

# CESRec: Constructing Pseudo Interactions for Sequential Recommendation via Conversational Feedback

Anonymous ACL submission

## Abstract

Sequential Recommendation Systems (SRS) have become essential in many real-world applications. However, existing SRS methods often rely on collaborative filtering signals and fail to capture real-time user preferences, while Conversational Recommendation Systems (CRS) excel at eliciting immediate interests through natural language interactions but neglect historical behavior. To bridge this gap, we propose CESRec, a novel framework that integrates the long-term preference modeling of SRS with the real-time preference elicitation of CRS. We introduce semantic-based pseudo interaction construction, which dynamically updates users' historical interaction sequences by analyzing conversational feedback, generating a pseudo-interaction sequence that seamlessly combines long-term and real-time preferences. Additionally, we reduce the impact of outliers in historical items that deviate from users' core preferences by proposing dual alignment outlier items masking, which identifies and masks such items using semantic-collaborative aligned representations. Extensive experiments demonstrate that CESRec achieves state-of-the-art performance by boosting strong SRS models, validating its effectiveness in integrating conversational feedback into SRS<sup>1</sup>.

## 1 Introduction

Sequential Recommendation Systems (SRS) are pivotal in various applications, such as e-commerce (Zhou et al., 2018) and streaming platforms (Pan et al., 2023), by providing personalized item recommendations based on users' historical interaction sequences (Fang et al., 2020). Recently, large language models (LLMs) have demonstrated remarkable reasoning capabilities (Mann et al., 2020; Zhang et al., 2022), making them promising method for enhancing recommendation tasks.

<sup>1</sup>Code is available at <https://anonymous.4open.science/r/NLESR-4342>

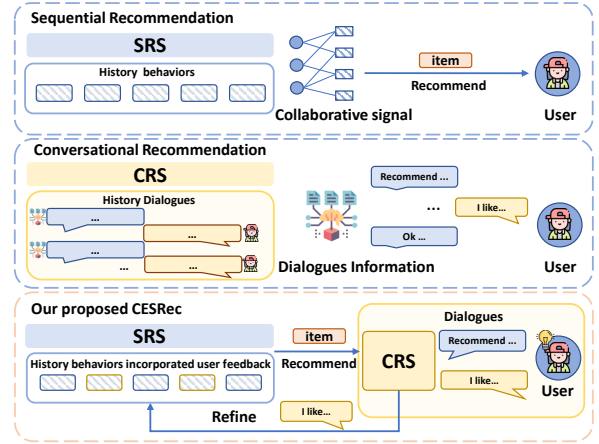


Figure 1: Comparison of sequential recommendation, conversational recommendation, and our CESRec, which combines the advantages of both sequential and conversational recommendation systems.

Several studies (Liao et al., 2024; Bao et al., 2023) have demonstrated the superiority of directly applying LLMs to sequential recommendation tasks. In contrast, Conversational Recommendation Systems (CRS) employ natural language interactions to inquire about user preferences and predict personalized item recommendations (Friedman et al., 2023; Mysore et al., 2023). However, existing SRS methods usually rely on collaborative filtering signals while neglecting the rich semantic information associated with items. A significant limitation of these approaches is their inability to capture users' real-time interests, as immediate preferences are not dynamically reflected in the behavior sequence. Conversely, while CRS methods excel at capturing immediate interests through natural language conversations, they typically fail to incorporate historical interaction sequences into their frameworks. Consequently, the first challenge lies in dynamically integrating the *long-term preference modeling* of SRS with the *real-time interests modeling* facilitated by natural language interactions in CRS.

In this paper, we propose **Conversation Enhanced Sequential Recommendation (CESRec)**. To address the first challenge, we introduce **semantic-based pseudo interaction construction**, a novel method that directly updates the historical interaction sequence based on users’ conversational feedback. Specifically, this approach analyzes users’ natural language inputs to model their current preferences and refines their historical interaction sequence, generating a *pseudo-interaction sequence* that seamlessly integrates both long-term and real-time preferences. Next, we use the pseudo-interaction sequence as input to SRS, which effectively combines the collaborative filtering signals of SRS with the semantic signals derived from conversational feedback. This enables accurate recommendations based on natural language interactions without requiring extensive modifications to existing SRS-based systems, ensuring seamless integration and enhanced user experience.

Since historical interaction sequences often contain items that deviate substantially from users’ main preferences, such as mistakenly clicked items or transient interests, as observed in many recent studies (Lin et al., 2023; Wang et al., 2021), these outliers can adversely affect the modeling of user behavior. These items can negatively influence the LLM’s modeling of user behavior, potentially misleading the construction of the pseudo-interaction sequence. For example, if a user’s primary preference is horror films, the inclusion of a comedy movie in the interaction sequence might lead the LLMs to utilize “horror-comedy” films to construct the pseudo-interaction sequence, rather than a pure horror film. In this work, we refer to such items as **outlier items**. Therefore, the second challenge is how to accurately identify these outlier items and mask them in the interaction sequence to minimize their impact on the generation of the pseudo-interaction sequence.

To address this, we propose dual alignment outlier items masking, a method that accurately identifies outlier items from the user’s historical interaction sequence based on semantic-collaborative aligned representations and subsequently masks these items. Specifically, we leverage LLMs to obtain semantic embeddings of items and extract collaborative representations from the SRS model. We then introduce a dual alignment mechanism to derive hybrid item representations, which simultaneously capture co-occurrence relationships and semantic information among items. Based on these

hybrid representations, we identify items that substantially deviate from the user’s core preferences, ensuring precise masking while preserving the integrity of the user’s historical behavior sequence. The experimental results demonstrate that our CESRec can boost the performance of several state-of-the-art SRS models in terms of HR and NDCG, which verifies that our CESRec effectively integrates the conversational feedback into the SRS.

The main contributions of this work are as follows:

- We propose CESRec, which combines the advantage of real-time conversational feedback with the efficiency of learning user preferences from historical behavior.
- We introduce semantic-based pseudo interaction construction method to refine user historical interaction sequences by leveraging user conversational feedback.
- We propose dual alignment outlier items masking method to optimize item selection during the sequence refinement process.
- Extensive experiments demonstrate that our proposed CESRec achieves state-of-the-art performance by boosting the performance of several strong SRS models.

## 2 Related Work

**Sequential Recommendation** Sequential recommendation aims to predict the next item that aligns with a user’s preferences based on their historical interaction sequence (Fang et al., 2020; Li et al., 2023a,b). Traditional sequential recommendation models capture user preferences by leveraging item co-occurrence relationships. To model complex sequential patterns, CNN-based (Tang and Wang, 2018) and GNN-based (He et al., 2020) methods have been introduced. Additionally, transformer-based approaches, such as SASRec (Kang and McAuley, 2018) and BERT4Rec (Sun et al., 2019), have been developed to capture long-term dependencies between arbitrary items. However, most of these methods primarily model user preferences based on long-term interaction histories, making it challenging to effectively capture dynamic shifts in user interests. As a result, they struggle to reflect users’ real-time preferences within interaction sequences, leading to recommendations that may not accurately align with users’ immediate interests.

**Conversational Recommendation** Conversational Recommendation System (CRS) aims to provide recommendations via natural language con-

versations (Zhou et al., 2020; Lei et al., 2020; He et al., 2023). Feng et al. (2023) propose an LLM-based CR method that utilizes LLMs for sub-task management, expert collaboration, and response generation. Fang et al. (2024) propose a multi-agent collaborative system that optimizes dialogue flow and recommendation accuracy, incorporating a user feedback-aware reflection mechanism to enhance the user interaction experience. While CRS methods excel at capturing immediate user interests through natural language conversations, they often fail to effectively integrate historical interaction sequences into their frameworks.

**LLMs for Recommendation** Large Language Models (LLMs) have demonstrated remarkable capabilities across various domains. By encoding extensive world knowledge during pretraining, LLMs have increasingly been utilized to enhance recommendation systems (Dai et al., 2023; Geng et al., 2022; Hou et al., 2024). LLaRA (Liao et al., 2024) utilizes a hybrid prompting approach, combining ID-based item embedding learned by traditional recommendation models with textual item features as input to predict the next item. Rajput et al. (2023) propose a generative retrieval approach in which the retrieval model decodes semantic IDs of target candidates. (Liu et al., 2024) propose leveraging LLMs to generate item embeddings, which can be seamlessly incorporated into sequential recommendation models to improve their performance. Hu et al. (2024) introduce a method for learning semantically aligned item ID embeddings from textual descriptions, using a projector module to map item IDs to embedding vectors, which are then transformed into descriptive text tokens by the LLM. (Bao et al., 2023) introduces a method that converts collaborative embeddings into binary sequences for LLM interpretability. While these approaches leverage LLMs to process textual information, they primarily focus on transforming item content into embedding representations. However, they do not fully exploit the rich semantic information contained in users’ conversational feedback, limiting their ability to dynamically adapt recommendation strategies based on real-time user preferences.

### 3 Problem Definition

In this paper, we follow the problem definition commonly used in sequential recommendation tasks (Hu et al., 2024). Given a user  $u \in \mathcal{U}$ , where

$\mathcal{U}$  represents the set of all users, and a historical interaction sequence  $\mathcal{I}(u) = \{v_1^{(u)}, v_2^{(u)}, \dots, v_{N_u}^{(u)}\}$ , the model aims to predict the next item the user is likely to interact with based on  $\mathcal{I}(u)$ . Here,  $v_i^{(u)}$  denotes the  $i$ -th item interacted by user  $u$ , and all items belong to the item set  $\mathcal{V}$ . The sequence length of  $\mathcal{I}(u)$  is denoted by  $N_u$ .

## 4 CESRec

### 4.1 Overview

In this section, we show the details of the Conversation Enhanced Sequential Recommendation (CESRec), which is illustrated in Figure 2. The proposed model consists of two main components: Semantic Pseudo Sequence Construction and Dual Alignment Outlier Items Masking. The **Semantic Pseudo Sequence Construction** module is designed to construct the pseudo-interaction sequence by refining the historical interaction sequence via users’ conversational feedback. Subsequently, the **Dual Alignment Outlier Items Masking** module further enhances the refinement process by identifying and masking items that deviate from the user’s core preferences.

### 4.2 Dual Alignment Outlier Items Masking

In the process of constructing a semantic-based pseudo interaction sequence, the model leverages the user historical sequences to capture their core preferences and selects appropriate replacement items based on conversational feedback. However, during the modification of the original interaction sequence, items in the historical sequence that deviate from the user’s core preferences can interfere with the LLM’s modeling of user behavior. This misalignment can introduce bias, potentially leading to the inappropriate replacement of items. In this work, we refer to such items as outlier items. To address this issue, we propose a dual-alignment outlier items masking method to ensure that such deviating items are appropriately masked.

According to a recent study (Sheng et al., 2024), LLMs can implicitly encode user preference information, and items sharing similar content tend to exhibit similar semantic embeddings. Based on this observation, we extract item embeddings from LLMs, which are rich in semantic information. Given an item  $v_i^{(u)}$  with content information  $c_i$  such as title, we employ an LLM to obtain the

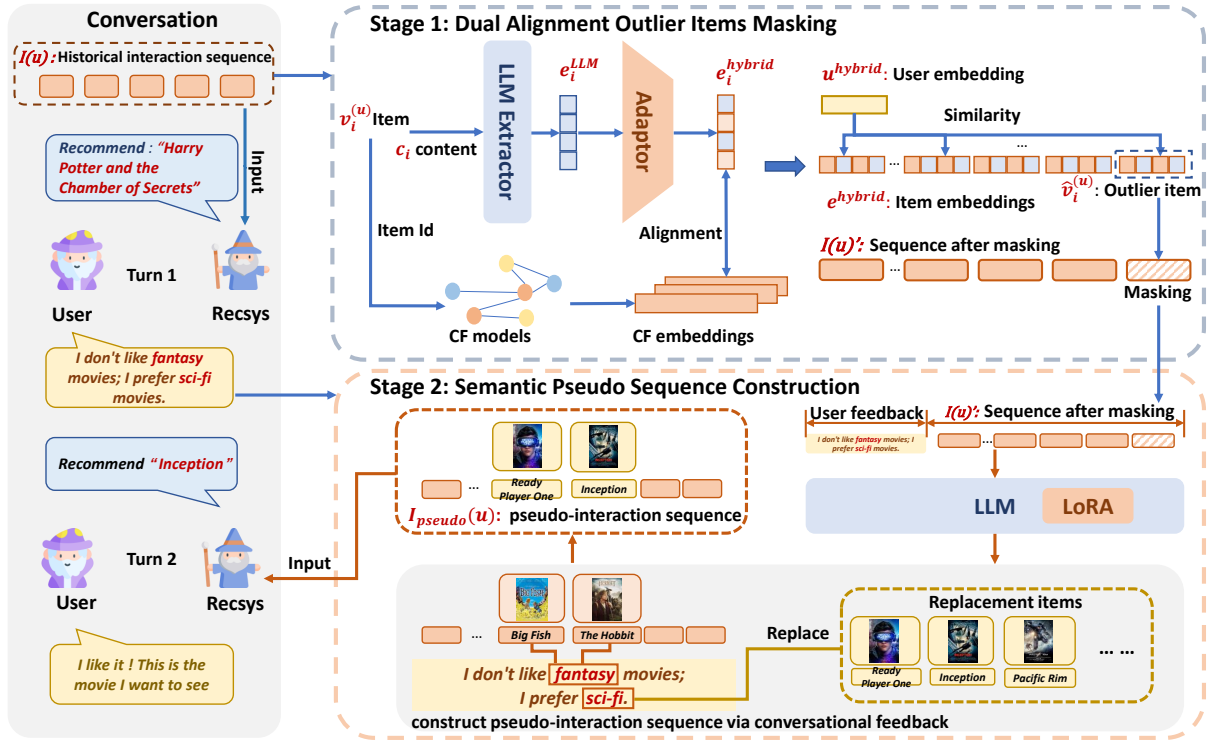


Figure 2: Overview of CESRec. In our proposed framework, we first employ the conventional sequential recommendation method (*a.k.a.*, Recsys) to predict an item based on the user’s historical interaction sequence. Next, our CESRec refines the interaction sequence by constructing the pseudo-interaction sequence and masking the outlier items. Finally, we employ Recsys to give a new recommendation by using the refined sequence.

semantic embeddings  $e_i^{LLM}$ :

$$e_i^{LLM} = \text{Extractor}(c_i), \quad (1)$$

where  $\text{Extractor}(\cdot)$  refers to the LLM tokenizer and encoder layers, and we utilize the output of the last hidden layer  $e_i^{LLM}$  as the semantic embedding.

Relying solely on semantic embeddings to identify outlier items may compromise the integrity of the user’s historical behavior sequence, thereby limiting the effectiveness of SRS in accurately modeling user preferences. We introduce a trainable adapter to align the semantic embeddings derived from LLMs with the collaborative signals typically used in SRS. This adapter is specifically trained to fuse the positional influence and co-occurrence information while utilizing semantic embeddings for masking:

$$e_i^{\text{hybrid}} = \text{Adapter}(\theta_{\text{collab}}; e_i^{LLM}), \quad (2)$$

where  $e_i^{\text{hybrid}}$  represents the hybrid embedding that integrates both semantic and collaborative information. The Adapter is a two-layer perception with trainable parameters  $\theta_{\text{collab}}$ .

Finally, to identify outlier items in interactions, we rank items based on the similarity between user

representation and each item. We first obtain all the hybrid embeddings of all the user interacted items in  $\mathcal{I}(u)$ , and fuse all the item representation as the user embedding  $u^{\text{hybrid}}$ :

$$u^{\text{hybrid}} = \text{Fuse}(\{e_1^{\text{hybrid}}, e_2^{\text{hybrid}}, \dots, e_{N_u}^{\text{hybrid}}\}), \quad (3)$$

where the  $\text{Fuse}(\cdot)$  denotes the mean-pooling operator. Then, we calculate the similarity between each item representation  $e_i^{\text{hybrid}}$  and user representation  $u^{\text{hybrid}}$ .

$$s_i = \text{Similarity}(e_i^{\text{hybrid}}, u^{\text{hybrid}}), \quad (4)$$

where  $s_i \in [0, 1]$  denotes the similarity score, and we employ the cosine similarity as the  $\text{Similarity}(\cdot)$  function to measure the semantic gap between  $e_i^{\text{hybrid}}$  and  $u^{\text{hybrid}}$ . To identify outlier items in interaction sequence, we rank items based on their similarity scores  $s_i$ . The top  $k$  items with the lowest similarity scores are considered as the outlier items and will be subsequently masked from the user interaction sequence.

The input and output format of the final dual alignment outlier items masking is as follows:

$$I(u)' = \text{Dual-Alignment}(I(u)), \quad (5)$$



where  $I(u)' = \{v_1^{(u)}, \dots, v_{N_u-k}^{(u)}, \hat{v}_1^{(u)}, \dots, \hat{v}_k^{(u)}\}$  represents interaction sequence after masking,  $\hat{v}_i^{(u)}$  represents the top  $k$  items with the lowest similarity scores. Using these hybrid representations, we identify and mask the outlier items that deviate from the user’s core preferences while preserving the integrity of their historical behavior sequence. This optimization enables the CESRec to better concentrate on core preferences when constructing a semantic-based pseudo interaction sequence.

### 4.3 Semantic Pseudo Sequence Construction

To address the challenge of dynamically integrating long-term preference modeling of SRS with real-time interest modeling driven by natural language interactions in CRS, we propose a semantic-based pseudo sequence construction approach. This method leverages natural language interaction with users to directly capture their current preferences, and generates semantic-based pseudo sequence by incorporating current preferences to historical interaction sequence. Specifically, we introduce a *constructor* that constructs semantic-based pseudo interaction sequences based on user-provided feedback. Following the previous conversational recommendation works (Fang et al., 2024), we ask the user for preference about the target target item attributes.

$$\text{feedback} = \text{User-Interaction}(v_{rec}^{(u)}, \text{Attr}_{\text{target}}) \quad (6)$$

where  $v_{rec}^{(u)}$  represents the recommended item generated by an SRS with input  $I(u)$ ,  $\text{Attr}_{\text{target}}$  refers to attributes of the target item, and feedback denotes a conversational feedback derived from the user that describes the user preference of the item attributes. For instance, if the SRS recommends <Avatar> to the user, but the user prefers films directed by Christopher Nolan, the user may respond with feedback such as: “*I don’t like film directed by James Cameron; I prefer Christopher Nolan.*”.

Next, the Constructor integrates user feedback to iterative refine the historical interaction sequence  $I'(u)$  and generate the pseudo-interaction sequence  $I_{\text{pseudo}}(u)$ :

$$I_{\text{pseudo}}(u) = \text{Constructor}(I'(u), \text{feedback}), \quad (7)$$

where  $I_{\text{pseudo}}(u)$  represents the pseudo-interaction sequence generated by the Constructor, dynamically adjusted based on user conversational feedback.

Dataset	#User	#Item	#Review	#Density
Video Games	55,223	17,408	496,315	0.051628%
Toys	208,180	78,772	1,826,430	0.011138%
MovieLens	6,040	3,883	1,000,209	4.264680%

Table 1: Statistics of three datasets.

To construct the training data for the constructor module, we construct a semantic pseudo sequence by replacing items that no longer align with the user’s current preference, considering both historical behavior and current preferences for the replacements. We construct training data by randomly selecting an item from the sequence as an “outdated” item. The target item, which reflects the user’s updated preference, serves as the ground truth, while the feedback generated between the outdated and target items is used as input for the model. The training instruction is as follows:

**Instruction:** Based on the preferences mentioned in the user feedback and the information about <items> contained in the historical interaction sequence, replace the <items> the user dislikes with <items> user may currently prefer.  
**Input:** historical interaction sequence: <sequence>; user feedback: <feedback>.  
**Output:** pseudo-interaction sequence:<pseudo sequence>

Finally, after refining the interaction sequence of the user by the Constructor, we use the semantic pseudo interaction sequence  $I_{\text{pseudo}}(u)$  as the input to the SRS to regenerate recommended items.

$$v_{N_u+1}^{(u)} = \text{SRS}(I_{\text{pseudo}}(u)), \quad (8)$$

where SRS represents sequential recommendation models,  $v_{N_u+1}^{(u)}$  represents the regenerated recommended item based on the semantic pseudo interaction sequence. Since our proposed CESRec is model-agnostic, it can be seamlessly integrated with existing sequential recommendation models.

## 5 Experimental Setup

### 5.1 Dataset and Evaluation Metric

We conduct experiments on two commonly used recommendation datasets, Video Games and Toys, constructed from the Amazon review datasets (Ni et al., 2019). We also employ the MovieLens datasets (Harper and Konstan, 2015) which is a widely adopted dataset for sequential recommendation tasks, which contains user interactions with movies. Statistics are shown in Table 1.

We adopt two widely used metrics to evaluate the performance: Normalized Discounted Cumulative Gain (NDCG@K) and Hit Rate (HR@K) with  $K=5, 10$ . We select 100 non-interacted items to construct the candidate set, ensuring the inclusion of the correct subsequent item.

## 5.2 Implementation Detail

For the sequential recommendation method, SASRec (Kang and McAuley, 2018), we train the model on all three datasets using the Adam optimizer (Kingma, 2014) for 200 epochs, with a learning rate of 0.001 and a batch size of 256. For the LLM-based recommendation method, LLaRA (Liao et al., 2024), the original configuration selects the top-ranked item from the candidate set as the recommendation result. To ensure consistency with our experimental setup, we adopt the ranking method from (Wang et al., 2024), which ranks the candidate items based on the cosine similarity between item embeddings and the output embeddings of LLaRA. In our CESRec, we mask 1 item in three datasets. We implement our CESRec using two LLMs as the backbone: LLaMA-2-7b (Touvron et al., 2023) and LLaMA-3-8b (Dubey et al., 2024). And we use the same user simulator as the previous conversational recommendation studies Fang et al. (2024) when training and evaluating the models.

## 5.3 Baselines

We conducted experiments using two strong SRS backbones: (1) **SASRec** (Kang and McAuley, 2018) is a widely used sequential recommendation model that employs a self-attention mechanism to effectively capture relationships between items within a user’s interaction sequence. (2) **LLaRA** (Liao et al., 2024) is an LLM-based recommendation model that utilizes a hybrid prompting approach, combining ID-based and text-based representations of items as input. This model aims to enhance recommendation accuracy by integrating both structured and unstructured data sources.

# 6 Experimental Results

## 6.1 Main Results

We evaluate the performance of our proposed CESRec and baseline methods on three datasets using four evaluation metrics. As shown in Table 2, SASRec+CESRec and LLaRA+CESRec consistently outperform their corresponding base SRS model

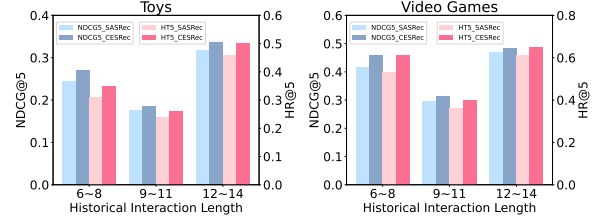


Figure 3: Performance of using different lengths of the historical interaction sequence.

(*a.k.a.*, SASRec and LLaRA) across all datasets and metrics. This demonstrates that the semantic-based pseudo interaction sequences, which incorporate users’ current feedback, enable recommendation models to more effectively capture users’ real-time preferences. Secondly, CESRec demonstrates improved performance when leveraging larger LLMs as the backbone, suggesting that more powerful LLMs possess the stronger capability to accurately model user preferences and select relevant replacement items.

## 6.2 Ablation Study

To validate the effectiveness of each module, we compare the performance of the following variants of CESRec-LLaMA3 on the SASRec backbone: (1) **CESRec w/o d.a.**: we solely employ user conversational feedback to construct pseudo interaction sequences and remove the **dual alignment** from CESRec. (2) **CESRec w/o c.**: we only leverage dual alignment method to mask outlier items and do not construct pseudo sequence. The results, as shown in Table 3, demonstrate that all modules in the model contribute to enhancing sequential recommendation. The superior performance of CESRec-LLaMA3 over CESRec w/o d.a. indicates that the dual alignment outlier items masking method enables CESRec to concentrate on user’s main preference, and construct semantic pseudo sequences that better align with user preferences. By employing the dual alignment and masking module to mask items that deviate from the user’s core preferences, SASRec+CESRec w/o c. demonstrates improved performance over SASRec. This indicates that our dual alignment method does not interfere with the SRS method’s ability to effectively capture user preferences.

## 6.3 The Impact of Historical Interaction Sequence Length

To investigate the impact of historical interaction sequence length, we evaluate model performance

Dataset	Model	HR@5	NDCG@5	HR@10	NDCG@10	Model	HR@5	NDCG@5	HR@10	NDCG@10
Video Games	SASRec	0.590	0.4629	0.717	0.5042	LLaRA	0.270	0.2277	0.360	0.2558
	+CESRec-LLaMA2	0.633	0.4847	0.725	0.5144	+CESRec-LLaMA2	0.380	0.3097	<b>0.450</b>	0.3316
	+CESRec-LLaMA3	<b>0.646</b>	<b>0.4923</b>	<b>0.745</b>	<b>0.5242</b>	+CESRec-LLaMA3	<b>0.380</b>	<b>0.3254</b>	0.440	<b>0.3445</b>
Movielens	SASRec	0.757	0.5688	0.866	0.6045	LLaRA	0.170	0.1416	0.210	0.1542
	+CESRec-LLaMA2	<b>0.824</b>	<b>0.6076</b>	0.882	<b>0.6264</b>	+CESRec-LLaMA2	0.260	0.2192	0.310	0.2347
	+CESRec-LLaMA3	0.810	0.5996	<b>0.886</b>	0.6244	+CESRec-LLaMA3	<b>0.280</b>	<b>0.2348</b>	<b>0.330</b>	<b>0.2508</b>
Toys	SASRec	0.431	0.3173	0.537	0.3509	LLaRA	0.420	0.3957	0.430	0.3986
	+CESRec-LLaMA2	0.472	0.3376	0.557	0.3647	+CESRec-LLaMA2	0.500	0.4671	0.590	0.4955
	+CESRec-LLaMA3	<b>0.478</b>	<b>0.3408</b>	<b>0.557</b>	<b>0.3659</b>	+CESRec-LLaMA3	<b>0.500</b>	<b>0.4671</b>	<b>0.600</b>	<b>0.4993</b>

Table 2: Performance on three datasets. We apply our proposed CESRec on two strong SRS: SASRec and LLaRA, and we implement CESRec based on two LLM: LLaMA2 and LLaMA3.

Dataset	Method	HR@5	NDCG@5	HR@10	NDCG@10
Video Games	+CESRec-LLaMA3	<b>0.646</b>	<b>0.4923</b>	<b>0.745</b>	<b>0.5242</b>
	+CESRec w/o d.a.	0.634	0.4849	0.723	0.5136
	+CESRec w/o c.	0.610	0.4711	0.723	0.5077
	SASRec	0.590	0.4629	0.717	0.5042
Movielens	+CESRec-LLaMA3	<b>0.810</b>	<b>0.5996</b>	<b>0.886</b>	<b>0.6244</b>
	+CESRec w/o d.a.	0.805	0.5940	0.880	0.6186
	+CESRec w/o c.	0.774	0.5766	0.866	0.6061
	SASRec	0.757	0.5688	0.866	0.6045
Toys	+CESRec-LLaMA3	<b>0.478</b>	<b>0.3408</b>	<b>0.557</b>	<b>0.3659</b>
	+CESRec w/o d.a.	0.468	0.3354	0.557	0.3638
	+CESRec w/o c.	0.443	0.3222	0.530	0.3501
	SASRec	0.431	0.3173	0.537	0.3509

Table 3: Performance of ablation models. We conduct ablation study on SASRec+CESRec.

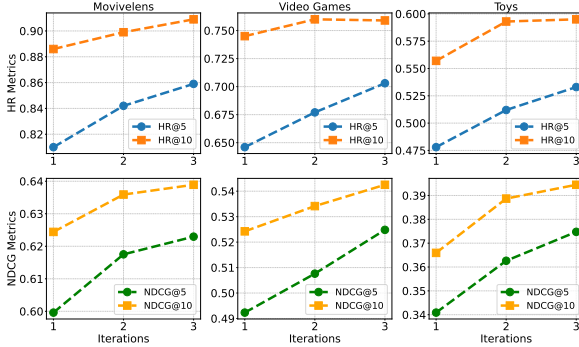


Figure 4: Performance of using different interaction numbers. We evaluate the impact of the number of conversational interactions between CESRec and users.

using different sequence lengths in terms of HR@5 and NDCG@5 on the Toys and Video Games datasets. As shown in Figure 3, the results demonstrate that our proposed CESRec consistently outperforms the baseline SASRec across all three sequence length ranges. This demonstrates the robustness of our model in effectively handling historical interaction sequences of varying lengths, further confirming its adaptability in diverse recommendation scenarios.

## 6.4 Analysis of Interaction Numbers

We further investigate the impact of the number of conversational interactions of CESRec-LLaMA3,

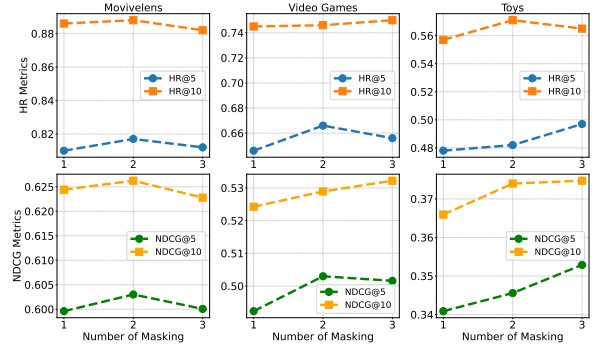


Figure 5: The impact of masking different numbers of outlier items.

based on SASRec. As illustrated in Figure 4, as the number of interactions between users and the CESRec-LLaMA3 increases, the performance of the recommendation system consistently improves. The HR@K and NDCG@K metrics (with K=5, 10) demonstrate a steady upward trend across all three real-world datasets. This indicates that as users provide more feedback, the recommendation system becomes increasingly effective at capturing users' real-time interests. By constructing semantic-based pseudo interaction sequences that reflect these interests, the system generates recommendations that better align with users' current preferences. Moreover, the improvement in both HR and NDCG metrics suggests that the recommendation system not only predicts items that users are more likely to engage with but also ranks relevant items higher in the recommendation list, thereby delivering more accurate and user-centric ranking results.

## 6.5 Analysis of Masking Outlier Items

We further investigated the impact of the number of masked outlier items on the performance of CESRec. The results show that for the MovieLens and Video Games datasets, the model achieves optimal performance when the number of masked items is set to 2. Beyond this threshold, performance



Figure 6: A case study of CESRec, where SASRec first recommends an item based on user historical interaction sequence, and then the user gives feedback. Next, our CESRec refines the interaction sequence and employs the SASRec to give a new recommendation by using the updated sequence.

begins to decline as the number of masked items increases. This decline can be attributed to the fact that excessive masking reduces the length of the user’s historical sequence, leading to a loss of valuable information regarding user preferences. Consequently, the model struggles to accurately capture user behavior and predict items that align with these preferences. In contrast, for the Toys dataset, the model’s performance improves as the number of masked items increases. This trend can be attributed to the higher sparsity of the Toys dataset compared to other two datasets, as shown in Table 1. With greater sparsity, the items in the constructed sequences exhibit more variability, and as the model adjusts these sequences based on user feedback expressed in natural language, the impact on the recommendation outcomes becomes more notable. Therefore, by masking items that deviate from the user’s preferences, the model can concentrate on the most relevant interactions, resulting in improved performance.

## 6.6 Case Study

To intuitively validate the effectiveness of our proposed CESRec, we randomly select an example from MovieLens dataset, as shown in Figure 6. The user’s historical interactions with movies include:

“I Still Know What You Did Last Summer”, “Jungle 2 Jungle”, “Two if by Sea, M. Butterfly”, “Super Mario Bros”, “Blank Check”, “Repossessed”, “The Evening Star”, “The Beautician and the Beast”, “Mr. Wrong”, “A Night at the Roxbury”, “Halloween: The Curse of Michael Myers”, “Stop! Or My Mom Will Shoot”, “Cops and Robbersons”. Given this sequence as input, SASRec generates “Jack Frost” as a recommended item by capturing the co-occurrence relationships between movies. However, “Jack Frost” is a comedy film, which does not align with the user’s current preference for horror films. To encourage the model’s focus on the user’s core interests, we employ the dual alignment outlier items masking method. This method masks the “Super Mario Bros.”, which belongs to the action/animation genre and deviates from the user’s core preference for horror films. Thus, the model can better align with the user’s primary interests and improve recommendation accuracy. This masking process enables the CESRec to better concentrate on the user’s core preferences. Since “Jack Frost” is inconsistent with the user’s preference, CESRec constructs a semantic-based pseudo-interaction sequence incorporating the user’s conversational feedback: “I don’t like comedy; I prefer horror.”. During this process, CESRec replaces “Cops and Robbersons (comedy)” with “Carnosaur 2 (horror)” to reinforce the user’s stated preference. Ultimately, based on this refined interaction sequence, CESRec predicts “Halloween: H20” as the recommended item.

## 7 Conclusion

In this paper, we proposed **Conversation Enhanced Sequential Recommendation (CESRec)**, a novel framework that seamlessly integrates the long-term preference modeling of SRS with the real-time preference elicitation of CRS. By leveraging users’ conversational feedback, CESRec dynamically refines historical interaction sequences to generate pseudo-interaction sequences that capture both long-term preferences and real-time interests. Additionally, the dual alignment outlier items masking method addresses the challenge of outlier items in historical sequences by accurately identifying and masking items that deviate from users’ core preferences. Extensive experiments on three real-world datasets demonstrate that CESRec enhances the performance of SOTA SRS models, achieving superior results in terms of HR and NDCG metrics.



## Limitations

Our method relies on user conversational feedback to dynamically refine the historical interaction sequence, aiming to better align with the user’s real-time preferences. However, if the user’s feedback is expressed in a vague, ambiguous, or unclear manner, the model may fail to capture the user’s real-time preferences accurately, leading to the generation of an imprecise pseudo-interaction sequence, which in turn affects the recommendation performance. In future work, we will investigate more sophisticated dialogue mechanisms that can effectively guide users to articulate their latent preferences.

## Ethical Considerations

The research conducted in this paper centers on investigating the effectiveness of leveraging LLMs to bridge the gap between conversational recommendation and sequential recommendation. Our work systematically benchmarks LLMs under various real-world scenarios and evaluates their performance. In the process of conducting this research, we have adhered to ethical standards to ensure the integrity and validity of our work. To minimize potential bias and ensure fairness, we employ the same prompts and experimental setups as those used in existing publicly accessible and freely available studies. We have made every effort to ensure that our research does not harm individuals or groups and does not involve any form of deception or misuse of information.

## References

Keqin Bao, Jizhi Zhang, Yang Zhang, Wenjie Wang, Fuli Feng, and Xiangnan He. 2023. Tallrec: An effective and efficient tuning framework to align large language model with recommendation. In *Proceedings of the 17th ACM Conference on Recommender Systems*, pages 1007–1014.

Sunhao Dai, Ninglu Shao, Haiyuan Zhao, Weijie Yu, Zihua Si, Chen Xu, Zhongxiang Sun, Xiao Zhang, and Jun Xu. 2023. Uncovering chatgpt’s capabilities in recommender systems. In *Proceedings of the 17th ACM Conference on Recommender Systems*, pages 1126–1132.

Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha Letman, Akhil Mathur, Alan Schelten, Amy Yang, Angela Fan, et al. 2024. The llama 3 herd of models. *arXiv preprint arXiv:2407.21783*.

Hui Fang, Danning Zhang, Yiheng Shu, and Guibing Guo. 2020. Deep learning for sequential recommendation: Algorithms, influential factors, and evaluations. *ACM Transactions on Information Systems (TOIS)*, 39(1):1–42.

Jiabao Fang, Shen Gao, Pengjie Ren, Xiuying Chen, Suzan Verberne, and Zhaochun Ren. 2024. A multi-agent conversational recommender system. *arXiv preprint arXiv:2402.01135*.

Yue Feng, Shuchang Liu, Zhenghai Xue, Qingpeng Cai, Lantao Hu, Peng Jiang, Kun Gai, and Fei Sun. 2023. A large language model enhanced conversational recommender system. *arXiv preprint arXiv:2308.06212*.

Luke Friedman, Sameer Ahuja, David Allen, Zhenning Tan, Hakim Sidahmed, Changbo Long, Jun Xie, Gabriel Schubiner, Ajay Patel, Harsh Lara, et al. 2023. Leveraging large language models in conversational recommender systems. *arXiv preprint arXiv:2305.07961*.

Shijie Geng, Shuchang Liu, Zuohui Fu, Yingqiang Ge, and Yongfeng Zhang. 2022. Recommendation as language processing (rlp): A unified pretrain, personalized prompt & predict paradigm (p5). In *Proceedings of the 16th ACM Conference on Recommender Systems*, pages 299–315.

F Maxwell Harper and Joseph A Konstan. 2015. The movielens datasets: History and context. *Acm transactions on interactive intelligent systems (tiis)*, 5(4):1–19.

Xiangnan He, Kuan Deng, Xiang Wang, Yan Li, Yongdong Zhang, and Meng Wang. 2020. Lightgcn: Simplifying and powering graph convolution network for recommendation. In *Proceedings of the 43rd International ACM SIGIR conference on research and development in Information Retrieval*, pages 639–648.

Zhankui He, Zhouhang Xie, Rahul Jha, Harald Steck, Dawen Liang, Yesu Feng, Bodhisattwa Prasad Majumder, Nathan Kallus, and Julian McAuley. 2023. Large language models as zero-shot conversational recommenders. In *Proceedings of the 32nd ACM international conference on information and knowledge management*, pages 720–730.

Yupeng Hou, Junjie Zhang, Zihan Lin, Hongyu Lu, Ruobing Xie, Julian McAuley, and Wayne Xin Zhao. 2024. Large language models are zero-shot rankers for recommender systems. In *European Conference on Information Retrieval*, pages 364–381. Springer.

Jun Hu, Wenwen Xia, Xiaolu Zhang, Chilin Fu, Weichang Wu, Zhaoxin Huan, Ang Li, Zuoli Tang, and Jun Zhou. 2024. Enhancing sequential recommendation via llm-based semantic embedding learning. In *Companion Proceedings of the ACM on Web Conference 2024*, pages 103–111.

Wang-Cheng Kang and Julian McAuley. 2018. Self-attentive sequential recommendation. In *2018 IEEE*

696	<i>international conference on data mining (ICDM)</i> ,	Yunzhu Pan, Chen Gao, Jianxin Chang, Yanan Niu,	749
697	pages 197–206. IEEE.	Yang Song, Kun Gai, Depeng Jin, and Yong Li. 2023.	750
698	Diederik P Kingma. 2014. Adam: A method for stochas-	Understanding and modeling passive-negative feed-	751
699	tic optimization. <i>arXiv preprint arXiv:1412.6980</i> .	back for short-video sequential recommendation. In	752
700	Wenqiang Lei, Gangyi Zhang, Xiangnan He, Yisong	<i>Proceedings of the 17th ACM conference on recom-</i>	753
701	Miao, Xiang Wang, Liang Chen, and Tat-Seng Chua.	<i>mender systems</i> , pages 540–550.	754
702	2020. Interactive path reasoning on graph for conver-	Shashank Rajput, Nikhil Mehta, Anima Singh, Raghu-	755
703	sational recommendation. In <i>Proceedings of the 26th</i>	nandan Hulikal Keshavan, Trung Vu, Lukasz Heldt,	756
704	<i>ACM SIGKDD international conference on knowl-</i>	Lichan Hong, Yi Tay, Vinh Tran, Jonah Samost, et al.	757
705	<i>edge discovery &amp; data mining</i> , pages 2073–2083.	2023. Recommender systems with generative re-	758
706	Chengxi Li, Yejing Wang, Qidong Liu, Xiangyu Zhao,	trieval. <i>Advances in Neural Information Processing</i>	759
707	Wanyu Wang, Yiqi Wang, Lixin Zou, Wenqi Fan, and	<i>Systems</i> , 36:10299–10315.	760
708	Qing Li. 2023a. Strec: Sparse transformer for sequen-	Leheng Sheng, An Zhang, Yi Zhang, Yuxin Chen, Xiang	761
709	tial recommendations. In <i>Proceedings of the 17th</i>	Wang, and Tat-Seng Chua. 2024. Language models	762
710	<i>ACM Conference on Recommender Systems</i> , pages	encode collaborative signals in recommendation.	763
711	101–111.	Fei Sun, Jun Liu, Jian Wu, Changhua Pei, Xiao Lin,	764
712	Muyang Li, Zijian Zhang, Xiangyu Zhao, Wanyu Wang,	Wenwu Ou, and Peng Jiang. 2019. Bert4rec: Se-	765
713	Minghao Zhao, Runze Wu, and Ruocheng Guo.	quential recommendation with bidirectional encoder	766
714	2023b. Automlp: Automated mlp for sequential rec-	representations from transformer. In <i>Proceedings of</i>	767
715	ommendations. In <i>Proceedings of the ACM Web</i>	<i>the 28th ACM international conference on informa-</i>	768
716	<i>Conference 2023</i> , pages 1190–1198.	<i>tion and knowledge management</i> , pages 1441–1450.	769
717	Jiayi Liao, Sihang Li, Zhengyi Yang, Jiancan Wu,	Jiaxi Tang and Ke Wang. 2018. Personalized top-n se-	770
718	Yancheng Yuan, Xiang Wang, and Xiangnan He.	quential recommendation via convolutional sequence	771
719	2024. Llara: Large language-recommendation assis-	embedding. In <i>Proceedings of the eleventh ACM</i>	772
720	tant. In <i>Proceedings of the 47th International ACM</i>	<i>international conference on web search and data</i>	773
721	<i>SIGIR Conference on Research and Development in</i>	<i>mining</i> , pages 565–573.	774
722	<i>Information Retrieval</i> , pages 1785–1795.	Hugo Touvron, Louis Martin, Kevin Stone, Peter Al-	775
723	Yujie Lin, Chenyang Wang, Zhumin Chen, Zhaochun	bert, Amjad Almahairi, Yasmine Babaei, Nikolay	776
724	Ren, Xin Xin, Qiang Yan, Maarten de Rijke, Xiuzhen	Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti	777
725	Cheng, and Pengjie Ren. 2023. A self-correcting	Bhosale, et al. 2023. Llama 2: Open founda-	778
726	sequential recommender. In <i>Proceedings of the ACM</i>	tion and fine-tuned chat models. <i>arXiv preprint</i>	779
727	<i>Web Conference 2023</i> , pages 1283–1293.	<i>arXiv:2307.09288</i> .	780
728	Qidong Liu, Xian Wu, Wanyu Wang, Yejing Wang,	Bohao Wang, Feng Liu, Jiawei Chen, Yudi Wu, Xingyu	781
729	Yuanshao Zhu, Xiangyu Zhao, Feng Tian, and Yefeng	Lou, Jun Wang, Yan Feng, Chun Chen, and Can	782
730	Zheng. 2024. Large language model empowered	Wang. 2024. Llm4dsr: Leveraing large language	783
731	embedding generator for sequential recommendation.	model for denoising sequential recommendation.	784
732	<i>arXiv preprint arXiv:2409.19925</i> .	<i>arXiv preprint arXiv:2408.08208</i> .	785
733	Ben Mann, N Ryder, M Subbiah, J Kaplan, P Dhari-	Wenjie Wang, Fuli Feng, Xiangnan He, Liqiang Nie,	786
734	wal, A Neelakantan, P Shyam, G Sastry, A Askill,	and Tat-Seng Chua. 2021. Denoising implicit feed-	787
735	S Agarwal, et al. 2020. Language models are few-	back for recommendation. In <i>Proceedings of the</i>	788
736	shot learners. <i>arXiv preprint arXiv:2005.14165</i> , 1.	<i>14th ACM international conference on web search</i>	789
737	Sheshera Mysore, Andrew McCallum, and Hamed Za-	<i>and data mining</i> , pages 373–381.	790
738	mani. 2023. Large language model augmented nar-	Zhuosheng Zhang, Aston Zhang, Mu Li, and Alex	791
739	rative driven recommendations. In <i>Proceedings of</i>	Smola. 2022. Automatic chain of thought prompt-	792
740	<i>the 17th ACM Conference on Recommender Systems</i> ,	ing in large language models. <i>arXiv preprint</i>	793
741	pages 777–783.	<i>arXiv:2210.03493</i> .	794
742	Jianmo Ni, Jiacheng Li, and Julian McAuley. 2019. Jus-	Guorui Zhou, Xiaoqiang Zhu, Chenru Song, Ying Fan,	795
743	tifying recommendations using distantly-labeled re-	Han Zhu, Xiao Ma, Yanghui Yan, Junqi Jin, Han Li,	796
744	views and fine-grained aspects. In <i>Proceedings of</i>	and Kun Gai. 2018. Deep interest network for click-	797
745	<i>the 2019 conference on empirical methods in natural</i>	through rate prediction. In <i>Proceedings of the 24th</i>	798
746	<i>language processing and the 9th international joint</i>	<i>ACM SIGKDD international conference on knowl-</i>	799
747	<i>conference on natural language processing (EMNLP-</i>	<i>edge discovery &amp; data mining</i> , pages 1059–1068.	800
748	<i>IJCNLP)</i> , pages 188–197.	Kun Zhou, Yuanhang Zhou, Wayne Xin Zhao, Xiaoke	801
		Wang, and Ji-Rong Wen. 2020. Towards topic-guided	802
		conversational recommender system. <i>arXiv preprint</i>	803
		<i>arXiv:2010.04125</i> .	804