

Revisiting scalable sequential recommendation with Multi-Embedding Approach and Mixture-of-Experts

Qiushi Pan

pqs@mail.ustc.edu.cn

University of Science and Technology
of China
Hefei, Anhui, China

Hao Wang

wanghao3@ustc.edu.cn

University of Science and Technology
of China
Hefei, Anhui, China

Guoyuan An

anguoyuan@h-partners.com

Huawei Noah's Ark Lab
Singapore, Singapore, Singapore

Luankang Zhang

zhanglk5@mail.ustc.edu.cn

University of Science and Technology
of China
Hefei, Anhui, China

Wei Guo

guowei67@huawei.com

Huawei Noah's Ark Lab
Singapore, Singapore, Singapore

Yong Liu

liu.yong6@huawei.com

Huawei Noah's Ark Lab
Singapore, Singapore, Singapore

Abstract

In recommendation systems, how to effectively scale up recommendation models has been an essential research topic. While significant progress has been made in developing advanced and scalable architectures for sequential recommendation (SR) models, there are still challenges due to items' multi-faceted characteristics and dynamic item relevance in the user context. To address these issues, we propose Fuxi-MME, a framework that integrates a multi-embedding strategy with a Mixture-of-Experts (MoE) architecture. Specifically, to efficiently capture diverse item characteristics in a decoupled manner, we decompose the conventional single embedding matrix into several lower-dimensional embedding matrices. Additionally, by substituting relevant parameters in the Fuxi Block with an MoE layer, our model achieves adaptive and specialized transformation of the enriched representations. Empirical results on public datasets show that our proposed framework outperforms several competitive baselines.

CCS Concepts

• Information systems → Recommender systems.

Keywords

Sequential Recommendation, Multi-embedding, Mixture-of-Experts

ACM Reference Format:

Qiushi Pan, Hao Wang, Guoyuan An, Luankang Zhang, Wei Guo, and Yong Liu. 2018. Revisiting scalable sequential recommendation with Multi-Embedding Approach and Mixture-of-Experts. In *Proceedings of Make sure to enter the correct conference title from your rights confirmation email (Conference acronym 'XX)*. ACM, New York, NY, USA, 10 pages. <https://doi.org/XXXX.XXX.XXXXXX>

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than the author(s) must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from permissions@acm.org.
Conference acronym 'XX, Woodstock, NY

© 2018 Copyright held by the owner/author(s). Publication rights licensed to ACM.
ACM ISBN 978-1-4503-XXXX-X/2018/06
<https://doi.org/XXXXXXX.XXXXXXX>

1 Introduction

Recommendation systems [23, 72, 16, 32, 54, 55, 53, 52, 51, 50, 46, 74, 62] aim to analyze users' historical behaviors to identify and present items that align with their potential interests. In industrial applications, Sequential Recommendation (SR) [78, 71, 73, 67, 69, 49, 13, 14] has become an indispensable component, as it focuses on predicting a user's next action by modeling the temporal dynamics of their interaction history. This subfield has consequently attracted substantial attention from both academia and industry.

In recent years, the field has embraced the autoregressive paradigm [23, 7], drawing inspiration from advances in natural language processing [47, 64]. Early implementations, such as GRU4Rec [18], utilized Recurrent Neural Networks (RNNs) [19, 35] to capture short-term user preferences but struggled with long-range patterns due to recursive state updates and gradient issues. Subsequently, SASRec [23] marked a significant leap forward by incorporating a Transformer-based [48] architecture that leverages self-attention to capture more complex and dynamic user preferences. Despite these advancements, scalability remains a persistent challenge—simply enlarging model size does not guarantee improved performance [44]. To overcome this, contemporary research has turned to scaling laws [24, 10, 43], revealing that model performance can grow predictably with increases in parameters, data, and compute. For instance, HSTU [71] modified the attention mechanism for large and non-stationary vocabulary and proposes the M-FALCON algorithm to speed up model inference, which exhibits scaling-up effects in recommendation. Similarly, Fuxi-Alpha [65] presents a novel architecture that disentangles temporal, positional, and semantic features, offering a more scalable and modular design. Together, these works demonstrate that autoregressive frameworks can successfully harness scaling laws for considerable performance enhancements in recommendation.

Despite the considerable progress achieved in designing scalable architectures for SR, most existing studies emphasize model structure while **overlooking the embedding layer**—a component that often constitutes the majority of parameters and represents a major bottleneck in expressive capacity [22, 55]. Consequently, the ability to capture rich item features is often limited by the constraints of a single, monolithic embedding vector. While some approaches attempt to mitigate this by increasing embedding dimensionality,

such brute-force scaling yields diminishing returns [12]. Drawing inspiration from Large Language Models (LLMs) [42, 34, 41, 70, 9] and their capacity to represent multifaceted semantic spaces [76, 21], we argue for a more principled rethinking of the embedding layer itself. By rearchitecting the embedding space, we can more effectively model complex item attributes and alleviate the representational bottleneck that limits current SR systems.

However, advancing this direction introduces new challenges. Research on embedding structures for sequential recommendation remains limited, and two critical challenges: (1) An item's identity is inherently multi-faceted—encompassing attributes such as category, brand, and style—yet a single dense embedding struggles to disentangle these heterogeneous signals, constraining the model's expressiveness. **Designing a mechanism that can represent such diverse features in a decoupled and structured manner is therefore essential.** (2) The relevance of item attributes varies with user context: one user may prioritize brand affinity, while another focuses on visual style. **Capturing these differences requires a dynamic, input-aware mechanism that can adaptively transform and weigh features according to sequence context.**

To tackle these challenges, we propose **Fuxi-MME**, a framework that integrates a multi-embedding strategy with a Mixture-of-Experts (MoE) architecture built upon the high-performing Fuxi-Alpha model. Our approach decomposes the conventional single embedding into multiple lower-dimensional sub-embeddings, each responsible for capturing a distinct facet of an item's identity. This design enhances expressiveness without inflating parameter count and can be seamlessly integrated as a plug-and-play component. Furthermore, to enable adaptive and specialized processing of these enriched representations, we introduce an MoE layer within the Fuxi block. This layer serves as a dynamic routing mechanism, directing sequence inputs to specialized expert networks based on their contextual characteristics. Extensive experiments show that the synergy between expressive embeddings and adaptive experts enables end-to-end learning that effectively models heterogeneous user behaviors and nuanced feature interactions. The main contributions of this paper are summarized as follows:

- We present Fuxi-MME, a framework that synergizes multi-embedding and expert routing to achieve more context-aware and adaptive sequence modeling.
- We propose a multi-faceted embedding strategy for sequential recommendation, designed to overcome the representational bottleneck of conventional single-vector embeddings in a parameter-efficient manner.
- We develop a MoE module tailored for Transformer-based architectures, enabling dynamic, input-dependent feature transformation and enhancing interaction modeling.
- Comprehensive experiments on three public benchmark datasets demonstrate that Fuxi-MME establishes new state-of-the-art results. Further studies confirm the framework's effectiveness.

2 Related Work

2.1 Sequential Recommendation

Sequential recommendation [1, 36, 61] aims to predict a user's next interaction based on their historical behavior. Early approaches predominantly relied on Markov chains [15, 39] to model transition probabilities between items. With the rise of deep learning, subsequent research adopted neural architectures such as recurrent neural networks (RNNs) [17, 18] and convolutional neural networks (CNNs) [3, 63] to better capture complex temporal dependencies.

More recently, Transformer-based architectures have achieved state-of-the-art performance in this domain. For instance, SASRec [23] formulates the task as an autoregressive prediction problem, employing self-attention to model user behavior sequences. In contrast, BERT4Rec [45] utilizes a bidirectional self-attention mechanism with a masked item prediction objective, allowing the model to leverage both past and future context for richer item representations.

Beyond attention-based models, generative paradigms such as GANs [2, 56], VAEs [77], and diffusion models [58, 59] have been explored to improve sequence modeling and uncertainty estimation. Furthermore, large language models (LLMs) [57] have recently been incorporated into recommendation systems to generate semantic input embeddings [29, 30, 40, 42], thereby introducing external knowledge and improving representation quality.

2.2 Scaling Recommendation Models

Building on the generative perspective, recent studies [6, 20, 33, 54, 66] have reformulated recommendation as a sequence generation problem, where the model autoregressively produces tokens corresponding to item identifiers. A common research trend has been to scale model parameters to enhance representational and generative capacity [10].

Following this direction, TIGER [38] employs RQ-VAE [27] to construct semantically meaningful discrete item codes and trains a Transformer to predict them in sequence. HSTU [71] extends this idea by scaling to trillion-parameter models that unify heterogeneous user behaviors into a single generative sequence, demonstrating consistent performance gains with model size. Wukong [72] introduces a stackable layer design for scalable architecture expansion, while another line of work [12] mitigates embedding collapse using a multi-embedding strategy.

This generative paradigm has recently been extended across the entire recommendation pipeline [11, 75]. For example, OneRec [7] replaces traditional ID-based item representations with semantic encodings, integrates Delayed Propagation Optimization (DPO) [37] into a Mixture-of-Experts (MoE) [5, 79] Transformer framework, and unifies multi-stage learning into a single end-to-end training process. Although these methods achieve strong performance and scalability, they primarily focus on architectural scaling, while the potential of scaling input representations remains underexplored.

3 Problem Definition

Let $\mathcal{U} = \{u_1, u_2, \dots, u_{|\mathcal{U}|}\}$ represent the set of all users and $\mathcal{I} = \{i_1, i_2, \dots, i_{|\mathcal{I}|}\}$ represent the set of all unique items in the system.

The historical interactions of each user are captured as a chronologically ordered sequence. For a given user $u \in \mathcal{U}$, their interaction history is denoted as $S^u = (i_1^u, i_2^u, \dots, i_{n_u}^u)$, where $i_t^u \in \mathcal{I}$ is the item that user u interacted with at time step t , and n_u is the total length of the sequence for that user.

For practical implementation and efficient batch processing in deep learning models, these variable-length sequences are typically converted to a fixed length, n . If a user's sequence is shorter than n ($n_u < n$), the sequence is padded with a special $\langle \text{PAD} \rangle$ token at the beginning. If the sequence is longer than n ($n_u > n$), only the most recent n interactions are retained, as these are generally the most indicative of the user's current interests.

The fundamental task of sequential recommendation is to predict the next item, i_{n+1}^u , that a user u is most likely to interact with, given their historical interaction sequence $S^u = (i_1^u, \dots, i_n^u)$. Formally, this is a probabilistic task of finding the item that maximizes the conditional probability:

$$\hat{i}_{n+1}^u = \arg \max_{i \in \mathcal{I}} P(i_{n+1} = i | S^u) \quad (1)$$

A sequential recommendation model, parameterized by θ , is trained to learn this probability distribution $P(\cdot | S^u; \theta)$. This is typically achieved by designing a model f_θ that performs two main functions:

- (1) **Sequence Encoding:** It encodes the input sequence S^u into a high-dimensional vector representation, $h_n^u \in \mathbb{R}^d$, which is intended to capture the user's current preferences and intent.
- (2) **Item Scoring:** It computes a relevance score between the user's representation h_n^u and the embedding vector $e_i \in \mathbb{R}^d$ of every candidate item $i \in \mathcal{I}$. The probability is then often estimated via a softmax function over these scores.

In practice, the goal is not just to predict a single item but to generate a ranked list of the Top-K most probable items for recommendation.

4 Methodology

In this section, we introduce the proposed **Fuxi-MME** framework. We begin by providing a concise overview of the Fuxi- α model, which serves as the foundational architecture for our work. Subsequently, we detail our two primary contributions: the **Multi-embedding Approach** designed to capture the multifaceted nature of items, and the **Mixture-of-Experts (MoE) enhanced Fuxi Block**, which enables dynamic and specialized feature processing. Finally, we describe the model's training objective.

4.1 The Foundational Backbone Fuxi- α

We select Fuxi- α [65] as our base model due to its demonstrated state-of-the-art performance and its scalable architecture, which is specifically designed to handle large-scale sequential recommendation tasks. The core of Fuxi- α is a stack of L identical decoder layers, referred to as *Fuxi blocks*. Each Fuxi block contains two main sub-layers: an Adaptive Multi-channel Self-attention (AMS) layer and a Multi-stage Feed-Forward Network (MFFN).

4.1.1 Adaptive Multi-channel Self-attention (AMS). Unlike standard self-attention, the AMS layer models user sequences through

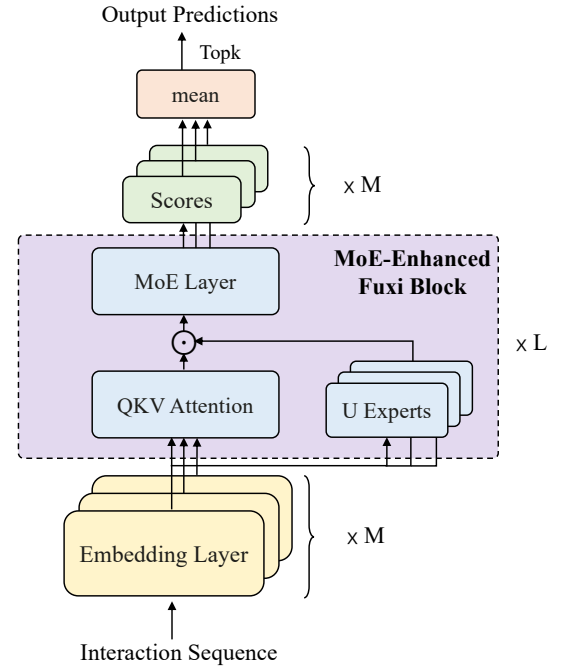


Figure 1: Overview of framework.

three distinct channels to capture different aspects of item-to-item relationships:

- (1) **Semantic Channel:** This channel captures the inherent feature-based relationships between items, akin to the standard self-attention mechanism in Transformers. For an input sequence representation x^{l-1} from the previous layer, the computation is as follows:

$$\hat{x}^l = \text{RMSN}(x^{l-1}) \quad (2)$$

$$q^l = \phi(\hat{x}^l W_q^l), \quad k^l = \phi(\hat{x}^l W_k^l), \quad v^l = \phi(\hat{x}^l W_v^l) \quad (3)$$

$$a_h^l = \frac{1}{\sqrt{d_k}} \phi(q^l (k^l)^T) \quad (4)$$

where $W_q^l, W_k^l, W_v^l \in \mathbb{R}^{d \times d_h}$ are learnable projection matrices for the attention mechanism. RMSN denotes Root Mean Square Layer Normalization, and ϕ is a non-linear activation function, typically SiLU.

- (2) **Temporal and Positional Channels:** To explicitly model the sequence order and time intervals, Fuxi- α incorporates two additional channels whose attention scores are calculated separately:

$$(a_t^l)_{i,j} = \alpha(t_j - t_i), \quad (a_p^l)_{i,j} = \beta_{j-i} \quad (5)$$

Here, $(a_t^l)_{i,j}$ represents the temporal bias based on the timestamp difference between items i and j , while $(a_p^l)_{i,j}$ represents the relative positional bias. The parameters α and β are learnable mappings that encode these biases.

The outputs from these three channels are then aggregated and combined with a gated projection, similar to the mechanism in

HSTU [71]:

$$h^l = \text{RMSN}(\text{concat}(a_h^l v^l, a_i^l v^l, a_p^l v^l) \otimes \phi(x^{l-1} W_u^l)) \quad (6)$$

where \otimes denotes element-wise multiplication, which acts as a dynamic gate to control the information flow from the attention module.

4.1.2 Multi-stage Feed-Forward Network (MFFN). The MFFN processes the output from the AMS layer through a two-stage process. First, a linear projection is applied, followed by a residual connection:

$$o^l = h^l W_o^l + x^{l-1} \quad (7)$$

Second, this intermediate representation is passed through a well-designed FFN structure that uses gating to control information flow and model complex feature interactions:

$$x^l = \text{FFN}_l(\text{RMSN}(o^l)) + o^l \quad (8)$$

$$\text{FFN}_l(x) = (\phi(x W_1^l) \otimes (x W_2^l)) W_3^l \quad (9)$$

where $W_1^l, W_2^l \in \mathbb{R}^{d \times d_{\text{FFN}}}$ and $W_3^l \in \mathbb{R}^{d_{\text{FFN}} \times d}$ are learnable weight matrices. This gated structure has been shown to be more effective than standard FFNs.

4.2 Multi-embedding Approach

A primary limitation of existing sequential models is their reliance on a single, monolithic embedding vector for each item. This design forces a single vector to encode all of an item's multifaceted characteristics (e.g., category, brand, style, price), creating a representational bottleneck and potentially leading to embedding collapse [12], where distinct features become entangled.

Inspired by [12], we propose a multi-embedding approach to address this. Instead of a single embedding matrix, we introduce M independent embedding matrices, $\{E_1, E_2, \dots, E_M\}$, where each $E_i \in \mathbb{R}^{|I| \times (d/M)}$ and $|I|$ is the total number of items. This partitions the total embedding dimension d into M smaller, specialized "sub-spaces," with the crucial benefit of maintaining the same total parameter count as a standard embedding layer.

Given a user's interaction sequence $S_u = (i_1, i_2, \dots, i_n)$, we map each item i_t to M corresponding sub-embeddings. This results in M distinct input representation matrices $\{X_1^0, X_2^0, \dots, X_M^0\}$, where $X_i^0 \in \mathbb{R}^{n \times (d/M)}$.

Critically, and in contrast to prior work that might process each representation stream independently, we feed all M representation matrices into a single, shared sequential decoder (i.e., the stacked Fuxi blocks):

$$\{X_1^L, \dots, X_M^L\} = \text{Decoder}(\{X_1^0, \dots, X_M^0\}; \Theta) \quad (10)$$

where Θ represents the shared parameters of the Fuxi- α decoder. This design choice is deliberate: it allows the model to learn universal sequential patterns (e.g., temporal dynamics, general user behavior) across all embedding spaces through the shared parameters of Θ , while still allowing each of the M streams to preserve its specialized, disentangled feature information.

After processing through all L decoder layers, we obtain M final output representations. During prediction, these are used to compute scores for candidate items, which are then aggregated via mean pooling to produce the final recommendation score.

4.3 MoE-enhanced Fuxi Block

While the multi-embedding approach enriches the input representations, the model must also be able to process these diverse and multifaceted features adaptively. A standard, dense network applies the same transformations to all inputs, which is suboptimal when dealing with varied signals. To address this, we enhance the Fuxi block by incorporating a Mixture-of-Experts (MoE) architecture.

The MoE paradigm decouples model capacity from computational cost by activating only a sparse subset of the network for any given input. We integrate MoE layers to replace the dense FFN and the gating projection matrix U within each Fuxi block. This allows the model to learn specialized "expert" networks and dynamically route different parts of the input sequence to the most appropriate experts.

Specifically, we replace the dense FFN from Equation 9 with a Sparse-Gated MoE layer. For a given input x , a lightweight gating network $G(x)$ calculates routing weights to select the top- k experts:

$$\text{MoE-Output}(x) = \sum_{i=1}^N G(x)_i \cdot \text{Expert}_i(x) \quad (11)$$

where N is the total number of experts (e.g., individual FFNs) and $G(x)_i$ is the weight assigned to the i -th expert for input x . The gating mechanism is designed to be sparse:

$$G(x) = \text{Softmax}(\text{KeepTopK}(H(x), k)) \quad (12)$$

where k is hyperparameter, ensuring only a small subset of experts are activated. The routing logits $H(x)$ are computed as:

$$H(x) = x \cdot W_g + \text{StandardNormal}() \cdot \text{Softplus}(x \cdot W_{\text{noise}}) \quad (13)$$

Here, W_g and W_{noise} are learnable parameters of the gating network. The noise term is added during training to improve load balancing across experts, preventing a few experts from being consistently chosen. The KeepTopK function sets the logits of all non-top- k experts to $-\infty$, effectively zeroing out their weights after the Softmax operation:

$$\text{KeepTopK}(v, k)_i = \begin{cases} v_i & \text{if } v_i \text{ is in the top } k \text{ elements of } v \\ -\infty & \text{otherwise} \end{cases} \quad (14)$$

By integrating this MoE structure, Fuxi-MME can dynamically allocate its parameters based on the input, allowing for a significant increase in model capacity while keeping the computational cost per forward pass nearly constant. This is particularly effective for modeling the diverse patterns emerging from our multi-embedding representations.

4.4 Model Training and Optimization Objective

We train the Fuxi-MME framework in an end-to-end fashion using an autoregressive next-item prediction task. For a given user sequence $(i_1, i_2, \dots, i_{n-1})$, the model's objective is to accurately predict the next item, i_n .

The final similarity score between the user's sequence representation at time t and a candidate item i is calculated by aggregating the scores from all M embedding spaces:

$$r(t, i) = \frac{1}{M} \sum_{k=1}^M \text{cosine_similarity}(x_{t,k}^L, e_{i,k}) \quad (15)$$

Table 1: Dataset statistics.

Dataset	#Users	#Items	#Interactions	Avg.Length	Sparsity
Amazon-Books	694897	695761	10053086	14.47	99.99%
Amazon-Beauty	40226	57288	353962	8.80	99.85%
Yelp	30431	20033	255492	8.40	99.96%

where $x_{t,k}^L$ is the final output representation for the t -th item in the sequence from the k -th embedding space, and $e_{i,k}$ is the corresponding k -th sub-embedding of the candidate item i .

Due to the massive size of the item vocabulary in real-world scenarios, computing a full softmax over all items is computationally prohibitive. Therefore, we employ a sampled-softmax loss [26] for efficient training. For each positive instance (the true next item), we randomly sample a set of N negative items from the item catalog. The model is then trained to maximize the score of the positive item relative to the negative ones.

The overall training loss is the negative log-likelihood over all sequences in the training set:

$$\mathcal{L}_{\text{rec}} = - \sum_{u \in U} \sum_{t=1}^{n_u} \log \frac{\exp(r(t, i_t))}{\exp(r(t, i_t)) + \sum_{j \in I_{\text{neg}}} \exp(r(t, j))} \quad (16)$$

where I_{neg} is the set of N randomly sampled negative items. This objective function effectively trains all components of the Fuxi-MME model, including the multiple embedding matrices, the shared Fuxi blocks, and the MoE layers, to work in concert for sequential recommendation.

5 Experiments

In this section, we present a comprehensive empirical evaluation of our proposed **Fuxi-MME** model. Our goal is to answer the following research questions:

- **(RQ1)** Does Fuxi-MME outperform state-of-the-art sequential recommendation baselines across different datasets?
- **(RQ2)** What are the individual contributions of the multi-embedding approach and the Mixture-of-Experts (MoE) architecture to the model's overall performance?
- **(RQ3)** How sensitive is Fuxi-MME to its key hyperparameters, namely the number of embeddings (M) and the number of experts (n)?
- **(RQ4)** Is the placement of the MoE layer within the attention block optimal, and what does this imply about its function?

5.1 Experimental Setup

5.1.1 Datasets. To evaluate the efficacy of our proposed method, we conduct extensive experiments on three public datasets. These datasets are derived from real-world online user interactions and are commonly adopted in sequential recommendation research. The statistics of these datasets are shown in 1. Detailed descriptions are provided as follows:

- *Amazon-Books* and *Amazon-Beauty*. The Amazon-Books and Amazon-Beauty datasets are subsets of the Amazon-Reviews dataset¹. It comprises real-world user interactions collected

from the Amazon website, including both product metadata and user review records. The product metadata encompasses fields such as product ID, title, and category classifications. The user review records contain information such as user IDs, product IDs, review texts, ratings, and timestamps.

- *Yelp*². The Yelp dataset is collected from the popular business website Yelp, comprises information about businesses, users, and reviews. In this work, we obtain a subset of the full dataset, following the preprocessing strategy in [60, 4, 68].

These three datasets are specifically chosen to provide a comprehensive testbed for our model. The two Amazon datasets represent dense e-commerce scenarios with distinct product feature distributions, while the Yelp dataset is known for its greater sparsity and different user behavior patterns related to local businesses. Success across these varied domains demonstrates the robustness and generalizability of a given approach.

5.1.2 Compared Methods. We compare our proposed method with several classic and state-of-the-art baselines for sequential recommendation. This includes classic RNN-based models that set early benchmarks, modern Transformer-based architectures that represent the current state-of-the-art in modeling complex dependencies, and emerging State-Space Models (SSMs) which offer an efficient alternative for sequence modeling. This diverse set ensures a thorough and challenging comparison. The specific baselines are as follows:

- GRU4Rec[18]. GRU4Rec applies GRU layers to capture user preferences within interaction sequences.
- NARM[28]. NARM combines RNN and attention mechanism to consider both user sequential behaviors and main purpose in sessions.
- SASRec[23]. SASRec proposes a self-attention based sequential model to capture long-term semantics in users' interaction sequences.
- Mamba4Rec[31]. Mamba4Rec is the first work to utilize state space models for efficient sequential recommendation.
- Tim4Rec[8]. Tim4Rec proposes a time-aware Mamba for sequential recommendation.
- HSTU[71]. HSTU designs a new self-attention based sequential recommendation architecture for higher efficiency and scaling up effects.
- Fuxi- α [65]. Fuxi- α introduces an Adaptive Multi-channel Self-attention mechanism to distinctly model temporal, positional and semantic features.

5.1.3 Evaluation Protocols. We employ a leave-one-out strategy for evaluation: the last item is for testing, the second-to-last is for validation, and all others are for training. The partitioning of the

¹<http://jmcauley.ucsd.edu/data/amazon/>

²<https://www.yelp.com/dataset>

Table 2: Overall Performance across three datasets. For each dataset, the best result is bolded while the second-best result is underlined.(p-value<0.05)

Model	Amazon-books				Amazon-beauty				Yelp			
	NG@10	NG@50	HR@10	HR@50	NG@10	NG@50	HR@10	HR@50	NG@10	NG@50	HR@10	HR@50
NARM	0.0109	0.0188	0.0209	0.0580	0.0201	0.0305	0.0361	0.0840	0.0172	0.0314	0.0338	0.1003
SASRec	0.0202	0.0324	0.0378	0.0944	0.0313	0.0445	0.0553	0.1160	0.0186	0.0333	0.0363	0.1052
GRU4Rec	0.0188	0.0302	0.0352	0.0877	0.0308	0.0442	0.0548	0.1165	0.0202	<u>0.0365</u>	0.0394	<u>0.1158</u>
Tim4Rec	0.0226	0.0358	0.0421	0.1030	0.0291	0.0421	0.0520	0.1120	0.0173	0.0328	0.0347	0.1072
Mamba4Rec	0.0221	0.0344	0.0405	0.0970	0.0290	0.0418	0.0508	0.1097	0.0173	0.0322	0.0344	0.1039
HSTU	0.0255	0.0391	0.0470	0.1096	<u>0.0337</u>	<u>0.0485</u>	<u>0.0587</u>	<u>0.1270</u>	<u>0.0205</u>	0.0356	<u>0.0399</u>	0.1102
Fuxi-alpha	<u>0.0279</u>	<u>0.0426</u>	<u>0.0509</u>	<u>0.1191</u>	0.0320	0.0459	0.0561	0.1198	0.0195	0.0352	0.0378	0.1109
Fuxi-MME	0.0294	0.0441	0.0534	0.1213	0.0358	0.0508	0.0626	0.1312	0.0225	0.0396	0.0443	0.1242

dataset is consistent with [71, 65]. We use two widely adopted metrics for Top- K recommendation lists: Hit Rate (HR@ K) and Normalized Discounted Cumulative Gain (NDCG@ K), with $K \in \{10, 50\}$. Crucially, we rank the ground-truth item against the entire item pool to ensure a fair and unbiased evaluation, avoiding biases from negative sampling during evaluation.

5.1.4 Implementation Details. To ensure fair evaluation of model performance, all models are implemented using PyTorch and trained on GPUs. We use the Adam [25] optimizer with a learning rate of $1e^{-3}$. The sampled-softmax loss used a batch size of 256 with 128 negative samples per positive instance. We train all models for 100 epochs and use an early stop strategy, which is quitting training without improving the metrics for 10 epochs. For self-attention-based models, the number of layers is fixed at 2. The embedding dimension of baselines is searched in [128, 256, 512] to compare with the multi-embedding approach. The embedding dimension of our proposed method is set to 128, and the number of embeddings is varied in [1, 2, 4].

5.2 Overall Performance (RQ1)

The main results of our comparative evaluation are presented in Table 2. The experimental findings clearly and consistently demonstrate the effectiveness of our proposed **Fuxi-MME** framework across all three benchmark datasets. We provide a detailed analysis below:

- The results confirm the strength of Transformer-based models in sequential recommendation. Self-attention models generally outperform RNN-based methods and the State-Space Models. This performance gap underscores the self-attention mechanism’s superior ability to capture complex, long-range dependencies and model dynamic user preferences. While SSM-based models offer impressive efficiency, these results suggest that on standard benchmarks, their recommendation performance has not yet surpassed that of state-of-the-art Transformer architectures.
- Our proposed model, Fuxi-MME, sets a new state-of-the-art by a significant margin across all datasets and evaluation

Table 3: Ablation Study.

Model	NG@10	NG@50	HR@10	HR@50	MRR
(a) w/o multi-embedding	0.0306	0.0444	0.0543	0.1178	0.0277
(b) ensemble	0.0341	0.0489	0.0588	0.1267	0.0311
(c) w/o MoE	0.0344	0.0489	0.0598	0.1262	0.0311
Fuxi-MME	0.0358	0.0508	0.0626	0.1312	0.0322

metrics. It consistently surpasses even the strongest and most recent baselines, including its own backbone, Fuxi- α . For instance, on the challenging and sparse Amazon-Books dataset, Fuxi-MME improves upon the best baseline (Fuxi- α) by 5.37% in NDCG@10 and 4.91% in HR@10. This substantial and consistent performance gain provides strong empirical evidence for our central thesis. The improvement is attributed to our novel architectural design: the multi-embedding approach provides a richer, more expressive representation of items by disentangling their multifaceted characteristics, and the MoE layer enables the model to process these nuanced features in an adaptive, input-dependent manner.

- Furthermore, comparing HSTU and Fuxi- α , Fuxi- α performs better on the Amazon-Books dataset, while HSTU performs better on the other two datasets. This is because the embedding size is searched within the range of [128, 256, 512] when evaluating the two models, while for Fuxi- α , its performance tends to decline when the embedding size increases to 512. This result empirically validates the core motivation of our paper. It suggests that a single, high-dimensional vector struggles to effectively organize diverse item features, leading to a representational bottleneck. Fuxi-MME addresses this by structuring the embedding space, leading to more consistent and scalable performance gains.

5.3 Ablation Study (RQ2)

To rigorously evaluate the individual contributions of our two primary architectural innovations, i.e. the multi-embedding approach

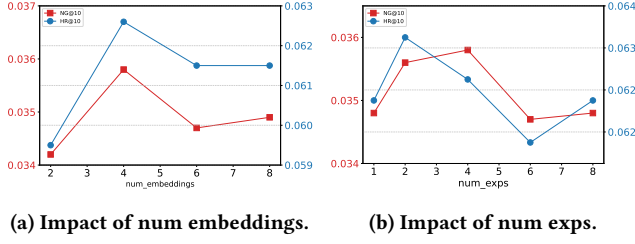


Figure 2: Analysis of Hyper-parameter.

and the Mixture-of-Experts (MoE) layer—we conduct a comprehensive ablation study. We design three distinct variants of our model and tested them on the Amazon-Beauty dataset. The results, presented in Table 3, allow us to systematically dissect the sources of Fuxi-MME’s performance gains.

- (a) **w/o multi-embedding:** This variant reverts to a standard, monolithic embedding layer with the same total dimension but removes the multi-embedding structure. It isolates the impact of our core representational hypothesis.
- (b) **ensemble:** This variant trains M separate Fuxi- α models, each with its own embedding space and its own decoder. The final predictions are generated by averaging their output scores. This tests whether Fuxi-MME’s performance is a naive ensemble effect.
- (c) **w/o MoE:** This variant retains the multi-embedding input but replaces the MoE-enhanced Fuxi blocks with the standard, dense Fuxi blocks. This quantifies the specific benefit of the adaptive, expert-based processing.

Our analysis of the experimental results in Table 3 leads to the following key insights:

- Comparing with (a), our proposed model achieves a significant improvement, highlighting the advantage of the multi-embedding approach when keeping the same parameter size. This is the most critical result, as it directly validates our primary hypothesis. By simply structuring the embedding space into multiple decoupled representations, our model gains significant expressive power without increasing the total number of parameters. The inferior performance of the monolithic embedding variant suggests it struggles with feature entanglement and cannot efficiently capture the multifaceted nature of items. This result effectively confirms the existence of the representational bottleneck we hypothesized in the introduction. A single, dense vector, even with the same total dimensionality, lacks the structural inductive bias to effectively disentangle and organize diverse item attributes, leading to a less effective optimization landscape.
- The ensemble variant (b) performs worse than our unified Fuxi-MME model while being far less parameter-efficient (it requires M times the number of decoder parameters). This highlights the benefit of shared learning. It also explains that there is a part of sharing information between decomposed representations within the sequential recommendation model. In the ensemble model, each network is

isolated and must learn these universal patterns independently, which is less efficient and effective. This proves that enabling input-adaptive parameter specialization within our proposed architecture is a key architectural advantage.

- The variant (c) shows a slight decrease in performance compared to our model, indicating that the multi-embedding design is the primary driver of the performance gains. However, the full Fuxi-MME model with MoE still achieves superior results. This demonstrates the synergistic relationship between our two contributions. The multi-embedding layer provides the model with a rich, disentangled set of features, while the MoE layer provides the mechanism for adaptive processing. The MoE-enhanced Fuxi block acts as a dynamic routing network, adaptively sending different feature facets to specialized expert networks for transformation. Without MoE, the model is forced to apply the same dense transformation to all parts of the rich input, which leads to suboptimal performance. The MoE layer provides the final, crucial step of specialization that unlocks the full potential of the multi-faceted representations.

5.4 Impact of Hyper-parameter (RQ3)

To provide practical guidance for deploying Fuxi-MME and to better understand the model’s behavior, we conducted a sensitivity analysis on its two most critical hyperparameters: the number of embeddings (M) and the number of experts (n). All experiments were performed on the Amazon-Beauty dataset, with results presented in Figure 2.

5.4.1 Impact of the Number of Embeddings (M). To validate the impact of the number of embeddings, we vary the number of embeddings M in the set $\{2, 4, 6, 8\}$ while keeping the total embedding dimension d fixed. This means the dimension of each sub-embedding (d/M) varies, respectively, isolating the effect of representational structure.

As shown in the left subfigure in Figure 2, the model’s performance exhibits a clear inverted U-shaped trend. Moving from $M = 2$ to $M = 4$ results in a significant performance improvement, peaking at $M = 4$. This strongly supports our core hypothesis: decomposing the monolithic embedding space allows the model to dedicate different sub-spaces to learning more disentangled and specialized representations of item facets, providing greater representational power.

Beyond $M = 4$, we observe a notable decline in performance. This reveals a critical trade-off. While increasing M promotes feature disentanglement, it simultaneously reduces the dimensionality (capacity) of each individual sub-embedding. When M is large, each sub-embedding has a small dimension, which may be insufficient to capture the complexity of its specialized feature facet. The model suffers from a lack of expressive power within each specialized vector. Therefore, $M = 4$ strikes the optimal balance for this dataset.

This result highlights a key design consideration: the trade-off between representational diversity (which increases with M) and the representational capacity of each sub-embedding (which decreases with M). For a given total dimension, there exists an optimal point where the model has enough distinct feature channels without making any single channel too simplistic to be useful.

Table 4: Model Performance of Different MoE placement.

Model	NG@10	NG@50	HR@10	HR@50
Fuxi-MME w/o attentionMoE	0.0344	0.0489	0.0598	0.1262
w/ qMoE	0.0326	0.0469	0.0566	0.1224
w/ kMoE	0.0332	0.0477	0.0581	0.1252
w/ vMoE	0.0305	0.0437	0.0534	0.1140
w/ uMoE	0.0358	0.0508	0.0626	0.1312
w/o MoE	0.0331	0.0474	0.0578	0.1232

5.4.2 Impact of the Number of Experts (n). We vary the number of experts n within the MoE layer from 1 to 8. $n = 1$ is equivalent to a dense FFN, serving as a non-MoE baseline. The subfigure in the right of the Figure 2 illustrates the impact of the number of experts. Increasing the number of experts from 1 to 4 leads to a steady improvement in performance, with the optimum reached at $n = 4$. This demonstrates the value of the MoE architecture: having multiple specialized expert networks allows the model to learn a richer set of transformations and apply the most appropriate one based on the input, handling diverse data patterns more effectively than a single, dense network.

However, a sharp performance drop occurs when n is increased to 6 and beyond, with $n = 8$ performing worse than the non-MoE baseline ($n = 1$). When we set more experts, the gating network may struggle to assign tokens meaningfully to many experts, leading to underutilized experts and increased noise. Moreover, too many experts dilute the shared knowledge across experts, reducing their collective effectiveness and even underperforming a non-MoE baseline.

5.5 An Analysis of MoE Placement (RQ4)

The Fuxi block’s attention mechanism involves several learnable projection matrices: for queries (W_q), keys (W_k), values (W_v), and the final output gate (W_u). While our primary model applies MoE to the Feed-Forward Network and the W_u gate, a natural question arises: is this the optimal placement? To investigate this and better understand the interplay between sparsity and self-attention, we conducted a controlled experiment to determine which projection is most amenable to being replaced by a Mixture-of-Experts layer.

We created several variants of our model. In each variant, we replaced exactly one of the dense projection matrices (W_q , W_k , W_v , or W_u) within the attention block with an MoE layer, keeping all other components standard. These variants were then evaluated on the Amazon-Beauty dataset, with the results presented in Table 4.

The variant where the output projection W_u is replaced with an MoE layer is the only one that shows a significant performance improvement. This validates our final architectural choice. The W_u matrix operates on the aggregated, context-aware output of the self-attention mechanism. Placing the MoE layer here allows the model to perform context-dependent transformations. For instance, a specific expert can be activated to process a context vector representing focused brand loyalty, while another handles a context vector representing broad category exploration. This aligns

perfectly with our goal of applying specialized, adaptive transformations to the rich outputs generated from our multi-embedding inputs.

In stark contrast, replacing the query, key, or value projection matrices with MoE layers leads to a significant degradation in performance. We attribute this to the fundamental requirements of the self-attention mechanism. The core of self-attention relies on a stable semantic space where queries and keys can be reliably compared. Replacing W_q or W_k with MoE means the projection into this crucial space is handled by different, sparse experts for different tokens. This disrupts the semantic consistency required for meaningful similarity calculations (i.e., Query \times Key). The model needs a single, shared, linear transformation to ensure all items are mapped into a common, comparable space. Besides, replacing W_v with MoE introduces noise and inconsistency into the final aggregated context vector, as the value representations that are weighted and summed are generated by different experts, making it harder for subsequent layers to interpret the output.

In summary, this analysis provides strong evidence that the design of Fuxi-MME preserves the integrity and stability of the self-attention mechanism while introducing powerful, adaptive, and context-aware modulation on its output, thereby validating the design of Fuxi-MME.

6 Conclusion and Future Work

In this paper, we addressed a critical yet often overlooked limitation in sequential recommendation: the representational bottleneck imposed by the single embedding layer. While contemporary research has focused on scaling architectural components like attention mechanisms, we argued that the true potential of these models is constrained by the inability of a single vector to capture the multi-faceted nature of items and the lack of adaptive mechanisms to process these features in a context-dependent manner. To overcome these challenges, we proposed Fuxi-MME, a framework that synergizes a **multi-embedding strategy** with a **Mixture-of-Experts (MoE) architecture** upon the powerful Fuxi- α backbone. Specifically, our multi-embedding approach directly tackles the representation issue by decomposing each item’s identity into multiple, lower-dimensional sub-embeddings, enabling a more structured and expressive capture of its diverse attributes. Complementing this, the MoE architecture provides the necessary machinery for dynamic, input-aware processing, routing sequence information to specialized expert networks that can apply the most appropriate transformations based on the context. Our comprehensive experiments on three public benchmarks confirmed the effectiveness of our method. Fuxi-MME not only established a new state-of-the-art by significantly outperforming strong baselines but also demonstrated through ablation studies the critical and complementary roles of both its core components. These results offer a new perspective on the principle of scaling laws in recommendation: instead of relying solely on "brute-force" increases in parameter count, significant performance gains can be unlocked through more principled, structured, and adaptive architectural design.

Looking forward, this work opens several promising avenues for future research. While Fuxi-MME provides a robust framework, the interpretability of its disentangled embeddings and specialized

experts remains an exciting direction. Investigating what specific item facets or user behavior patterns are learned by each component could pave the way for more explainable and controllable recommendation systems. Furthermore, the principles of structured representation and adaptive computation can be extended to other domains, such as multi-modal recommendation, where different embedding spaces could naturally correspond to different data modalities (e.g., text, image). By continuing to explore these directions, we believe the architectural concepts pioneered in Fuxi-MME will contribute to building the next generation of more intelligent, expressive, and effective sequential recommendation systems.

References

- [1] Tesfaye Fenta Boka, Zhendong Niu, and Rama Bastola Neupane. 2024. A survey of sequential recommendation systems: techniques, evaluation, and future directions. *Information Systems*, 125, 102427.
- [2] Junshu Chen and GuangCong Liu. 2024. Collaborative filtering algorithm based on generative adversarial networks. In *2024 IEEE 14th International Conference on Electronics Information and Emergency Communication (ICEIEC)*. IEEE, 1–6.
- [3] Qi Chen, Guohui Li, Quan Zhou, Si Shi, and Deqing Zou. 2022. Double attention convolutional neural network for sequential recommendation. *ACM Transactions on the Web*, 16, 4, 1–23.
- [4] Yongjun Chen, Zhiwei Liu, Jia Li, Julian McAuley, and Caiming Xiong. 2022. Intent contrastive learning for sequential recommendation. In *Proceedings of the ACM web conference 2022*, 2172–2182.
- [5] Damai Dai et al. 2024. Deepseekmoe: towards ultimate expert specialization in mixture-of-experts language models. *arXiv preprint arXiv:2401.06066*.
- [6] Yashar Deldjoo et al. 2024. Recommendation with generative models. *arXiv preprint arXiv:2409.15173*.
- [7] Jiaxin Deng, Shiyao Wang, Kuo Cai, Lejian Ren, Qigen Hu, Weifeng Ding, Qiang Luo, and Guorui Zhou. 2025. Onerec: unifying retrieve and rank with generative recommender and iterative preference alignment. *arXiv preprint arXiv:2502.18965*.
- [8] Hao Fan, Mengyi Zhu, Yanrong Hu, Hailin Feng, Zhijie He, Hongjiu Liu, and Qingyang Liu. 2025. Tim4rec: an efficient sequential recommendation model based on time-aware structured state space duality model. *Neurocomputing*, 131270.
- [9] Hongchao Gu, Dexun Li, Kuicai Dong, Hao Zhang, Hang Lv, Hao Wang, Defu Lian, Yong Liu, and Enhong Chen. 2025. Rapid: efficient retrieval-augmented long text generation with writing planning and information discovery. *arXiv preprint arXiv:2503.00751*.
- [10] Wei Guo et al. 2024. Scaling new frontiers: insights into large recommendation models. *arXiv preprint arXiv:2412.00714*.
- [11] Xian Guo, Ben Chen, Siyuan Wang, Ying Yang, Chenyi Lei, Yuqing Ding, and Han Li. 2025. Onesug: the unified end-to-end generative framework for e-commerce query suggestion. *arXiv preprint arXiv:2506.06913*.
- [12] Xingzhuo Guo, Junwei Pan, Ximei Wang, Baixu Chen, Jie Jiang, and Mingsheng Long. 2023. On the embedding collapse when scaling up recommendation models. *arXiv preprint arXiv:2310.04400*.
- [13] Yongqiang Han, Hao Wang, Kefan Wang, Likang Wu, Zhi Li, Wei Guo, Yong Liu, Defu Lian, and Enhong Chen. 2024. Efficient noise-decoupling for multi-behavior sequential recommendation. In *Proceedings of the ACM Web Conference 2024*, 3297–3306.
- [14] Yongqiang Han, Likang Wu, Hao Wang, Guifeng Wang, Mengdi Zhang, Zhi Li, Defu Lian, and Enhong Chen. 2023. Guesr: a global unsupervised data-enhancement with bucket-cluster sampling for sequential recommendation. In *International conference on database systems for advanced applications*. Springer, 286–296.
- [15] Ruining He and Julian McAuley. 2016. Fusing similarity models with markov chains for sparse sequential recommendation. In *2016 IEEE 16th international conference on data mining (ICDM)*. IEEE, 191–200.
- [16] Xinran He et al. 2014. Practical lessons from predicting clicks on ads at facebook. In *Proceedings of the eighth international workshop on data mining for online advertising*, 1–9.
- [17] Balázs Hidasi and Alexandros Karatzoglou. 2018. Recurrent neural networks with top-k gains for session-based recommendations. In *Proceedings of the 27th ACM international conference on information and knowledge management*, 843–852.
- [18] Balázs Hidasi, Alexandros Karatzoglou, Linas Baltrunas, and Domonkos Tikk. 2015. Session-based recommendations with recurrent neural networks. *arXiv preprint arXiv:1511.06939*.
- [19] Sepp Hochreiter and Jürgen Schmidhuber. 1997. Long short-term memory. *Neural computation*, 9, 8, 1735–1780.
- [20] Yupeng Hou, An Zhang, Leheng Sheng, Zhengyi Yang, Xiang Wang, Tat-Seng Chua, and Julian McAuley. 2025. Generative recommendation models: progress and directions. In *Companion Proceedings of the ACM on Web Conference 2025*, 13–16.
- [21] Yuqing Huang et al. 2024. Chemeval: a comprehensive multi-level chemical evaluation for large language models. *arXiv preprint arXiv:2409.13989*.
- [22] Gangwei Jiang, Hao Wang, Jin Chen, Haoyu Wang, Defu Lian, and Enhong Chen. 2021. Xlightfm: extremely memory-efficient factorization machine. In *Proceedings of the 44th International ACM SIGIR Conference on Research and Development in Information Retrieval*, 337–346.
- [23] Wang-Cheng Kang and Julian McAuley. 2018. Self-attentive sequential recommendation. In *2018 IEEE international conference on data mining (ICDM)*. IEEE, 197–206.
- [24] Jared Kaplan et al. 2020. Scaling laws for neural language models. *arXiv preprint arXiv:2001.08361*.
- [25] Diederik P Kingma and Jimmy Ba. 2014. Adam: a method for stochastic optimization. *arXiv preprint arXiv:1412.6980*.
- [26] Anton Klenitskiy and Alexey Vasilev. 2023. Turning dross into gold loss: is bert4rec really better than sasrec? In *Proceedings of the 17th ACM Conference on Recommender Systems*, 1120–1125.
- [27] Doyup Lee, Chiheon Kim, Saehoon Kim, Minsu Cho, and Wook-Shin Han. 2022. Autoregressive image generation using residual quantization. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, 11523–11532.
- [28] Jing Li, Pengjie Ren, Zhumin Chen, Zhaochun Ren, Tao Lian, and Jun Ma. 2017. Neural attentive session-based recommendation. In *Proceedings of the 2017 ACM on Conference on Information and Knowledge Management*, 1419–1428.
- [29] Ruyi Li, Wenhao Deng, Yu Cheng, Zheng Yuan, Jiaqi Zhang, and Fajie Yuan. 2023. Exploring the upper limits of text-based collaborative filtering using large language models: discoveries and insights. *arXiv preprint arXiv:2305.11700*.
- [30] Xiangyang Li, Bo Chen, Lu Hou, and Ruiming Tang. 2023. Ctrl: connect tabular and language model for ctr prediction. *CoRR*.
- [31] Chengkai Liu, Jianghao Lin, Jianling Wang, Hanzhou Liu, and James Caverlee. 2024. Mamba4rec: towards efficient sequential recommendation with selective state space models. *arXiv preprint arXiv:2403.03900*.
- [32] Weiwen Liu, Wei Guo, Yong Liu, Ruiming Tang, and Hao Wang. 2023. User behavior modeling with deep learning for recommendation: recent advances. In *Proceedings of the 17th ACM Conference on Recommender Systems*, 1286–1287.
- [33] Zihan Liu, Yupeng Hou, and Julian McAuley. 2024. Multi-behavior generative recommendation. In *Proceedings of the 33rd ACM International Conference on Information and Knowledge Management*, 1575–1585.
- [34] Hang Lv, Sheng Liang, Hao Wang, Hongchao Gu, Yaxiong Wu, Wei Guo, Defu Lian, Yong Liu, and Enhong Chen. 2025. Costeer: collaborative decoding-time personalization via local delta steering. (2025). *arXiv: 2507.04756 [cs.CL]*.
- [35] Larry R Medsker, Lakhmi Jain, et al. 2001. Recurrent neural networks. *Design and applications*, 5, 64–67, 2.
- [36] Liwei Pan, Weiwei Pan, Meiyang Wei, Hongzhi Yin, and Zhong Ming. 2024. A survey on sequential recommendation. *arXiv preprint arXiv:2412.12770*.
- [37] Rafael Rafailov, Archit Sharma, Eric Mitchell, Christopher D Manning, Stefano Ermon, and Chelsea Finn. 2023. Direct preference optimization: your language model is secretly a reward model. *Advances in neural information processing systems*, 36, 53728–53741.
- [38] Shashank Rajput et al. 2023. Recommender systems with generative retrieval. *Advances in Neural Information Processing Systems*, 36, 10299–10315.
- [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*, 811–820.
- [40] Tianshu Shen, Jiaru Li, Mohamed Reda Bouadjenek, Zheda Mai, and Scott Sanner. 2023. Towards understanding and mitigating unintended biases in language model-driven conversational recommendation. *Information Processing & Management*, 60, 1, 103139.
- [41] Tingjia Shen, Hao Wang, Chuan Qin, Ruijun Sun, Yang Song, Defu Lian, Hengshu Zhu, and Enhong Chen. 2025. Genki: enhancing open-domain question answering with knowledge integration and controllable generation in large language models. (2025). <https://arxiv.org/abs/2505.19660> *arXiv: 2505.19660 [cs.CL]*.
- [42] Tingjia Shen, Hao Wang, Jiaqing Zhang, Sirui Zhao, Liangyue Li, Zulong Chen, Defu Lian, and Enhong Chen. 2024. Exploring user retrieval integration towards large language models for cross-domain sequential recommendation. *arXiv preprint arXiv:2406.03085*.
- [43] Tingjia Shen et al. 2025. Optimizing sequential recommendation models with scaling laws and approximate entropy. (2025). <https://arxiv.org/abs/2412.00430> *arXiv: 2412.00430 [cs.AI]*.
- [44] Tingjia Shen et al. 2024. Optimizing sequential recommendation models with scaling laws and approximate entropy. *arXiv preprint arXiv:2412.00430*.
- [45] 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*, 1441–1450.
- [46] Junxiong Tong, Mingjia Yin, Hao Wang, Qiushi Pan, Defu Lian, and Enhong Chen. 2024. Mdap: a multi-view disentangled and adaptive preference learning framework for cross-domain recommendation. In *International Conference on Web Information Systems Engineering*. Springer, 164–178.
- [47] Hugo Touvron et al. 2023. Llama: open and efficient foundation language models. *arXiv preprint arXiv:2302.13971*.
- [48] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. 2017. Attention is all you need. *Advances in neural information processing systems*, 30.
- [49] Hao Wang, Yongqiang Han, Kefan Wang, Kai Cheng, Zhen Wang, Wei Guo, Yong Liu, Defu Lian, and Enhong Chen. 2024. Denoising pre-training and customized prompt learning for efficient multi-behavior sequential recommendation. *arXiv preprint arXiv:2408.11372*.
- [50] Hao Wang, Defu Lian, Hanghang Tong, Qi Liu, Zhenya Huang, and Enhong Chen. 2021. Decoupled representation learning for attributed networks. *IEEE Transactions on Knowledge and Data Engineering*, 35, 3, 2430–2444.
- [51] Hao Wang, Defu Lian, Hanghang Tong, Qi Liu, Zhenya Huang, and Enhong Chen. 2021. Hypersorec: exploiting hyperbolic user and item representations with multiple aspects for social-aware recommendation. *ACM Transactions on Information Systems (TOIS)*, 40, 2, 1–28.
- [52] Hao Wang, Tong Xu, Qi Liu, Defu Lian, Enhong Chen, Dongfang Du, Han Wu, and Wen Su. 2019. Mcne: an end-to-end framework for learning multiple conditional network representations of social network. In *Proceedings of the 25th ACM SIGKDD international conference on knowledge discovery & data mining*, 1064–1072.
- [53] Hao Wang, Mingjia Yin, Luankang Zhang, Sirui Zhao, and Enhong Chen. 2025. Mf-gslae: a multi-factor user representation pre-training framework for dual-target cross-domain recommendation. *ACM Transactions on Information Systems*, 43, 2, 1–28.
- [54] Hao Wang et al. 2025. Generative large recommendation models: emerging trends in llms for recommendation. In *Companion Proceedings of the ACM on Web Conference 2025*, 49–52.
- [55] Kefan Wang, Hao Wang, Kenan Song, Wei Guo, Kai Cheng, Zhi Li, Yong Liu, Defu Lian, and Enhong Chen. 2025. A universal framework for compressing embeddings in ctr prediction. *arXiv preprint arXiv:2502.15355*.
- [56] Qingxian Wang, Renjian Zhang, Kangkang Ma, Bo Chen, Jiufang Chen, and Xiaoyu Shi. 2021. Siamese generative adversarial predicting network for extremely sparse data in recommendation system. In *2021 IEEE Intl Conf on Parallel & Distributed Processing with Applications, Big Data & Cloud Computing, Sustainable Computing & Communications, Social Computing & Networking (ISPA/BDCloud/SocialCom/SustainCom)*. IEEE, 1236–1241.
- [57] Likang Wu et al. 2024. A survey on large language models for recommendation. *World Wide Web*, 27, 5, 60.
- [58] Tengqing Wu. 2024. A diffusion data enhancement retentive model for sequential recommendation. In *2024 7th International Conference on Computer Information Science and Application Technology (CISAT)*. IEEE, 114–118.
- [59] Wenjia Xie, Hao Wang, Luankang Zhang, Rui Zhou, Defu Lian, and Enhong Chen. 2024. Breaking determinism: fuzzy modeling of sequential recommendation using discrete state space diffusion model. *Advances in Neural Information Processing Systems*, 37, 22720–22744.
- [60] Xu Xie, Fei Sun, Zhaoyang Liu, Shiwen Wu, Jinyang Gao, Jiandong Zhang, Bolin Ding, and Bin Cui. 2022. Contrastive learning for sequential recommendation. In *2022 IEEE 38th international conference on data engineering (ICDE)*. IEEE, 1259–1273.
- [61] Mingming Xu, Fangai Liu, and Weizhi Xu. 2019. A survey on sequential recommendation. In *2019 6th international conference on information science and control engineering (ICISCE)*. IEEE, 106–111.
- [62] Xiang Xu, Hao Wang, Wei Guo, Luankang Zhang, Wanshan Yang, Runlong Yu, Yong Liu, Defu Lian, and Enhong Chen. 2024. Multi-granularity interest retrieval and refinement network for long-term user behavior modeling in ctr prediction. *arXiv preprint arXiv:2411.15005*.
- [63] An Yan, Shuo Cheng, Wang-Cheng Kang, Mengting Wan, and Julian McAuley. 2019. Cosrec: 2d convolutional neural networks for sequential recommendation. In *Proceedings of the 28th ACM international conference on information and knowledge management*, 2173–2176.
- [64] An Yang et al. 2025. Qwen3 technical report. *arXiv preprint arXiv:2505.09388*.
- [65] Yufei Ye et al. 2025. Fuxi- α : scaling recommendation model with feature interaction enhanced transformer. In *Companion Proceedings of the ACM on Web Conference 2025*, 557–566.
- [66] Yufei Ye et al. 2025. Fuxi- β : towards a lightweight and fast large-scale generative recommendation model. *arXiv preprint arXiv:2508.10615*.
- [67] Mingjia Yin, Hao Wang, Wei Guo, Yong Liu, Zhi Li, Sirui Zhao, Zhen Wang, Defu Lian, and Enhong Chen. 2024. Learning partially aligned item representation for cross-domain sequential recommendation. *arXiv preprint arXiv:2405.12473*.
- [68] Mingjia Yin, Hao Wang, Wei Guo, Yong Liu, Suojuan Zhang, Sirui Zhao, Defu Lian, and Enhong Chen. 2024. Dataset regeneration for sequential recommendation. In *Proceedings of the 30th ACM SIGKDD Conference on Knowledge Discovery and Data Mining*, 3954–3965.
- [69] Mingjia Yin et al. 2023. Appl4sr: a generic framework with adaptive and personalized global collaborative information in sequential recommendation. In *Proceedings of the 32nd ACM international conference on information and knowledge management*, 3009–3019.
- [70] Mingjia Yin et al. 2024. Entropy law: the story behind data compression and llm performance. *arXiv preprint arXiv:2407.06645*.
- [71] Jiaqi Zhai et al. 2024. Actions speak louder than words: trillion-parameter sequential transducers for generative recommendations. *arXiv preprint arXiv:2402.17152*.
- [72] Buyun Zhang et al. 2024. Wukong: towards a scaling law for large-scale recommendation. *arXiv preprint arXiv:2403.02545*.
- [73] Jiaqing Zhang, Mingjia Yin, Hao Wang, Yawen Li, Yuyang Ye, Xingyu Lou, Junping Du, and Enhong Chen. 2025. Td3: tucker decomposition based dataset distillation method for sequential recommendation. *arXiv preprint arXiv:2502.02854*.
- [74] Luankang Zhang, Hao Wang, Suojuan Zhang, Mingjia Yin, Yongqiang Han, Jiaqing Zhang, Defu Lian, and Enhong Chen. 2024. A unified framework for adaptive representation enhancement and inversed learning in cross-domain recommendation. In *International Conference on Database Systems for Advanced Applications*. Springer, 115–130.
- [75] Luankang Zhang et al. 2025. Killing two birds with one stone: unifying retrieval and ranking with a single generative recommendation model. In *Proceedings of the 48th International ACM SIGIR Conference on Research and Development in Information Retrieval*, 2224–2234.
- [76] Xin Zhang et al. 2024. Gme: improving universal multimodal retrieval by multimodal llms. *arXiv preprint arXiv:2412.16855*.
- [77] Chung-Han Zhou and Yi-Ling Chen. 2023. Vcgan: variational collaborative generative adversarial network for recommendation systems. In *ICC 2023-IEEE International Conference on Communications*. IEEE, 6324–6330.
- [78] 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 *Proceedings of the 29th ACM international conference on information & knowledge management*, 1893–1902.
- [79] Barret Zoph. 2022. Designing effective sparse expert models. In *2022 IEEE International Parallel and Distributed Processing Symposium Workshops (IPDPSW)*. IEEE, 1044–1044.

Received 20 February 2007; revised 12 March 2009; accepted 5 June 2009