$L_{\max}$	64, 128, 256, 512, 1024
$ \bar{y} $	32, 64, 128, 256, 512

Table 7: Hyperparameter tuning grid used for training and optimizing ClusterRAG.

lar users provides the most informative collaborative signals. Beyond  $k = 3$ , gains saturate or slightly decline, suggesting that adding more users introduces weaker or noisier preferences that dilute personalization benefits. This trend highlights the importance of selectively leveraging collaborative information rather than aggregating large numbers of loosely related users.

**Impact of the Number of Retrieved Profile Documents ( $m$ ).** As shown in Table 9, increasing the number of retrieved profile documents leads to steady performance improvements up to a moderate range ( $m \approx 6-7$ ), after which the gains plateau. Larger profile sizes consistently reduce prediction error (MAE/RMSE) and improve generation quality (ROUGE scores), reflecting richer contextual grounding. At the same time, marginal benefits diminish for larger  $m$ , indicating that excessively long profiles provide limited additional signal while increasing prompt complexity. Overall, these results suggest that ClusterRAG is robust to the choice of  $m$  and performs best when balancing sufficient contextual coverage with prompt efficiency.

## F Cluster Cohesion for Collaborative User Retrieval

To evaluate whether ClusterRAG forms meaningful user clusters for collaborative filtering, we measure the Silhouette score (Rousseeuw, 1987) of user clusters produced by HDBSCAN and  $k$ -means and report results in Table 10. The Silhouette score is an internal clustering metric that jointly captures intra-cluster cohesion and inter-cluster separation, with values ranging from  $-1$  to  $1$ , where higher scores indicate better-defined clusters.

As shown in Table 10, ClusterRAG consistently achieves higher Silhouette scores with HDBSCAN than with  $k$ -means, whose scores remain below  $0.5$  across all tasks. In high-dimensional embedding spaces, moderately positive Silhouette scores are common and still reflect meaningful structure. Therefore, the scores above  $0.5$  obtained with HDBSCAN indicate that ClusterRAG learns cohesive and well-separated user clusters, enabling the retrieval of more similar users for collaborative filtering in personalized RAG. These findings further corroborate the superior performance of ClusterRAG when combined with HDBSCAN.

## G Case Study

We randomly sample a case from **LaMP\_2 (Personalized Movie Tagging)** in Table 11 to illustrate the effectiveness of ClusterRAG in leveraging both target and similar user profiles for personalized generation. In this task, the *User Query* corresponds to a movie description, and the objective is to generate an appropriate movie tag based on this description and the user’s historical tagging behavior. Due to space constraints, we include three profile documents per user. The target user’s historical tags are *twist ending* and *action*, while the similar user’s historical tags include *true story*, *action*, and *violence*. The gold label for the given query is *violence*.

When we use a non-personalized prompt or the target user profile only, ClusterRAG incorrectly pre-

$k$ value	LaMP-1		LaMP-2		LaMP-3		LaMP-4		LaMP-5		LaMP-7	
	Acc.↑	F1↑	Acc.↑	F1↑	MAE↓	RMSE↓	R-1↑	R-L↑	R-1↑	R-L↑	R-1↑	R-L↑
1	0.690	0.690	0.661	0.620	0.270	0.594	0.190	0.176	0.490	0.440	0.521	<b>0.470</b>
2	0.692	0.697	0.665	0.631	<b>0.269</b>	<b>0.582</b>	0.193	0.179	0.496	0.444	0.524	0.469
3	<b>0.700</b>	<b>0.700</b>	<b>0.668</b>	<b>0.642</b>	0.270	0.595	<b>0.196</b>	<b>0.180</b>	0.495	<b>0.445</b>	<b>0.528</b>	0.468
4	0.690	0.688	0.660	0.621	0.272	0.595	0.192	0.178	<b>0.497</b>	0.441	0.523	0.469
5	0.689	0.690	0.658	0.619	0.273	0.596	0.191	0.177	0.492	0.438	0.521	0.467

Table 8: Effect of the number of similar users ( $k$ ) on ClusterRAG performance in hybrid mode across LaMP tasks.

$m$ value	LaMP-1		LaMP-2		LaMP-3		LaMP-4		LaMP-5		LaMP-7	
	Acc.↑	F1↑	Acc.↑	F1↑	MAE↓	RMSE↓	R-1↑	R-L↑	R-1↑	R-L↑	R-1↑	R-L↑
1	0.674	0.673	0.648	0.612	0.289	0.608	0.182	0.166	0.480	0.431	0.515	0.464
2	0.690	0.690	0.661	0.620	0.270	0.594	0.190	0.176	0.490	0.440	0.521	0.470
3	0.691	0.691	0.665	0.632	0.269	0.591	0.191	0.178	0.501	0.465	0.528	0.474
4	0.693	0.692	0.669	0.639	0.268	0.590	0.196	0.179	0.510	0.471	0.530	0.473
5	0.695	0.695	<b>0.675</b>	<b>0.650</b>	0.267	0.586	0.198	0.180	0.512	0.473	0.533	0.475
6	<b>0.707</b>	<b>0.707</b>	0.673	0.649	0.262	0.584	<b>0.200</b>	<b>0.182</b>	<b>0.514</b>	<b>0.479</b>	0.534	0.476
7	0.703	0.700	0.672	0.647	<b>0.260</b>	<b>0.581</b>	0.199	0.180	0.511	0.474	<b>0.538</b>	<b>0.481</b>
8	0.704	0.701	0.670	0.645	0.263	0.583	0.197	0.178	0.509	0.472	0.536	0.478
9	0.707	0.701	0.671	0.643	0.263	0.585	0.196	0.177	0.511	0.474	0.534	0.476
10	0.706	0.705	0.671	0.644	0.261	0.582	0.191	0.174	0.508	0.471	0.533	0.476
11	0.701	0.700	0.669	0.641	0.263	0.584	0.191	0.173	0.504	0.469	0.534	0.475
12	0.700	0.700	0.670	0.664	0.264	0.584	0.192	0.174	0.504	0.470	0.532	0.473

Table 9: Effect of the number of retrieved profile documents ( $m$ ) on ClusterRAG performance in hybrid mode across LaMP tasks.

Task	HDBSCAN Score	$k$ -means Score
LaMP-1	0.601	0.389
LaMP-2	0.535	0.326
LaMP-3	0.551	0.328
LaMP-4	0.570	0.274
LaMP-5	0.562	0.347
LaMP-7	0.537	0.323

Table 10: Silhouette scores of user clusters produced by ClusterRAG using HDBSCAN and  $k$ -means on the LaMP benchmark.

## H AI Assistance Usage

In this work, ChatGPT has been used solely as a writing assistant. Specifically, draft passages were provided to the tool for paraphrasing and language refinement, after which we manually reviewed, edited, and finalized the text.

1216  
1217  
1218  
1219  
1220  
1221

dicts the movie tag as *action*, driven by its higher frequency in the current user’s history. However, when similar-user information is incorporated via hybrid profile retrieval, the model correctly predicts *violence*, as the movie tag. This occurs because the top-ranked retrieved profile, originating from a similar user, emphasizes personal and retaliatory violence that closely aligns with the query, whereas the lower-ranked *action*-tagged profile reflects more stylized narratives less relevant to the description.

1205  
1206  
1207  
1208  
1209  
1210  
1211  
1212  
1213  
1214  
1215

---

**User Query:**

**Movie description:** When the Davison family comes under attack during their wedding anniversary getaway, the gang of mysterious killers soon learns that one of their victims harbors a secret talent for fighting back.

---

**Gold Output:** Violence

---

**Target User Profile:**

(1) **Movie tag:** “twist ending”, **Movie description:** “Soon after his insufferably arrogant father wins the Nobel Prize for chemistry, Barkley Michaelson is kidnapped by Thaddeus James, a young genius who claims to be Barkley’s illegitimate half-brother. Motivated not so much by money as revenge, Thaddeus tries to convince Barkley to help him carry out a multimillion-dollar extortion plot against their patriarch.”

(2) **Movie tag:** “action”, **Movie description:** “*The Bride unwaveringly continues on her roaring rampage of revenge against the band of assassins who had tried to kill her and her unborn child. She visits each of her former associates one-by-one, checking off the victims on her Death List Five until there’s nothing left to do . . . but kill Bill.*”

(3) **Movie tag:** “action”, **Movie description:** “NYPD cop John McClane’s plan to reconcile with his estranged wife is thrown for a serious loop when, minutes after he arrives at her office, the entire building is overtaken by a group of terrorists. With little help from the LAPD, wisecracking McClane sets out to single-handedly rescue the hostages and bring the bad guys down.”

---

**Similar User Profile:**

(1) **Movie tag:** “true story”, **Movie description:** “The mostly true story of the legendary ‘worst director of all time’, who, with the help of his strange friends, filmed countless B-movies without ever becoming famous or successful.”

(2) **Movie tag:** “action”, **Movie description:** “Liu Jian, an elite Chinese police officer, comes to Paris to arrest a Chinese drug lord. When Jian is betrayed by a French officer and framed for murder, he must go into hiding and find new allies.”

(3) **Movie tag:** “violence”, **Movie description:** “*An elderly ex-serviceman and widower looks to avenge his best friend’s murder by doling out his own form of justice.*”

---

**Top Ranked Documents:** Doc # (3) from similar user then doc # (2) from current user.

---

**Personalized Prompt:**

Given the user previous movie tag pairs: The tag for the movie description: “*An elderly ex-serviceman and widower looks to avenge his best friend’s murder by doling out his own form of justice*” is “violence”, and the tag for the movie description: “*The Bride unwaveringly continues on her roaring rampage of revenge against the band of assassins who had tried to kill her and her unborn child. She visits each of her former associates one-by-one, checking off the victims on her Death List Five until there’s nothing left to do . . . but kill Bill.*” is “action”, which tag does the movie description: “*When the Davison family comes under attack during their wedding anniversary getaway, the gang of mysterious killers soon learns that one of their victims harbors a secret talent for fighting back.*” relate to among the following tags? Just answer with the tag name without further explanation. Movie tags: “sci-fi, based on a book, comedy, action, twist ending, dystopia, dark comedy, classic, psychology, fantasy, romance, thought-provoking, social commentary, violence, and true story”

---

**Generated Output:** Violence

---

Table 11: A case study illustrating how ClusterRAG leverages user profile and similar-user information for personalized RAG.