# RaSeRec: Retrieval-Augmented Sequential Recommendation

Xinping Zhao[1], Baotian Hu[1*], Yan Zhong[2], Shouzheng Huang[1], Zihao Zheng[1], Meng Wang[3],
Haofen Wang[3*], and Min Zhang[1]

[1]Harbin Institute of Technology (Shenzhen), Shenzhen, China
[2]Peking University, Beijing, China, [3]Tongji University, Shanghai, China
zhaoxinping@stu.hit.edu.cn,mengwangtj@tongji.edu.cn
carter.whfcarter@gmail.com,{hubaotian,zhangmin2021}@hit.edu.cn

## Abstract

Although prevailing supervised and self-supervised learning augmented sequential recommendation (SeRec) models have achieved improved performance with powerful neural network architectures, we argue that they still suffer from two limitations: **(1) Preference Drift**, where models trained on past data can hardly accommodate evolving user preference; and **(2) Implicit Memory**, where head patterns dominate parametric learning, making it harder to recall long tails. In this work, we explore retrieval augmentation in SeRec, to address these limitations. Specifically, we propose a Retrieval-Augmented Sequential Recommendation framework, named **RaSeRec**, the main idea of which is to maintain a dynamic memory bank to accommodate preference drifts and retrieve relevant memories to augment user modeling explicitly. It consists of two stages: (i) collaborative-based pre-training, which learns to recommend and retrieve; (ii) retrieval-augmented fine-tuning, which learns to leverage retrieved memories. Extensive experiments on three datasets fully demonstrate the superiority and effectiveness of RaSeRec. The implementation code is available at https://github.com/HITsz-TMG/RaSeRec.

## CCS Concepts

• **Information systems** → **Recommender systems**; **Information retrieval**.

## Keywords

Sequential Recommendation, Retrieval Augmentation, Preference Drift, Implicit Memory

---

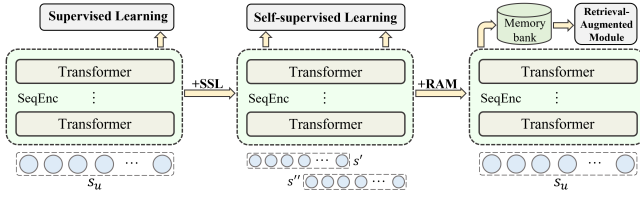*Baotian Hu and Haofen Wang are the corresponding authors.

---

## 1 Introduction

Sequential Recommendation (SeRec) aims to predict the next items that users would like to adopt, by calculating the matching probabilities between their historical interactions and the candidate items [15, 21, 58]. As such, it is essential to learn high-quality user representations from the user's historical interactions for better recommendation performance [29, 34, 53]. However, with enormous parameters to be optimized, existing SeRec models commonly suffer from the data sparsity issue, *i.e.*, users only interacted with a tiny fraction of the total items [16], making it difficult to learn high-quality user representations. Recently, there has been a surge of research interest in leveraging self-supervised learning (SSL) within the realm of SeRec to alleviate the data sparsity issue caused by sparse supervised signals [7, 14, 29, 34, 53]. Generally, it involves two procedures: **(1) data augmentation**, which generates multiple views for each user sequence, and **(2) self-supervised learning**, which maximizes the agreement between different views of the same (or similar) user sequence. In this way, SSL-Augmented models endow their model backbones with additional parametric knowledge to generate more high-quality user representations than those without SSL [48]. Despite their effectiveness, we argue that current SSL-Augmented ones still suffer from some severe limitations:

- **Preference Drift.** In real-world recommender systems, user preference frequently drifts as time goes on [3, 30] due to many reasons, *e.g.*, life changes [1]. As such, models trained on past user preferences can hardly handle preference drifts and may recommend undesirable items to users, which leads to sub-optimal performance.
- **Implicit Memory.** Existing methods encode user sequential patterns into implicit memory (model parameters), where the distribution of observed patterns usually obeys a power law [5, 32]. As such, long-tailed patterns that lack supervised signals may be overwhelmed by head ones during training, making it hard to recall long-tailed ones during inference.

In this work, we explore Retrieval-Augmented Generation (RAG) [12] in SeRec, to address the above limitations. Though being well-studied in Large Language Models (LLMs) [2, 17, 18, 24, 60], RAG is rarely explored in SeRec. The main idea of RAG is to retrieve relevant documents from an external knowledge base to facilitate LLMs' generation, especially for knowledge-intensive tasks [24]. For example, RAG [24] uses the input sequence as the query to retrieve text documents from Wikipedia and then adopts them as additional context to generate the target sequence; Self-RAG [2] learns to retrieve passages on demand adaptively, and interleaves generation and retrieval with reflection tokens by unifying them as the next token prediction. Compared to SSL (similar to pre-training),

**Figure 1: Comparison between Vanilla, SSL-Augmented, and Retrieval-Augmented SeRec paradigms, where their system working flows are illustrated from left to right, respectively. $s_u$ denotes the user sequence, while $s'_u, s''_u$ represent two augmented views. SeqEnc denotes the sequence encoder, referring to §2.2 for more technical details.**

RAG allows us to accommodate preference drift rapidly and recall long tails via looking up an external memory bank directly.

Here, we wish to bring the RAG's superiority into SeRec to enhance user representation learning, which differs from NLP tasks since they model *language semantics* while recommender systems model *collaborative semantics* [61]. To address the above limitations of SSL-Augmented SeRec models, we design a retrieval-augmented mechanism to improve the quality of user representations. Specifically, it consists of two key components: **(1) memory retrieval**, which recalls collaborative memories for the input user, and **(2) representation augmentation**, which leverages retrieved memories to augment user representation. On the one hand, *memory retrieval* adapts well to dynamic preference drifts with a real-time update memory bank. On the other hand, *representation augmentation* explicitly leverages retrieved memories, rendering it easier to recall long tails. Figure 1 illustrates the difference between the three SeRec paradigms and the continuation between them. The vanilla SeRec is trained on supervised signals to predict the next item based on the input user sequence [10]. The SSL-Augmented SeRec further enriches the parametric knowledge of the model backbone with SSL signals [55]. Standing on their shoulders, the retrieval-augmented SeRec handles dynamic preference drifts and relieves the heavy reliance on implicit memory by *memory retrieval* and *representation augmentation*. Note that the retrieval-augmented SeRec is designed to further perfect existence, rather than completely replacing them.

To this end, we propose a novel Retrieval-Augmented Sequential Recommendation paradigm, named **RaSeRec**, which learns to refine user representations with retrieved memories. Specifically, RaSeRec's training process consists of two key stages: (i) collaborative-based pre-training and (ii) retrieval-augmented fine-tuning. In the **collaborative-based pre-training** stage, the model is trained with two objectives built upon collaborative signals: (1) recommendation learning, which learns to predict the next item based on input user sequence, and (2) retrieval training, which learns to retrieve memories with the same (or similar) preference as the input user sequence. During the **retrieval-augmented fine-tuning** stage, we treat <user sequence, target item> pairs as the *reference set* except the current input user sequence. We then use the pre-trained model to encode the reference set into a memory bank. Given a user sequence, our model first retrieves similar collaborative memories from the memory bank. Then, a Retrieval-Augmented Module

(RAM) is built upon the pre-trained encoder, which learns to augment the representation of the input user with retrieved memories. An intuitive explanation of our idea is that it can be seen as an open-book exam. Specifically, with retrieval augmentation, RaSeRec does not have to memorize all sequential patterns whereas it learns to use retrieved ones (cheat sheets) from the memory bank. Note that RaSeRec is model-agnostic and can be applied to any ID-based SeRec backbones. Here, we adopt SASRec [21] as the backbone since it is simple yet effective. We also implement it on different backbones (see §4.3). In summary, our contributions are:

- We unveil the major issues existing models suffer from and shed light on the potential of RAG in SeRec.
- We propose a new SeRec paradigm, RaSeRec, which explicitly retrieves collaborative memories and then learns to use them for better user representation modeling.
- Extensive experiments on three benchmark datasets fully demonstrate that RaSeRec improves overall performance, enhances long-tailed recommendation, and alleviates preference drift.

## 2 Preliminaries

### 2.1 Problem Formulation

In this section, we first introduce the symbols and then formalize the problem of SeRec. Let $\mathcal{U}$ and $\mathcal{V}$ denote the set of users and items, respectively, where $|\mathcal{U}|$ and $|\mathcal{V}|$ denote the number of users and items. The historically interacted items of a user are sorted in chronological order: $s_u = [v_1^{(u)}, v_2^{(u)}, ..., v_{|s_u|}^{(u)}]$, where $v_t^{(u)} \in \mathcal{V}, 1 \le t \le |s_u|$ denotes the item interacted by user $u \in \mathcal{U}$ at the time step $t$ and $|s_u|$ denotes the length of interaction sequence of user $u$. Besides, $s_{u,t} = [v_1^{(u)}, v_2^{(u)}, ..., v_t^{(u)}]$ is a subsequence, where items are interacted by user $u$ before the time step $t + 1$. The goal of sequential recommendation is to predict the next item $v_{|s_u|+1}^{(u)}$ given the interaction sequence $s_u$, which can be formulated as follows:

$$v_* = \arg \max_{v_i \in \mathcal{V}} P\left(v_{|s_u|+1}^{(u)} = v_i | s_u\right), \tag{1}$$

where $v_*$ denotes the predicted item with the largest matching probability. In what follows, we employ bold lowercase and uppercase symbols to represent vectors and matrices, respectively.

### 2.2 Transformer for SeRec

In this section, we introduce how to model users' historical interactions to get their representations. Specifically, we adopt SASRec [21] as the model backbone, whose encoding module is based on the Transformer encoder [44]. To leverage the transformer's strong encoding ability, we first convert items to embeddings. Then, we apply the Transformer encoder to generate the user representation.

*2.2.1 Embedding Layer.* Formally, an embedding table $\mathbf{V} \in \mathbb{R}^{|\mathcal{V}| \times d}$ is created upon the whole item set $\mathcal{V}$ to project each item (one-hot representation) into a low-dimensional dense vector for better sequential modeling, where $d$ is the embedding dimension. Additionally, to perceive the time order of the sequence, a learnable position encoding matrix $\mathbf{P} \in \mathbb{R}^{T \times d}$ is constructed, where $T$ is the maximum sequence length. Owing to the constraint of the maximum sequence length $T$, when inferring user $u$ representation at time step $t + 1$, we truncate the user sequence $s_{u,t}$ to the last $T$ items

**Figure 2: The overall system framework of the proposed RaSeRec. The upper layer illustrates the workflow of collaborative-based pre-training while the bottom layer shows the workflow of retrieval-augmented fine-tuning.**

if $t > T$, otherwise left unchanged:

$$s_{u,t} = [v_{t-T+1}, v_{t-T+2}, ..., v_t], \tag{2}$$

where we omit the superscript $(u)$ for brevity. Subsequently, the item embedding and position encoding are added together to make up the input vector:

$$\mathbf{h}_i^0 = \mathbf{v}_i + \mathbf{p}_i, \tag{3}$$

where $\mathbf{h}_i^0$ is the input vector at position $i$. All input vectors collectively form the hidden representation of user sequence $s_{u,t}$ as $\mathbf{H}^0 = [\mathbf{h}_{t-T+1}^0, ..., \mathbf{h}_t^0]$.

*2.2.2 Transformer Encoder.* A Transformer encoder is composed of a multi-head self-attention module and a position-wise feed-forward Network. More technical details about these two modules can be found in [21] and [44]. We stack multiple blocks to capture more complex sequential patterns for high-quality user representation. Specifically, an $L$-layer Transformer encoder (unidirectional attention) is applied to refine the hidden representation of each item in $\mathbf{H}^0$, which can be formulated as follows:

$$\mathbf{H}^L = \text{Trm}(\mathbf{H}^0), \tag{4}$$

where $\text{Trm}(\cdot)$ denotes the $L$-layer Transformer encoder, $\mathbf{H}^L = [\mathbf{h}_{t-T+1}^L, ..., \mathbf{h}_t^L]$ is the refined hidden representations of the sequence. The last refined hidden representation $\mathbf{h}_t^L$ is selected as the representative of the user sequence. Here, we omit the superscript $L$ and subscript $t$ for brevity, *i.e.,* $\mathbf{h}$. In addition, we use $\text{SeqEnc}(\cdot)$ to denote the sequence encoder, that is $\mathbf{h} = \text{SeqEnc}(s_{u,t})$ for simplicity.

## 3 Methodology

We propose the Retrieval-Augmented Sequential Recommendation (RaSeRec) paradigm, which endows SeRec to accommodate preference drift and recall long-tailed patterns, as shown in Figure 2. Specifically, it consists of a two-stage training strategy: **(1) collaborative-based pre-training**, which learns to recommend and retrieve based on collaborative signals (§3.1); **(2) retrieval-augmented fine-tuning**, which learns to leverage retrieved memories to augment user representation (§3.2). Lastly, we introduce the inference process and analyze the complexity of RaSeRec (§3.3). The notation table can be found in Appendix A.

### 3.1 Collaborative-based Pre-training

In this section, we pre-train the model backbone $\text{SeqEnc}(\cdot)$ with two learning objections: **(1) recommendation learning**, which endows the backbone with the ability to generate the next item based on the input user sequence, and **(2) retrieval training**, which endows the backbone with the ability to retrieve memories with the same (or similar) preference as the input user sequence.

*3.1.1 Recommendation Learning.* After computing the user representation $\mathbf{h} \in \mathbb{R}^d$ (§2.2), a prediction layer is built upon it to predict how likely user $u$ would adopt item $v_i$ at time step $t + 1$. A simple yet effective solution is to calculate the inner product, *i.e.,* $\mathbf{h}^T \mathbf{v}_i$. After that, we employ the next item prediction over the whole item set $\mathcal{V}$ as the recommendation learning objective. To be specific, we compute the negative log-likelihood with the softmax as the recommendation learning objective:

$$\mathcal{L}_{rec} = -\log\left(e^{\mathbf{h}^T \mathbf{v}_{t+1}} / \sum_{v_i \in \mathcal{V}} e^{\mathbf{h}^T \mathbf{v}_i}\right), \tag{5}$$
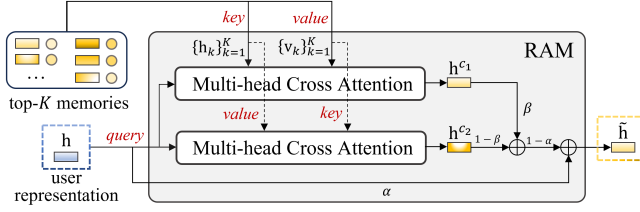
where $\mathbf{h}, \mathbf{v}_{t+1}, \mathbf{v}_i$ denote the last refined hidden representation $\mathbf{h}_t^L$, the embedding vector of the next item $v_{t+1}$, and the embedding vector of item $v_i \in \mathcal{V}$, respectively. We train the model to predict the next item auto-regressively similar to the language modeling's objective [35].

*3.1.2 Retrieval Training.* To retrieve memories that share the same (or similar) preference as the input user sequence, a retrieval training objective is developed to maximize the agreement of positive pairs and minimize that of negative pairs. The positive pair consists of two user sequences that share the same (or similar) preference. In SeRec, the goal is to predict the next item. Therefore, if the next item of two user sequences is the same, these two sequences can be thought to have similar collaborative semantics [34]. Particularly, the positive pair can be defined as $s_{u',t^i} = [v_{t^i-T+1}, ..., v_{t^i}]$, and $s_{u'',t^j} = [v_{t^j-T+1}, ..., v_{t^j}]$, where $v_{t^i+1}$ and $v_{t^j+1}$ are the same item. For brevity, we mark $s_{u',t^i}, s_{u'',t^j}$ as $s', s''$, respectively. The user representation of them can be defined as:

$$\mathbf{h}' = \text{SeqEnc}(s'), \quad \mathbf{h}'' = \text{SeqEnc}(s''). \tag{6}$$

In retrieval training settings, negative samples usually need to be selected from a large pool [22]. To efficiently construct negatives, we treat the other samples in the same training batch as negatives. Assume that we have a training batch $\mathcal{B}$ that consists of $|\mathcal{B}|$ positive pairs:

$$\mathcal{B} = \{\langle \mathbf{h}_1', \mathbf{h}_1'' \rangle, \langle \mathbf{h}_2', \mathbf{h}_2'' \rangle, ..., \langle \mathbf{h}_{|\mathcal{B}|}', \mathbf{h}_{|\mathcal{B}|}'' \rangle\}. \tag{7}$$

**Figure 3: The architecture of Retrieval-Augmented Module, where $\tilde{\mathbf{h}}$ is the augmented user representation.**

As such, for each positive pair, there are $2(|\mathcal{B}|-1)$ negatives forming the negative set $\mathcal{S}^-$. For example, for the positive pair $\langle \mathbf{h}'_1, \mathbf{h}''_1 \rangle$, its corresponding negative set is $\mathcal{S}_1^- = \{\mathbf{h}'_2, \mathbf{h}''_2, ..., \mathbf{h}'_{|\mathcal{B}|}, \mathbf{h}''_{|\mathcal{B}|}\}$. Having established the positive and negative pairs, we adopt InfoNCE [13] as the loss function, which can be defined as follows:

$$\mathcal{L}_{ret} = - \Big( \log \frac{e^{s(\mathbf{h}'_i, \mathbf{h}''_i)/\tau}}{e^{s(\mathbf{h}'_i, \mathbf{h}''_i)/\tau} + \sum_{s^- \in \mathcal{S}_i^-} e^{s(\mathbf{h}'_i, \mathbf{h}^-)/\tau}}$$
$$+ \log \frac{e^{s(\mathbf{h}''_i, \mathbf{h}'_i)/\tau}}{e^{s(\mathbf{h}''_i, \mathbf{h}'_i)/\tau} + \sum_{s^- \in \mathcal{S}_i^-} e^{s(\mathbf{h}''_i, \mathbf{h}^-)/\tau}} \Big), \tag{8}$$

where we take the $i$-th positive pair as an example; $\mathcal{S}_i^-$ denotes the negative set; $s(\cdot)$ measures the similarity between two vectors, which is implemented as the inner product; $\tau$ is the temperature, which is used to control the smoothness of softmax. We pre-train the model backbone with the training objectives of $\mathcal{L}_{rec}$ as well as $\mathcal{L}_{ret}$, simultaneously.

## 3.2 Retrieval-Augmented Fine-tuning

After pre-training, the model backbone has been endowed with the ability to recommend the next item and retrieve relevant memories. However, it still suffers from two limitations: **(i)** It can hardly accommodate dynamic *preference drifts* as it is trained on past user preferences, and **(ii)** It may fail to recall long-tailed patterns from *implicit memory* due to skewed data distribution. Given the above issues, we propose Retrieval-Augmented Fine-Tuning, RAFT, making the model quicker to adapt to preference drifts and easier to recall long-tailed patterns. It mainly consists of two components: **(1) memory retrieval**, which recalls useful memories from the memory bank, and **(2) representation augmentation**, which leverages the retrieved memories to augment user representation.

*3.2.1 Memory Retrieval.* Collaborative <user sequence $s_{u,t}$, target item $v_{t+1}^{(u)}$> pairs explicitly reveal the sequential patterns, where pairs sharing the same (or similar) preference as the input user sequence can act as important references. Inspired by the above idea, we construct a reference set $\mathcal{R}$ by auto-regressively enumerating all <$s_{u,t}, v_{t+1}^{(u)}$> pairs from the training data except the current input user sequence. We then employ the sequence encoder SeqEnc$(\cdot)$ to encode all pairs in $\mathcal{R}$ so that we can get a memory bank $\mathcal{M}$, where each entry in $\mathcal{M}$ consists of a user representation and the corresponding target item embedding. Given the input representation $\mathbf{h}$, RaSeRec retrieves top-$K$ similar user representations $\{\mathbf{h}_k\}_{k=1}^K$ as well as their corresponding target item embedding $\{\mathbf{v}_k\}_{k=1}^K$ from

the memory bank $\mathcal{M}$, where we employ the Faiss library [20] to speed up the retrieval process.

*3.2.2 Representation Augmentation.* With the retrieved memories, we design a Retrieval-Augmented Module (RAM), which learns to leverage the retrieved memories to augment the current user representation. Its overall architecture is illustrated in Figure 3. Inspired by [52], RAM employs dual-channel multi-head cross attention (MHCA) which takes $\mathbf{h}$ as *query*, $\{\mathbf{h}_k\}_{k=1}^K$ (or $\{\mathbf{v}_k\}_{k=1}^K$) as *key* and $\{\mathbf{v}_k\}_{k=1}^K$ (or $\{\mathbf{h}_k\}_{k=1}^K$) as *value* to get the augmented user representation. In the first channel, RAM learns to weight and aggregate $\{\mathbf{v}_k\}_{k=1}^K$ via modeling the relationship between $\mathbf{h}$ and $\{\mathbf{h}_k\}_{k=1}^K$:

$$\mathbf{h}^{c_1} = \text{MHCA}(\mathbf{h}, \{\mathbf{h}_k\}_{k=1}^K, \{\mathbf{v}_k\}_{k=1}^K), \tag{9}$$

where MHCA$(\cdot)$ is the multi-head cross attention module; $\mathbf{h}^{c_1}$ is the first augmented representation. Analogously, RAM compute the second augmented representation $\mathbf{h}^{c_2}$ through aggregating $\{\mathbf{h}_k\}_{k=1}^K$:

$$\mathbf{h}^{c_2} = \text{MHCA}(\mathbf{h}, \{\mathbf{v}_k\}_{k=1}^K, \{\mathbf{h}_k\}_{k=1}^K). \tag{10}$$

Lastly, RAM weights them together and obtains the final retrieval-augmented user representation:

$$\tilde{\mathbf{h}} = \alpha \mathbf{h} + (1-\alpha)(\beta \mathbf{h}^{c_1} + (1-\beta)\mathbf{h}^{c_2}), \tag{11}$$

where $0 \le \alpha, \beta \le 1$ control the strength of each representation. A possible extension is to learn $\alpha$ and $\beta$, *e.g.,* designing attention mechanism [50]. However, it is not the focus of this work, and we leave it for future exploration. We further fine-tune the model with RAFT to improve recommendations with the retrieval-augmented representation $\tilde{\mathbf{h}}$. Specifically, we freeze the model backbone SeqEnc$(\cdot)$ and only the RAM is being updated during fine-tuning. Formally, we adopt the same objective as Equ (5), which learns to leverage the retrieved memories for better recommendation:

$$\mathcal{L}_{raft} = -\log \Big( e^{\tilde{\mathbf{h}}^T \mathbf{v}_{t+1}} / \sum_{v_i \in \mathcal{V}} e^{\tilde{\mathbf{h}}^T \mathbf{v}_i} \Big). \tag{12}$$

In this way, RaSeRec does not have to memorize all long-tailed patterns but learns to leverage retrieved explicit memories to accommodate preference drift.

## 3.3 Inference and Complexity Analyses

In this section, we first introduce the model inference process and then analyze the model complexity in terms of both time and space in detail.

*3.3.1 Model Inference.* During inference, RaSeRec first computes the user representation $\mathbf{h}$ via SeqEnc$(\cdot)$ for the input user sequence $s_u$. Then, RaSeRec retrieves top-$K$ similar memories (*i.e.,* $\{\mathbf{h}_k\}_{k=1}^K$ and $\{\mathbf{v}_k\}_{k=1}^K$) from the memory bank $\mathcal{M}$. After that, RaSeRec employs RAM that consists of dual-channel MHCA to weight and aggregate valuable patterns from both $\{\mathbf{h}_k\}_{k=1}^K$ and $\{\mathbf{v}_k\}_{k=1}^K$. We obtain the final retrieval-augmented user representation $\tilde{\mathbf{h}}$ by combining both the original user representation $\mathbf{h}$ and the augmented user ones of the dual channels *i.e.,* $\mathbf{h}^{c_1}$ and $\mathbf{h}^{c_2}$. Lastly, we recommend the top-$N$ items based on the dot product scores between $\tilde{\mathbf{h}}$ and item embedding matrix $\mathbf{V}$.

**Table 1: The space and time complexity comparison between the Vanilla and SSL/Retrieval-Augmented SeRec.**

| Paradigms | Space | | Time | | | |
|---|---|---|---|---|---|---|
| | Backbone | RAM | Encode | Retrieve | Augment | Recommend |
| Vanilla SeRec | $O(\|\theta\|)$ | - | $O(L(Td^2 + T^2d))$ | - | - | $O(\|\mathcal{V}\|d)$ |
| SSL-Augmented SeRec | $O(\|\theta\|)$ | - | $O(L(Td^2 + T^2d))$ | - | - | $O(\|\mathcal{V}\|d)$ |
| Retrieval-Augmented SeRec | $O(\|\theta\|)$ | $O(\|\phi\|)$ | $O(L(Td^2 + T^2d))$ | $O(kd + \frac{\|\mathcal{V}\|}{k}d)$ | $O(Kd)$ | $O(\|\mathcal{V}\|d)$ |

*3.3.2 Complexity Analyses.* Table 1 summarizes the space and time complexity of three SeRec paradigms, where we compare the complexity of RaSeRec with the Vanilla and SSL-Augmented SeRec to manifest the efficiency of RaSeRec. For space complexity, RaSeRec needs only $O(\|\phi\| + \|\theta\|)$, where $\phi$ and $\theta$ represent the trainable parameters of the RAM and sequence encoder, respectively. Generally, $\|\phi\|$ is far below than $\|\theta\|$. As such, the analytical space complexity of RaSeRec is in the same magnitude as the Vanilla and SSL-Augmented one ($O(\|\theta\|)$). In practice, taking the Amazon Sports dataset as an example, the space occupation of RaSeRec is merely 0.05 times larger than theirs.

As recommender systems mainly focus on the time cost of online serving, we analyze the time complexity during inference. The time complexity for encoding user representations is $O(L(Td^2 + T^2d))$, where $T$ is the maximum sequence length and $d$ is the hidden size. In addition to encoding, RaSeRec also needs to retrieve memories and augment user representation. For retrieval, RaSeRec instantiates Faiss using Inverted File Indexing (IVF), whose retrieval complexity is $O(kd + \frac{\|\mathcal{V}\|}{k}d)$, where $k$ denotes the number of clusters and we set the number of clusters accessed as 1. On the other hand, the time complexity of augmenting user representation is $O(Kd)$. Finally, these three SeRec paradigms need to compute dot product scores over all items to make recommendations, resulting in a time complexity $O(\|\mathcal{V}\|d)$. The time complexity of additional retrieval and augmentation processes (i.e., $O(kd + \frac{\|\mathcal{V}\|}{k}d + Kd)$) is much lower than the encoding and recommendation processes (i.e., $O(L(Td^2 + T^2d) + \|\mathcal{V}\|d)$). Specifically, $k + \frac{\|\mathcal{V}\|}{k} + K$ is on the order of thousands, while $\|\mathcal{V}\|$ is on the order of tens of thousands. Therefore, the analytical time complexity of RaSeRec is of the same magnitude as the Vanilla and SSL-Augmented SeRec.

## 4 Experiments

To verify the effectiveness of RaSeRec, we conduct extensive experiments on three benchmark datasets to answer the following research questions (**RQs**):

- **RQ1:** How does RaSeRec perform compared with *state-of-the-art* SeRec models?
- **RQ2:** How much gain can RaSeRec bring to the existing base backbones?
- **RQ3:** Can RaSeRec improve long-tailed recommendation?
- **RQ4:** Can RaSeRec alleviate preference drift?
- **RQ5:** How do different partitions of the memory bank contribute to RaSeRec's performance?
- **RQ6:** Can RaSeRec perform robustly against the data noise issue?
- **RQ7:** How does the performance of RaSeRec vary with different hyper-parameter values? (Appendix F)

- **RQ8:** Whether RaSeRec benefits both high-frequency and low-frequency users? (Appendix G)

### 4.1 Experimental Settings

*4.1.1 Datasets.* We conduct extensive experiments on three benchmark datasets collected from Amazon [31]. We adopt three subcategories, *i.e.*, Beauty, Sports, and Clothing[1], with different sparsity degrees. Following [21, 41], we discard users and items with fewer than 5 interactions. The statistics of datasets after preprocessing are provided in Appendix B.

*4.1.2 Baselines.* To verify the effectiveness of RaSeRec, we compare it with the following three groups of competitive models: **(1) Non-sequential recommendation models (Non-SeRec)** include PopRec and BPR-MF [38]; **(2) Vanilla sequential recommendation models (Vanilla SeRec)** include GRU4Rec [15], Caser [42], SASRec [21], BERT4Rec [41], and $S^3Rec_{MIP}$ [63]; and **(3) SSL-Augmented sequential recommendation models (SSL-Augmented SeRec)** include CL4SRec [53], CoSeRec [29], ICLRec [4], DuoRec [34], and MCLRec [33]. Detailed descriptions of each baseline model can be seen in Appendix C.

*4.1.3 Implementation and Evaluation.* We use Recbole [57] to implement baselines. For models with learnable item embedding, we set the hidden size $d = 64$. For each baseline, we follow the suggested settings reported in their original paper to set hyperparameters. More implementation details can be seen in Appendix E. Following the leave-one-out strategy, we hold out the last interacted item of each user sequence for testing, the second last item for validation, and all the earlier ones for training. We measure the recommendation performance via HR and NDCG (see Appendix D), where the cut-off position (@$N$) is set as 5 and 10.

### 4.2 Comparison with Baselines (RQ1)

Table 2 and 8 summarize the performance of all baselines and RaSeRec *w.r.t.* @5 and @10, respectively. From the experimental results, we observe that: **(1)** The performance of the Non-SeRec models is unsurprisingly worse than the vanilla SeRec models, indicating the necessity of modeling sequential patterns for capturing more accurate user preferences. On the other hand, SSL-Augmented SeRec models perform better than those without SSL, indicating the superiority of supplementing the recommendation task with self-supervised learning. **(2)** Benefiting from the retrieval-augmented mechanism, RaSeRec considerably outperforms all baseline models in 16 out of 18 cases with $p$-value less than 0.005, demonstrating the improvement of RaSeRec over the SOTA baselines is statistically significant. The reasons are mainly twofold: **(i)** By sustainably

---

[1]http://jmcauley.ucsd.edu/data/amazon/

**Table 2: Performance comparison (@5), where the best results are boldfaced and the second-best ones are underlined.**

| Dataset | Beauty | | Sports | | Clothing | | Average | |
|---|---|---|---|---|---|---|---|---|
| Method | HR@5 | NDCG@5 | HR@5 | NDCG@5 | HR@5 | NDCG@5 | HR@5 | NDCG@5 |
| PopRec | 0.0072 | 0.0040 | 0.0055 | 0.0040 | 0.0030 | 0.0018 | 0.0052 | 0.0033 |
| BPR-MF | 0.0120 | 0.0065 | 0.0092 | 0.0053 | 0.0067 | 0.0052 | 0.0093 | 0.0057 |
| GRU4Rec | 0.0164 | 0.0086 | 0.0137 | 0.0096 | 0.0095 | 0.0061 | 0.0132 | 0.0081 |
| Caser | 0.0259 | 0.0127 | 0.0139 | 0.0085 | 0.0108 | 0.0067 | 0.0169 | 0.0093 |
| SASRec | 0.0365 | 0.0236 | 0.0218 | 0.0127 | 0.0168 | 0.0091 | 0.0250 | 0.0151 |
| BERT4Rec | 0.0193 | 0.0187 | 0.0176 | 0.0105 | 0.0125 | 0.0075 | 0.0165 | 0.0122 |
| $S^3Rec_{MIP}$ | 0.0327 | 0.0175 | 0.0157 | 0.0098 | 0.0163 | 0.0101 | 0.0216 | 0.0125 |
| CL4SRec | 0.0394 | 0.0246 | 0.0238 | 0.0146 | 0.0163 | 0.0098 | 0.0265 | 0.0163 |
| CoSeRec | 0.0463 | 0.0303 | 0.0278 | 0.0188 | 0.0161 | 0.0108 | 0.0301 | 0.0200 |
| ICLRec | 0.0469 | 0.0305 | 0.0271 | 0.0182 | 0.0165 | 0.0104 | 0.0302 | 0.0197 |
| DuoRec | 0.0541 | 0.0337 | 0.0315 | 0.0196 | 0.0191 | 0.0109 | 0.0349 | 0.0214 |
| MCLRec | 0.0552 | 0.0347 | 0.0294 | 0.0187 | **0.0197** | 0.0110 | 0.0348 | 0.0215 |
| RaSeRec (Ours) | **0.0569** | **0.0369** | **0.0331** | **0.0211** | 0.0194 | **0.0118** | **0.0365** | **0.0233** |
| %Improv. | 3.08% | 6.34% | 5.08% | 7.65% | -0.02% | 7.27% | 4.58% | 8.37% |
| $p$-value. | $8e^{-4}$ | $3e^{-5}$ | $1e^{-4}$ | $1e^{-5}$ | - | $2e^{-3}$ | $4e^{-4}$ | $5e^{-6}$ |

**Table 3: Performance comparison in terms of different sequential recommendation model backbones.**

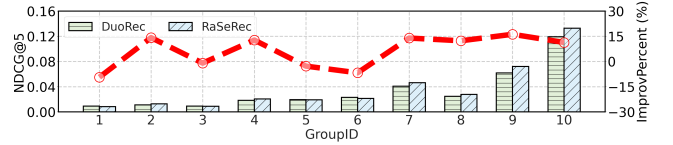| Backbone | Model | Beauty | | | | Sports | | | |
|---|---|---|---|---|---|---|---|---|---|
| | | HR@5 | HR@10 | NDCG@5 | NDCG@10 | HR@5 | HR@10 | NDCG@5 | NDCG@10 |
| GRU4Rec | Base | 0.0164 | 0.0365 | 0.0086 | 0.0142 | 0.0137 | 0.0274 | 0.0096 | 0.0137 |
| | CL4SRec | 0.0314 | 0.0506 | 0.0196 | 0.0257 | 0.0196 | 0.0317 | 0.0126 | 0.0165 |
| | DuoRec | 0.0443 | 0.0691 | 0.0296 | 0.0377 | 0.0248 | 0.0386 | 0.0162 | 0.0206 |
| | RaSeRec | **0.0469** | **0.0718** | **0.0322** | **0.0403** | **0.0261** | **0.0400** | **0.0173** | **0.0218** |
| BERT4Rec | Base | 0.0193 | 0.0401 | 0.0187 | 0.0254 | 0.0176 | 0.0326 | 0.0105 | 0.0153 |
| | CL4SRec | 0.0322 | 0.0499 | 0.0195 | 0.0253 | 0.0226 | 0.0352 | 0.0132 | 0.0173 |
| | DuoRec | 0.0505 | 0.0759 | 0.0318 | 0.0400 | 0.0275 | 0.0436 | 0.0170 | 0.0222 |
| | RaSeRec | **0.0534** | **0.0787** | **0.0344** | **0.0426** | **0.0288** | **0.0451** | **0.0182** | **0.0236** |

maintaining a dynamic memory bank, RaSeRec can quickly adapt to preference drifts. **(ii)** By retrieving memories, RaSeRec can explicitly recall long-tailed patterns that may be overwhelmed by head ones. These above two innovative designs contribute to considerable improvements together.

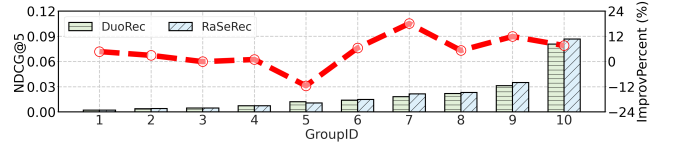## 4.3 Improving Base Backbones (RQ2)

RaSeRec is a model-agnostic retrieval-augmented SeRec paradigm that can enhance various base models. To verify this, instead of using the default Self-attentive backbone SASRec [21], we evaluate RaSeRec using other backbones, *i.e.*, RNN-based GRU4Rec [15] and Bert-based BERT4Rec [41]. As the average performance of DuoRec and MCLRec is close, we mainly compare with DuoRec in the following. As shown in Table 3, RaSeRec performs best with different backbones, indicating that RaSeRec has a good generalization ability. The impressive performance gain demonstrates the effectiveness of our retrieval-augmented mechanism, even though we just add a very small set of parameters (about 5%). In summary, RaSeRec possesses the virtue of the plug-and-play property, which can enhance various given base models with retrieval augmentation while leaving the base model backbone unchanged.

## 4.4 Long-tailed Recommendation (RQ3)

As stated in Section 1, existing SeRec models make it hard to recall long-tailed patterns via implicit memories due to skewed supervised signal distribution. To verify whether RaSeRec is of promise in solving this issue, we split the test user sequences into 10 groups based on the target item's popularity measured by the number



(a) Results on the Beauty dataset *w.r.t.* NDCG@5.



(b) Results on the Sports dataset *w.r.t.* NDCG@5.

**Figure 4: Performance comparison over different item groups between RaSeRec and DuoRec. The bar represents NDCG@5, while the line represents the performance improvement percentage of RaSeRec compared with DuoRec.**

of interactions, meanwhile keeping the total number of test user sequences in each group the same. And the larger the GroupId, the more popular the target item. The results are shown in Figure 4. From the results, we observe that: **(1)** Compared with DuoRec, the performance improvements of RaSeRec mostly come from accurately recommending long-tailed and head items. This verifies that retrieval-augmented user representation well recalls both long-tailed and head patterns. **(2)** We observe that RaSeRec performs relatively worse in recommending general items, *i.e.*, 5-th and 6-th

(a) Results on the Beauty dataset *w.r.t.* HR@5 and NDCG@5.

(b) Results on the Sports dataset *w.r.t.* HR@5 and NDCG@5.

(c) Results on the Clothing dataset *w.r.t.* HR@5 and NDCG@5.

**Figure 5: Performance comparison when ablating different partitions of the memory bank $\mathcal{M}$. The bar represents HR@5 or NDCG@5, while the line represents the percentage of performance degradation compared to the "All".**

**Table 4: Model performance *w.r.t.* preference drifts.**

| Metrics | | HR | | NDCG | |
|---|---|---|---|---|---|
| **Datasets** | **Drifts** | **@5** | **@10** | **@5** | **@10** |
| **Beauty** | (A) Full | **0.0569** | **0.0860** | **0.0369** | **0.0463** |
| | (B) 10% | 0.0560 | 0.0842 | 0.0360 | 0.0450 |
| | (C) 20% | 0.0550 | 0.0843 | 0.0355 | 0.0449 |
| | (D) 30% | 0.0542 | 0.0843 | 0.0350 | 0.0447 |
| **Sports** | (A) Full | **0.0331** | **0.0497** | **0.0211** | **0.0264** |
| | (B) 10% | 0.0326 | 0.0492 | 0.0208 | 0.0261 |
| | (C) 20% | 0.0322 | 0.0490 | 0.0204 | 0.0258 |
| | (D) 30% | 0.0328 | 0.0490 | 0.0208 | 0.0260 |

**Table 5: Statistics of each partition.**

| Dataset | #len.range | #avg.length | %proportion |
|---|---|---|---|
| Beauty | <3 | 1.50 | 34.03% |
| | 3-6 | 4.15 | 31.18% |
| | >6 | 18.88 | 34.79% |
| Sports | <3 | 1.50 | 37.56% |
| | 3-6 | 4.14 | 33.98% |
| | >6 | 15.31 | 28.46% |
| Clothing | <2 | 1.00 | 24.54% |
| | 2-3 | 2.38 | 39.38% |
| | >3 | 8.06 | 36.08% |

groups. We think the main reason is that the relevant memories are not retrieved while the noisy memories are introduced. How to effectively retrieve relevant knowledge is also an open challenge in RAG [27, 46, 59]. In summary, we suggest performing retrieval augmentation when recommending the long-tailed and head items.

## 4.5 Alleviating Preference Drift (RQ4)

As stated in Section 1, existing SeRec models trained on past data may recommend undesirable items due to preference drift. To verify whether RaSeRec is promising in alleviating this issue, we simulate this scenario by removing a certain ratio (*i.e.,* 10%, 20%, and 30%) of the preference data with the latest timestamp from the memory bank. And, "Full" denotes using all preference data in the memory bank to remedy preference drift. From the results shown in Table 4, we find that (1) In general, the larger the preference drifts, the worse the model performance; (2) "Full" performs best as it can remedy preference drift by maintaining a dynamic memory bank with the latest preference data; and (3) In a few cases, large preference drift
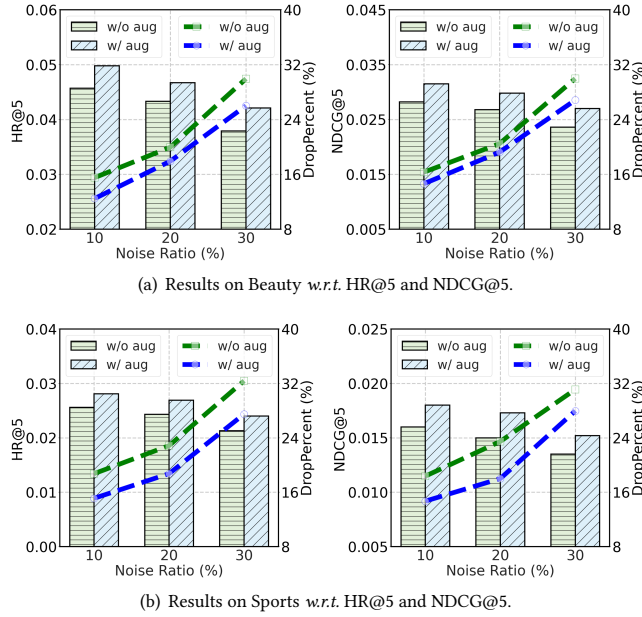
is not necessarily worse performance, *e.g.,* (D) on the sports dataset. The main reason is different datasets have different preferences towards short-, medium-, and long-term memories, which also greatly affects performance (more details can be seen in §4.6). In a nutshell, the above observations verify that RaSeRec can alleviate preference drift effectively.

## 4.6 Study on Memory Bank (RQ5)

To investigate the influence of different partitions of the memory bank on the model performance, we divide the memory bank $\mathcal{M}$ into three partitions according to the length of user sequences. The statistics of each partition are summarized in Table 5. For convenience, we use S, M, and L to represent the <u>S</u>hort-term, <u>M</u>edium-term, and <u>L</u>ong-term memory bank, respectively. For example, on the beauty dataset, we treat user sequences with length less than 3, those between 3 and 6, and those greater than 6 as S, M, and L memory bank respectively, where the proportion of different partitions accounts for about 1/3. The experimental results are shown in Figure 5. "All" denotes using all parts of the memory bank, including

(a) Results on Beauty *w.r.t.* HR@5 and NDCG@5.



(b) Results on Sports *w.r.t.* HR@5 and NDCG@5.

**Figure 6: Model performance *w.r.t.* noise ratio on Beauty and Sports datasets. The bar represents HR@5 or NDCG@5, while the line represents the percentage of performance degradation compared to the test user sequences without added noise.**

S, M, and L, while {S, L} denotes using parts of the memory bank, *i.e.,* S and L, and so on. From the experimental results, we observe that different datasets show different preferences toward the combinations of the partitions. Specifically, the beauty dataset prefers "All"; the sports dataset prefers {S, M}; the clothing dataset prefers {S, L}. We think the main reason is associated with #avg.length.

As for the beauty dataset, its average length is longer than the other two datasets and it prefers the L partition. Specifically, the models with L partition largely outperform those without L partition, as shown in Figure 5(a). As for the sports and clothing datasets, their average length is shorter than that of the beauty dataset, and they perform relatively worse when only provided with L partition, as shown in Figure 5(b) and 5(c). Besides, we observe that the model equipped with parts of the memory bank (*e.g.,* {S, M} is the best for the sports dataset and {S, L} is the best for the clothing dataset) even outperforms that with all parts of the memory bank. The above observations demonstrate that not all memories are beneficial and some ones are harmful to recommendation performance. This is also a challenging issue [6, 9, 51] in RAG that warrants further research. In a nutshell, we suggest trying different combinations or directly using all of the memories, because "All" usually leads to satisfactory performance. In this work, we use "All" by default.

### 4.7 Robustness against Noise Data (RQ6)

In real-world scenarios, recommender systems usually suffer from the noisy interaction issue [11, 43], *e.g.,* a user clicked an item by mistake. As such, it is necessary to verify RaSeRec's robustness against noisy data. To simulate this scenario, we train models with full training data and randomly add a certain ratio (*i.e.,* 10%, 20%,

and 30%) of negative items into each test user sequence. The experimental results are shown in Figure 6, where models with retrieval augmentation mention "w/ aug" otherwise "w/o aug". From the experimental results, we observe that adding noises significantly degrades the performance of models with or without retrieval augmentation. However, "w/ aug" consistently outperforms "w/o aug" under any noise ratio in terms of HR@5 and NDCG@5. Moreover, the blue line is always below the green line, as shown in Figure 6, demonstrating that "w/ aug" can endow the model backbone with more robustness against the noisy interactions during inference. For example, on the Beauty and Sports datasets, "w/ aug" with a 20% noise ratio outperforms "w/o aug" with a 10% noise ratio. Moreover, on the Sports dataset, "w/ aug" drops 14.7% of its original performance, while "w/o aug" drops 18.4%, in terms of NDCG@5, when both have a 10% noise ratio. This fully implies that with retrieval augmentation, RaSeRec can retrieve useful sequential patterns as references to alleviate the influence of noisy interactions and make a suitable recommendation. In a nutshell, we suggest performing retrieval augmentation to alleviate the negative influence of noisy interactions.

## 5 Related Works

### 5.1 Retrieval-Augmented Generation

Recently, RAG prevails in LLMs, introducing non-parametric knowledge into LLMs to supplement their incomplete, incorrect, or outdated parametric knowledge [12]. It has been shown highly effective in many research fields, not limited to natural language, including computer vision [36, 37, 40, 52], speech [45, 54], code [28, 47], and graph [8]. Although being well-studied in many fields, very limited works exploit RAG in recommender systems, let alone SeRec. A very recent one is RUEL [49]. However, it differs from our work in **(1)** it lacks retrieval training, and **(2)** it falls short in decoupling implicit and explicit memories, while our work proposes a model-agnostic collaborative-based pre-training and retrieval-augmented fine-tuning paradigm for SeRec that can be applied to any ID-based SeRec backbones and improve their performance.

### 5.2 Sequential Recommendation

Pioneering attempts on SeRec are based on Markov Chain (MC) to model item-item transition relationships, *e.g.,* FPMC [39]. To capture long-term dependencies, some studies first adopt RNN for SeRec by modeling sequence-level item transitions, *e.g.,* GRU4Rec [15]. Simultaneously, CNN has been explored for SeRec by modeling the union-level sequential patterns [42]. Recently, the success of self-attention models in NLP [44] has spawned a series of Transformer-based sequential recommendation models, *e.g.,* SASRec [21] and BERT4Rec [41]. Considering sparse supervised signals, SSL has been widely adopted in SeRec to alleviate this issue [33, 34, 53]. Enlightened by generative models, some studies have attempted to apply LLMs to solve SeRec [25, 26, 61, 62]. Albeit studied for ages, almost all existing methods focus on exploiting implicit memories hidden in the model. This work is in an orthogonal direction. It exploits explicit memories in a retrieval-augmented manner and opens a new research direction for SeRec.

# 6 Conclusion and Future Works

In this work, we unveiled the limitations of the Vanilla SeRec and SSL-Augmented SeRec paradigms and explored the potential of RAG to solve them. In particular, we propose a retrieval-augmented SeRec paradigm, which learns to refine user representation with explicit retrieved memories and adapt to preference drifts by maintaining a dynamic memory bank. Extensive experiments on three datasets manifest the advantage of RaSeRec in terms of recommendation performance, long-tailed recommendation, and alleviating preference drift. Despite our innovations and improvements, we recognize some issues that warrant further study. In this work, we fixed the number of retrieved memories $K$ during training and inference. In Appendix F, we have studied the impact of different $K$ values on the model performance. We find both too large and small $K$ values hurt the model performance. Furthermore, in Section 4.4 and Appendix G, we have studied the impact of retrieval augmentation on different item groups and user groups. The experimental results show that retrieval augmentation may hurt the performance of some groups. These observations motivate us to develop an active retrieval augmentation mechanism for RaSeRec, *i.e.,* retrieval when needed.

# References

[1] Mehdi Hosseinzadeh Aghdam, Negar Hariri, Bamshad Mobasher, and Robin D. Burke. 2015. Adapting Recommendations to Contextual Changes Using Hierarchical Hidden Markov Models. In *Proceedings of the 9th ACM Conference on Recommender Systems, RecSys 2015, Vienna, Austria, September 16-20, 2015*, Hannes Werthner, Markus Zanker, Jennifer Golbeck, and Giovanni Semeraro (Eds.). ACM, 241–244. doi:10.1145/2792838.2799684

[2] 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, ICLR 2024, Vienna, Austria, May 7-11, 2024*. OpenReview.net. https://openreview.net/forum?id=hSyW5go0v8

[3] Chao Chen, Dongsheng Li, Junchi Yan, and Xiaokang Yang. 2022. Modeling Dynamic User Preference via Dictionary Learning for Sequential Recommendation. *IEEE Trans. Knowl. Data Eng.* 34, 11 (2022), 5446–5458. doi:10.1109/TKDE.2021.3050407

[4] Yongjun Chen, Zhiwei Liu, Jia Li, Julian J. McAuley, and Caiming Xiong. 2022. Intent Contrastive Learning for Sequential Recommendation. In *WWW '22: The ACM Web Conference 2022, Virtual Event, Lyon, France, April 25 - 29, 2022*, Frédérique Laforest, Raphaël Troncy, Elena Simperl, Deepak Agarwal, Aristides Gionis, Ivan Herman, and Lionel Médini (Eds.). ACM, 2172–2182. doi:10.1145/3485447.3512090

[5] Aaron Clauset, Cosma Rohilla Shalizi, and Mark E. J. Newman. 2009. Power-Law Distributions in Empirical Data. *SIAM Rev.* 51, 4 (2009), 661–703. doi:10.1137/070710111

[6] Florin Cuconasu, Giovanni Trappolini, Federico Siciliano, Simone Filice, Cesare Campagnano, Yoelle Maarek, Nicola Tonellotto, and Fabrizio Silvestri. 2024. The Power of Noise: Redefining Retrieval for RAG Systems. In *Proceedings of the 47th International ACM SIGIR Conference on Research and Development in Information Retrieval, SIGIR 2024, Washington DC, USA, July 14-18, 2024*, Grace Hui Yang, Hongning Wang, Sam Han, Claudia Hauff, Guido Zuccon, and Yi Zhang (Eds.). ACM, 719–729. doi:10.1145/3626772.3657834

[7] Yizhou Dang, Enneng Yang, Guibing Guo, Linying Jiang, Xingwei Wang, Xiaoxiao Xu, Qinghui Sun, and Hong Liu. 2023. Uniform Sequence Better: Time Interval Aware Data Augmentation for Sequential Recommendation. In *Thirty-Seventh AAAI Conference on Artificial Intelligence, AAAI 2023, Thirty-Fifth Conference on Innovative Applications of Artificial Intelligence, IAAI 2023, Thirteenth Symposium on Educational Advances in Artificial Intelligence, EAAI 2023, Washington, DC, USA, February 7-14, 2023*, Brian Williams, Yiling Chen, and Jennifer Neville (Eds.). AAAI Press, 4225–4232. doi:10.1609/AAAI.V37I4.25540

[8] Darren Edge, Ha Trinh, Newman Cheng, Joshua Bradley, Alex Chao, Apurva Mody, Steven Truitt, and Jonathan Larson. 2024. From Local to Global: A Graph RAG Approach to Query-Focused Summarization. *CoRR* abs/2404.16130 (2024). doi:10.48550/ARXIV.2404.16130 arXiv:2404.16130

[9] Feiteng Fang, Yuelin Bai, Shiwen Ni, Min Yang, Xiaojun Chen, and Ruifeng Xu. 2024. Enhancing Noise Robustness of Retrieval-Augmented Language Models with Adaptive Adversarial Training. In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), ACL 2024, Bangkok, Thailand, August 11-16, 2024*, Lun-Wei Ku, Andre Martins, and Vivek Srikumar (Eds.). Association for Computational Linguistics, 10028–10039. doi:10.18653/V1/2024.ACL-LONG.540

[10] Hui Fang, Danning Zhang, Yiheng Shu, and Guibing Guo. 2020. Deep Learning for Sequential Recommendation: Algorithms, Influential Factors, and Evaluations. *ACM Trans. Inf. Syst.* 39, 1 (2020), 10:1–10:42. doi:10.1145/3426723

[11] Yunjun Gao, Yuntao Du, Yujia Hu, Lu Chen, Xinjun Zhu, Ziquan Fang, and Baihua Zheng. 2022. Self-Guided Learning to Denoise for Robust Recommendation. In *SIGIR '22: The 45th International ACM SIGIR Conference on Research and Development in Information Retrieval, Madrid, Spain, July 11 - 15, 2022*, Enrique Amigó, Pablo Castells, Julio Gonzalo, Ben Carterette, J. Shane Culpepper, and Gabriella Kazai (Eds.). ACM, 1412–1422. doi:10.1145/3477495.3532059

[12] 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. *CoRR* abs/2312.10997 (2023). doi:10.48550/ARXIV.2312.10997 arXiv:2312.10997

[13] Michael Gutmann and Aapo Hyvärinen. 2010. Noise-contrastive estimation: A new estimation principle for unnormalized statistical models. In *Proceedings of the Thirteenth International Conference on Artificial Intelligence and Statistics, AISTATS 2010, Chia Laguna Resort, Sardinia, Italy, May 13-15, 2010 (JMLR Proceedings, Vol. 9)*, Yee Whye Teh and D. Mike Titterington (Eds.). JMLR.org, 297–304. http://proceedings.mlr.press/v9/gutmann10a.html

[14] Yongjing Hao, Pengpeng Zhao, Junhua Fang, Jianfeng Qu, Guanfeng Liu, Fuzhen Zhuang, Victor S. Sheng, and Xiaofang Zhou. 2024. Meta-Optimized Joint Generative and Contrastive Learning for Sequential Recommendation. In *40th IEEE International Conference on Data Engineering, ICDE 2024, Utrecht, The Netherlands, May 13-16, 2024*. IEEE, 705–718. doi:10.1109/ICDE60146.2024.00060

[15] Balázs Hidasi, Alexandros Karatzoglou, Linas Baltrunas, and Domonkos Tikk. 2016. Session-based Recommendations with Recurrent Neural Networks. In *4th International Conference on Learning Representations, ICLR 2016, San Juan, Puerto Rico, May 2-4, 2016, Conference Track Proceedings*, Yoshua Bengio and Yann LeCun (Eds.). http://arxiv.org/abs/1511.06939

[16] Folasade Olubusola Isinkaye, Yetunde O Folajimi, and Bolande Adefowoke Ojokoh. 2015. Recommendation systems: Principles, methods and evaluation. *Egyptian informatics journal* 16, 3 (2015), 261–273. https://www.sciencedirect.com/science/article/pii/S1110866515000341?via%3Dihub

[17] Gautier Izacard and Edouard Grave. 2021. Distilling Knowledge from Reader to Retriever for Question Answering. In *9th International Conference on Learning Representations, ICLR 2021, Virtual Event, Austria, May 3-7, 2021*. OpenReview.net. https://openreview.net/forum?id=NTEz-6wysdb

[18] Gautier Izacard, Patrick S. H. Lewis, Maria Lomeli, Lucas Hosseini, Fabio Petroni, Timo Schick, Jane Dwivedi-Yu, Armand Joulin, Sebastian Riedel, and Edouard Grave. 2023. Atlas: Few-shot Learning with Retrieval Augmented Language Models. *J. Mach. Learn. Res.* 24 (2023), 251:1–251:43. https://jmlr.org/papers/v24/23-0037.html

[19] Bowen Jin, Jinsung Yoon, Jiawei Han, and Sercan O Arik. 2024. Long-Context LLMs Meet RAG: Overcoming Challenges for Long Inputs in RAG. *arXiv preprint arXiv:2410.05983* (2024). https://arxiv.org/pdf/2410.05983

[20] Jeff Johnson, Matthijs Douze, and Hervé Jégou. 2021. Billion-Scale Similarity Search with GPUs. *IEEE Trans. Big Data* 7, 3 (2021), 535–547. doi:10.1109/TBDATA.2019.2921572

[21] Wang-Cheng Kang and Julian J. McAuley. 2018. Self-Attentive Sequential Recommendation. In *IEEE International Conference on Data Mining, ICDM 2018, Singapore, November 17-20, 2018*. IEEE Computer Society, 197–206. doi:10.1109/ICDM.2018.00035

[22] Vladimir Karpukhin, Barlas Oguz, Sewon Min, Patrick S. H. Lewis, Ledell Wu, Sergey Edunov, Danqi Chen, and Wen-tau Yih. 2020. Dense Passage Retrieval for Open-Domain Question Answering. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing, EMNLP 2020, Online, November 16-20, 2020*, Bonnie Webber, Trevor Cohn, Yulan He, and Yang Liu (Eds.). Association for Computational Linguistics, 6769–6781. doi:10.18653/V1/2020.EMNLP-MAIN.550

[23] Diederik P. Kingma and Jimmy Ba. 2015. Adam: A Method for Stochastic Optimization. In *3rd International Conference on Learning Representations, ICLR 2015, San Diego, CA, USA, May 7-9, 2015, Conference Track Proceedings*, Yoshua Bengio and Yann LeCun (Eds.). http://arxiv.org/abs/1412.6980

[24] Patrick S. H. Lewis, Ethan Perez, Aleksandra Piktus, Fabio Petroni, Vladimir Karpukhin, Naman Goyal, Heinrich Küttler, Mike Lewis, Wen-tau Yih, Tim Rocktäschel, Sebastian Riedel, and Douwe Kiela. 2020. Retrieval-Augmented Generation for Knowledge-Intensive NLP Tasks. In *Advances in Neural Information Processing Systems 33: Annual Conference on Neural Information Processing Systems 2020, NeurIPS 2020, December 6-12, 2020, virtual*, Hugo Larochelle, Marc'Aurelio Ranzato, Raia Hadsell, Maria-Florina Balcan, and Hsuan-Tien Lin (Eds.). https://proceedings.neurips.cc/paper/2020/hash/6b493230205f780e1bc26945df7481e5-Abstract.html

[25] Jiacheng Li, Ming Wang, Jin Li, Jinmiao Fu, Xin Shen, Jingbo Shang, and Julian J. McAuley. 2023. Text Is All You Need: Learning Language Representations for

Sequential Recommendation. In *Proceedings of the 29th ACM SIGKDD Conference on Knowledge Discovery and Data Mining, KDD 2023, Long Beach, CA, USA, August 6-10, 2023*, Ambuj K. Singh, Yizhou Sun, Leman Akoglu, Dimitrios Gunopulos, Xifeng Yan, Ravi Kumar, Fatma Ozcan, and Jieping Ye (Eds.). ACM, 1258–1267. doi:10.1145/3580305.3599519

[26] Jinming Li, Wentao Zhang, Tian Wang, Guanglei Xiong, Alan Lu, and Gerard Medioni. 2023. GPT4Rec: A Generative Framework for Personalized Recommendation and User Interests Interpretation. In *Proceedings of the 2023 SIGIR Workshop on eCommerce co-located with the 46th International ACM SIGIR Conference on Research and Development in Information Retrieval (SIGIR 2023), Taipei, Taiwan, July 27, 2023 (CEUR Workshop Proceedings, Vol. 3589)*, Surya Kallumadi, Yubin Kim, Tracy Holloway King, Shervin Malmasi, Maarten de Rijke, and Jacopo Tagliabue (Eds.). CEUR-WS.org. https://ceur-ws.org/Vol-3589/paper_2.pdf

[27] Mingda Li, Xinyu Li, Yifan Chen, Wenfeng Xuan, and Weinan Zhang. 2024. Unraveling and Mitigating Retriever Inconsistencies in Retrieval-Augmented Large Language Models. In *Findings of the Association for Computational Linguistics, ACL 2024, Bangkok, Thailand and virtual meeting, August 11-16, 2024*, Lun-Wei Ku, Andre Martins, and Vivek Srikumar (Eds.). Association for Computational Linguistics, 4833–4850. doi:10.18653/V1/2024.FINDINGS-ACL.288

[28] Xinze Li, Zhenghao Liu, Chenyan Xiong, Shi Yu, Yu Gu, Zhiyuan Liu, and Ge Yu. 2023. Structure-Aware Language Model Pretraining Improves Dense Retrieval on Structured Data. In *Findings of the Association for Computational Linguistics: ACL 2023, Toronto, Canada, July 9-14, 2023*, Anna Rogers, Jordan L. Boyd-Graber, and Naoaki Okazaki (Eds.). Association for Computational Linguistics, 11560–11574. doi:10.18653/V1/2023.FINDINGS-ACL.734

[29] Zhiwei Liu, Yongjun Chen, Jia Li, Philip S. Yu, Julian J. McAuley, and Caiming Xiong. 2021. Contrastive Self-supervised Sequential Recommendation with Robust Augmentation. *CoRR* abs/2108.06479 (2021). arXiv:2108.06479 https://arxiv.org/abs/2108.06479

[30] Yung-Yin Lo, Wanjiun Liao, Cheng-Shang Chang, and Ying-Chin Lee. 2018. Temporal Matrix Factorization for Tracking Concept Drift in Individual User Preferences. *IEEE Trans. Comput. Soc. Syst.* 5, 1 (2018), 156–168. doi:10.1109/TCSS.2017.2772295

[31] Julian J. McAuley, Christopher Targett, Qinfeng Shi, and Anton van den Hengel. 2015. Image-Based Recommendations on Styles and Substitutes. In *Proceedings of the 38th International ACM SIGIR Conference on Research and Development in Information Retrieval, Santiago, Chile, August 9-13, 2015*, Ricardo Baeza-Yates, Mounia Lalmas, Alistair Moffat, and Berthier A. Ribeiro-Neto (Eds.). ACM, 43–52. doi:10.1145/2766462.2767755

[32] Stasa Milojevic. 2010. Power law distributions in information science: Making the case for logarithmic binning. *J. Assoc. Inf. Sci. Technol.* 61, 12 (2010), 2417–2425. doi:10.1002/ASI.21426

[33] Xiuyuan Qin, Huanhuan Yuan, Pengpeng Zhao, Junhua Fang, Fuzhen Zhuang, Guanfeng Liu, Yanchi Liu, and Victor S. Sheng. 2023. Meta-optimized Contrastive Learning for Sequential Recommendation. In *Proceedings of the 46th International ACM SIGIR Conference on Research and Development in Information Retrieval, SIGIR 2023, Taipei, Taiwan, July 23-27, 2023*, Hsin-Hsi Chen, Wei-Jou (Edward) Duh, Hen-Hsen Huang, Makoto P. Kato, Josiane Mothe, and Barbara Poblete (Eds.). ACM, 89–98. doi:10.1145/3539618.3591727

[34] Ruihong Qiu, Zi Huang, Hongzhi Yin, and Zijian Wang. 2022. Contrastive Learning for Representation Degeneration Problem in Sequential Recommendation. In *WSDM '22: The Fifteenth ACM International Conference on Web Search and Data Mining, Virtual Event / Tempe, AZ, USA, February 21 - 25, 2022*, K. Selcuk Candan, Huan Liu, Leman Akoglu, Xin Luna Dong, and Jiliang Tang (Eds.). ACM, 813–823. doi:10.1145/3488560.3498433

[35] Alec Radford, Karthik Narasimhan, Tim Salimans, and Ilya Sutskever. 2018. Improving Language Understanding by Generative Pre-Training. In *OpenAI*. https://cdn.openai.com/research-covers/language-unsupervised/language_understanding_paper.pdf

[36] Jun Rao, Liang Ding, Shuhan Qi, Meng Fang, Yang Liu, Li Shen, and Dacheng Tao. 2023. Dynamic contrastive distillation for image-text retrieval. *IEEE Transactions on Multimedia* 25 (2023), 8383–8395.

[37] Jun Rao, Fei Wang, Liang Ding, Shuhan Qi, Yibing Zhan, Weifeng Liu, and Dacheng Tao. 2022. Where Does the Performance Improvement Come From? -A Reproducibility Concern about Image-Text Retrieval. In *Proceedings of the 45th international ACM SIGIR conference on research and development in information retrieval*. 2727–2737.

[38] Steffen Rendle, Christoph Freudenthaler, Zeno Gantner, and Lars Schmidt-Thieme. 2009. BPR: Bayesian Personalized Ranking from Implicit Feedback. In *UAI 2009, Proceedings of the Twenty-Fifth Conference on Uncertainty in Artificial Intelligence, Montreal, QC, Canada, June 18-21, 2009*, Jeff A. Bilmes and Andrew Y. Ng (Eds.). AUAI Press, 452–461. https://www.auai.org/uai2009/papers/UAI2009_0139_48141db02b9f0b02bc7158819ebfa2c7.pdf

[39] Steffen Rendle, Christoph Freudenthaler, and Lars Schmidt-Thieme. 2010. Factorizing personalized Markov chains for next-basket recommendation. In *Proceedings of the 19th International Conference on World Wide Web, WWW 2010, Raleigh, North Carolina, USA, April 26-30, 2010*, Michael Rappa, Paul Jones, Juliana Freire, and Soumen Chakrabarti (Eds.). ACM, 811–820. doi:10.1145/1772690.1772773

[40] Sahel Sharifymoghaddam, Shivani Upadhyay, Wenhu Chen, and Jimmy Lin. 2024. UniRAG: Universal Retrieval Augmentation for Multi-Modal Large Language Models. *CoRR* abs/2405.10311 (2024). doi:10.48550/ARXIV.2405.10311 arXiv:2405.10311

[41] Fei Sun, Jun Liu, Jian Wu, Changhua Pei, Xiao Lin, Wenwu Ou, and Peng Jiang. 2019. BERT4Rec: Sequential Recommendation with Bidirectional Encoder Representations from Transformer. In *Proceedings of the 28th ACM International Conference on Information and Knowledge Management, CIKM 2019, Beijing, China, November 3-7, 2019*, Wenwu Zhu, Dacheng Tao, Xueqi Cheng, Peng Cui, Elke A. Rundensteiner, David Carmel, Qi He, and Jeffrey Xu Yu (Eds.). ACM, 1441–1450. doi:10.1145/3357384.3357895

[42] Jiaxi Tang and Ke Wang. 2018. Personalized Top-N Sequential Recommendation via Convolutional Sequence Embedding. In *Proceedings of the Eleventh ACM International Conference on Web Search and Data Mining, WSDM 2018, Marina Del Rey, CA, USA, February 5-9, 2018*, Yi Chang, Chengxiang Zhai, Yan Liu, and Yoelle Maarek (Eds.). ACM, 565–573. doi:10.1145/3159652.3159656

[43] Changxin Tian, Yuexiang Xie, Yaliang Li, Nan Yang, and Wayne Xin Zhao. 2022. Learning to Denoise Unreliable Interactions for Graph Collaborative Filtering. In *SIGIR '22: The 45th International ACM SIGIR Conference on Research and Development in Information Retrieval, Madrid, Spain, July 11 - 15, 2022*, Enrique Amigó, Pablo Castells, Julio Gonzalo, Ben Carterette, J. Shane Culpepper, and Gabriella Kazai (Eds.). ACM, 122–132. doi:10.1145/3477495.3531889

[44] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N. Gomez, Lukasz Kaiser, and Illia Polosukhin. 2017. Attention is All you Need. In *Advances in Neural Information Processing Systems 30: Annual Conference on Neural Information Processing Systems 2017, December 4-9, 2017, Long Beach, CA, USA*, Isabelle Guyon, Ulrike von Luxburg, Samy Bengio, Hanna M. Wallach, Rob Fergus, S. V. N. Vishwanathan, and Roman Garnett (Eds.). 5998–6008. https://proceedings.neurips.cc/paper/2017/hash/3f5ee243547dee91fbd053c1c4a845aa-Abstract.html

[45] Mingqiu Wang, Izhak Shafran, Hagen Soltau, Wei Han, Yuan Cao, Dian Yu, and Laurent El Shafey. 2024. Retrieval Augmented End-to-End Spoken Dialog Models. In *IEEE International Conference on Acoustics, Speech and Signal Processing, ICASSP 2024, Seoul, Republic of Korea, April 14-19, 2024*. IEEE, 12056–12060. doi:10.1109/ICASSP48485.2024.10447448

[46] Zhiruo Wang, Jun Araki, Zhengbao Jiang, Md. Rizwan Parvez, and Graham Neubig. 2023. Learning to Filter Context for Retrieval-Augmented Generation. *CoRR* abs/2311.08377 (2023). doi:10.48550/ARXIV.2311.08377 arXiv:2311.08377

[47] Zora Zhiruo Wang, Akari Asai, Xinyan Velocity Yu, Frank F. Xu, Yiqing Xie, Graham Neubig, and Daniel Fried. 2024. CodeRAG-Bench: Can Retrieval Augment Code Generation? *CoRR* abs/2406.14497 (2024). doi:10.48550/ARXIV.2406.14497 arXiv:2406.14497

[48] Jiancan Wu, Xiang Wang, Fuli Feng, Xiangnan He, Liang Chen, Jianxun Lian, and Xing Xie. 2021. Self-supervised Graph Learning for Recommendation. In *SIGIR '21: The 44th International ACM SIGIR Conference on Research and Development in Information Retrieval, Virtual Event, Canada, July 11-15, 2021*, Fernando Diaz, Chirag Shah, Torsten Suel, Pablo Castells, Rosie Jones, and Tetsuya Sakai (Eds.). ACM, 726–735. doi:10.1145/3404835.3462862

[49] Ning Wu, Ming Gong, Linjun Shou, Jian Pei, and Daxin Jiang. 2023. RUEL: Retrieval-Augmented User Representation with Edge Browser Logs for Sequential Recommendation. In *Proceedings of the 32nd ACM International Conference on Information and Knowledge Management, CIKM 2023, Birmingham, United Kingdom, October 21-25, 2023*, Ingo Frommholz, Frank Hopfgartner, Mark Lee, Michael Oakes, Mounia Lalmas, Min Zhang, and Rodrygo L. T. Santos (Eds.). ACM, 4871–4878. doi:10.1145/3583780.3615498

[50] Shu Wu, Yuyuan Tang, Yanqiao Zhu, Liang Wang, Xing Xie, and Tieniu Tan. 2019. Session-Based Recommendation with Graph Neural Networks. In *The Thirty-Third AAAI Conference on Artificial Intelligence, AAAI 2019, The Thirty-First Innovative Applications of Artificial Intelligence Conference, IAAI 2019, The Ninth AAAI Symposium on Educational Advances in Artificial Intelligence, EAAI 2019, Honolulu, Hawaii, USA, January 27 - February 1, 2019*. AAAI Press, 346–353. doi:10.1609/AAAI.V33I01.3301346

[51] Chong Xiang, Tong Wu, Zexuan Zhong, David A. Wagner, Danqi Chen, and Prateek Mittal. 2024. Certifiably Robust RAG against Retrieval Corruption. *CoRR* abs/2405.15556 (2024). doi:10.48550/ARXIV.2405.15556 arXiv:2405.15556

[52] Chen-Wei Xie, Siyang Sun, Xiong Xiong, Yun Zheng, Deli Zhao, and Jingren Zhou. 2023. RA-CLIP: Retrieval Augmented Contrastive Language-Image Pre-Training. In *IEEE/CVF Conference on Computer Vision and Pattern Recognition, CVPR 2023, Vancouver, BC, Canada, June 17-24, 2023*. IEEE, 19265–19274. doi:10.1109/CVPR52729.2023.01846

[53] Xu Xie, Fei Sun, Zhaoyang Liu, Shiwen Wu, Jinyang Gao, Jiandong Zhang, Bolin Ding, and Bin Cui. 2022. Contrastive Learning for Sequential Recommendation. In *38th IEEE International Conference on Data Engineering, ICDE 2022, Kuala Lumpur, Malaysia, May 9-12, 2022*. IEEE, 1259–1273. doi:10.1109/ICDE53745.2022.00099

[54] Jinlong Xue, Yayue Deng, Yingming Gao, and Ya Li. 2024. Retrieval Augmented Generation in Prompt-based Text-to-Speech Synthesis with Context-Aware Contrastive Language-Audio Pretraining. *CoRR* abs/2406.03714 (2024). doi:10.48550/ARXIV.2406.03714 arXiv:2406.03714

[55] Junliang Yu, Hongzhi Yin, Xin Xia, Tong Chen, Jundong Li, and Zi Huang. 2024. Self-Supervised Learning for Recommender Systems: A Survey. *IEEE Trans. Knowl. Data Eng.* 36, 1 (2024), 335–355. doi:10.1109/TKDE.2023.3282907

[56] Zhenrui Yue, Honglei Zhuang, Aijun Bai, Kai Hui, Rolf Jagerman, Hansi Zeng, Zhen Qin, Dong Wang, Xuanhui Wang, and Michael Bendersky. 2024. Inference Scaling for Long-Context Retrieval Augmented Generation. *CoRR* abs/2410.04343 (2024). doi:10.48550/ARXIV.2410.04343 arXiv:2410.04343

[57] Wayne Xin Zhao, Yupeng Hou, Xingyu Pan, Chen Yang, Zeyu Zhang, Zihan Lin, Jingsen Zhang, Shuqing Bian, Jiakai Tang, Wenqi Sun, Yushuo Chen, Lanling Xu, Gaowei Zhang, Zhen Tian, Changxin Tian, Shanlei Mu, Xinyan Fan, Xu Chen, and Ji-Rong Wen. 2022. RecBole 2.0: Towards a More Up-to-Date Recommendation Library. In *Proceedings of the 31st ACM International Conference on Information & Knowledge Management, Atlanta, GA, USA, October 17-21, 2022*, Mohammad Al Hasan and Li Xiong (Eds.). ACM, 4722–4726. doi:10.1145/3511808.3557680

[58] Xinping Zhao, Chaochao Chen, Jiajie Su, Yizhao Zhang, and Baotian Hu. 2024. Enhancing Attributed Graph Networks with Alignment and Uniformity Constraints for Session-based Recommendation. In *IEEE International Conference on Web Services, ICWS 2024, Shenzhen, China, July 7-13, 2024*. IEEE, 247–257. doi:10.1109/ICWS62655.2024.00047

[59] Xinping Zhao, Dongfang Li, Yan Zhong, Boren Hu, Yibin Chen, Baotian Hu, and Min Zhang. 2024. SEER: Self-Aligned Evidence Extraction for Retrieval-Augmented Generation. In *Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing, EMNLP 2024, Miami, FL, USA, November 12-16, 2024*, Yaser Al-Onaizan, Mohit Bansal, and Yun-Nung Chen (Eds.). Association for Computational Linguistics, 3027–3041. https://aclanthology.org/2024.emnlp-main.178

[60] Xinping Zhao, Yan Zhong, Zetian Sun, Xinshuo Hu, Zhenyu Liu, Dongfang Li, Baotian Hu, and Min Zhang. 2024. FunnelRAG: A Coarse-to-Fine Progressive Retrieval Paradigm for RAG. *arXiv preprint arXiv:2410.10293* (2024). https://arxiv.org/pdf/2410.10293

[61] Bowen Zheng, Yupeng Hou, Hongyu Lu, Yu Chen, Wayne Xin Zhao, Ming Chen, and Ji-Rong Wen. 2024. Adapting Large Language Models by Integrating Collaborative Semantics for Recommendation. In *40th IEEE International Conference on Data Engineering, ICDE 2024, Utrecht, The Netherlands, May 13-16, 2024*. IEEE, 1435–1448. doi:10.1109/ICDE60146.2024.00118

[62] Zhi Zheng, Wenshuo Chao, Zhaopeng Qiu, Hengshu Zhu, and Hui Xiong. 2024. Harnessing Large Language Models for Text-Rich Sequential Recommendation. In *Proceedings of the ACM on Web Conference 2024, WWW 2024, Singapore, May 13-17, 2024*, Tat-Seng Chua, Chong-Wah Ngo, Ravi Kumar, Hady W. Lauw, and Roy Ka-Wei Lee (Eds.). ACM, 3207–3216. doi:10.1145/3589334.3645358

[63] Kun Zhou, Hui Wang, Wayne Xin Zhao, Yutao Zhu, Sirui Wang, Fuzheng Zhang, Zhongyuan Wang, and Ji-Rong Wen. 2020. S3-Rec: Self-Supervised Learning for Sequential Recommendation with Mutual Information Maximization. In *CIKM '20: The 29th ACM International Conference on Information and Knowledge Management, Virtual Event, Ireland, October 19-23, 2020*, Mathieu d'Aquin, Stefan Dietze, Claudia Hauff, Edward Curry, and Philippe Cudré-Mauroux (Eds.). ACM, 1893–1902. doi:10.1145/3340531.3411954

## A Notations

Table 6 summarizes the main notations used in this paper, including some essential hyper-parameters.

**Table 6: Notation table.**

| Notation | Meaning |
| --- | --- |
| $\mathcal{U}, \mathcal{V}$ | set of users and items |
| $\theta, \phi$ | parameters of the backbone and RAM |
| $\mathbf{V}$ | item embedding table |
| $\mathbf{P}$ | position encoding matrix |
| $\mathbf{h}, \mathbf{h}', \mathbf{h}''$ | user representation |
| $\tilde{\mathbf{h}}, \mathbf{h}^{c_1}, \mathbf{h}^{c_2}$ | augmented user representation |
| $\mathbf{v}_i$ | item embedding of $i$-th item |
| $s_u$ | interaction sequence of user $u$ |
| $v_t^{(u)}$ | item interacted by user $u$ at timestamp $t$ |
| $L$ | number of Transformer layers |
| $N$ | number of items recommended |
| $T$ | maximum input sequence length |
| $d$ | hidden size |
| $K$ | number of memories to retrieve |
| $k$ | number of clusters when building Faiss |
| $\alpha, \beta$ | control coefficient |

## B Datasets

Table 7 summarizes the detailed statistics of three benchmark datasets after preprocessing, where #users, #items, and #inters represent the number of users, items, and interactions, #avg.length denotes the average length of all user sequences, and 'sparsity' calculates the proportion of entries in the user-item matrix without any interaction records.

**Table 7: Statistics of the datasets.**

| Dataset | #users | #items | #inters | #avg.length | sparsity |
| --- | --- | --- | --- | --- | --- |
| Beauty | 22363 | 12101 | 198502 | 8.88 | 99.93% |
| Sports | 35598 | 18357 | 296337 | 8.32 | 99.95% |
| Clothing | 39387 | 23033 | 278677 | 7.08 | 99.97% |

## C Baselines

We compare RaSeRec with the following three groups of competitive baseline models:

**(1) Non-sequential recommendation models (Non-SeRec):**

- PopRec is a non-personalized model based on the popularity of items to recommend.
- BPR-MF [38] uses matrix factorization to model users and items with BPR loss.

**(2) Vanilla sequential recommendation models (Vanilla SeRec):**

- GRU4Rec [15] uses the Gate Recurrent Unit (GRU) to model user sequence.
- Caser [42] applies Convolutional Neural Networks (CNN) to model union-level sequential patterns for SeRec.
- SASRec [21] models user sequences via self-attention mechanism.
- BERT4Rec [41] applies the bi-directional self-attention mechanism with masked item training to model user sequence.
- $S^3Rec_{MIP}$ [63] learns the correlation among items, attributes and etc. Here, we adopt the Mask Item Prediction (MIP) variant.

**(3) SSL-Augmented sequential recommendation models (SSL-Augmented SeRec):**

- CL4SRec [53] proposes item crop, mask, and reorder to construct different views of user sequences.
- CoSeRec [29] propose two robust data augmentations (*i.e.,* item substitute and insert), to create high-quality views for SSL.
- ICLRec [4] first mines the latent users' intents and creates SSL signals between sequences and intents.
- DuoRec [34] proposes model-level augmentation with different sets of dropout masks on the model backbone.
- MCLRec [33] designs learnable augmenters to generate views and applies meta-learning to guide the training of augmenters.

## D Metrics

We rank the prediction results over the whole item set $\mathcal{V}$ and employ Hit Rate (HR) and Normalized Discounted Cumulative Gain (NDCG) to measure recommendation performance. Assuming the cut-off position is $N$, HR@$N$ and NDCG@$N$ can be defined as:

**Table 8: Performance comparison (@10), where the best results are boldfaced and the second-best ones are underlined.**

| Dataset | Beauty | | Sports | | Clothing | | Average | |
|---|---|---|---|---|---|---|---|---|
| Method | HR@10 | NDCG@10 | HR@10 | NDCG@10 | HR@10 | NDCG@10 | HR@10 | NDCG@10 |
| PopRec | 0.0114 | 0.0053 | 0.0090 | 0.0051 | 0.0052 | 0.0025 | 0.0085 | 0.0043 |
| BPR-MF | 0.0299 | 0.0122 | 0.0188 | 0.0083 | 0.0094 | 0.0069 | 0.0194 | 0.0091 |
| GRU4Rec | 0.0365 | 0.0142 | 0.0274 | 0.0137 | 0.0165 | 0.0083 | 0.0268 | 0.0121 |
| Caser | 0.0418 | 0.0253 | 0.0231 | 0.0126 | 0.0174 | 0.0098 | 0.0274 | 0.0159 |
| SASRec | 0.0627 | 0.0281 | 0.0336 | 0.0169 | 0.0272 | 0.0124 | 0.0412 | 0.0191 |
| BERT4Rec | 0.0401 | 0.0254 | 0.0326 | 0.0153 | 0.0208 | 0.0102 | 0.0312 | 0.0170 |
| $S^3Rec_{MIP}$ | 0.0591 | 0.0268 | 0.0265 | 0.0135 | 0.0237 | 0.0132 | 0.0364 | 0.0178 |
| CL4SRec | 0.0667 | 0.0334 | 0.0392 | 0.0195 | 0.0274 | 0.0134 | 0.0444 | 0.0221 |
| CoSeRec | 0.0705 | 0.0381 | 0.0422 | 0.0234 | 0.0248 | 0.0136 | 0.0458 | 0.0250 |
| ICLRec | 0.0729 | 0.0389 | 0.0433 | 0.0234 | 0.0256 | 0.0134 | 0.0473 | 0.0252 |
| DuoRec | 0.0834 | 0.0431 | <u>0.0479</u> | <u>0.0248</u> | 0.0293 | 0.0142 | 0.0535 | <u>0.0274</u> |
| MCLRec | <u>0.0844</u> | <u>0.0437</u> | 0.0462 | 0.0241 | **0.0306** | <u>0.0144</u> | <u>0.0537</u> | <u>0.0274</u> |
| RaSeRec (Ours) | **0.0860** | **0.0463** | **0.0497** | **0.0264** | <u>0.0301</u> | **0.0152** | **0.0553** | **0.0293** |
| %Improv. | 1.90% | 5.95% | 3.76% | 6.45% | -0.02% | 5.56% | 2.98% | 6.93% |
| $p$-value. | $9e^{-3}$ | $5e^{-5}$ | $4e^{-4}$ | $1e^{-4}$ | - | $5e^{-4}$ | $2e^{-4}$ | $1e^{-6}$ |

**Table 9: Hyperparameter settings, where * denotes different hyperparameter settings for different datasets.**

| Configuration | Value |
|---|---|
| number of Transformer layers $L$ | 2 |
| number of heads | 2 |
| hidden size $d$ | 64 |
| maximum sequence length $L$ | 50 |
| number of clusters $k$ | 128 |
| number of clusters accessed | 1 |
| control coefficient $\alpha^*$ | [0.0, 1.0] |
| control coefficient $\beta^*$ | [0.0, 1.0] |
| number of retrieved memories $K^*$ | [5, 55] |
| learning rate | 0.001 |
| optimizer | Adam |
| mini-batch size | 1024 |

- **HR@N** measures whether the target item of the test user sequence appears in the first $N$ recommended items:

$$HR@N = \frac{1}{|\mathcal{U}|} \sum_{u \in \mathcal{U}} \delta(v_{t+1} \in R(u)), \qquad (13)$$

where $\delta(\cdot)$ is an indicator function; $v_{t+1}$ denotes the target item; $R(u)$ denotes the top-$N$ recommended list of user $u$.
- **NDCG@N** further measures the ranking quality. It logarithmically discounts the position:

$$NDCG@N = \frac{1}{|\mathcal{U}|} \sum_{u \in \mathcal{U}} \sum_{j=1}^{N}$$
$$\delta(v_{t+1} = R(u)_j) \frac{1}{\log_2(j+1)}, \qquad (14)$$

where $R(u)_j$ denotes the $j$-th recommended item for user $u$.
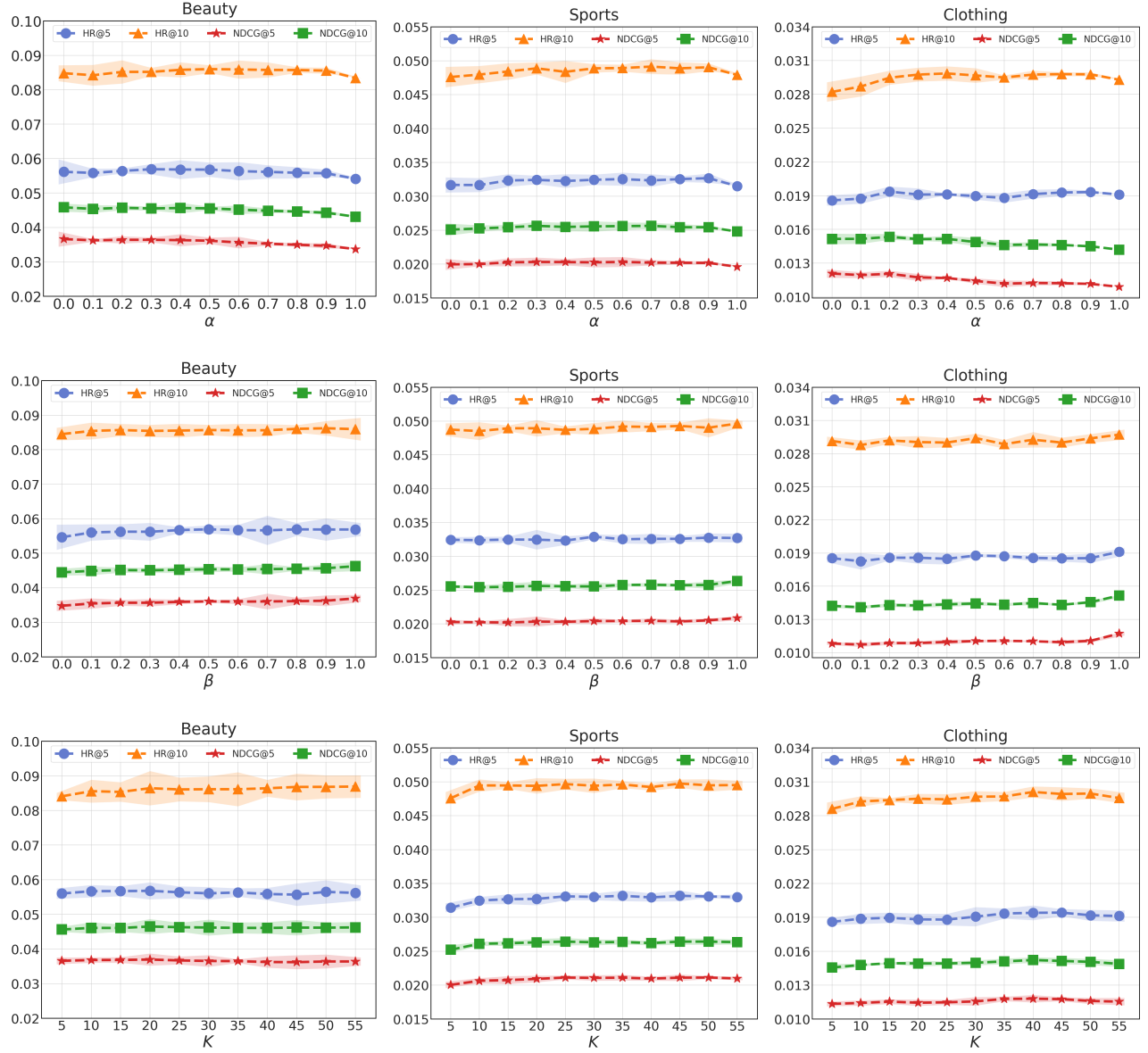
## E  Implementation Details

We use the Adam optimizer [23] to optimize model parameters with the learning rate of 0.001, the mini-batch size of 1024, $\beta_1 = 0.9$, and $\beta_2 = 0.999$, where the maximum number of epochs is set to 100 for both pre-training and fine-tuning stages. We train models with the early stopping strategy that interrupts training if the HR@10 result on the validation set continues to drop until 10 epochs. For

RAM, we tune $\alpha$, $\beta$ and $K$ within the ranges of $\{0.0, 0.1, 0.2, ..., 1.0\}$, $\{0.0, 0.1, 0.2, ..., 1.0\}$ and $\{5, 10, 15, ..., 55\}$, respectively, referring to Appendix F for more details. As for Faiss, we set the number of clusters $k$ to 128, set the number of clusters accessed as 1, and adopt the cosine similarity as the measurement to retrieve relevant memories. As for the sequence encoder, we adopt a two-layer Transformer and set the head number as 2 for each self-attention block, where the hidden size $d$ and the maximum sequence length $L$ are set to 64 and 50 respectively. Additionally, for Non-SeRec and Vanilla SeRec methods, we use results reported by DuoRec [34] as these methods already have generally acknowledged experimental results on the three benchmark datasets used in this paper. We ran each experiment five times and reported the average. Table 9 summarizes the detailed hyper-parameter settings.

## F  Parameter Sensitivity (RQ7)

We study the effects of RaSeRec's key hyper-parameters on its performance, including the control coefficients $\alpha$, $\beta$ (Equ. (11)) and the number of retrieved memories $K$ (§3.2). For $\alpha$ and $\beta$, we vary them from 0 to 1.0, incrementing by 0.1; for $K$, we vary it from 5 to 55, incrementing by 5. We ran each experiment five times and computed the average as well as standard deviation results. The experimental results are shown in Figure 7, where the line denotes the average results and the error bar denotes the scaled-up standard deviation results. From the results, we mainly have the following observations:

Firstly, we can observe that setting $\alpha$ to a too-big value (*e.g.,* 1.0) or to a too-small value (*e.g.,* 0.1) will significantly hurt the model performance. When alpha is high, implicit memories dominate the user representation, whereas if alpha is low, explicit memories dominate the user representation. The above observation indicates the user representation should not be dominated by one side, *i.e.,* implicit memories or explicit memories. In a nutshell, we suggest tuning $\alpha$ in the range of $0.3 \sim 0.7$. For example, setting $\alpha$ as 0.5 will lead to satisfactory performance on the Beauty dataset. Besides, we observe that the performance volatility of RaSeRec generally increases as $\alpha$ decreases. The main reason is that retrieved explicit memories are more diverse than implicit ones since the model backbone is frozen. In conclusion, we suggest setting $\alpha$ at around
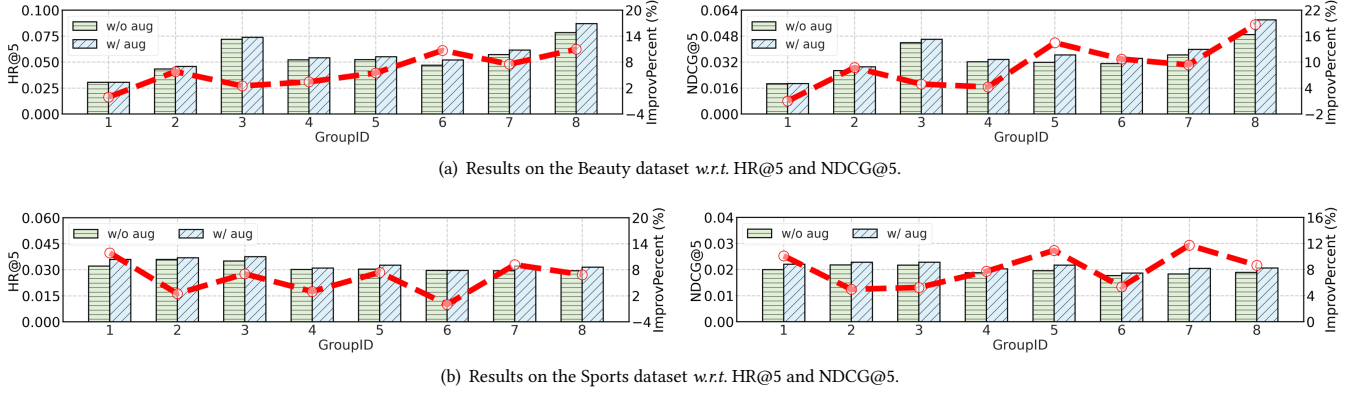
**Figure 7: Parameter sensitivity *w.r.t.* the control coefficients $\alpha, \beta$ and the number of retrieved memories $K$, where we plot the error bar with the mean $\pm \gamma \times$standard deviation and $\gamma$ is set as $500 \times$the scale interval of the y-axis.**

0.5 to achieve stable recommendation performance and relatively diverse results.

Secondly, we can observe the results *w.r.t.* different values of $\beta$ on all three datasets following similar trends. As $\beta$ increases, the recommendation performance and the performance volatility gradually increase. Generally, $\beta = 0.9$ or $\beta = 1.0$ leads to the best performance. It demonstrates that augmenting user representation via modeling the relationship between user representation and the retrieved user representation (Equ. (9)) can produce more informative and accurate user representation. We think the main reason is the model backbone learns to recommend items based on the

user sequence instead of recommending users based on the target item. As such, users with similar behaviors lead to similar target items, whereas users who like the same item do not necessarily have similar behaviors. For the above reasons, augmenting user representation via modeling the relationship between user representation and the retrieved user representation (Equ. (9)) performs better than that between user representation and the corresponding target item embedding of the retrieved user representation (Equ. (10)). In conclusion, we suggest tuning $\beta$ in the range of $0.7 \sim 1.0$ carefully.

(a) Results on the Beauty dataset *w.r.t.* HR@5 and NDCG@5.



(b) Results on the Sports dataset *w.r.t.* HR@5 and NDCG@5.

**Figure 8: Performance comparison over different user groups between the model with retrieval augmentation and without retrieval augmentation. The bar represents HR@5 and NDCG@5, while the line denotes the performance improvement percentage of "w/ aug" compared to "w/o aug".**

Thirdly, we observe that the performance of RaSeRec reaches the peak to different $K$ across three datasets, *e.g.,* 20 is the best value for the Beauty dataset. This demonstrates the effectiveness of augmenting user representation via retrieved explicit memories and manifests that setting a suitable number of retrieved memories can considerably enhance performance. The results on the three datasets follow similar trends, where the performance gradually improves as $K$ increases and then starts to degrade when further increasing $K$. This phenomenon is also observed in long-context RAG [19, 56], where the generative modeling performance initially increases and then declines when consistently increasing the number of retrieved passages. The main reason is that increasing the number of retrieved memories will inevitably introduce irrelevant memories (*i.e.,* "noise") that can mislead the retrieval-augmented module. In conclusion, we suggest tuning $K$ in the range of 10 ~ 30 for the Beauty and Sports datasets, and 30 ~ 50 for the Clothing dataset.

## G  Study on User Interaction Frequency (RQ8)

To investigate the robustness of RaSeRec against low-frequency (short user sequence) and high-frequency users (long user sequence), we split the test user sequences into 8 groups based on their length and ensure the number of test user sequences within each group is the same. The larger the GroupID, the longer the test user sequences. The experimental results are shown in Figure 8, where models with retrieval augmentation mention "w/ aug" otherwise "w/o aug". From the results, we mainly have the following observations:

Firstly, it is surprising that augmenting user representations with retrieved memories is highly effective for high-frequency users. For example, on the 8-th group of the beauty dataset, the performance improvements of "w/ aug" over "w/o aug" are 10.96% and 18.67% in terms of HR@5 and NDCG@5, respectively. This demonstrates that even for high-frequency users, implicit memories may fall short of capturing user preferences accurately, while explicit memories can remedy this issue to a large extent.

Secondly, we find the tendency of improvements *w.r.t.* different groups vary in datasets. For example, on the Beauty dataset, the improvement of high-frequency users is larger than that of low-frequency users. While on the sports dataset, different groups have similar performance improvements. We think the main reason is that the #avg.length of the sports dataset is shorter than that of the beauty dataset. As such, there are a lot of short- and medium-term memories in the memory bank of the sports dataset, as shown in Table 5. These memories can largely facilitate the performance of low- and medium-frequency users.

Thirdly, we are surprised to find that there is no significant difference in the performance of each group on the sports dataset. However, on the beauty dataset, high-frequency users perform better than low-frequency ones, which is in accordance with people's intuition. We think the main reason is that the sports dataset is sparser than the beauty dataset. As such, there are not enough supervised signals for the model to learn high-quality user representation for long user sequences. This also explains why the model performs poorly when only the L partition is retained, as shown in Figure 5(b).

In a nutshell, we suggest performing retrieval augmentation for both low- and high-frequency users.