

AutoMLP: Automated MLP for Sequential Recommendations

Muyang Li*
City University of Hong Kong
University of Sydney
muli0371@uni.sydney.edu.au

Zijian Zhang*
City University of Hong Kong
Jilin University
zhangzj2114@mails.jlu.edu.cn

Xiangyu Zhao†
City University of Hong Kong
xianzhao@cityu.edu.hk

Wanyu Wang
City University of Hong Kong
wanyuwang4-c@my.cityu.edu.hk

Minhao Zhao
Fuxi AI Lab, NetEase
zhaominghao@corp.netease.com

Runze Wu
Fuxi AI Lab, NetEase
wurunze1@corp.netease.com

Ruocheng Guo
Bytedance AI Lab
rguo.asu@gmail.com

ABSTRACT

Sequential recommender systems aim to predict users' next interested item given their historical interactions. However, a long-standing issue is how to distinguish between users' long/short-term interests, which may be heterogeneous and contribute differently to the next recommendation. Existing approaches usually set predefined short-term interest length by exhaustive search or empirical experience, which is either highly inefficient or yields subpar results. The recent advanced transformer-based models can achieve state-of-the-art performances despite the aforementioned issue, but they have a quadratic computational complexity to the length of the input sequence. To this end, this paper proposes a novel sequential recommender system, AutoMLP, aiming for better modeling users' long/short-term interests from their historical interactions. In addition, we design an automated and adaptive search algorithm for preferable short-term interest length via end-to-end optimization. Through extensive experiments, we show that AutoMLP has competitive performance against state-of-the-art methods, while maintaining linear computational complexity.

CCS CONCEPTS

• Information systems → Recommender systems.

KEYWORDS

Sequential Recommendations, MLP, AutoML

ACM Reference Format:

Muyang Li, Zijian Zhang, Xiangyu Zhao, Wanyu Wang, Minhao Zhao, Runze Wu, and Ruocheng Guo. 2023. AutoMLP: Automated MLP for Sequential Recommendations. In *Proceedings of the ACM Web Conference 2023 (WWW '23, May 1–5, 2023, Austin, TX, USA)*

*Both authors contributed equally to this research.

†Xiangyu Zhao is the corresponding author.

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](https://www.acm.org).

WWW '23, May 1–5, 2023, Austin, TX, USA

© 2023 Copyright held by the owner/author(s). Publication rights licensed to ACM.

ACM ISBN 978-1-4503-9416-1/23/04...\$15.00

<https://doi.org/10.1145/3543507.3583440>



Figure 1: Illustrative examples of users' long/short-term interests. The target items of user 1 and 2 are related to their long- and short-term interests respectively, and user 3 is influenced by both.

'23), May 1–5, 2023, Austin, TX, USA. Austin, Texas, USA, 9 pages. <https://doi.org/10.1145/3543507.3583440>

1 INTRODUCTION

Sequential recommender systems have been playing an essential role in real-world service providers, such as purchase prediction [34, 36, 41], web page recommendations [7], and next point-of-interest recommendations [3, 9, 24]. It models the sequential dependency based on users' historical behaviors, which consistently outperforms static recommendation models due to the utilization of sequential information [4, 5, 21].

Essential information for making informed sequential recommendations could be summarized into three-fold. The first fold is the long-term user sequential dependency, which relates to the entire sequence of interacted items. It illustrates the users' relatively static long-term interests. The second fold is the short-term user sequential dependency, which may deviate from the long-term interest in a dynamic recommendation environment. It describes the users' most recent interest, which could be temporarily independent. In addition, item features contribute significantly to illustrating the

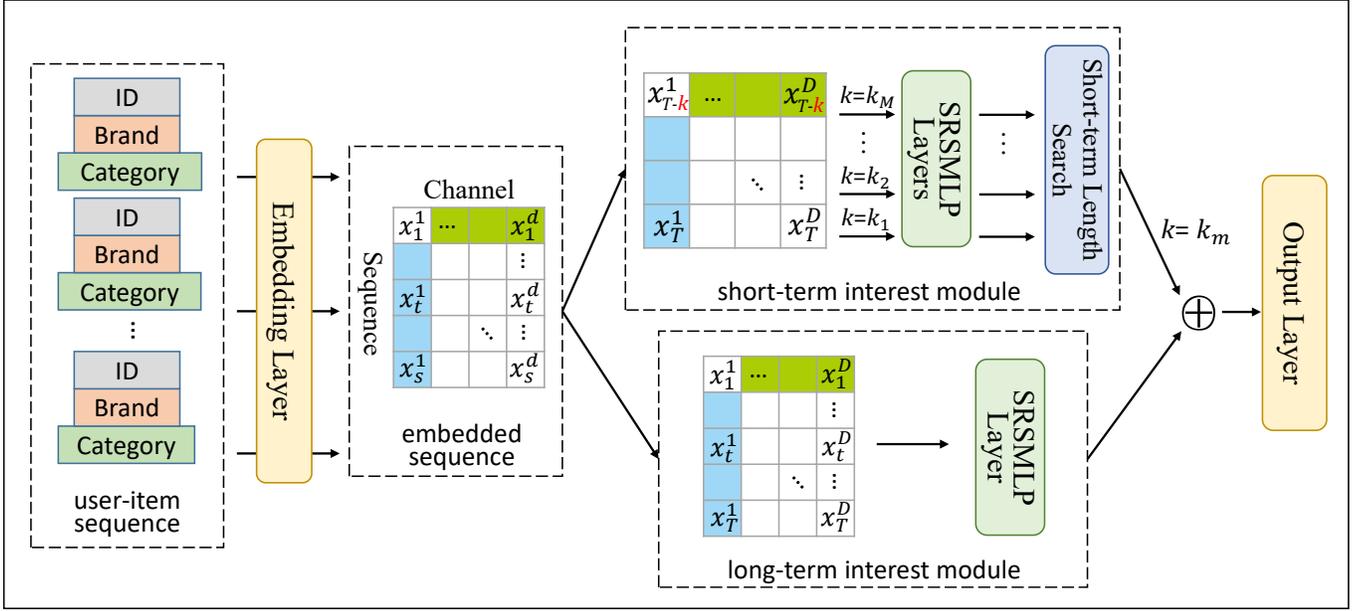


Figure 2: Framework overview of AutoMLP.

users' short-term dependency [19], which is the third-fold critical information for sequential recommendations. Obviously, users' behaviors in sequential recommendations are usually affected by both their long- and short-term interests, which describe users' preferences from different perspectives, as illustrated in Figure 1. However, while the majority of existing works are adept at capturing long-term interest, only a few specifically address the relatively dynamic short-term interest.

Existing methods in sequential recommendation usually tend to model the sequential dependencies among interacted items using Recurrent Neural Network (RNN) [5, 25, 32] and RNN with attention mechanism [14, 33]. RNN-related methods can successfully capture the short-term dependency, but it is well-known that RNN is prone to degenerate when facing long sequences, even-though techniques such as LSTM and GRU were proposed to mitigate this problem. Besides, recent advances show that Transformer-based methods can outperform previous methods by a large margin [10, 23, 37]. However, due to the nature of self-attention, they are insensitive to the order of the input sequence, so they must rely on auxiliary positional embedding to learn the sequential information. Despite it has been shown that such an approach sometimes may hurt the performance [10, 11]. In addition, Transformer's computational complexity is quadratic to the input sequence, so the computational cost for handling a long user interaction sequence is not neglectable. Furthermore, since users' behaviors could be heterogeneous in the long-term and short-term, *i.e.*, the learned short-term dependencies might not accord with the long-term or vice-versa, most existing works train separate models to capture long- and short-term interests, respectively. However, they restrict the short-term sequence length through empirical practices, such as selecting a fixed amount of most recent interactions, or from a fixed time

interval [19, 26, 33]. These approaches are obviously not adaptive or automated enough, because the optimal short-term length varies across different recommendation tasks or scenarios.

To mitigate aforementioned issues, we propose a simple yet efficient method called **Automated Long-term Short-term Multi-Layer Perceptron** for sequential recommendation (**AutoMLP**). AutoMLP only consists of MLP blocks, therefore maintaining linear time and space complexity. We devise a long-term interest module and a short-term interest module to capture long- and short-term dependencies respectively. To automatically adapt the short-term interest window across different tasks, we leverage continuous relaxation to convert the discrete sequence length to a continuous and differentiable representation via AutoML techniques, which could be optimized via gradient descent [18]. To show the effectiveness of our proposed method, we test it on commonly used benchmarks. To summarize, our paper has the following contributions:

- We propose a new long/short-term sequential recommendation model which only consists of MLP blocks. Despite its simple architecture and linear-complexity, it presents competitive performance against state-of-the-art methods.
- We devise an automated short-term session length learner to search the local-optimal short-term interest length adaptive to the given context, which enhances the model's generality.
- We test our proposed method through extensive experiments on both commonly used datasets and real-world private datasets from the industry.

2 FRAMEWORK

In this section, we will make a detailed introduction to our proposed method - AutoMLP, which can adaptively choose users' short-term

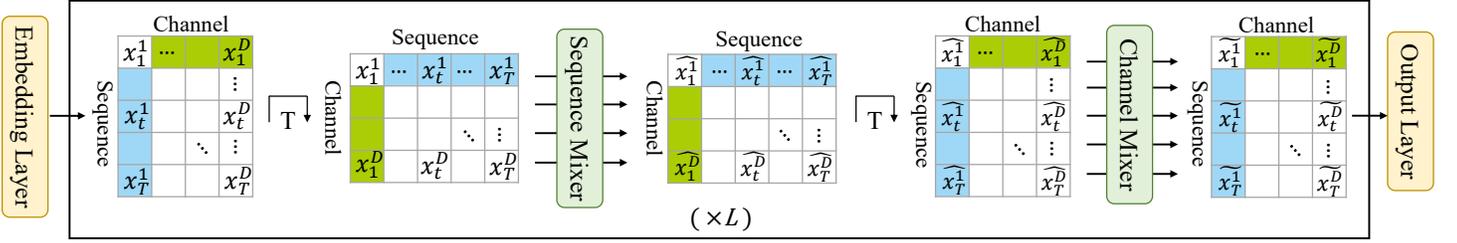


Figure 3: Architecture of Sequential Recommender System MLP Layer (SRMLP Layer).

interest length, and use MLP blocks to capture long-term and short-term interests efficiently. We will first present the overview of the AutoMLP framework, and then discuss the components in detail.

2.1 Problem Formulation

The goal of sequential recommendation is to predict the user’s next possible interacted item given their historical interaction sequences. Thus, for each user u in the user set U , there is a corresponding sequence from $S = \{S_1, \dots, S_u, \dots, S_U\}$, where S_u indicates the interaction sequence for user u . Within each sequence, user-interacted items will be sorted chronologically, *i.e.*, $S_u = \{i_1, \dots, i_t, \dots, i_T\}$, where i_t is the item interacted at timestamp t , and T is the sequence length. Then, the objective of sequential recommendation can be formally defined as: *given user u ’s historical interaction sequence S_u , the goal is to find a function $f : S_u \rightarrow i_{T+1}$, where i_{T+1} is the next item to be recommended to the user u .*

2.2 Framework Overview

Now we introduce the overview of AutoMLP - a sequential recommender system based solely on MLP architectures, which can automatically learn the appropriate short-term user interest length for different sequential recommendation applications. The main body of AutoMLP is made up of two separate MLP-based networks, namely the long-term user interest module and short-term interest module as shown in Figure 2. To be specific, the long-term user interest module takes the entire user historical behavior sequence for prediction and therefore is more leaning towards long-term sequential dependencies.

On the other hand, the short-term user interest module takes a certain number of latest interactions before time T , tending to model the short-term sequential dependencies. The number of recent interactions k will be determined by a Neural Architecture Search (NAS) algorithm, DARTS [18], which leverages continuous relaxation to make the neural architecture search space differentiable and thus can be optimized by gradient descent. Finally, the outputs of separated modules will be fused by a fully-connected layer to predict the next item interacted.

2.3 Detailed Architecture

In this section, we will introduce the macro-structures as well as the methodology behind our proposed method in detail.

2.3.1 Embedding Layer. Following the commonly used strategy in sequential recommender system [10, 23, 36], we use a lookup table to transform the item ID and other features into embedding

vectors. Assuming each item i_t has C features, we then use a fully-connected layer to fuse C embedding vectors from different features into one embedding vector $x_t = [x_t^1, \dots, x_t^D]$ with predefined dimension D . For a user interaction sequence with length T , we then attain an embedding table with size $T \times D$, *i.e.*, the embedded sequence in Figure 2.

2.3.2 Long-term Interest Module. Now we present the key component for modeling the long-term sequential dependencies, *i.e.*, long-term interest module. As we mentioned above, the long-term interest module takes the entire user interaction sequence as input and encoded the interest as the hidden representation for prediction. The encoding process is carried out by a sequential recommender system MLP layer, *i.e.*, SRMLP Layer in Figure 3, which is comprised of two MLP blocks, namely the *sequence-mixer* and *channel-mixer*. SRMLP Layer alternately employs both MLP blocks to capture the sequential correlations within the user interaction sequence and the cross-channel correlations within each item’s embedding vector.

Sequence-mixer. The illustration of the structure of the sequence-mixer is shown in Figure 4. The objective of the sequence-mixer is to learn the sequential information within the user interaction sequence. The input dimension of the sequence-mixer is the length of the input sequence, which is T . Recall that after the embedding layer, we have an embedding table with a size of $T \times D$, so sequence-mixer is applied to the transposed embedding table with a size of $D \times T$, as shown in Figure 3. Sequence-mixer performs non-linear projection that maps input x^d to output \hat{x}^d of the same shapes, *i.e.*, $\mathbb{R}^{D \times T} \mapsto \mathbb{R}^{D \times T}$, fusing the cross-item (sequential) information within the sequence to each item’s embedding vector. More specifically, since the inputs are the same embedding dimension of the entire sequence, *i.e.*, $|[x_1^d, \dots, x_T^d]| = T, \forall d \in [1, D]$, we can observe that sequence-mixer can learn the sequential correlation throughout the sequence, which often reveals users’ long-term interest evolution. We can denote the output of the sequence-mixer at layer l as:

$$\hat{x}^d = x^d + W^2 g^l(W^1 \text{LayerNorm}(x^d)) \quad (1)$$

where $d = 1, 2, \dots, D$. x^d is the d th example from the embedding table. \hat{x}^d is the output of the sequence-mixer block, g^l is the non-linear activation function at layer l , $W^1 \in \mathbb{R}^{R_s \times T}$ denotes the learnable weights representing the first fully connected layer in the sequence-mixer, $W^2 \in \mathbb{R}^{T \times R_s}$ signifies the learnable weights of the

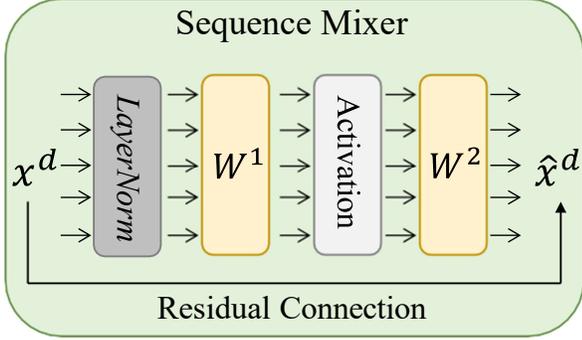


Figure 4: Sequence mixer. Channel mixer has a similar architecture, except for the input x^d and output \hat{x}^d .

second fully connected layer in the sequence-mixer, R_s is the tunable hidden size of sequence-mixer. We employ layer normalization (*LayerNorm*) and residual connection as in MLP-mixer [28].

Channel-mixer. Sharing a similar macro-structure to sequence-mixer in Figure 4, the key distinction between channel-mixer and sequence-mixer is their purpose. The objective of the channel-mixer is to learn the correlation within each item’s embedding vector. Since each dimension of an embedding vector usually expresses different semantics, collectively learning their representation is vital for making an informed prediction [13]. Also, after the sequence-mixer, the sequential information is fused in the individual embedding dimension of each embedding vector, channel-mixer can also communicate the sequential information from different dimensions and thus coherently learn the hidden representation of sequential information [28]. Therefore, the input dimension of the channel-mixer must match the size of the embedding vector, which is D . So for the embedding table, we obtained from the embedding layer, after passing through the sequence-mixer and encoding with sequential information, we first transpose it back to the initial shape, *i.e.*, $T \times D$. Then, channel-mixer maps input x^t to output \hat{x}^t of the same shapes, *i.e.*, $\mathbb{R}^{T \times D} \mapsto \mathbb{R}^{T \times D}$, fusing the cross-channel correlation to each item’s embedding vector. We can denote the output of the channel-mixer at layer l as:

$$\hat{x}^t = x^t + W^4 g^l(W^3 \text{LayerNorm}(x^t)) \quad (2)$$

where $t = 1, 2, \dots, T$, x^t is the input, which is the t^{th} example from the embedding table. \hat{x}^t is the output of the channel-mixer that contains cross-channel correlations. $W^3 \in \mathbb{R}^{R_c \times D}$ is learnable weights of the first fully connected layer in the channel-mixer, and $W^4 \in \mathbb{R}^{D \times R_c}$ is that of the second fully connected layer, where R_c denotes the tunable size of hidden layer in channel-mixer.

2.3.3 Short-term Interest Module. Now we introduce the short-term interest module for capturing the short-term sequential dependencies in users’ historical sequences. The short-term interest module follows a similar framework as the long-term interest module, meaning it also contains a sequence-mixer and channel-mixer. The key difference is the short-term interest length search section

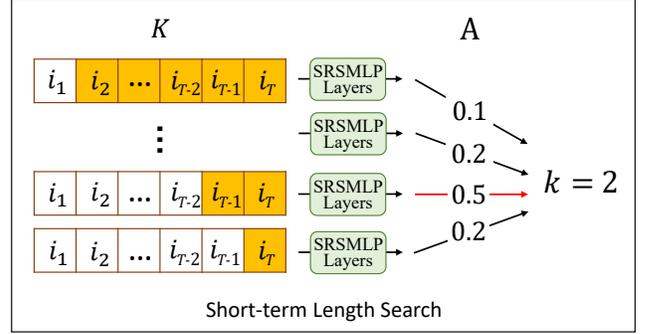


Figure 5: Short-term session length search process. K is the set of candidate lengths, represented by the highlight part.

that can automatically select optimal short-term interest length k in a data-driven manner.

Session Length Search. The session length search part is a Neural Architecture Search (NAS) algorithm based on DARTS [18]. To be specific, we first define some candidate short-term lengths before the last item i_T , representing the possible range of users’ short-term interest length, *i.e.*, $K = \{k_1, \dots, k_m, \dots, k_M\}$, as shown in Figure 5. Then we create M SRSMLPs, each for individual candidate k_m . Next, as shown in Figure 5, we assign a set of learnable *architectural weights* $A = \{\alpha_1, \dots, \alpha_m, \dots, \alpha_M\}$ to the outputs from different SRSMLPs respectively. After that, we apply softmax to transform weights as continuous and differentiable approximations:

$$p_m = \frac{\exp(\alpha_m)}{\sum_{m=1}^M \exp(\alpha_j)} \quad (3)$$

where $m = 1, \dots, M$, and α_m is the learnable weight for m^{th} candidate short-term sequence. After Equation (3), we now have a set of continuous and differential weights $\{p_1, \dots, p_m, \dots, p_M\}$ that could adaptively decide the length of the short-term sequence.

Discussion. Since the effect between the value of short-term interest length and the model performance is not monotonic, thus to determine a local optimal value must apply an exhaustive search, which is extremely computationally expensive for a long user-item interaction sequence, since there are more possible candidates. Therefore, the main upside of such an approach is to learn a local optimal user short-term interest length without enumerating every possible model architecture and train them repetitively, thus making the decision process of selecting short-term interest length highly efficient and adaptive.

2.3.4 Output Layer. Once we find the optimal length of the short-term sequence, we can obtain the outputs from the long-term interest module and the optimal short-term interest module. Then, we use a fully-connected layer to learn their joint representation as follows:

$$h_T = W^o \text{LayerNorm}(x_T^s; x_T^l) + b^o \quad (4)$$

where h_T is the final hidden representation of time step T , x_T^s and x_T^l are the output from the short-term interest module and long-term interest module at time step T , respectively. W^o is learnable

weights that projects the combined x_T^s and x_T^l to a lower dimension representation, hence $\mathbf{W}^o \in \mathbb{R}^{D \times 2D}$, and $\mathbf{b}^o \in \mathbb{R}^D$ is learnable bias vector. Finally, h_T will be used for predicting users' next interaction, which we will later introduce in Section 2.4.2 of *model inference*.

2.4 Training and Inference

Following the common setting in the sequential recommendation, our training loss is defined by a Cross-Entropy loss function:

$$\mathcal{L} = - \sum_{S_u \in \mathcal{S}} \sum_{t \in [1, \dots, T]} [\log(\sigma(r_{i_t, t})) + \sum_{j \notin S_n} \log(1 - \sigma(r_{i_j, t}))] \quad (5)$$

where σ demotes sigmoid function. $r_{i_t, t}$ is model's predicted similarity to ground-truth item i_t , and $r_{i_j, t}$ is the predicted similarity to sampled items at time step t , where j is the negative sampled items. S is the superset that contains all users' interaction sequences.

2.4.1 Training. The training process contains two phases. The first phase is the search phase, aiming to find local optimal A^* representing preferable short-term length. The second phase is the retraining phase, where after A^* is found, we retrain the AutoMLP framework with the optimal short-term length.

Search Phase. The previous section illustrates how continuous relaxation can make the search space \mathcal{K} differentiable and thus could be optimized via back-propagation. Now our objective is to jointly learn the architectural weights $A = \{\alpha_m\}$, and all other learnable parameters of AutoMLP \mathbf{W} . Note that each p_m is corresponding with each α_m , *i.e.*, learning $\{p_m\}$ is the same to learning $\{\alpha_m\}$. While A is a subset of the learnable parameters of AutoMLP, literature shows that simply updating \mathbf{W} and A altogether will cause overfitting problems in the training process, since they are highly dependent from each other [39, 40]. Therefore, we will use the training dataset to optimize \mathbf{W} following common practice, while using the validation dataset to optimize A . To be specific, we formulate this as a bilevel optimization [2], with A as upper-level variable and \mathbf{W} as lower-level variable [18]. Formally, we denote that as:

$$\begin{aligned} \min_A \mathcal{L}_{val}(\mathbf{W}^*(A), A) \\ \text{s.t. } \mathbf{W}^*(A) = \arg \min_{\mathbf{W}} \mathcal{L}_{train}(\mathbf{W}, A) \end{aligned} \quad (6)$$

In addition, it is notable that the inner optimization (lower-level variable optimization) for solving Equation (6) is usually computationally expensive, since every time we intend to optimize A , we have to first train \mathbf{W} till converged. Alternatively, we apply a one-step approximation:

$$\mathbf{W}^*(A) \approx \mathbf{W} - \xi \nabla_{\mathbf{W}} \mathcal{L}_{train}(\mathbf{W}, A) \quad (7)$$

where ξ is the learning rate. The motivation of this approximation is to approximate $\mathbf{W}^*(A)$ by a single training step instead of from training thoroughly. The detailed training algorithm of the search phase is shown in Algorithm 1.

Retraining Phase. In the search stage, we incorporate all candidate short-term sequences into the AutoMLP framework, where the suboptimal short-term sequences may lead to a negative impact on model performance. Thus, we need to retrain AutoMLP with only the optimal short-term sequence.

Algorithm 1 Optimization for AutoMLP in Search Phase

Input: embedding table \mathbf{x} and ground-truth interaction sequences S

Output: well-learned parameters \mathbf{W}^* and A

- 1: **while** not converged **do**
 - 2: Sample a mini-batch of validation data examples
 - 3: Estimate the approximation of $\mathbf{W}^*(A)$ via Eq.(7)
 - 4: Update \mathbf{V} by descending $\nabla_A \mathcal{L}_{val}(\mathbf{W}^*(A), A)$
 - 5: Sample a mini-batch of training data examples
 - 6: Update \mathbf{W} by descending $\nabla_{\mathbf{W}} \mathcal{L}_{train}(\mathbf{W}, A)$
 - 7: **end while**
-

To be specific, after the optimization of the search stage, we obtain well-trained A^* . We select the optimal short-term length with the highest α_m , and discard all other candidate lengths. Then, we retrain the AutoMLP framework based on the outputs from the long-term interest module and the optimal short-term interest module following common deep recommender system training paradigm, *i.e.*, optimizing loss $\min_{\mathbf{W}} \mathcal{L}_{train}(\mathbf{W})$ by gradient descent.

2.4.2 Inference. Our proposed method employs a commonly used inference process in sequential recommender systems [10, 23, 36, 42], which is to compute the cosine similarity between model output and all candidate items. More specifically, once we have the hidden representation h_T from the output layer as the final hidden representation of the entire sequence, we then calculate the dot product of h_T and the item embeddings of all candidate items I :

$$p_i = \text{softmax}(h_T \cdot E_i^T) \quad (8)$$

where $i = 1, 2, \dots, |I|$, and $|I|$ is the total number of candidate items. E_i is the item embedding of i , p_i is the predicted probability for i being the next possible interaction. Then, candidate item with highest p_i is the predicted next interaction i_{T+1} .

2.5 Complexity Analysis

This section will show the important strength of AutoMLP, which is the linear computational complexity. The complexity of AutoMLP is primarily determined by two MLP blocks, sequence-mixer, and channel-mixer. Since the time complexity of an MLP with one hidden layer is $O(\text{input units} \times \text{hidden units} \times \text{output units})$, thus we can know that the complexity of sequence-mixer is $O(T \times R_s + R_s \times T)$. Since the determinant variable here is maximum sequence length T , and R_s is a constant that is independent of T , we can rewrite that as $O(T)$. Similarly, we have the complexity of the channel-mixer as $O(D)$. And since short-term modules have the exact same structure except for the sequence length, we have the complexity of the short-term sequence-mixer as $O(k)$, and the complexity of the short-term channel-mixer as $O(D)$ as well. Hence the complexity of AutoMLP is $O(T + D + k)$, and it is linear to the three determinant variables.

3 EXPERIMENTS

In this section, we will evaluate our proposed method - AutoMLP through extensive experiments on two commonly used benchmark datasets in the recommender system.

Table 1: Statistics of the datasets.

Data	MovieLens	Amazon Beauty
# Interactions	1,000,209	2,023,070
# Items	3,952	249,274
# Users	6,040	1,210,271

3.1 Datasets

MovieLens¹: MovieLens is a commonly used benchmark dataset for recommender systems, it contains users' rating history on a wide range of movies with corresponding timestamps. In this study we choose MovieLens-1M for experiments, which contains over a million interactions, detailed statistics are available in Table 1.

Amazon Beauty²: This dataset contains reviews and ratings of various items that are crawled from Amazon [20], it is divided into multiple categories, and in this paper, we will use "Beauty" category, statistics of this category are available in Table 1.

In order to model the long-term dependencies in a more realistic scenario, we filter out users who have less than 10 interactions in MovieLens and Amazon Beauty. More statistics about the datasets can be found in Table 1.

3.2 Experimental Setting

We follow the common scenario in the sequential recommendation, which is the next-item prediction (leave-one-out evaluation). In this setting, we use the last interactions from user interaction sequences as the test set, the second last interactions as the validation set, and all of the previous interactions as training set [10, 15, 23, 36, 42]. We will also pair 100 negative samples for each ground-truth item in evaluation, sampling based on their popularity [8, 23].

3.3 Evaluation Metrics

We adopt the three most commonly used metrics for evaluation, namely Hit Ratio (HR), Normalized Discounted Cumulative Gain (NDCG), and Mean Reciprocal Rank (MRR). HR is the accuracy for the ground-truth item appearing in top N recommendations, NDCG is a ranking loss that measures the position of the ground-truth item in the top N recommendations, and MRR is the reciprocal of the ground-truth item's ranking in top N recommendations.

3.4 Baselines

To illustrate the effectiveness of our proposed method, we compare it against several popular sequential recommendation models, where some of them achieve state-of-the-art performances. Of their design motivation and adapted techniques, they could be divided into three groups:

General sequential recommender systems refers to the representative sequential recommendation models of some important categories, namely MC-based, RNN-based, and Transformer-based.

- **FPMC** Factorizing Personalized Markov Chains (FPMC) [22] represents the earlier works in sequential recommender systems,

which is using Markov Chains to model the sequential dependencies between items.

- **GRU4Rec** Hidasi *et al.* proposed using RNN for the sequential recommendation, and improved vanilla RNN with Gated Recurrent Unit (GRU) and point-wise ranking loss [5].
- **BERT4Rec** Sun *et al.* utilized Bi-directional self-attention to improve transformer-based sequential recommendation [23], and several studies have shown that BERT4Rec surpasses vanilla self-attention approach such as SASRec constantly [13, 23, 42]

Long-term short-term sequential recommender systems focus on modeling the short-term sequential dependencies while maintaining long-term sequential dependencies as well:

- **NextItNet** Yuan *et al.* pointed out some existing limitations of Caser [26], especially its drawback in modeling long-term sequential dependencies and proposed a generative model which does not require pooling operation for a bigger receptive field to improve CNN-based sequential recommender systems [35].

Feature-level sequential recommender systems refer to the sequential recommendation models that can learn information within item features to aid next-item prediction.

- **GRU4Rec+** is the improved version of GRU4Rec, which leverages item features for better prediction under a feature-rich scenario and parallel RNN for higher efficiency [6].
- **FDSA** Feature-level Deeper Self-Attention Network for Sequential Recommendation (FDSA) [36] utilize item features for the sequential recommendation, and has been considered to have state-of-the-art performance in supervised-methods [42].

3.5 Implementation Details

The implementations of baseline methods as well as AutoMLP are based on RecBole framework [38], an open-sourced library for recommender systems, which offers an unbiased environment to evaluate the performance of our proposed methods.

For baseline methods, we use hyper-parameters as suggested in the original papers. For those unspecified hyper-parameters, we use grid search for hyper-parameter tuning. The searching strategy is to select the hyper-parameters that yield the best performance on the validation set. We use the early-stop strategy, *i.e.*, if, in the next 10 epochs, the validation performance does not improve, we stop the training and use the model from the current best validation performance. For some larger models such as BERT4Rec, when dealing with large datasets, their space complexity could be too large to fit into GPU memory, therefore we define a "space efficient" setting, referring to the search best hyper-parameters within our GPU memory capacity, which is 32GB. We use Adam [12] as the optimizer for all implementations, the learning rate is set as $1e-3$, β_1 as 0.9, β_2 as 0.999.

3.6 Overall Performance

In Table 2, we show the overall comparison of our proposed methods against a wide range of popular baseline methods.

From the results, we can make the following observations:

¹<https://grouplens.org/datasets/movielens/1m/>

²<http://jmcauley.ucsd.edu/data/amazon/>

Table 2: Overall performance comparison. Best performances are bold, next best performances are underlined

Methods	MovieLens			Beauty			Param
	MRR@10	NDCG@10	HR@10	MRR@10	NDCG@10	HR@10	
FPMC	0.2453	0.3088	0.5156	0.0991	0.1251	0.2098	100 M
GRU4Rec	<u>0.3893</u>	<u>0.4553</u>	0.6666	0.1162	0.1435	0.2324	33.5 M
BERT4Rec	0.3535	0.4289	<u>0.6695</u>	0.0907	0.1198	0.2154	17 M
NextItNet	0.2085	0.2642	0.4455	0.1087	0.1393	0.2385	16.8 M
GRU4Rec ⁺	0.3736	0.4412	0.6578	<u>0.1325</u>	<u>0.1638</u>	<u>0.2657</u>	34.1 M
FDSA	0.3725	0.4409	0.6594	0.1305	0.1595	0.2536	34.3 M
AutoMLP	0.3912*	0.4593*	0.6767*	0.1438*	0.1754*	0.2779*	16.8 M

“**” indicates the statistically significant improvements (i.e., two-sided t-test with $p < 0.05$) over the original model.

“Param” refers to the number of trainable model parameters on Beauty dataset, and M = million.

- When item features are available, feature-level models usually obtain better performances than those who cannot process features, indicating the importance of item features in the sequential recommendation.
- AutoMLP maintains competitive performance across all datasets. Specifically, AutoMLP exceeds existing long/short-term sequential recommender systems such as NextItNet significantly, showing it is not only an efficient choice but also a strong alternative in terms of accuracy.
- Finally, we can observe that AutoMLP consumes the least trainable model parameters, which indicates its superior space efficiency in real-world sequential recommender systems.

In summary, AutoMLP can achieve comparable performances against popular baselines on public datasets with lower space complexity, which validates its effectiveness.

3.7 Efficiency Analysis

In this section, we will study the time/space efficiency of AutoMLP, which are critical metrics to launch a recommendation model in industrial recommender systems.

As one of the fundamental motivations, we argue that via automated short-term interest length search, AutoMLP is more efficient than the exhaustive search for optimal short-term interest length. In Figure 6, we compare AutoMLP against AutoMLP-s, which has identical architecture as AutoMLP but a different short-term interest length search algorithm. Specifically, AutoMLP-s exhaustively searches for a short-term interest length that yields the best performance as existing methods [26, 33]. In Figure 6(a), we compare the search time of AutoMLP and AutoMLP-s when the size of candidate search space (i.e., candidate short-term length) is 5. And in Figure 6(b), we show that as the size of the search space increases, the increment ratio of AutoMLP is much smaller, suggesting better scalability of AutoMLP.

In addition, from Figure 7, we can also observe that AutoMLP shows higher training efficiency in terms of time usage and memory consumption against state-of-the-art Transformer-based methods. Figure 7(a) shows the total training time used for training AutoMLP, FDSA, and BERT4Rec till convergence in seconds. We can observe that, because BERT4Rec implements complex training techniques such as Cloze objective and bi-directional self-attention [23], its

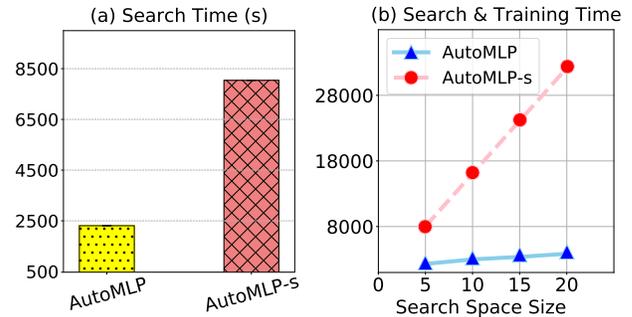


Figure 6: Searching efficiency on Beauty dataset

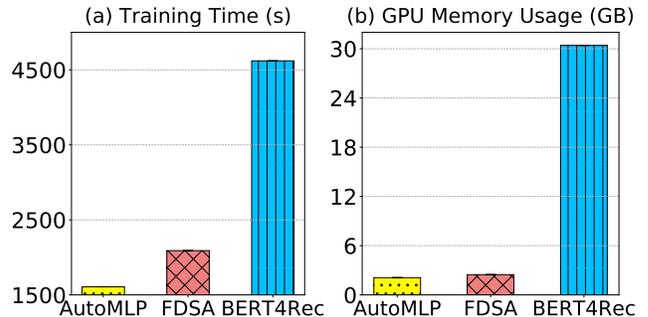


Figure 7: Training efficiency on Beauty dataset

training efficiency is relatively low. While FDSA obtained a more balanced training efficiency by combining self-attention and vanilla attention, its training cost is still more expensive than AutoMLP.

3.8 Hyper-parameters Analysis

In this section, we present how the essential hyper-parameters in AutoMLP affect the model performance. Unlike transformer-based methods [10, 23, 36] and CNN-based methods [26, 35], which have bigger sets of hyper-parameters and larger search space, our proposed methods only have a few key hyper-parameters that can significantly affect its performance. We investigate how will the number of layers L and embedding size D affect the recommendation performance of AutoMLP.

As shown in Figure 8 (a), we can observe that the optimal number of layers is 8. In other words, fewer layers (e.g., $L = 4$) can

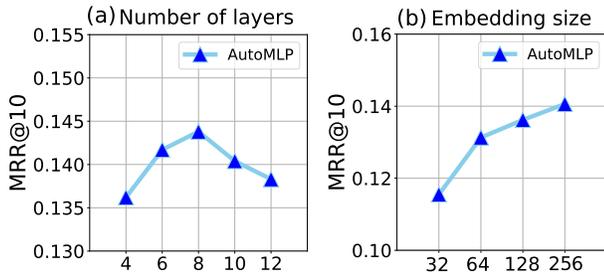


Figure 8: Influence of hyper-parameters on Beauty dataset degenerate model representation ability, while more layers (e.g., $L = 12$) may lead to the overfitting issue. From Figure 8 (b), it can be summarized that small embedding sizes will downgrade the model performance; with the increase of embedding size, AutoMLP possesses representation capacity to better learn the user-item interaction representations.

3.9 Ablation Study

In previous subsections, we have illustrated the effectiveness and efficiency of AutoMLP. An intuitive question behind such effectiveness is how we assign credits to each module of our proposed model. To validate the importance of each part of our model, we devise simplified alternatives architecture, each with one key component removed. Performances are reported in Table 3. Where w/o Sequence Mixer is removing Sequence Mixer from AutoMLP, and w/o Channel Mixer is to remove Channel Mixer from AutoMLP.

As shown in Table 3, we can observe that without Sequence-Mixer, the performance of AutoMLP drops most significantly. As we introduced in the Framework section, Sequence-Mixer aims to capture the sequential correlation of input sequence, without it the model will simply make predictions assuming input sequences are independent. The significant decrease indicates the core importance of sequential information in our recommendation task, and Sequence-Mixer can effectively capture such information.

When we only remove Channel-Mixer, we can also observe a drop in performance, but less significantly compared with removing Sequence-Mixer. The performance drop suggests that exchanging cross-channel information after each layer can also improve the quality of prediction by better learning the joint correlations of different implicit attributes of items.

4 RELATED WORK

In this section, we briefly summarize several lines of works related to us, consisting of sequential recommendation and MLP-mixer.

4.1 Sequential Recommendation

Sequential recommendation plays an essential part in recommender systems. It aims at depicting users’ successive preferences based on sequential interactions between users and items. RNN enjoys a natural strength in modeling sequential dependencies and takes the major part of the deep learning-based sequential recommendation. GRU4Rec [5] employs the Gated Recurrent Unit (GRU) in the session-based recommendation. It’s improved version [4] further boosts the performance with a ranking loss function as well as an improved sampling strategy. However, RNN suffers from a

Table 3: Ablation study comparison on Beauty dataset

Model	MRR@10	NDCG@10	HR@10
w/o Sequence Mixer	0.1132	0.1315	0.2169
w/o Channel Mixer	0.1408	0.1720	0.2700
AutoMLP	0.1438	0.1754	0.2779

high training cost, especially when modeling long-term sequences. Furthermore, both RNN-based models can not capture long-term dependencies well due to the vanishing gradient problem [31].

Since the successful application in [1], attention has been widely applied to the sequential recommendation. NARM [14] captures sequential user-item interactions as well as the user’s main purpose in the current session through an encoder-decoder framework. Yu *et al.* [34] propose to maintain user preference through jointly learning the individual- and union-level item interactions.

Lately, thanks to the advances in related areas [30], the transformer has achieved better performance than RNN structure. To mitigate the deficiency of RNN-based models, SASRec [10] proposes to use an attention mechanism to model long-term sequential dependency. Bert4Rec [23] employs deep bidirectional self-attention to learn user preferences. FDSA [36] models transition patterns by integrating heterogeneous features with different weights. Besides, transformer-based models do not account for the order of input sequence without a heuristic process[11].

4.2 MLP-mixer

Recent advances in MLP architectures [17, 27–29] show that with simple alternations in design, MLP can serve as a strong alternative to Transformer-based models and attain competitive performances against state-of-the-art methods while maintaining both smaller time and space complexity. Among them, MLP-Mixer [28] is usually considered the most representative work, by using separate MLP blocks to learn the per-location correlations and cross-location correlations from images, MLP-Mixer shows great scalability on large datasets and comparable performance against Transformer-based and Convolution Neural Network (CNN) based methods.

We note that our proposed method is not the first endeavor to try to apply MLP-only architecture in a recommender system. MOI-Mixer [13] first investigates the possibility of using MLP architecture as a substitute for Transformer-based methods. Subsequently, MLP4Rec [16] proposed a tri-directional information fusion scheme, to coherently capture higher-order interactions across different attribute levels under a feature-rich scenario.

5 CONCLUSION

In this paper, we proposed a long-term short-term sequential recommender system named AutoMLP. By only leveraging MLP architectures, AutoMLP shows competitive performances against state-of-the-art methods on both open-sourced benchmark datasets and real-world dataset from industrial application. Moreover, we devised an automated short-term interest length search algorithm that can efficiently learn an optimal short-term interest length, together with the linear complexity of MLP architectures, our method shows better efficiency and more promising improvement space compared with the existing methods.

ACKNOWLEDGEMENT

This research was partially supported by APRC - CityU New Research Initiatives (No.9610565, Start-up Grant for New Faculty of City University of Hong Kong), SIRG - CityU Strategic Interdisciplinary Research Grant (No.7020046, No.7020074), HKIDS Early Career Research Grant (No.9360163), Huawei Innovation Research Program and Ant Group (CCF-Ant Research Fund).

REFERENCES

- [1] Dzmitry Bahdanau, Kyunghyun Cho, and Yoshua Bengio. 2014. Neural machine translation by jointly learning to align and translate. *arXiv preprint arXiv:1409.0473* (2014).
- [2] Benoit Colson, Patrice Marcotte, and Gilles Savard. 2007. An overview of bilevel optimization. *Annals of operations research* 153 (2007), 235–256.
- [3] Jie Feng, Yong Li, Chao Zhang, Funing Sun, Fanchao Meng, Ang Guo, and Depeng Jin. 2018. Deepmove: Predicting human mobility with attentional recurrent networks. In *Proceedings of the 2018 world wide web conference*. 1459–1468.
- [4] 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.
- [5] Balázs Hidasi, Alexandros Karatzoglou, Linas Baltrunas, and Domonkos Tikk. 2015. Session-based recommendations with recurrent neural networks. *arXiv preprint arXiv:1511.06939* (2015).
- [6] Balázs Hidasi, Massimo Quadrana, Alexandros Karatzoglou, and Domonkos Tikk. 2016. Parallel recurrent neural network architectures for feature-rich session-based recommendations. In *Proceedings of the 10th ACM conference on recommender systems*. 241–248.
- [7] Linmei Hu, Chen Li, Chuan Shi, Cheng Yang, and Chao Shao. 2020. Graph neural news recommendation with long-term and short-term interest modeling. *Information Processing & Management* 57, 2 (2020), 102142.
- [8] Jin Huang, Wayne Xin Zhao, Hongjian Dou, Ji-Rong Wen, and Edward Y Chang. 2018. Improving sequential recommendation with knowledge-enhanced memory networks. In *The 41st international ACM SIGIR conference on research & development in information retrieval*. 505–514.
- [9] Md Ashrafur Islam, Mir Mahathir Mohammad, Sarkar Snigdha Sarathi Das, and Mohammed Eunus Ali. 2022. A survey on deep learning based Point-of-Interest (POI) recommendations. *Neurocomputing* 472 (2022), 306–325.
- [10] Wang-Cheng Kang and Julian McAuley. 2018. Self-attentive sequential recommendation. In *2018 IEEE international conference on data mining (ICDM)*. IEEE, 197–206.
- [11] Guolin Ke, Di He, and Tie-Yan Liu. 2020. Rethinking positional encoding in language pre-training. *arXiv preprint arXiv:2006.15595* (2020).
- [12] Diederik P Kingma and Jimmy Ba. 2014. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980* (2014).
- [13] Hojoon Lee, Dongyoon Hwang, Sunghwan Hong, Changyeon Kim, Seungryong Kim, and Jaegul Choo. 2021. Moi-mixer: Improving mlp-mixer with multi order interactions in sequential recommendation. *arXiv preprint arXiv:2108.07505* (2021).
- [14] 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 conference on information and knowledge management*. 1419–1428.
- [15] Jiacheng Li, Yujie Wang, and Julian McAuley. 2020. Time interval aware self-attention for sequential recommendation. In *Proceedings of the 13th international conference on web search and data mining*. 322–330.
- [16] Muyang Li, Xiangyu Zhao, Chuan Lyu, Minghao Zhao, Runze Wu, and Ruocheng Guo. 2022. MLP4Rec: A Pure MLP Architecture for Sequential Recommendations. In *31st International Joint Conference on Artificial Intelligence and the 25th European Conference on Artificial Intelligence (IJCAI-ECAI 2022)*. International Joint Conferences on Artificial Intelligence, 2138–2144.
- [17] Hanxiao Liu, Zihang Dai, David So, and Quoc V Le. 2021. Pay attention to mlps. *Advances in Neural Information Processing Systems* 34 (2021), 9204–9215.
- [18] Hanxiao Liu, Karen Simonyan, and Yiming Yang. 2018. Darts: Differentiable architecture search. *arXiv preprint arXiv:1806.09055* (2018).
- [19] Chen Ma, Peng Kang, and Xue Liu. 2019. Hierarchical gating networks for sequential recommendation. In *Proceedings of the 25th ACM SIGKDD international conference on knowledge discovery & data mining*. 825–833.
- [20] 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.
- [21] Massimo Quadrana, Paolo Cremonesi, and Dietmar Jannach. 2018. Sequence-aware recommender systems. *ACM Computing Surveys (CSUR)* 51, 4 (2018), 1–36.
- [22] 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.
- [23] 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.
- [24] Ke Sun, Tiejun Qian, Tong Chen, Yile Liang, Quoc Viet Hung Nguyen, and Hongzhi Yin. 2020. Where to go next: Modeling long-and short-term user preferences for point-of-interest recommendation. In *Proceedings of the AAAI Conference on Artificial Intelligence*, Vol. 34. 214–221.
- [25] Yong Kiam Tan, Xinxing Xu, and Yong Liu. 2016. Improved recurrent neural networks for session-based recommendations. In *Proceedings of the 1st workshop on deep learning for recommender systems*. 17–22.
- [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] Yuki Tatsunami and Masato Taki. 2021. Raftmlp: Do mlp-based models dream of winning over computer vision. *arXiv preprint arXiv:2108.04384* 3 (2021).
- [28] Ilya O Tolstikhin, Neil Houlsby, Alexander Kolesnikov, Lucas Beyer, Xiaohua Zhai, Thomas Unterthiner, Jessica Yung, Andreas Steiner, Daniel Keysers, Jakob Uszkoreit, et al. 2021. Mlp-mixer: An all-mlp architecture for vision. *Advances in neural information processing systems* 34 (2021), 24261–24272.
- [29] Hugo Touvron, Piotr Bojanowski, Mathilde Caron, Matthieu Cord, Alaeldin El-Nouby, Edouard Grave, Gautier Izacard, Armand Joulin, Gabriel Synnaeve, Jakob Verbeek, et al. 2022. Resmlp: Feedforward networks for image classification with data-efficient training. *IEEE Transactions on Pattern Analysis and Machine Intelligence* (2022).
- [30] 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).
- [31] 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).
- [32] Chao-Yuan Wu, Amr Ahmed, Alex Beutel, Alexander J Smola, and How Jing. 2017. Recurrent recommender networks. In *Proceedings of the tenth ACM international conference on web search and data mining*. 495–503.
- [33] Haochao Ying, Fuzhen Zhuang, Fuzheng Zhang, Yanchi Liu, Guandong Xu, Xing Xie, Hui Xiong, and Jian Wu. 2018. Sequential recommender system based on hierarchical attention network. In *IJCAI International Joint Conference on Artificial Intelligence*.
- [34] Lu Yu, Chuxu Zhang, Shangsong Liang, and Xiangliang Zhang. 2019. Multi-order attentive ranking model for sequential recommendation. In *Proceedings of the AAAI conference on artificial intelligence*, Vol. 33. 5709–5716.
- [35] Fajie Yuan, Alexandros Karatzoglou, Ioannis Arapakis, Joemon M Jose, and Xiangnan He. 2019. A simple convolutional generative network for next item recommendation. In *Proceedings of the twelfth ACM international conference on web search and data mining*. 582–590.
- [36] Tingting Zhang, Pengpeng Zhao, Yanchi Liu, Victor S Sheng, Jiajie Xu, Deqing Wang, Guanfeng Liu, Xiaofang Zhou, et al. 2019. Feature-level Deeper Self-Attention Network for Sequential Recommendation. In *IJCAI*. 4320–4326.
- [37] Kesen Zhao, Xiangyu Zhao, Zijian Zhang, and Muyang Li. 2022. MAE4Rec: Storage-saving Transformer for Sequential Recommendations. In *Proceedings of the 31st ACM International Conference on Information & Knowledge Management*. 2681–2690.
- [38] Wayne Xin Zhao, Shanlei Mu, Yupeng Hou, Zihan Lin, Yushuo Chen, Xingyu Pan, Kaiyuan Li, Yujie Lu, Hui Wang, Changxin Tian, et al. 2021. Recbole: Towards a unified, comprehensive and efficient framework for recommendation algorithms. In *Proceedings of the 30th ACM International Conference on Information & Knowledge Management*. 4653–4664.
- [39] Xiangyu Zhao, Haochen Liu, Wenqi Fan, Hui Liu, Jiliang Tang, and Chong Wang. 2021. Autoloss: Automated loss function search in recommendations. In *Proceedings of the 27th ACM SIGKDD Conference on Knowledge Discovery & Data Mining*. 3959–3967.
- [40] Xiangyu Zhao, Haochen Liu, Hui Liu, Jiliang Tang, Weiwei Guo, Jun Shi, Sida Wang, Huiji Gao, and Bo Long. 2021. Autodim: Field-aware embedding dimension search in recommender systems. In *Proceedings of the Web Conference 2021*. 3015–3022.
- [41] Guorui Zhou, Na Mou, Ying Fan, Qi Pi, Weijie Bian, Chang Zhou, Xiaoqiang Zhu, and Kun Gai. 2019. Deep interest evolution network for click-through rate prediction. In *Proceedings of the AAAI conference on artificial intelligence*, Vol. 33. 5941–5948.
- [42] 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.