

Facet-Aware Multi-Head Mixture-of-Experts Model for Sequential Recommendation

Mingrui Liu
Nanyang Technological University
Singapore, Singapore
mingrui001@e.ntu.edu.sg

Sixiao Zhang
Nanyang Technological University
Singapore, Singapore
sixiao001@e.ntu.edu.sg

Cheng Long
Nanyang Technological University
Singapore, Singapore
c.long@ntu.edu.sg

ABSTRACT

Sequential recommendation (SR) systems excel at capturing users' dynamic preferences by leveraging their interaction histories. Most existing SR systems assign a single embedding vector to each item to represent its features, and various types of models are adopted to combine these item embeddings into a sequence representation vector to capture the user intent. However, we argue that this representation alone is insufficient to capture an item's multi-faceted nature (e.g., movie genres, starring actors). Besides, users often exhibit complex and varied preferences within these facets (e.g., liking both action and musical films in the facet of genre), which are challenging to fully represent. To address the issues above, we propose a novel structure called *Facet-Aware Multi-Head Mixture-of-Experts Model for Sequential Recommendation (FAME)*. We leverage sub-embeddings from each head in the last multi-head attention layer to predict the next item separately. This approach captures the potential multi-faceted nature of items without increasing model complexity. A gating mechanism integrates recommendations from each head and dynamically determines their importance. Furthermore, we introduce a Mixture-of-Experts (MoE) network in each attention head to disentangle various user preferences within each facet. Each expert within the MoE focuses on a specific preference. A learnable router network is adopted to compute the importance weight for each expert and aggregate them. We conduct extensive experiments on four public sequential recommendation datasets and the results demonstrate the effectiveness of our method over existing baseline models.

CCS CONCEPTS

• Information systems → Recommender Systems.

KEYWORDS

Recommender System; Sequential Recommendation

ACM Reference Format:

Mingrui Liu, Sixiao Zhang, and Cheng Long. 2025. Facet-Aware Multi-Head Mixture-of-Experts Model for Sequential Recommendation. In *Proceedings of Make sure to enter the correct conference title from your rights confirmation*

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, June 03–05, 2018, Woodstock, NY

© 2018 Copyright held by the owner/author(s). Publication rights licensed to ACM.
ACM ISBN 978-1-4503-XXXX-X/18/06...\$15.00
<https://doi.org/XXXXXXX.XXXXXXX>

emai (Conference acronym 'XX). ACM, New York, NY, USA, 9 pages. <https://doi.org/XXXXXXX.XXXXXXX>

1 INTRODUCTION

The explosion of information online presents users with a vast and ever-growing sea of items, from products [11] and apps [4] to videos [7, 40]. With limited time to explore everything, recommender systems (RS) have become crucial tools for helping users make efficient and satisfying choices. However, user interests are inherently dynamic, evolving over time and making it challenging for platforms to deliver consistently relevant recommendations [29]. To address these challenges, sequential recommendation (SR) has emerged as a powerful technique. This approach leverages the sequential nature of user interactions, typically captured as sessions containing a series of recent item interactions, to predict the user's next action [9, 30].

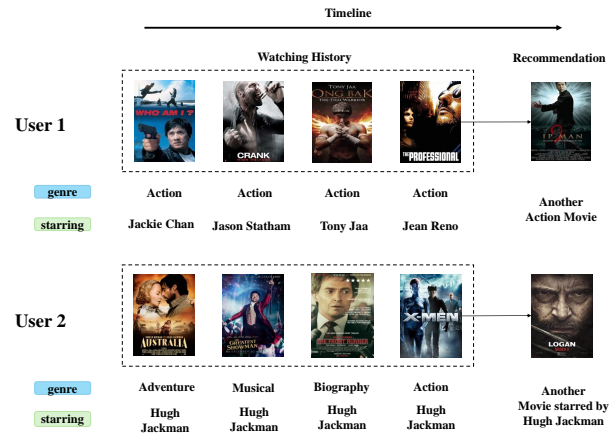


Figure 1: A motivation example.

Mainstream SR systems assign a single embedding vector to each item, capturing its features. Recurrent neural networks (RNNs) [14, 15], attention-based models [17, 24, 25, 39], graph-based models [2, 12, 32, 35], and others combine these item embeddings into a sequence representation vector to capture the user intent. This representation is used to predict the next item (e.g., selecting the item with the highest inner product with the sequence representation vector). However, a single embedding cannot well capture an item's multifaceted nature (e.g., movie genres and starring) [6, 38]. This is particularly problematic when different facets of an item can influence user intent. As illustrated in Figure 1, User 1's watch history suggests a strong preference for action movies. In this case,

recommending another action movie might be appropriate. Conversely, User 2's movie choices range across genres but all feature Hugh Jackman. In this case, recommending another movie starring Hugh Jackman might be more relevant. These examples highlight how user interests can be dominated by a single facet (genre or actor) within a category (movie). Furthermore, in more realistic scenarios, users can have multiple preferences within a single facet. For example, a user might enjoy both action and musical movies in the facet of genre. Recognizing and addressing these diverse preferences within a sequence is crucial for generating more effective recommendations that cater to specific user interests [6]. Failing to capture the dominant facet and the specific preferences within each facet can lead to suboptimal recommendations. This highlights the need for recommender systems that can effectively model the dynamic and multi-faceted nature of user interests.

Existing research addresses user intent complexity by using hierarchical windows [12, 37] to capture multi-level user intents from recent items, or by utilizing item representations from multiple items in the sequence instead of using the last item's representation only to recommend the next interacted item [6]. However, these methods still neglect the multi-faceted nature of items themselves.

To solve the aforementioned issues, we propose a novel structure called *Facet-Aware Multi-Head Mixture-of-Experts Model for Sequential Recommendation (FAME)*. We leverage sub-embeddings from each head in the last multi-head attention layer to predict the next item separately. This approach captures the potential multi-faceted nature of items (e.g., genres, starring) without increasing model complexity. A gating mechanism integrates recommendations from each head and dynamically determines their importance. Furthermore, we introduce a Mixture-of-Experts (MoE) network: this network replaces the query matrix in the self-attention layer, enabling the model to disentangle various user preferences within each facet. Each expert within the MoE focuses on a specific preference within each facet. A learnable router network is adopted to compute the importance weight for each expert and aggregate them. (e.g., whether action or musical movies are the stronger preference).

To summarize, our contributions in this paper are as follows:

- We propose a Multi-Head Prediction Mechanism to enhance the recommendation quality. This design facilitates capturing the potential multi-facet features of the items without increasing space requirements or the number of parameters.
- We propose a Mixture-of-Experts (MoE) network that improves user preference modeling by disentangling multiple preferences in each facet within a sequence. This module seamlessly integrates with existing attention-based models.
- Our model demonstrates significant effectiveness compared to various baseline categories (sequential, pre-trained, multi-intent) on four public datasets.

2 RELATED WORK

2.1 Sequential Recommendation

Recent advancements in neural networks and deep learning have spurred the development of various models to extract rich latent semantics from user behavior sequences and generate accurate recommendations. Convolutional Neural Networks (CNNs) [26],

Recurrent Neural Networks (RNNs) [14], Transformer-based models [17, 24, 25, 39], and Graph Neural Networks (GNNs) [2, 12, 32, 35] have been widely employed to enhance representation learning and recommendation performance. Self-supervised learning (SSL) has emerged as a promising technique for sequential recommendation [3, 20, 23, 33, 35, 38], with methods like CL4SRec [33] and ICLRec [3] employing data augmentation and contrastive learning to improve sequence representations and capture user intents. Additionally, research has focused on modeling multiple user intents, such as the hierarchical window approach in MSGIFSR [12] and Atten-Mixer [37], or the multi-item-based representation in MiasRec [6]. Furthermore, incorporating auxiliary information like item categories or attributes [1, 41] and textual descriptions [21] has been explored to enrich item representations. The integration of large language models (LLMs) is another emerging trend in the field [28, 31, 36]. However, these approaches are beyond the scope of this paper.

2.2 Sparse Mixtures of Experts (SMoE)

The Mixture-of-Experts (MoE) architecture has emerged as a powerful tool for handling complex tasks by distributing computations across multiple specialized models, or experts. While MoE models can significantly enhance model capacity, their computational overhead due to routing data to all experts can be prohibitive. To address this, Sparse Mixture of Experts (SMoE) was introduced, enabling each data point to be processed by a carefully selected subset of experts [5, 8, 10]. This approach offers the potential for substantial computational savings without compromising performance.

While SMoE has shown promise in various domains, its application in sequential recommendation remains relatively under-explored. Leveraging SMoE in this context could unlock new opportunities to enhance recommendation quality by effectively capturing and modeling diverse user preferences within a sequence.

3 PRELIMINARIES

3.1 Notations and Problem Statement

Let \mathcal{U} and \mathcal{V} represent the user set and item set, where $u \in \mathcal{U}$ (resp. $v \in \mathcal{V}$) denotes an individual user (resp. item). Consequently, $|\mathcal{U}|$ and $|\mathcal{V}|$ denote the sizes of user set and item set, respectively. For each user u , we define their interaction sequence $\mathcal{S}_u = \{v_1^{(u)}, \dots, v_i^{(u)}, \dots, v_t^{(u)}\}$ as a chronologically ordered list of items. Here, $v_i^{(u)} \in \mathcal{V}$ represents the item that user u interacted with at time step i , and t denotes the length of the interaction sequence for user u . Given a user interaction sequence \mathcal{S}_u , the goal of sequential recommendation is to predict the item that user u will interact with at the next time step, $t + 1$. Formally, we can define the problem as:

$$v_u^{(*)} = \arg \max_{v_i \in \mathcal{V}} P(v_{t+1}^{(u)} = v_i | \mathcal{S}_u) \quad (1)$$

3.2 Multi-Head Self-Attention

- (1) **Item Embeddings:** The model first obtains embeddings for each item in the sequence (denoted as $x \in \mathbb{R}^d$).
- (2) **Query, Key, Value Vectors:** Each item embedding (x) is then projected into three vectors:

- **Query Vector (q):** Represents what the model is currently looking for in the sequence.
- **Key Vector (k):** Captures the content of the current item.
- **Value Vector (v):** Contains the actual information associated with the item.

These projections are calculated using three trainable weight matrices (denoted by $W_Q, W_K \in \mathbb{R}^{d \times d_k}, W_V \in \mathbb{R}^{d \times d_v}$):

$$q = x^T \cdot W_Q, \quad k = x^T \cdot W_K, \quad v = x^T \cdot W_V, \quad (2)$$

where d_k is the dimension of query and key vector, and d_v is the dimension of value vector.

- (3) **Attention Scores:** The model calculates an attention score (α_{ij}) for each pair of items (i, j) in the sequence. This score reflects the similarity between the current item's query vector (q_i) and the key vector (k_j) of each other item. A normalization term (\sqrt{d}) is used to account for the vector dimension. The attention scores are then normalized using a softmax function (denoted by $\tilde{\alpha}_{ij}$) to create a probability distribution across all items, indicating the relative importance of each item to the current one.

$$\alpha_{ij} = \frac{q_i^T \cdot k_j}{\sqrt{d}}, \quad \tilde{\alpha}_{ij} = \frac{\exp(\alpha_{ij})}{\sum_{j=1}^t \exp(\alpha_{ij})} \quad (3)$$

- (4) **Item representation:** The item representation f_i is calculated based on weighted sum of value vectors in the sequence. The weights for this summation are derived from the previously calculated attention scores ($\tilde{\alpha}_{ij}$):

$$f_i = \sum_{j=1}^t \tilde{\alpha}_{ij} \cdot v_j \quad (4)$$

4 METHODS

4.1 Overview

This section introduces our proposed framework with a high-level overview, which is displayed in Figure 2. The framework incorporates two key components: the **Facet-Aware Multi-Head Prediction Mechanism** (detailed in Section. 4.2), which learns to represent each item with multiple sub-embedding vectors, each capturing a specific facet of the item; and the **Mixture-of-Experts Self-Attention Layer** (detailed in Section. 4.3), which employs a Mixture-of-Experts (MoE) network within each subspace to capture the users' specific preferences within each facet. Our framework can be seamlessly integrated to any attention-based recommendation model. In this paper, we incorporate our framework to SASRec for illustration.

4.2 Facet-Aware Multi-Head Prediction Mechanism

4.2.1 Original SASRec Prediction Process. In the original SASRec model, the final prediction for the next item is based on the last item's representation (f_t , calculated by Equation 4, which can also be regarded as the sequence representation) obtained from the last self-attention layer. This representation is processed through a feed-forward network (FFN) with ReLU activation for non-linearity,

followed by layer normalization, dropout, and a residual connection:

$$\begin{aligned} \text{FFN}(f_t) &= \text{RELU}(f_t^T \cdot W_1 + b_1)^T \cdot W_2 + b_2, \\ F_t &= \text{LayerNorm}(f_t + \text{Dropout}(\text{FFN}(f_t))), \end{aligned} \quad (5)$$

Here, $W_1, W_2 \in \mathbb{R}^{d \times d}, b_1, b_2 \in \mathbb{R}^d$ are all learnable parameters. The final user preference score for item v at step ($t + 1$) is then calculated as the dot product between the item embedding (x_v) and the sequence representation (F_t):

$$P(v_{t+1} = v | \mathcal{S}_u) = x_v^T \cdot F_t, \quad (6)$$

Top- k items with the highest preference scores are recommended to the user.

4.2.2 Motivation for Our Approach. The multi-head self-attention mechanism splits the sequence representation and item embeddings into multiple subspaces (heads). Research suggests that these heads can allocate different attention distributions so as to perform different tasks [27]. We hypothesize that these heads could also capture different facets of items (e.g., genre and starring actors in the context of movie recommendation). This ability to capture multi-faceted information has the potential to improve recommendation quality.

4.2.3 Proposed Multi-Head Recommendation. Instead of performing a single attention function with d -dimensional keys, values and queries, it is found beneficial to linearly project the queries, keys and values H times with different, learned linear projections to d_k, d_k and d_v dimensions, respectively [27]. Here, H is the number of heads, and d_k, d_k and d_v are typically set to $d' = \frac{d}{H}$.

Leveraging the multi-head attention mechanism, we propose a novel approach where each head independently generates recommendations. The final item embedding from head h is denoted as $f_t^{(h)} \in \mathbb{R}^{d'}$. We then process this embedding similarly as we do for the original model:

$$\begin{aligned} \text{FFN}'(f_t^{(h)}) &= \text{RELU}(f_t^{(h)T} \cdot W_1' + b_1')^T \cdot W_2' + b_2', \\ F_t^{(h)} &= \text{LayerNorm}(f_t^{(h)} + \text{Dropout}(\text{FFN}'(f_t^{(h)}))), \end{aligned} \quad (7)$$

Unlike the original FFN (Equation 5), the feed-forward network applied to each head (FFN') operates on a reduced dimension of d' . The learnable parameters for FFN' are therefore adjusted accordingly: $W_1', W_2' \in \mathbb{R}^{d' \times d'}, b_1', b_2' \in \mathbb{R}^{d'}$. This adaptation aligns with the dimensionality of sub-embeddings within each head. To enhance parameter efficiency and improve performance, we adopt a shared feed-forward network (FFN') across all attention heads. Each head generates the preference score for each item independently, i.e.,

$$P^{(h)}(v_{t+1} = v | \mathcal{S}_u) = x_v^{(h)T} \cdot F_t^{(h)}, \quad (8)$$

where $x_v^{(h)} \in \mathbb{R}^{d'}$ is the sub-embedding of the item v , reflecting the features of the specific facet corresponding to the attention head h . Specifically, it is calculated by a linear transformation from its original embedding:

$$x_v^{(h)} = x_v^T \cdot W_f^{(h)}, \quad (9)$$

with $W_f^{(h)} \in \mathbb{R}^{d \times d'}$ being a learnable matrix.

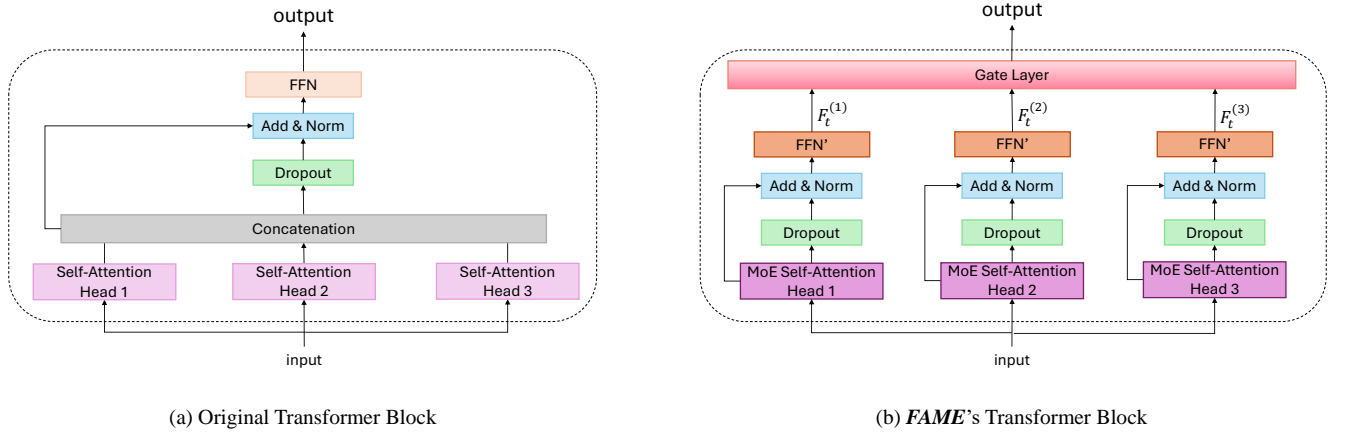


Figure 2: Overview of the proposed model: (a) illustrates the original Transformer block, while (b) depicts the architecture of our proposed *FAME* model. For simplicity, the LayerNorm and Dropout operations following the FFN (FFN') are omitted from the Figure

In order to integrate the recommendation results from each head, we employ a gate mechanism to determine the relative importance of each head's recommendations:

$$g = \left[F_t^{(1)} \parallel \dots \parallel F_t^{(H)} \right]^T \cdot W_g + b_g, \quad (10)$$

$$\tilde{g} = \text{softmax}(g)$$

Here, $[\cdot \parallel \cdot]$ denotes the concatenation operation. Each element $\tilde{g}^{(h)} \in [0, 1]$ within the vector \tilde{g} represents the importance of head h in determining the user's dominant interest or preference. For instance, a higher $\tilde{g}^{(h)}$ for a genre-focused head indicates a stronger preference for specific movie genres, while a higher value for an actor-focused head suggests a preference for movies starring particular actors. The gate mechanism, parameterized by $W_g \in \mathbb{R}^{d \times H}$ and $b_g \in \mathbb{R}^H$, learns to assign appropriate weights to each head based on the user's current context. Finally, we compute a unified preference score for each item by weighting the recommendations from each head:

$$P(v_{t+1} = v | \mathcal{S}_u) = \sum_{i=1}^H \tilde{g}^{(i)} \cdot P^{(i)}(v_{t+1} = v | \mathcal{S}_u) \quad (11)$$

This approach allows the model to exploit the strengths of each head while assigning appropriate weights based on their importance in the specific context.

4.3 Mixture-of-Experts Self-Attention Layer

While the Facet-Aware Multi-Head mechanism effectively captures item facets, users often exhibit more granular and diverse preferences within these facets. To address this, we introduce the Mixture-of-Experts Self-Attention Layer (MoE-Attention), as illustrated in Figure 3.

We assume that each facet can be decomposed into N distinct preferences. For instance, a genre facet might include preferences for action, comedy, musicals, etc. To capture the nuanced preferences within each facet of a sequence, we replace the standard query

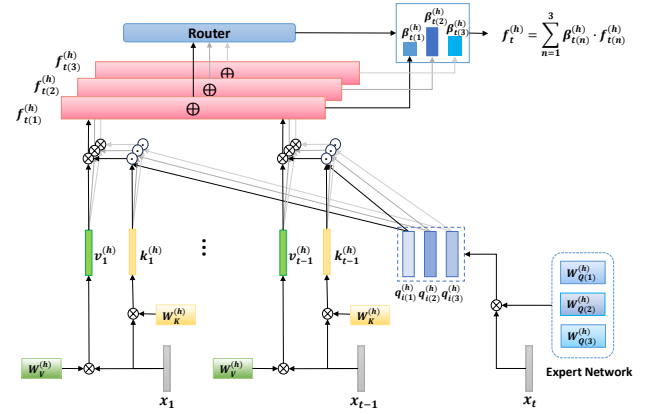


Figure 3: MoE Self-Attention Network: Integrated Item Representation Calculation. This diagram visualizes the computational process for determining the integrated item representation of the final item ($f_t^{(h)}$) within a specific head (h) of our proposed model.

generation mechanism in self-attention (Equation 2) with a Mixture-of-Experts (MoE) network in each head. This network consists of N experts, each represented by a trainable matrix $W_{Q(n)}^{(h)} \in \mathbb{R}^{d \times d'}$ (where $n \in [1, N]$). Each expert within a head is designed to capture one of these preferences by transforming an item embedding x_i (i.e., the embedding of the i^{th} item in the sequence) into an expert query vector $q_{i(n)}^{(h)} \in \mathbb{R}^{d'}$ as follows:

$$q_{i(n)}^{(h)} = x_i^T \cdot W_{Q(n)}^{(h)} \quad (12)$$

The key vector ($k_j^{(h)}$) and value vector ($v_j^{(h)}$) of the j^{th} sequence item in head h are computed using the same linear transformations

as in the original SASRec model:

$$k_j^{(h)} = x_j^T \cdot W_K^{(h)}, \quad v_j^{(h)} = x_j^T \cdot W_V^{(h)} \quad (13)$$

Then the attention score for the i^{th} item relative to the j^{th} item in head h by expert n is computed as:

$$\alpha_{ij(n)}^{(h)} = \frac{q_{i(n)}^{(h)T} \cdot k_j^{(h)T}}{\sqrt{d'}}, \quad (14)$$

$$\tilde{\alpha}_{ij(n)}^{(h)} = \text{softmax}(\alpha_{i1(n)}^{(h)}, \dots, \alpha_{it(n)}^{(h)})$$

The item representation of the i^{th} item for head h and expert n ($f_{i(n)}^{(h)}$) is then calculated as a weighted sum of value vectors, where the weights are the corresponding attention scores:

$$f_{i(n)}^{(h)} = \sum_{j=1}^t \tilde{\alpha}_{ij(n)}^{(h)} \cdot v_j^{(h)} \quad (15)$$

As illustrated in Figure 4, consider a genre-focused head with two experts: one for action movies and another for musical movies. As detailed in Section 4.2.1, the standard SASRec model treats the representation of the final item in a sequence as the overall sequence representation. To illustrate our MoE attention mechanism, we focus on the attention scores associated with the 4^{th} (and final) item in the sequence. The first expert's query vector of the 4^{th} item ($q_{4(1)}^{(h)}$) would assign higher attention scores to action movies (items 1 and 3), while the second expert's query vector ($q_{4(2)}^{(h)}$) would focus on musical movies (items 2 and 4). Consequently, the final item's representation ($f_{4(1)}^{(h)}$) generated by the first expert would lean towards recommending action movies, whereas the representation ($f_{4(2)}^{(h)}$) from the second expert would favor musical movies.

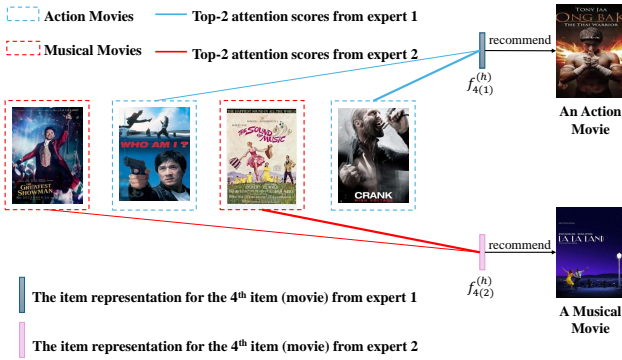


Figure 4: An example on attention scores distribution and recommendation results among different experts on genre-focused head

To dynamically determine the importance of each preference within a facet (e.g., whether action or musical is the preferred genre), we introduce a router network parameterized by $W_{exp}^{(h)} \in \mathbb{R}^{(n \cdot d') \times n}$. This network assigns an importance score $\beta_{i(n)}^{(h)} \in (0, 1)$ to each

Table 1: Dataset statistic

Dataset	#users	#items	#actions	avg.length	density
Beauty	22,363	12,101	198,502	8.8	0.07%
Sports	25,598	18,357	296,337	8.3	0.05%
Toys	19,412	19,392	167,597	8.6	0.04%
ML-20m	96,726	16,297	185,6746	19.2	0.11%

item representation generated by each expert $f_{i(n)}^{(h)}$. The importance scores are computed as follows:

$$\beta_i^{(h)} = \text{softmax}([f_{i(1)}^{(h)} | \dots | f_{i(n)}^{(h)}]^T \cdot W_{exp}^{(h)}) \quad (16)$$

The integrated item representation $f_i^{(h)}$ for the i^{th} item in head h is then computed as a weighted sum of the expert query vectors:

$$f_i^{(h)} = \sum_{n=1}^N \beta_{i(n)}^{(h)} \cdot f_{i(n)}^{(h)} \quad (17)$$

The integrated item representation ($f_i^{(h)}$) represents the overall preference at the i^{th} timestamp within head h . For instance, for the case in Figure 4, a higher weight for $f_{4(1)}^{(h)}$ (resp. $f_{4(2)}^{(h)}$) would push the model towards recommending action (resp. musical) movies.

4.4 Deployment and Training

4.4.1 Model Deployment. Our FAME model is built upon the SASRec (or any attention-based) framework, with the final Transformer layer replaced by our proposed architecture.

4.4.2 Training Pipeline. We initiate our model by pre-training an attention-based sequential recommendation model (e.g., SASRec). Subsequently, we replace Transformer block's query matrix at the final layer with our proposed MoE network (Section 4.3) while retaining the original key and value matrices. The newly introduced components, including the head-specific FFN' (Equation 7), gate mechanism (Equation 10), and router (Equation 16), are randomly initialized. The entire model is then fine-tuned end-to-end.

4.4.3 Training Objectives. We empirically found that using a global cross-entropy loss function leads to better performance in sequential recommendation. Therefore, we adopt the following loss function for training:

$$\mathcal{L}_{ce} = - \sum_{u \in \mathcal{U}} \log \left(\frac{\exp(x_{t+1}^{(u)T} \cdot f_t^{(u)})}{\sum_i \exp(x_i^T \cdot f_t^{(u)})} \right) \quad (18)$$

5 EXPERIMENTS

5.1 Datasets

We conduct experiments on four public datasets. *Beauty*, *Sports* and *Toys* are three subcategories of Amazon review data introduced in [19]. *ML-20m* is a subset of the MovieLens dataset [13], containing approximately 20 million ratings from 138,493 users on 27,278 movies. Following [33, 41], only "5-core" sequences are remained in the 4 datasets, in which all users and items have at least 5 interactions. The statistics of the prepared datasets are summarized in Table 1.

Table 2: Performance comparison of different methods on top- k recommendation

Dataset	Metric	GRU4Rec	SASRec	BERT4Rec	CORE	CL4SRec	ICLRec	DuoRec	A-Mixer	MSGIFSR	MiasRec	FAME	Improv.
Beauty	HR@5	0.0408	0.0508	0.0510	0.0331	0.0623	<u>0.0664</u>	0.0504	0.0507	0.0518	0.0524	0.0710	6.9%
	HR@10	0.0623	0.0761	0.0745	0.0664	0.0877	<u>0.0918</u>	0.0691	0.0752	0.0771	0.0795	0.0978	6.2%
	HR@20	0.0895	0.1057	0.1075	0.1071	0.1195	<u>0.1252</u>	0.0912	0.1033	0.1105	0.1125	0.1345	7.4%
	NDCG@5	0.0273	0.0318	0.0343	0.0164	0.0440	<u>0.0480</u>	0.0363	0.0350	0.0344	0.0362	0.0508	5.8%
	NDCG@10	0.0342	0.0400	0.0419	0.0271	0.0521	<u>0.0562</u>	0.0424	0.0421	0.0429	0.0449	0.0593	5.5%
	NDCG@20	0.0410	0.0474	0.0502	0.0373	0.0601	<u>0.0646</u>	0.0479	0.0504	0.0508	0.0532	0.0687	6.3%
Sports	HR@5	0.0210	0.0266	0.0252	0.0150	0.0338	<u>0.0384</u>	0.0225	0.0217	0.0268	0.0270	0.0400	4.2%
	HR@10	0.0339	0.0412	0.0395	0.0342	0.0498	<u>0.0543</u>	0.0327	0.0321	0.0425	0.0435	0.0580	6.8%
	HR@20	0.0527	0.0618	0.0607	0.0609	0.0723	<u>0.0753</u>	0.0476	0.0469	0.0634	0.0651	0.0820	8.9%
	NDCG@5	0.0136	0.0158	0.0166	0.0072	0.0235	<u>0.0266</u>	0.0161	0.0165	0.0171	0.0180	0.0277	4.1%
	NDCG@10	0.0178	0.0205	0.0212	0.0134	0.0287	<u>0.0317</u>	0.0193	0.0188	0.0221	0.0233	0.0337	6.3%
	NDCG@20	0.0225	0.0256	0.0265	0.0201	0.0344	<u>0.0370</u>	0.0231	0.0235	0.0279	0.0288	0.0402	8.6%
Toys	HR@5	0.0369	0.0489	0.0464	0.0338	0.0658	<u>0.0792</u>	0.0481	0.0565	0.0576	0.0581	0.0820	3.5%
	HR@10	0.0524	0.0676	0.0677	0.0699	0.0912	<u>0.1043</u>	0.0666	0.0819	0.0831	0.0828	0.1065	2.1%
	HR@20	0.076	0.0908	0.0968	0.1114	0.1209	<u>0.1382</u>	0.0879	0.1099	0.1150	0.1143	0.1409	1.9%
	NDCG@5	0.0247	0.0329	0.0322	0.0158	0.047	<u>0.0579</u>	0.0356	0.403	0.0407	0.0408	0.0603	4.1%
	NDCG@10	0.0296	0.0389	0.0391	0.0274	0.0552	<u>0.0660</u>	0.0415	0.481	0.0492	0.0488	0.0681	3.2%
	NDCG@20	0.0356	0.0448	0.0464	0.0378	0.0627	<u>0.0745</u>	0.0469	0.574	0.0577	0.0567	0.0759	1.9%
ML-20m	HR@5	0.1365	0.1305	0.1446	0.0655	0.1205	0.1380	<u>0.1458</u>	0.1325	0.1303	0.1367	0.1538	5.5%
	HR@10	0.2052	0.2016	0.2172	0.1312	0.1853	0.2070	0.2164	0.2022	0.2013	0.2071	0.2246	3.4%
	HR@20	0.2981	0.2996	<u>0.3132</u>	0.2251	0.2760	0.2997	0.3108	0.2994	0.2978	0.3021	0.3230	3.1%
	NDCG@5	0.0927	0.0858	0.0964	0.0347	0.0804	0.0927	<u>0.0986</u>	0.0899	0.0844	0.0918	0.1046	6.1%
	NDCG@10	0.1148	0.1086	0.1197	0.0558	0.1012	0.1149	<u>0.1212</u>	0.1121	0.1119	0.1144	0.1276	5.3%
	NDCG@20	0.1382	0.1333	0.1438	0.0794	0.1240	0.1382	<u>0.1450</u>	0.1334	0.1357	0.1383	0.1513	4.3%

5.2 Evaluation Metrics

We rank the prediction on the whole item set without negative sampling [18]. Performance is evaluated on a variety of evaluation metrics, including Hit Ratio@ k (HR@ k), and Normalized Discounted Cumulative Gain@ k (NDCG@ k) where $k \in \{5, 10, 20\}$. Following standard practice in sequential recommendation [17, 22, 25, 41], we employ a *leave-one-out* evaluation strategy: for each user sequence, the final item serves as the test data, the penultimate item as the validation data, and the remaining items as the training data.

5.3 Baselines

We compare our proposed method against a set of baseline models as follows:

- **GRU4Rec** [14]: it employs a GRU to encode sequences and incorporates a ranking-based loss.
- **SASRec** [17]: this method is a pioneering work utilizing self-attention to capture dynamic user interests.
- **BERT4Rec** [25]: this approach adapts the BERT architecture for sequential recommendation using a cloze task.
- **CORE** [16]: it proposes a representation-consistent encoder based on linear combinations of item embeddings to ensure that sequence representations are in the same space with item embeddings.
- **CL4SRec** [33]: this method combines contrastive learning with a Transformer-based model through data augmentation techniques (i.e., item crop, mask, and reorder).
- **ICLRec** [3]: this approach improves sequential recommendation by conducting clustering and contrastive learning on

user intentions represented by cluster centroids to enhance recommendation.

- **DuoRec** [20]: this research investigates the representation degeneration issue in sequential recommendation and offers solutions based on contrastive learning techniques.
- **MSGIFSR** [12]: it captures multi-level user intents using a Multi-granularity Intent Heterogeneous Session Graph.
- **Atten-Mixer** [37]: this method leverages concept-view and instance-view readouts for multi-level intent reasoning instead of using the GNN propagation.
- **MiasRec** [6]: this approach utilizes multiple item representations in the sequence instead of only using the last item’s representation as the sequence representation to capture diverse user intents.

5.4 Settings and Implementation Details

We employ original implementations for SASRec, ICLRec, MSGIFSR, Atten-Mixer, and MiasRec public in their papers. For GRU4Rec and CORE, we leverage the RecBole library¹ [34], while BERT4Rec, CL4SRec, and DuoRec are implemented using the SSLRec library² [23]. Hyperparameters for all models are set according to their respective papers. We experiment with embedding dimensions of 64 and 128 (as experimented, larger dimensions often lead to convergence issues) and select the configuration that yields the best performance for each model.

Our method is implemented in PyTorch. The model is optimized by Adam optimizer with a learning rate of 0.001, $\beta_1 = 0.9$, $\beta_2 =$

¹<https://github.com/RUCAIBox/RecBole>

²<https://github.com/HKUDS/SSLRec>

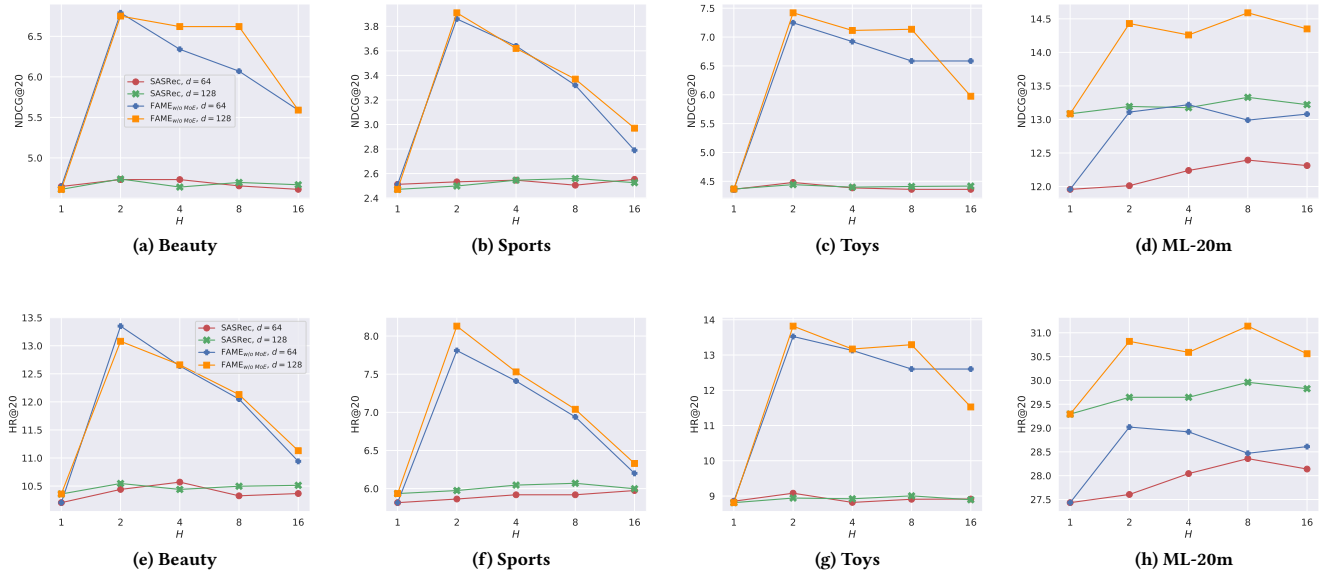


Figure 5: The performances comparison varying the number of heads in each dataset. The metric in (a)-(d) is NDCG@20, and the metric in (e)-(h) is HR@20.

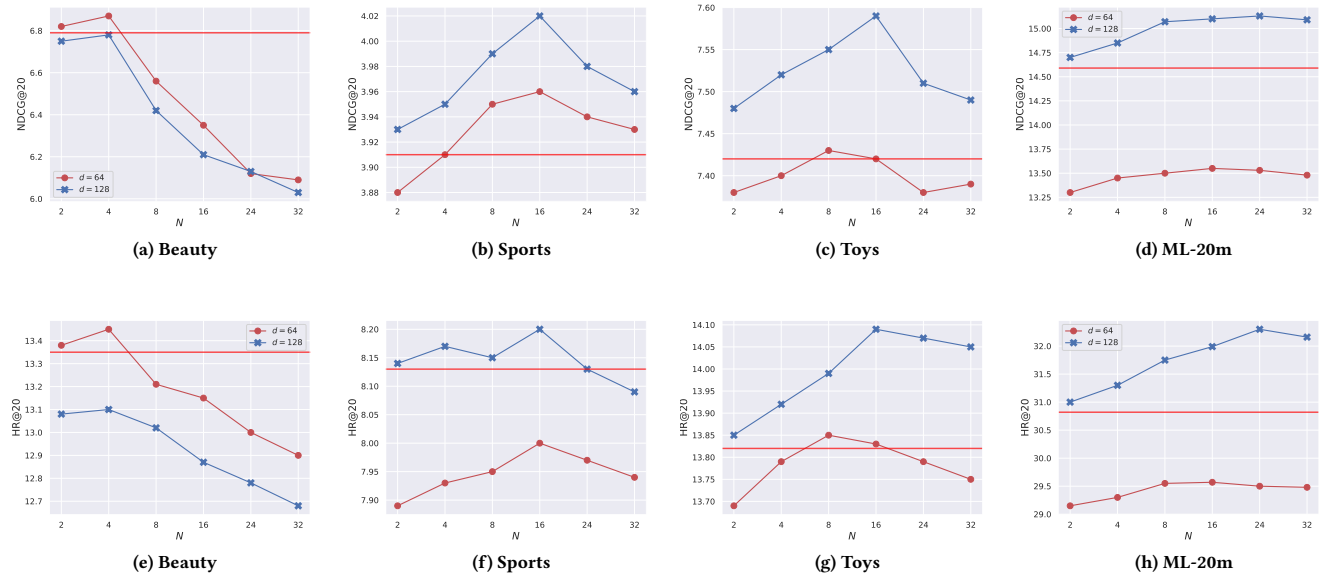


Figure 6: The performances comparison varying the number of experts in each dataset. The metric in (a)-(d) is NDCG@20, and the metric in (e)-(h) is HR@20. The red horizontal line in each subfigure indicates the peak performance (NDCG@20 or HR@20) achieved by FAME_{w/o MoE} within that dataset, as shown in Figure 5.

0.999. We employ a batch size of 256. For FAME hyperparameters, H and N are tuned within the ranges $\{1, 2, 4, 8, 16\}$ and $\{2, 4, 8, 16, 32\}$, respectively. All experiments were conducted on a single NVIDIA RTX A5000 GPU.

5.5 Overall Performance

Table 2 presents a comprehensive comparison of FAME against various baseline models. Our experimental findings reveal several key observations:

User 230

History Sequence: [1975,2018,1447,2019,2020,2021,518,2022,2023,1411,2024,2025,2026,509,2027,2028,523]

Item Representation	Top-20 Recommendation Results (IDs)					Gate Value	Integrated Results (IDs)				
Head 1	Results from Head 1					$g^{(1)}$	FAME				
	997	4608	3141	1390	8049		523	649	2356	490	403
recommend	3601	5832	8106	14486	3142	0.18	997	2028	4608	752	5173
$P_{17}^{(1)}$	523	4866	4049	6430	1783		386	371	825	1328	1278
	1782	404	5029	3225	996		996	4345	57	4936	1789
Head 2	Results from Head 2					$g^{(2)}$	Simple Concatenation (equal gate value)				
	523	2028	649	2362	57		523	997	4608	2356	490
recommend	1335	403	2356	551	1328	0.82	649	5173	752	386	371
$P_{17}^{(2)}$	490	2039	3558	825	371		386	996	825	1278	1789
	386	752	4345	4963	5173		3142	1390	3601	5214	404

Figure 7: Recommendation results for user 230 in the Sports dataset. User history is displayed at the top. The ground truth next item (item 2028) is highlighted.

- **Limitations of Traditional Models:** While RNN and Transformer-based models have shown success in sequential tasks, their direct application to recommendation often yields suboptimal results due to a lack of consideration for real-world user and item complexities (e.g., GRU4Rec, BERT4Rec).
- **Importance of Intent Modeling:** Models that explicitly capture user intents significantly outperform traditional sequential models. This improvement is attributed to their ability to: 1) handle noisy user behavior by focusing on underlying preferences rather than superficial interactions (e.g., ICLRec), or 2) disentangle multiple co-existing user intents within a sequence (e.g., MiasRec).
- **Superiority of our model FAME:** Our proposed FAME model consistently outperforms all baselines across datasets. This highlights the importance of considering item multi-faceted nature and disentangling user preferences within each facet for effective sequential recommendation.

5.6 Ablation Study and Parameters Study

This subsection presents the ablation study to evaluate the contributions of our proposed components and conduct corresponding hyperparameter tuning. We begin by examining $FAME_{w/o MoE}$, which excludes the MoE module, to assess the impact of the facet-aware multi-head prediction mechanism (introduced in Section. 4.2) and determine the optimal number of attention heads in Section. 5.6.1. Subsequently, using the optimized head configuration, we evaluate the complete FAME model to validate the effectiveness of the MoE module and identify the optimal number of experts in Section. 5.6.2.

5.6.1 Impact of the number of heads. Figure 5 illustrates the performance variation with different numbers of heads (H), treated as a hyperparameter. We compare the original SASRec model with $FAME_{w/o MoE}$ to isolate the impact of our multi-head prediction mechanism. We experiment with H values of $\{1, 2, 4, 8, 16\}$. When $H = 1$, our $FAME_{w/o MoE}$ is reduced to the original SASRec model with single head. As noted in [27], computational costs remain constant when varying the number of heads (H) while maintaining a fixed embedding dimension (d).

Benefits of multi-head attention: Both SASRec and $FAME_{w/o MoE}$

exhibit performance improvements with multiple heads, however, excessive heads can lead to diminishing returns, aligning with findings in Transformer [27] and SASRec [17].

Superiority of facet-aware architecture: $FAME_{w/o MoE}$ consistently outperforms SASRec, demonstrating the effectiveness of our facet-aware approach.

Dataset-specific optimal head count: The optimal number of heads varies across datasets. Beauty, Sports, and Toys benefit from fewer heads, suggesting simpler item facets, while ML-20m requires more heads to capture complex item characteristics.

5.6.2 Impact of the number of experts. Figure 6 illustrates the influence of the number of experts (N) within each attention head on model performance. We set H to the optimal value determined for $FAME_{w/o MoE}$ and compare its performance (red horizontal line in each subfigure) to that of FAME by varying N in $\{2, 4, 8, 16, 32\}$. FAME simplifies to $FAME_{w/o MoE}$ when N is set to 1.

FAME outperforms $FAME_{w/o MoE}$ across all datasets. This improvement can be attributed to the effectiveness of the MoE component, as evidenced by the existence of an optimal N value in each subfigure that surpasses the performance of $FAME_{w/o MoE}$. While the Beauty dataset exhibits diminishing returns for N greater than 4, suggesting simpler user preferences, the other datasets benefit from a larger number of experts. In particular, ML-20m show performance gains with increasing N , indicating the presence of more complex and diverse user preferences. However, excessive experts ($N = 32$) might lead to overfitting in the Sports and Toys dataset.

5.7 Case Study

To illustrate the effectiveness of our facet-aware mechanism, Figure 7 presents recommendation results for user 230, along with corresponding head importance scores (calculated using Equation 10). For simplicity, we set the number of heads to two and focus on the Sports dataset.

The figure clearly demonstrates the diversity of recommendations across different heads, highlighting the ability of our model to capture distinct item facets. The calculated head importance scores reveal that head 2 better aligns with user 230’s preferences (0.82 vs 0.18), as evidenced by the inclusion of the ground truth item (item 2028) in its recommendation list. The integrated recommendation, incorporating both heads with appropriate weights, successfully predicts the ground truth item.

In contrast, a traditional approach concatenating sub-embeddings from all heads without considering head importance fails to capture the user’s dominant preference, resulting in the omission of the ground truth item in the recommendation list.

6 CONCLUSION

In this paper, we propose a Facet-Aware Multi-Head Mixture-of-Experts Model for Sequential Recommendation (FAME), leveraging sub-embeddings from each head in the last multi-head attention layer to predict the next item separately. This approach captures the potential multi-faceted nature of items without increasing model complexity. A Mixture-of-Experts (MoE) network is adopted in each attention head to disentangle various user preferences within each facet. Each expert within the MoE focuses on a specific preference, and the importance score is calculated by a router network, which

is used to aggregate the overall preference. Extensive experiments demonstrate the effectiveness of our method over existing baseline models.

REFERENCES

- [1] Renqin Cai, Jibang Wu, Aidan San, Chong Wang, and Hongning Wang. 2021. Category-aware collaborative sequential recommendation. In *Proceedings of the 44th international ACM SIGIR conference on research and development in information retrieval*. 388–397.
- [2] Jianxin Chang, Chen Gao, Yu Zheng, Yiqun Hui, Yanan Niu, Yang Song, Depeng Jin, and Yong Li. 2021. Sequential recommendation with graph neural networks. In *Proceedings of the 44th international ACM SIGIR conference on research and development in information retrieval*. 378–387.
- [3] 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.
- [4] Heng-Tze Cheng, Levent Koc, Jeremiah Harmsen, Tal Shaked, Tushar Chandra, Hrishikesh Aradhya, Glen Anderson, Greg Corrado, Wei Chai, Mustafa Ispir, et al. 2016. Wide & deep learning for recommender systems. In *Proceedings of the 1st workshop on deep learning for recommender systems*. 7–10.
- [5] Zewen Chi, Li Dong, Shaohan Huang, Damai Dai, Shuming Ma, Barun Patra, Saksham Singhal, Payal Bajaj, Xia Song, Xian-Ling Mao, et al. 2022. On the representation collapse of sparse mixture of experts. *Advances in Neural Information Processing Systems* 35 (2022), 34600–34613.
- [6] Minjin Choi, Hye-young Kim, Hyunsouk Cho, and Jongwuk Lee. 2024. Multi-intent-aware Session-based Recommendation. *arXiv preprint arXiv:2405.00986* (2024).
- [7] Paul Covington, Jay Adams, and Emre Sargin. 2016. Deep neural networks for youtube recommendations. In *Proceedings of the 10th ACM conference on recommender systems*. 191–198.
- [8] Nan Du, Yanping Huang, Andrew M Dai, Simon Tong, Dmitry Lepikhin, Yuanzhong Xu, Maxim Krikun, Yanqi Zhou, Adams Wei Yu, Orhan Firat, et al. 2022. Glam: Efficient scaling of language models with mixture-of-experts. In *International Conference on Machine Learning*. PMLR, 5547–5569.
- [9] 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 (2020), 1–42.
- [10] William Fedus, Barret Zoph, and Noam Shazeer. 2022. Switch transformers: Scaling to trillion parameter models with simple and efficient sparsity. *Journal of Machine Learning Research* 23, 120 (2022), 1–39.
- [11] 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*. 299–315.
- [12] Jiayan Guo, Yaming Yang, Xiangchen Song, Yuan Zhang, Yujing Wang, Jing Bai, and Yan Zhang. 2022. Learning multi-granularity consecutive user intent unit for session-based recommendation. In *Proceedings of the fifteenth ACM International conference on web search and data mining*. 343–352.
- [13] F Maxwell Harper and Joseph A Konstan. 2015. The movielens datasets: History and context. *ACM transactions on interactive intelligent systems (tiis)* 5, 4 (2015), 1–19.
- [14] Balázs Hidasi, Alexandros Karatzoglou, Linas Baltrunas, and Domonkos Tikk. 2015. Session-based recommendations with recurrent neural networks. *arXiv preprint arXiv:1511.06939* (2015).
- [15] Balázs Hidasi, Alexandros Karatzoglou, Linas Baltrunas, and Domonkos Tikk. 2015. Session-based recommendations with recurrent neural networks. *arXiv preprint arXiv:1511.06939* (2015).
- [16] Yupeng Hou, Binbin Hu, Zhiqiang Zhang, and Wayne Xin Zhao. 2022. Core: simple and effective session-based recommendation within consistent representation space. In *Proceedings of the 45th international ACM SIGIR conference on research and development in information retrieval*. 1796–1801.
- [17] Wang-Cheng Kang and Julian McAuley. 2018. Self-attentive sequential recommendation. In *2018 IEEE international conference on data mining (ICDM)*. IEEE, 197–206.
- [18] Walid Krichene and Steffen Rendle. 2020. On sampled metrics for item recommendation. In *Proceedings of the 26th ACM SIGKDD international conference on knowledge discovery & data mining*. 1748–1757.
- [19] Julian McAuley, Christopher Targett, Qinfeng Shi, and Anton Van Den Hengel. 2015. Image-based recommendations on styles and substitutes. In *Proceedings of the 38th international ACM SIGIR conference on research and development in information retrieval*. 43–52.
- [20] Ruihong Qiu, Zi Huang, Hongzhi Yin, and Zijian Wang. 2022. Contrastive learning for representation degeneration problem in sequential recommendation. In *Proceedings of the fifteenth ACM international conference on web search and data mining*. 813–823.
- [21] Shashank Rajput, Nikhil Mehta, Anima Singh, Raghunandan Hulikal Keshavan, Trung Vu, Lukasz Heldt, Lichan Hong, Yi Tay, Vinh Tran, Jonah Samost, et al. 2024. Recommender systems with generative retrieval. *Advances in Neural Information Processing Systems* 36 (2024).
- [22] Ruiyang Ren, Zhaoyang Liu, Yaliang Li, Wayne Xin Zhao, Hui Wang, Bolin Ding, and Ji-Rong Wen. 2020. Sequential recommendation with self-attentive multi-adversarial network. In *Proceedings of the 43rd international ACM SIGIR conference on research and development in information retrieval*. 89–98.
- [23] Xubin Ren, Lianghao Xia, Yuhao Yang, Wei Wei, Tianle Wang, Xuheng Cai, and Chao Huang. 2024. Sslrec: A self-supervised learning framework for recommendation. In *Proceedings of the 17th ACM International Conference on Web Search and Data Mining*. 567–575.
- [24] Weiping Song, Chence Shi, Zhiping Xiao, Zhijian Duan, Yewen Xu, Ming Zhang, and Jian Tang. 2019. AutoInt: Automatic feature interaction learning via self-attentive neural networks. In *Proceedings of the 28th ACM international conference on information and knowledge management*. 1161–1170.
- [25] 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.
- [26] Jiayi Tang and Ke Wang. 2018. Personalized top-n sequential recommendation via convolutional sequence embedding. In *Proceedings of the eleventh ACM international conference on web search and data mining*. 565–573.
- [27] 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 (2017).
- [28] Jianling Wang, Haokai Lu, James Caverlee, Ed H Chi, and Minmin Chen. 2024. Large Language Models as Data Augmenters for Cold-Start Item Recommendation. In *Companion Proceedings of the ACM on Web Conference 2024*. 726–729.
- [29] Shoujin Wang, Longbing Cao, Yan Wang, Quan Z Sheng, Mehmet A Orgun, and Defu Lian. 2021. A survey on session-based recommender systems. *ACM Computing Surveys (CSUR)* 54, 7 (2021), 1–38.
- [30] Shoujin Wang, Liang Hu, Yan Wang, Longbing Cao, Quan Z Sheng, and Mehmet Orgun. 2019. Sequential recommender systems: challenges, progress and prospects. *arXiv preprint arXiv:2001.04830* (2019).
- [31] Xinyuan Wang, Liang Wu, Liangjie Hong, Hao Liu, and Yanjie Fu. 2024. LLM-Enhanced User-Item Interactions: Leveraging Edge Information for Optimized Recommendations. *arXiv preprint arXiv:2402.09617* (2024).
- [32] Shu Wu, Yuyuan Tang, Yanqiao Zhu, Liang Wang, Xing Xie, and Tieniu Tan. 2019. Session-based recommendation with graph neural networks. In *Proceedings of the AAAI conference on artificial intelligence*, Vol. 33. 346–353.
- [33] Xu Xie, Fei Sun, Zhaoyang Liu, Shihwen 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.
- [34] Lanling Xu, Zhen Tian, Gaowei Zhang, Junjie Zhang, Lei Wang, Bowen Zheng, Yifan Li, Jiakai Tang, Zeyu Zhang, Yupeng Hou, Xingyu Pan, Wayne Xin Zhao, Xu Chen, and Ji-Rong Wen. 2023. Towards a More User-Friendly and Easy-to-Use Benchmark Library for Recommender Systems. In *SIGIR*. ACM, 2837–2847.
- [35] Yuhao Yang, Chao Huang, Lianghao Xia, Chunzhen Huang, Da Luo, and Kangyi Lin. 2023. Debaised contrastive learning for sequential recommendation. In *Proceedings of the ACM web conference 2023*. 1063–1073.
- [36] Zhenrui Yue, Sara Rabhi, Gabriel de Souza Pereira Moreira, Dong Wang, and Even Oldridge. 2023. LlamaRec: Two-stage recommendation using large language models for ranking. *arXiv preprint arXiv:2311.02089* (2023).
- [37] Peiyang Zhang, Jiayan Guo, Chaozhuo Li, Yueqi Xie, Jae Boum Kim, Yan Zhang, Xing Xie, Haohan Wang, and Sunghun Kim. 2023. Efficiently leveraging multi-level user intent for session-based recommendation via atten-mixer network. In *Proceedings of the sixteenth ACM international conference on web search and data mining*. 168–176.
- [38] Shengyu Zhang, Lingxiao Yang, Dong Yao, Yujie Lu, Fuli Feng, Zhou Zhao, Tat-Seng Chua, and Fei Wu. 2022. Re4: Learning to re-contrast, re-attend, re-construct for multi-interest recommendation. In *Proceedings of the ACM Web Conference 2022*. 2216–2226.
- [39] Tingting Zhang, Pengpeng Zhao, Yanchi Liu, Victor S Sheng, Jiajie Xu, Deqing Wang, Guanfang Liu, Xiaofang Zhou, et al. 2019. Feature-level deeper self-attention network for sequential recommendation. In *IJCAI*. 4320–4326.
- [40] Zhe Zhao, Lichan Hong, Li Wei, Jilin Chen, Aniruddh Nath, Shawn Andrews, Aditee Kumbhakar, Maheswaran Sathiamoorthy, Xinyang Yi, and Ed Chi. 2019. Recommending what video to watch next: a multitask ranking system. In *Proceedings of the 13th ACM conference on recommender systems*. 43–51.
- [41] 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.