# TriMLP: Revenge of a MLP-like Architecture in Sequential Recommendation

Yiheng Jiang[1,2], Yuanbo Xu[1,2] ✉, Yongjian Yang[1,2], Funing Yang[1,2], Pengyang Wang[3], Hui Xiong[4,5]

[1]*Lab of Mobile Intelligence (MIC), Jilin University, Changchun, China*
[2]*College of Computer Science and Technology, Jilin University, Changchun China*
[3]*Department of Computer and Information Science, University of Macau, Macao SAR, China*
[4]*Thrust of Artificial Intelligence, Hong Kong University of Science and Technology (Guangzhou), Guangzhou, China*
[5]*Department of Computer Science and Engineering, Hong Kong University of Science and Technology, Hong Kong SAR, China*
[1,2]jiangyh22@mails.jlu.edu.cn; [1,2]{yuanbox, yyj, yfn}@jlu.edu.cn;
[3]pywang@um.edu.mo; [4,5]xionghui@ust.hk

*Abstract*—Sequential recommenders concentrate on modeling the transmission patterns shrouded in sequences of historical user-item interactive behaviors (or referred as token) and inferring dynamic preferences over candidate items. Fueled by diverse advanced neural network architectures like RNN, CNN and Transformer, existing methods have enjoyed rapid performance boost in the past years. Recent progress on MLP lights on an efficient method, token-mixing MLP, to establish contact among tokens. However, due to the unrestricted cross-token communications, i.e., information leakage issue, caused by the inherent fully-connection structure, we find that directly migrating these modern MLPs in recommendation task would neglect the chronological order of historical sequences and lead to subpar performances. In this paper, we present a MLP-like architecture for sequential recommendation, namely TriMLP, with a novel Triangular Mixer for cross-token communications. In designing Triangular Mixer, we simplify the cross-token operation in MLP as the basic matrix multiplication, and drop the lower-triangle neurons of the weight matrix to block the anti-chronological order connections from future tokens. Accordingly, the information leakage issue can be remedied and the prediction capability of MLP can be fully excavated under the standard auto-regressive mode. Take a step further, the mixer serially alternates two delicate MLPs with triangular shape, tagged as global and local mixing, to separately capture the long range dependencies and local patterns on fine-grained level, i.e., long and short-term preferences. Empirical study on 12 datasets of different scales (50K~10M user-item interactions) from 4 benchmarks (Amazon, MovieLens, Tenrec and LBSN) show that TriMLP consistently attains promising accuracy/efficiency trade-off, where the average performance boost against several state-of-the-art baselines achieves up to 14.88% with 8.65% less inference cost. Our code is available at https://github.com/jiangyiheng1/TriMLP.

*Index Terms*—sequential recommendation, data mining, multi-layer perceptron

## I. INTRODUCTION

Sequential recommendation processes sequences of historical user-item interactive behaviors (or referred as tokens in this paper), concentrates on mining dependencies among tokens and inferring preferences over time, and provides pleasing suggestions [1]. The performances of sequential recommenders are closely tied with the reliant neural network architectures, which serve as essential components for establishing contact among tokens and capturing transformation patterns.
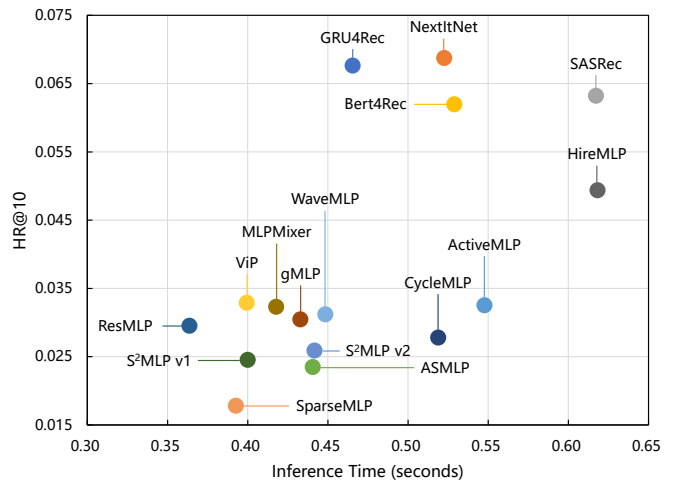


Fig. 1. Accuracy/Efficiency traded-off on `QB-Video`. Along the vertical axis, the higher, the better recommendation performance; along the horizontal axis, the more left, the less inference cost.

Innovations in neural network architectures have consistently played a major role in sequential recommendation. Recurrent neural network (RNN)-based sequential recommenders, represented by [2], [3], transmits the information in tokens step-by-step. Methods like [4], [5], adopting convolutional neural network (CNN), aggregate the local spatial features with sliding filters. Credited to the superb adaptability with sequential tasks, Transformer architecture [6] that dynamically re-weights and integrates tokens through the self-attention mechanism has become the *de-facto* backbone in modern sequential recommenders [7]–[9].

In most recent, "retrospective" research on purely multi-layer perceptron-based (MLP) models, pioneered by MLP-Mixer [10] and ResMLP [11], investigates an conceptually simple and computationally efficiency idea to realize the cross-token communication, namely token-mixing MLP, where tokens interact with each other independently and identically across dimensions [12]. In terms of structural properties, token-mixing MLP holds the global reception field as Transformer while reserves the learned sequential dependency as
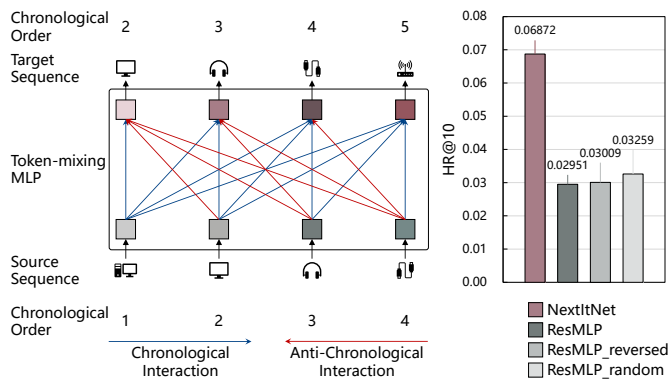
Fig. 2. The left part illustrates the unrestricted cross-token interactions in MLP. The blue arrows denote the interaction in chronological order where the current token can only attend to itself and previous tokens, while the red ones are in contrast which lead to the information leakage. The right histogram reveals that the fully-connected MLP is insensitive to the sequence order.

static weights like RNN and CNN. Despite token-mixing MLP is originally derived from the vision community, it is intuitive that such MLP owns promising potential for sequential tasks.

However, the empirical observation draws apart from the exception when we refer and explore token-mixing MLP in sequential recommendation. Following the implementations in [13], we reproduce RNN-based GRU4Rec [2], CNN-based NextItNet [4], Transformer-based SASRec [7] and Bert4Rec [8] and various modern MLP-based models [10], [11], [14]–[18], [18]–[22] on `QB-Video` [13]. As shown in Figure 1, albeit most MLPs inherit the merit of efficiency (along the horizontal axis), their recommendation performances stand far behind other neural network architecture-based sequential recommenders, e.g. the strongest HireMLP lags behind NextItNet by 28.19% on the metric of HR@10.

We argue such subpar performance is oriented from the inherent *Fully-Connection Design* in MLPs. As the example in the left part of Figure 2, a user has historically interacted with 1: Computer→2: Monitor→3: Headphone→4: USB, and the MLP aims at predicting the $i + 1$ th token at step $i$ under the standard auto-regressive training mode. Unfortunately, except for the cross-token interactions in chronological order (denoted as blue arrows in Figure 2) consistent with the natural behavior pattern [2], [4], [5], [7], the fully-connected MLP inevitably conducts the anti-chronological ones (red arrows) which would leak the future information to the current time step and suppress the prediction ability throughout the training procedure [1]. To further verify whether or not MLP is sensitive to the sequence order, we train ResMLP on `QB-Video` with different ordered sequences, i.e., chronological, reserved and random. As summarized in the histogram of Figure 2, the results of MLPs are sharing inferior performances to CNN-based NextItNet with negligible standard deviations. It supports our opinion about the incompatibility between MLP and auto-regressive manner, that the fully-connection nullifies the capacity of implicitly encoding and differentiating the position of each token [11].

---

[1] We also consider utilizing the bidirectional attribute of MLP and conduct the auto-encoding training mode [23]. See Section V-K for more details.

In this paper, we propose to build a MLP-like architecture for sequential recommendation, with the aim of persisting the computationally efficiency advantage, gearing to the auto-regressive training fashion and catching up to the performances obtained using advanced neural network architectures.

In designing the MLP-based token mixer, we present Triangular Mixer to remedy the issues brought by fully-connection. It is inspired by the use of masking strategy in Transformer-based methods [9]. In principle, since the cross-token interactions endowed by MLP can be simplified as the matrix multiplication, undesirable interactions can be forbidden by disabling specific neurons in MLP. In practice, we drop the lower-triangle elements in the weight matrix of MLP to block the connections from future tokens and ensure that each token can only attend to itself and previous ones. Naturally, the information leakage issue is avoided, and the potential of MLP can be fully excavated under the auto-regressive training.

Take a step further, since MLP with global reception field excels in modeling the long-range relations among tokens while fails in capturing local patterns [24], we derive two mixing layers based on the above delicate MLP with triangular shape, tagged as global mixing and local mixing. The global mixing follows the vanilla triangular shape and attaches importance to all tokens in sequences for inferring long-term preference. The local one further drops specific upper-triangle neurons of weight matrix and treats the input sequence as multiple non-overlapped independent sessions with equal length. Specifically, the shape of active neurons is converted as several isosceles right sub-triangles arranged along the main diagonal whose sides are equal to the session length. Each sub-triangle is responsible for capturing the short-term preference contained in the corresponding session. Triangular Mixer serially alternates global mixing and local mixing for the fine-grained sequential dependency modeling.

To this end, we present a MLP-like sequential recommender TriMLP based on the proposed Triangular Mixer. In summary, our major contributions can be listed as follows:

- We refer and explore the idea of all-MLP architecture in sequential recommendation. To the best of our knowledge, we are the first to empirically point out that the fully-connection in MLP is not compatible with the standard auto-regressive training mode.
- We present a MLP-like sequential recommender, namely TriMLP, with a novel Triangular Mixer which endows the chronological interactions among tokens.
- We put forward Triangular Mixer with the global mixing and local mixing, to reconcile the long-rang dependencies and local patterns in sequences.
- We evaluate TriMLP with 12 datasets of different scales from 4 benchmarks (MovieLens, Amazon, Tenrec and LBSN) which contain 50K~10M user-item interactive behaviors. The experimental results demonstrate that TriMLP attains stable and promising accuracy/efficiency trade-off over all validated datasets, i.e., averagely surpasses the performance of several state-of-the-art baselines by 14.88% with 8.65% less inference cost.

## II. RELATED WORK

### A. Sequential Recommendation

Sequential recommendation aims at capturing dynamic preferences from sequences of historical user-item interactive behaviors and providing pleasant suggestions [1]. Building upon technological breakthroughs in the past decade [25]–[28], this field has ushered a new era of deep learning. Hidasi et al. [2] leveraged RNN to model the sequential dependency which transmits the information contained in token step-by-step. The spatial local information aggregation in CNN also benefits sequential recommenders [4], [5]. SASRec [7] and Bert4Rec [8] separately employed unidirectional and bidirectional Transformer-base encoder [6], [23] to dynamically extract the relationship between target and relevant items. Towards all-MLP methods, FMLP4Rec [29] referred the learnable filter-based denoising manner and encoded sequential patterns with Fourier transform, and MLP4Rec [30] incorporated contextual information (e.g., item characteristics) into the MLPMixer architecture [10]. In separate lines of research, [31], [32] utilized graph neural network (GNN) to enhance item representations, [33] adopted hierarchical structures, [34], [35] introduced data augmentation and [36]–[38] exploited pre-training techniques.

TriMLP architecture focuses on improving the primary sequential modeling capacity of MLP under the essential auto-regressive training mode without assistance of any auxiliary information. Credited to the triangular design, TriMLP successfully merges the performance gap between MLP and other advanced neural network-based sequential recommenders.

### B. Toke-mixing MLP

Since the pioneering MLPMixer [10] and ResMLP [11] have been proposed in the early 2020s, all-MLP models are staging a comeback in vision community. These models rely on the novel MLP-based token mixer where tokens interact independently and identically across the channel dimension. Due to the simple concept with less intensive computation, such deep MLP models have stirred up a lot of interest [12] and derive a surge of variants. According to the dimensions of mixing tokens, these variants can be divided into three categories: (i) employing both the axial direction and the channel dimension [14], [16], [17], [39], [40] which proposes to orthogonally decompose the cross-token mixing, maintain long-range dependency and encode cross-token relation along axial directions, (ii) considering only the channel dimension [18]–[22], [41], [42] which aligns features at different spatial locations to the same channel by axial shifting, and then interacting with spatial information through channel projection. (iii) reserving the entire spatial and channel dimensions [10], [11], [15], [17], [24], [43] which retains the global reception field and channel projection.

In the pilot experiments, we empirically investigate that the inherent fully-connection design in MLP is incompatible with sequential tasks especially under the auto-regressive training fashion. In contrast, the proposed Triangular Mixer provides a simple, effective and efficient alternative to remedy the issue.

## III. PRELIMINARY

### A. Basic Definition

Let $\mathcal{U}$ and $\mathcal{I}$ denote the user and item set, respectively. Accordingly, we have the following basic definitions.

**Definition 1: (User-Item Interactive Behavior)** A user-item interactive behavior, or referred as token in this paper, is represented as a triplet $x = \langle u, i, t \rangle$, which denotes that the user $u \in \mathcal{U}$ interacted with the item $i \in \mathcal{I}$ at time $t$.

**Definition 2: (Historical Sequence)** A historical sequence, tagged as $X^u = x_1 \rightarrow x_2 \rightarrow \cdots \rightarrow x_{|X^u|}$, chronologically records $|X^u|$ user-item interactive behaviors of the user $u$.

### B. Problem Statement

**Sequential Recommendation** Given the specific user $u$ and his/her historical sequence $X^u$. Sequential recommendation problem infers the dynamic preferences and provides the top-K recommendation list, which contains $K$ items that the user might be most likely to interact in the next time step. It can be formulated as the following equation,

$$\mathcal{F}(X^u) \rightarrow TopK^u, \tag{1}$$

where $TopK^u$ denotes the top-K recommendation list and $\mathcal{F}(\cdot)$ is the abstract symbol of any sequential recommender.

## IV. METHODOLOGY

### A. Architecture Overview

The macro overview of TriMLP architecture is depicted in the Figure 3 (a). TriMLP takes a historical sequence of $n$ user-item interactive behaviors (or tokens) as input, where $n$ is the maximum sequence length. The tokens are independently pass through the Embedding layer to form the $d$-dimension sequence representation matrix. The resulting embedding is then fed into the Triangular Mixer to produce cross-token interactions. The Classifier takes these encoded representations as input, and predicts the probabilities over all candidate items.

### B. Embedding

Considering that the historical sequences of different users are inconsistent in length, we set the maximum sequence length to $n$. Following the operation in [7], we split the longer sequences into several non-overlapping sub-sequences of length $n$. For the shorter ones, we repeatedly add the "padding" token in the head until their lengths grow to $n$. For the clarity and conciseness, we omit the superscript that denotes the specific user $u$, and the embedding layer can be formulated with the following equation,

$$\text{Embed}(X) \rightarrow \boldsymbol{X} \in \mathbb{R}^{n \times d}, \tag{2}$$

where $\boldsymbol{X}$ is the sequence representation matrix. Note that the padding tokens are encoded with constant zero vectors [9] and excluded from the gradient update.

(a) TriMLP Architecture
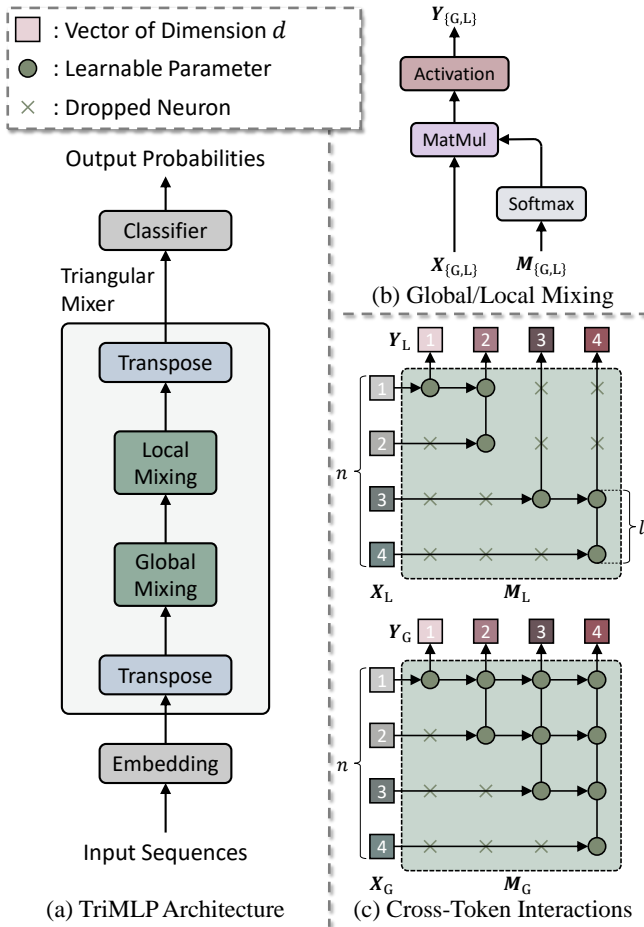
(b) Global/Local Mixing

(c) Cross-Token Interactions

Fig. 3. (a) depicts the proposed MLP-like architecture TriMLP. (b) reveals the details of global and local mixing in Triangular Mixer. (c) presents an illustrative example with sequence length $n = 4$ and session length $l = 2$ to explain the cross-token interactions in global and local mixing.

```
1  import torch
2  import torch.nn as nn
3  # n: input sequence length
4  # l: session length
5  # s: number of sessions
6  # n = l * s
7  def generate_kernel(n, l, s):
8    mask = torch.zeros([n, n])
9    for i in range(0, s):
10     mask[i*l: i*l+l, i*l: i*l+l] = torch.ones(l, l)
11   M_G = torch.triu(torch.ones([n, n]))
12   M_L = M_G.masked_fill(mask == 0.0, 0)
13   M_G = nn.parameter.Parameter(M_G, requires_grad=True)
14   M_L = nn.parameter.Parameter(M_L, requires_grad=True)
15   return M_G, M_L
16
17 class TriangularMixer(nn.Module):
18   def __init__(self, n, l, s, act):
19     super().__init__()
20     assert l * s == n
21     self.M_G, self.M_L = generate_kernel(n, l, s)
22     self.act = act
23
24   def forward(self, X):
25     # X: input sequence embedding, [b, n, d]
26     X_G = X.permute(0, 2, 1)
27     Y_G = self.act(torch.matmul(X_G, self.M_G).softmax(dim=-1))
28     Y_L = self.act(torch.matmul(Y_G, self.M_L).softmax(dim=-1))
29     Y = Y_L.permute(0, 2, 1)
30     return Y
```

### C. Triangular Mixer

Triangular Mixer endows the cross-token communication in strict compliance with chronological order. As shown in Figure 3 (a), the mixer takes as input the sequence representations $X$, and encodes the sequential dependency through the global mixing layer and local mixing layer, successively. Formulaically, it can be expressed as,

$$Y = \mathrm{TriMix}(X) = \mathrm{Mix_L}\left(\mathrm{Mix_G}\left(X^\top\right)\right)^\top, \quad (3)$$

where "$\top$" is the matrix transposition and $Y \in \mathbb{R}^{n \times d}$ is the encoded sequence representations. The global mixing $\mathrm{Mix_G}(\cdot)$ injects the long-range sequential dependency and the local mixing $\mathrm{Mix_L}(\cdot)$ further captures the local patterns. These two mixing layers share the identical structure (as Figure 3 (b)), and can be expressed as,

$$\begin{aligned} Y_{\{G,L\}} &= \mathrm{Mix_{\{G,L\}}}(X_{\{G,L\}}) \\ &= \mathrm{Act}(X_{\{G,L\}} \cdot \mathrm{Softmax}(M_{\{G,L\}})) \end{aligned} \quad (4)$$

where $M_{\{G,L\}} \in \mathbb{R}^{n \times n}$ are the mixing kernels in global and local mixing, separately, i.e., the learnable weight matrices in MLPs. Since the cross-token operations in MLP is actually re-weighting and integrating tokens based on the weight matrix,

i.e., linear combination, we utilize activation function $\mathrm{Act}(\cdot)$ to inject the non-linearity. We also adopt $\mathrm{Softmax}(\cdot)$ to convert the parameters in MLP as the probabilities over tokens. Note that we employ the unified notations in Eq. 4 for simplicity. Specifically, the input for global mixing $X_G \in \mathbb{R}^{d \times n}$ is the transposed sequence embedding, and the corresponding output $Y_G \in \mathbb{R}^{d \times n}$ also serves as the input for local mixing $X_L \in \mathbb{R}^{d \times n}$. The output of local mixing $Y_L \in \mathbb{R}^{d \times n}$ is transposed to form the final output of the mixer $Y$.

Next, we devote into the details of the delicate $M_{\{G,L\}}$ in global and local mixing, respectively.

*1) Global Mixing:* The triangular design in $M_G$ gets insight from the utilization of mask strategy in Transformer-based sequential recommenders [7], [9], which masks the lower-triangular elements in the attentive map to prevent the information leakage. Similarly, we drop the lower-triangular neurons in $M_G$ to cut off the contact from future tokens. The lower part of Figure 3 (c) provides an illustrative example with the input sequence of length $n = 4$ to explain the cross-token communication, where "1, 2, 3, 4" denote the chronological order. $\mathrm{Mix_G}(\cdot)$ compels that the 2 nd token can only attend to the previous 1 st token and itself. Mathematically, the $i$ th token in the global mixing interact with each other as

$$Y_{G_{*,i}} = \sum_{j=1}^{i} X_{G_{*,j}} M_{G_{j,i}}, \quad \text{for } i = 1, \cdots n, \quad (5)$$

where "$*$" denotes any dimension in $d$ and the upper bound $i$ of cumulative sum blocks the information from future tokens. In the premise of avoiding information leakage, global mixing $\mathrm{Mix_G}(\cdot)$ reserves the global reception field for the long-range sequential dependency.

*2) Local Mixing:* Based on the aforementioned global triangular design, local mixing $\mathrm{Mix_L}(\cdot)$ further selectively drops specific upper-triangle neurons to capture the local patterns, calling for the short-term preferences. In principle, it treats the input sequence as $s$ non-overlapping sessions of length $l$

where $n = s \times l$ and forbids the cross-session communications. In practice, the shape of active neurons in $\boldsymbol{M}_{\mathrm{L}}$ converts into $s$ isosceles right triangles of equal side length $l$, which arrange along the main diagonal. The resulting sub-triangles are responsible for capturing the local patterns contained in corresponding sessions. As the example in Figure 3 (c) where the session length $l = 2$, the 4 th token attaches importance to the 3 rd token and itself, while the information from the 1 st and 2 nd token is ignored. Generally, the $i$ th token in the local mixing interact with each other as

$$\boldsymbol{Y}_{\mathrm{L}*,i} = \sum_{j=\lceil i/l \rceil}^{i} \boldsymbol{X}_{\mathrm{L}*,j} \boldsymbol{M}_{\mathrm{L}j,i}, \quad \text{for } i = 1, \cdots n, \quad (6)$$

where the lower bound of cumulative sum, i.e., the round-up operation $j = \lceil i/l \rceil \in [1, s]$, further cuts off the connections from previous sessions.
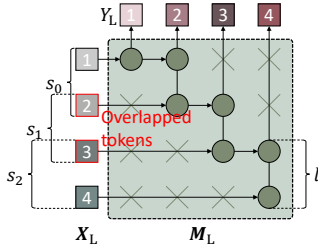


Fig. 4. Cross-token in another variant of local mixing that endows the cross-session connections. The input sequence length of $n = 4$ are divided 3 overlapped sessions $s_0$, $s_1$ and $s_2$. Session $s_1$ extracts the information from session $s_0$ based on the overlapped 2 nd token and so on.

We exploit another variant of local mixing which endows the cross-session communications[2]. It treats the input historical sequence of length $n$ into $s = n - l + 1$ sessions with length $l$. Accordingly, the active neurons in local mixing are reshaped as an isosceles trapezoid with the waist length of $l$. The cross-session connections are built upon the overlapped tokens, i.e., for adjacent sessions, there are $l - 1$ over lapped tokens. Along with the example in Figure 4, the input sequence of length $n = 4$ is split into $s = 3$ sessions $s_{\{0,1,2\}}$ of length $l = 2$, and the session $s_0$, $s_1$ are connected by the overlapped 2 nd token, et cetera. This variant interacts tokens as

$$\boldsymbol{Y}_{\mathrm{L}*,i} = \sum_{j=\max(1,i-l+1)}^{i} \boldsymbol{X}_{\mathrm{L}*,j} \boldsymbol{M}_{\mathrm{L}j,i}, \quad \text{for } i = 1, \cdots n. \quad (7)$$

The PyTorch-like pseudo-code of Triangular Mixer is presented in Algorithm 1. Since we convert the weights in $\boldsymbol{M}_{\{\mathrm{G,L}\}}$ to the probabilities over tokens by $\mathrm{Softmax}(\cdot)$, the dropping operation can be easily implemented by padding these neurons with "$-\infty$". Moreover, we find that initializing the active neurons with 1 which enforces each token contributes equally to the target during the early stage of training procedure[3].

---

[2]We block the cross-session interactions in TriMLP and Section V-J compares the performances of these two local mixing layers.

[3]Section V-H2 analyses the performances with different initialization.

*3) Discussion:* Triangular Mixer shares some similarities with the self-attention mechanism [6], including the global reception field and parallel processing capability.

Notably, our mixer departs from the self-attention mechanism with the following peculiarities:

- **Positional Sensitive**: Since Triangular Mixer compels the cross-token interactions in strict line with the chronological order, the extra positional information (e.g. Positional Encoding [6] or Embedding [23]) is no more necessary.
- **Independent and Static Weights**: Triangular Mixer reserves the sequential dependency as static weights in MLP which is agnostic to the input, rather than the attention map which is dynamically generated by the scale-dot product of query and key matrices.
- **Fewer Parameters and Higher Efficiency**: Triangular Mixer shrinks the parameter-scale by removing the query, key, value matrices mapping and Feed-Forward Network. Accordingly, our method is less computationally intensive than the self-attention mechanism.

### D. Classifier

Recall that the output of Triangular Mixer is denoted as $\boldsymbol{Y}$, the Classifier, i.e., implemented with the plain linear layer and Softmax function, converts the $d$-dimension representation vector of each token to the probabilities over all candidate items at each time step. As follows,

$$\boldsymbol{P} = \mathrm{Softmax}(\boldsymbol{Y} \cdot \boldsymbol{W} + \boldsymbol{b}), \quad (8)$$

where $\boldsymbol{W} \in \mathbb{R}^{d \times |\mathcal{I}|}$ and $\boldsymbol{b} \in \mathbb{R}^{|\mathcal{I}|}$ are learnable parameters. $\boldsymbol{P} \in \mathbb{R}^{n \times |\mathcal{I}|}$ is the calculated probability matrix where $p_{i,c} \in [0, 1]$ is the probability over candidate item $c$ at time step $i$.

### E. Model Training and Recommendation

During the training processing, we apply the standard autoregressive fashion. Specifically, TriMLP takes the historical sequence excluded the last token $X = x_1 \rightarrow x_2 \rightarrow \cdots \rightarrow x_{n-1}$ as source, and the sequence excluded the first token $X = x_2 \rightarrow x_3 \rightarrow \cdots \rightarrow x_n$ as target. At each time step $i$, TriMLP aims at predicting the $i + 1$ th token, i.e., maximizing the probability of the $i + 1$ th interacted item. We use the following cross entropy loss to optimize TriMLP,

$$\mathcal{L} = - \sum_{X^u \in X^{\mathcal{U}}} \sum_{i=1}^{n} \log(p_{i,t_i}), \quad (9)$$

where $X^{\mathcal{U}}$ is a training set of all users' historical sequences , $t_i$ is the target item at step $i$ and $p_{i,t_i}$ is the probability.

During the recommendation stage, TriMLP first extracts the last row $\boldsymbol{p}_n \in \mathbb{R}^{|\mathcal{I}|}$ from $\boldsymbol{P}$ which contains the information of all interacted items in the historical sequence. Then, it ranks all candidate items according to the probabilities and retrieves $K$ items as the top-K recommendation list.

TABLE I
DATA STATISTICS (AFTER PRE-PROCESSED)

| Dataset | Scale | # Users | # Items | # Interactions | Avg. Seq. Length | Max Seq. Length $n$ | Sparsity |
|---------|-------|---------|---------|----------------|------------------|--------------------|----------|
| Beauty | Tiny | 1,664 | 36,938 | 56,558 | 33.99 | 32 | 99.91% |
| Sports | | 1,958 | 55,688 | 58,844 | 30.05 | 32 | 99.95% |
| ML-100K | | 932 | 1,152 | 97,746 | 104.88 | 128 | 90.90% |
| NYC | Small | 1,031 | 5,135 | 142,237 | 137.96 | 128 | 97.31% |
| QB-Article | | 4,671 | 1,844 | 164,939 | 35.31 | 32 | 98.09% |
| TKY | | 2,267 | 7,873 | 444,183 | 195.93 | 128 | 97.51% |
| ML-1M | Base | 6,034 | 3,260 | 998,428 | 165.47 | 128 | 94.92% |
| QB-Video | | 19,047 | 15,608 | 1,370,577 | 71.96 | 64 | 99.54% |
| Brightkite | | 5,714 | 48,181 | 1,765,247 | 308.93 | 256 | 99.36% |
| Yelp | Large | 42,461 | 101,269 | 2,199,786 | 51.81 | 64 | 99.95% |
| Gowalla | | 32,439 | 131,329 | 2,990,783 | 92.20 | 64 | 99.93% |
| ML-10M | | 69,865 | 9,708 | 9,995,230 | 143.06 | 128 | 98.53% |

## V. EXPERIMENT AND DISCUSSION

In this section, we start from introducing the datasets, metrics, baselines and the implement details. Then, we analyze the experimental results, including the overall recommendation performance and ablation study. Take a step further, we explore various characteristics of our Triangular Mixer. In summary,we conduct a large amount of experiments to answer the following four research questions:

- **RQ 1** How is the recommendation performance and inference cost of TriMLP compared to other neural network-based state-of-the-art sequential recommenders?
- **RQ 2** How is the effectiveness of Triangular Mixer under the TriMLP architecture? Can global and local mixing both contribute to model the sequential dependency?
- **RQ 3** How is the property of Triangular Mixer with respect to various internal structures and micro designs?
- **RQ 4** How do the different settings of session number $s$ (or session length $l$) influence the performance of TriMLP and token-mixing manner in Triangular Mixer?

### A. Datasets

We evaluate our method on 12 publicly available datasets from 4 benchmarks. Specifically, we select `Beauty`, `Sports` from Amazon[4] [44], `ML-100K`, `ML-1M`, `ML-10M` from MovieLens[5] [45], `QB-Article`, `QB-Video` from Tenrec[6] [13], and `NYC`[7], `TKY`[7], `Brightkite`[8], `Yelp`[9], `Gowalla`[10] from the scenario of Location-based Social Network (LBSN). In accordance with [9], we remove the "inactive" users who interact with fewer than 20 items and the "unpopular" items which are interacted by less than 10 times. According to the number of interactions, we categorize the 12 datasets into 4 different scales: Tiny, Small, Base and Large which separately contain 50K~100K, 150K~500K, 1M~2M and 2M~10M interactions. Table I summarizes the statistics.

We set the maximum sequence length $n$ of each dataset according to the average sequence length. For each user, we take the last previously un-interacted item as the target and utilize all prior items for training during the data partition.

[4]http://snap.stanford.edu/data/web-Amazon-links.html
[5]https://grouplens.org/datasets/movielens/
[6]https://static.qblv.qq.com/qblv/h5/algo-frontend/tenrec_dataset.html
[7]https://sites.google.com/site/yangdingqi/home/foursquare-dataset
[8]https://snap.stanford.edu/data/loc-brightkite.html
[9]https://www.yelp.com/dataset
[10]https://snap.stanford.edu/data/loc-Gowalla.html

### B. Metric

We introduce the following 3 metrics to measure the efficiency and accuracy of sequential recommenders.

- **Inference Time (Infer. Time)** calculates the average time cost of finishing 100 rounds recommendation.
- **Hit Rate (HR@K)** [46] counts the fraction of times that target item is among the top-K recommendation list.
- **Normalized Discounted Cumulative Gain (NDCG@K)** [46] rewards the method that ranks the positive items in first few positions of the top-K recommendation list.

The smaller Infer. Time stands for the better efficiency, and the recommendation performance is positively correlated with the values of HR and NDCG. We report $K = \{5, 10\}$ in our experiments. To avoid the bias brought by different negative sampling strategies [47], we compare the probability of the target item with all other items in the dataset, and compute the HR, NDCG based on the ranking of all items.

### C. Baselines

Since TriMLP concentrates on improving the primary ability of MLP in encoding cross-token interactions, we compare it with the following four representative sequential recommenders developed from different neural networks:

- **GRU4Rec** [2] utilizes RNN to model historical sequences and dynamic preferences for sequential recommendation.
- **NextItNet** [4] is a state-of-the-art CNN-based generative model for recommendation, which learns high-level representation from both short and long-range dependencies.
- **SASRec** [7] employs Transformer-based encoder for recommendation where the self-attention mechanism dynamically models the sequential dependency.
- **FMLP4Rec** [29] is a state-of-the-art all-MLP sequential recommender, which follows the denoising manner and establishes cross-token contact by Fourier Transform.

### D. Implementation Details

For the rigorous comparison, we uniform the width and depth of TriMLP and baselines, and ensure that all the compared method differentiate only in the neural network architecture. Specifically, we set the embedding dimension $d$ to 128, and the intermediate dimension in Triangular Mixer to $n$. Since TriMLP contains 2 MLP layers for cross-token interactions, we stack 2 token-mixing encoders in each baseline.

During the training stage, we perform the standard auto-regressive manner, and adopt the identical gradient-updating strategy for all compared methods where dropout rate is 0.5 and the optimizer is Adam [48] with the learning rate of 0.001. The parameters keep updating until the performance no longer increases for consecutive 10 epochs. All experiments are conducted on a single server with 64GB RAM, AMD Ryzen 5900X CPU and NVIDIA RTX 3090 GPU.

### E. Overall Recommendation Performance (RQ 1)

The experimental results of all compared methods on 12 datasets are summarized in Table II. From the table, we have the following observations.

## TABLE II
OVERALL RECOMMENDATION PERFORMANCE. THE ARROW "↑" (OR "↓") DENOTES THAT THE HIGHER (OR LOWER) VALUE, THE BETTER METRIC. WE USE BOLDFACE AND UNDERLINE TO INDICATE THE BEST AND SECOND RESULTS IN EACH COLUMN, RESPECTIVELY. THE "COMPARISON" ROW REPORTS THE RELATIVE IMPROVEMENT OR DECLINE OF TRIMLP AGAINST THE STRONGEST BASELINE.

| Dataset-Tiny | | Beauty | | | | | Sports | | | | | ML-100K | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| Model | Mixer | Infer. Time (s)↓ | HR@5↑ | NDCG@5↑ | HR@10↑ | NDCG@10↑ | Infer. Time (s)↓ | HR@5↑ | NDCG@5↑ | HR@10↑ | NDCG@10↑ | Infer. Time (s)↓ | HR@5↑ | NDCG@5↑ | HR@10↑ | NDCG@10↑ |
| GRU4Rec | RNN | 0.3176 | 0.07752 | 0.05130 | 0.11839 | 0.06440 | 0.3361 | 0.01124 | 0.00665 | 0.01839 | 0.00900 | 0.3778 | 0.05472 | 0.03557 | 0.12124 | 0.05682 |
| NextItNet | CNN | 0.3196 | 0.08534 | 0.06376 | 0.10998 | 0.07164 | 0.3339 | 0.01532 | 0.01177 | 0.02043 | 0.01334 | 0.3861 | 0.06974 | 0.04463 | 0.12017 | 0.06083 |
| SASRec | Trans. | 0.3185 | 0.08113 | 0.06106 | 0.11358 | 0.07160 | 0.3293 | 0.00409 | 0.00191 | 0.00817 | 0.00322 | 0.3953 | 0.02682 | 0.01444 | 0.05365 | 0.02300 |
| FMLP4Rec | Four. | 0.3246 | 0.08353 | 0.05962 | 0.12019 | 0.07121 | 0.3323 | 0.01481 | 0.01167 | 0.02298 | 0.01420 | 0.3948 | 0.06760 | 0.04144 | 0.11373 | 0.05646 |
| TriMLP | MLP | 0.3103 | 0.09615 | 0.07070 | 0.12560 | 0.08003 | 0.3176 | 0.01839 | 0.01258 | 0.02451 | 0.01442 | 0.3763 | 0.08691 | 0.05848 | 0.15451 | 0.07988 |
| Comparison | | -2.30% | +11.24% | +9.82% | +4.50% | +11.71% | -3.55% | +20.04% | +6.88% | +6.66% | +1.55% | -0.40% | +24.62% | +31.03% | +27.44% | +31.32% |

| Dataset-Small | | NYC | | | | | QB-Article | | | | | TKY | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| Model | Mixer | Infer. Time (s)↓ | HR@5↑ | NDCG@5↑ | HR@10↑ | NDCG@10↑ | Infer. Time (s)↓ | HR@5↑ | NDCG@5↑ | HR@10↑ | NDCG@10↑ | Infer. Time (s)↓ | HR@5↑ | NDCG@5↑ | HR@10↑ | NDCG@10↑ |
| GRU4Rec | RNN | 0.3315 | 0.02037 | 0.01268 | 0.03589 | 0.01783 | 0.3315 | 0.06487 | 0.03811 | 0.13466 | 0.06052 | 0.3540 | 0.02603 | 0.01715 | 0.04279 | 0.234 |
| NextItNet | CNN | 0.3482 | 0.03686 | 0.02272 | 0.05141 | 0.02748 | 0.3389 | 0.06101 | 0.03565 | 0.11732 | 0.05373 | 0.3716 | 0.03264 | 0.01836 | 0.06308 | 0.02802 |
| SASRec | Trans. | 0.3568 | 0.04074 | 0.02370 | 0.05723 | 0.02889 | 0.3461 | 0.05973 | 0.03820 | 0.10662 | 0.05324 | 0.3909 | 0.02955 | 0.01821 | 0.05073 | 0.02498 |
| FMLP4Rec | Four. | 0.3551 | 0.03395 | 0.02077 | 0.05432 | 0.02722 | 0.3465 | 0.05331 | 0.03158 | 0.10026 | 0.47290 | 0.3854 | 0.04499 | 0.02566 | 0.08293 | 0.03794 |
| TriMLP | MLP | 0.3350 | 0.04365 | 0.02574 | 0.06111 | 0.03121 | 0.3225 | 0.07621 | 0.04751 | 0.13552 | 0.06639 | 0.3350 | 0.05293 | 0.03175 | 0.09043 | 0.04384 |
| Comparison | | -0.42% | +7.14% | +8.61% | +6.78% | +8.03% | -2.79% | +17.48% | +24.37% | +0.64% | +9.70% | -5.76% | +17.65% | +23.73% | +9.04% | +15.55% |

| Dataset-Base | | ML-1M | | | | | QB-Video | | | | | Brightkite | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| Model | Mixer | Infer. Time (s)↓ | HR@5↑ | NDCG@5↑ | HR@10↑ | NDCG@10↑ | Infer. Time (s)↓ | HR@5↑ | NDCG@5↑ | HR@10↑ | NDCG@10↑ | Infer. Time (s)↓ | HR@5↑ | NDCG@5↑ | HR@10↑ | NDCG@10↑ |
| GRU4Rec | RNN | 0.4204 | 0.14700 | 0.09583 | 0.20655 | 0.10151 | 0.4656 | 0.03560 | 0.02135 | 0.06762 | 0.03159 | 0.5385 | 0.01540 | 0.00948 | 0.02678 | 0.01314 |
| NextItNet | CNN | 0.4600 | 0.15291 | 0.09969 | 0.21113 | 0.12480 | 0.5227 | 0.03602 | 0.02149 | 0.06872 | 0.03200 | 0.6272 | 0.02468 | 0.01232 | 0.04705 | 0.02266 |
| SASRec | Trans. | 0.5279 | 0.15214 | 0.09286 | 0.20968 | 0.10877 | 0.6173 | 0.03418 | 0.02051 | 0.06631 | 0.02977 | 0.8505 | 0.03220 | 0.02054 | 0.04953 | 0.02615 |
| FMLP4Rec | Four. | 0.5110 | 0.08336 | 0.05189 | 0.12728 | 0.06602 | 0.6287 | 0.02426 | 0.01332 | 0.04415 | 0.01969 | 0.7262 | 0.02520 | 0.01330 | 0.04428 | 0.01951 |
| TriMLP | MLP | 0.3924 | 0.16390 | 0.11196 | 0.23434 | 0.13454 | 0.4176 | 0.04284 | 0.02550 | 0.07445 | 0.03563 | 0.4878 | 0.04848 | 0.03016 | 0.05863 | 0.03335 |
| Comparison | | -6.66% | +7.19% | +12.31% | +10.99% | +7.80% | -10.31% | +18.93% | +18.66% | +8.34% | +11.34% | -9.42% | +50.56% | +46.84% | +18.37% | +27.53% |

| Dataset-Large | | Yelp | | | | | Gowalla | | | | | ML-10M | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| Model | Mixer | Infer. Time (s)↓ | HR@5↑ | NDCG@5↑ | HR@10↑ | NDCG@10↑ | Infer. Time (s)↓ | HR@5↑ | NDCG@5↑ | HR@10↑ | NDCG@10↑ | Infer. Time (s)↓ | HR@5↑ | NDCG@5↑ | HR@10↑ | NDCG@10↑ |
| GRU4Rec | RNN | 0.6634 | 0.01347 | 0.00820 | 0.02508 | 0.01190 | 0.7060 | 0.02022 | 0.01075 | 0.05012 | 0.02025 | 1.1730 | 0.10844 | 0.07214 | 0.17025 | 0.09523 |
| NextItNet | CNN | 0.8083 | 0.01503 | 0.00901 | 0.02737 | 0.01295 | 0.8100 | 0.05012 | 0.03022 | 0.08752 | 0.04216 | 1.6706 | 0.12769 | 0.08992 | 0.19066 | 0.10748 |
| SASRec | Trans. | 0.9810 | 0.01006 | 0.00602 | 0.01948 | 0.00904 | 0.9352 | 0.04254 | 0.02595 | 0.07904 | 0.03761 | 2.5037 | 0.11720 | 0.08830 | 0.17869 | 0.09905 |
| FMLP4Rec | Four. | 1.0085 | 0.00794 | 0.00411 | 0.01557 | 0.00657 | 0.9492 | 0.02799 | 0.0144 | 0.04834 | 0.02095 | 2.2933 | 0.06319 | 0.03677 | 0.09812 | 0.04799 |
| TriMLP | MLP | 0.5241 | 0.01554 | 0.00963 | 0.02784 | 0.01357 | 0.5815 | 0.06175 | 0.03968 | 0.09997 | 0.05194 | 0.8946 | 0.13900 | 0.09636 | 0.20273 | 0.11685 |
| Comparison | | -20.80% | +3.39% | +6.88% | +1.72% | +4.79% | -17.63% | +23.20% | +31.30% | +14.23% | +23.20% | -23.73% | +8.86% | +7.16% | +6.33% | +8.72% |

*Observation 1: Consistent superior recommendation performance.* NextItNet with temporal CNN architecture is the strongest baseline with decent scores on most datasets, while the performances of GRU4Rec, SASRec and FMLP4Rec drifts sharply. Notably, TriMLP achieves the state-of-the-art performances on all validated 12 datasets. Specifically, TriMLP is substantial ahead of the strongest baseline averagely by 15.57%, 18.23%, 18.35% and 11.66% cross 4 different scales of datasets in terms of HR and NDCG. It demonstrates that our method equips MLP with the ample sequential modeling ability under the same training manner, which is competitive to RNN, CNN, Transformer and Fourier transform.

*Observation 2: Incremental ascendancy in efficiency.* According to the metric of Infer. Time, TriMLP saves the inference cost by 2.08%, 2.99%, 8.80% and 20.72% respectively on Tiny, Small, Base and Large datasets compared to the fastest competitor. The reduction increases with the scale of datasets which is in line with our expectation. Since the number of interactions are minor in Tiny and Small datasets, the computation is more intensive on the Embedding part and Classifier. Along with the number of interactions increases, i.e., mixing token occupies the main part of the computation, the efficiency advantages of TriMLP show up. It is credited to the plain structure in Triangular Mixer, which only contains 2 matrix transposition and multiplications. We provide a more detailed case study on the largest ML-10M in Section V-F to reveal the advantage of TriMLP in computational complexity.

*Observation 3: Surprisingly good accuracy/efficiency trade-off.* Throughout all 12 validated datasets, TriMLP can averagely provide 14.88% higher recommendation performance against SOTA and reduce 8.64% inference time. The proposed TriMLP architecture reveals the promising potential to be served as an alternative for sequential recommenders.

## TABLE III
COMPUTATIONAL COMPLEXITY COMPARISON ON ML-10M.

| Dataset-Large | ML-10M | | | | | |
|---|---|---|---|---|---|---|
| Settings | batch size: 512, sequence length $n = 128$, dimension $d = 128$ and kernel size $k = 3$ | | | | | |
| Model | Complexity | MACs↓ | Par. Scale↓ | GPU Mem.↓ | Infer. Time↓ | HR@5↑ |
| GRU4Rec | $O(nd^2)$ | 13.04 G | 0.19 M | 2,536 MB | 1.1730 s | 0.10844 |
| NextItNet | $O(k \cdot nd^2)$ | 26.03 G | 0.40 M | 1,810 MB | 1.6706 s | 0.12769 |
| SASRec | $O(n^2d)$ | 25.94 G | 0.40 M | 1,819 MB | 2.5037 s | 0.11720 |
| FMLP4Rec | $O(nd \cdot \log(nd))$ | 17.18 G | 0.26 M | 2,093 MB | 2.2933 s | 0.06319 |
| TriMLP | $O(n^2d)$ | 2.15 G | 0.03 M | 1,391 MB | 0.8946 s | 0.13900 |
| Comparison | - | -83.51% | -84.21% | -23.15% | -23.73% | +8.86% |

### F. Computational Complexity Analysis on ML-10M

We count the Multiply–Accumulate Operations (MACs), Parameter Scale (Para. Scale) and GPU Memory Occupation (GPU Mem.) on the largest dataset ML-10M. Since all the compared methods share the common Embedding layer and Classifier, we only calculate MACs and Par. Scale of the encoder in each model, i.e., the RNN (or CNN) layer in GRU4Rec (or NextItNet) and the self-attention (or denoised Fourier) layer in SASRec (or FMLP4Rec).

The experimental results are summarized in Table III. Albeit the computational complexity is quadratic correlated to the input sequence length $n$, TriMLP possesses the higher efficiency due to parallel and minimalist matrix multiplication. Notably, TriMLP shrinks 83.51% MACs, 84.21% Para. Scale, 23.15% GPU Mem. and 23.73% Infer. Time.

### G. Macro-Design of Triangular Mixer (RQ 2)

Recall our vanilla implementation of TriMLP contains the complete Triangular Mixer, including both global and local mixing (without cross-session interactions), we derive the following 4 variants to carry on the ablation study:

- **EyeMLP** replaces Triangular Mixer with the identity matrix where tokens no longer interact with each other.

| Dataset-Tiny | Beauty | | | | | Sports | | | | | ML-100K | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| Model | HR@5↑ | NDCG@5↑ | HR@10↑ | NDCG@10↑ | Avg. Impv.↑ | HR@5↑ | NDCG@5↑ | HR@10↑ | NDCG@10↑ | Avg. Impv.↑ | HR@5↑ | NDCG@5↑ | HR@10↑ | NDCG@10↑ | Avg. Impv.↑ |
| EyeMLP | 0.08053 | 0.05942 | 0.09976 | 0.06561 | - | 0.01685 | 0.01118 | 0.02043 | 0.01229 | - | 0.05579 | 0.03624 | 0.09227 | 0.04768 | - |
| SqrMLP | 0.00421 | 0.00362 | 0.00781 | 0.00473 | -93.41% | 0.00051 | 0.00022 | 0.00153 | 0.00053 | -95.80% | 0.01180 | 0.00759 | 0.02253 | 0.01100 | -77.60% |
| TriMLP_G | 0.08233 | 0.06177 | 0.11599 | 0.07261 | +8.28% | 0.01583 | 0.01063 | 0.02298 | 0.01287 | +1.56% | 0.06545 | 0.03819 | 0.11803 | 0.05517 | +16.58% |
| TriMLP_L | 0.08173 | 0.06218 | 0.09916 | 0.06754 | +2.12% | 0.01733 | 0.01176 | 0.02331 | 0.01399 | +9.73% | 0.06009 | 0.03973 | 0.10622 | 0.05451 | +11.70% |
| TriMLP | **0.09615** | **0.07070** | **0.12560** | **0.08003** | **+21.57%** | **0.01839** | **0.01258** | **0.02451** | **0.01442** | **+14.74%** | **0.08691** | **0.05848** | **0.15451** | **0.07988** | **+63.03%** |

| Dataset-Small | NYC | | | | | QB-Article | | | | | TKY | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| Model | HR@5↑ | NDCG@5↑ | HR@10↑ | NDCG@10↑ | Avg. Impv.↑ | HR@5↑ | NDCG@5↑ | HR@10↑ | NDCG@10↑ | Avg. Impv.↑ | HR@5↑ | NDCG@5↑ | HR@10↑ | NDCG@10↑ | Avg. Impv.↑ |
| EyeMLP | 0.01940 | 0.01297 | 0.02813 | 0.01570 | - | 0.04389 | 0.02650 | 0.08885 | 0.04074 | - | 0.04801 | 0.02715 | 0.08255 | 0.03318 | - |
| SqrMLP | 0.00291 | 0.00128 | 0.00582 | 0.00221 | -85.09% | 0.01841 | 0.00885 | 0.05523 | 0.02058 | -53.00% | 0.00706 | 0.00431 | 0.01279 | 0.00612 | -83.87% |
| TriMLP_G | 0.02813 | 0.02030 | **0.06111** | 0.03065 | +78.49% | 0.07193 | 0.04421 | 0.13273 | 0.06356 | +59.03% | 0.04955 | 0.02722 | 0.08640 | 0.03364 | +2.38% |
| TriMLP_L | 0.04171 | 0.02337 | 0.05626 | 0.02797 | +93.33% | 0.07129 | 0.04363 | 0.13702 | 0.06456 | +59.94% | 0.05026 | 0.02924 | 0.08907 | 0.04163 | +11.48% |
| TriMLP | **0.04365** | **0.02574** | **0.06111** | **0.03121** | **+109.87%** | **0.07621** | **0.04751** | **0.13552** | **0.06639** | **+64.10%** | **0.05293** | **0.03175** | **0.09043** | **0.04384** | **+17.22%** |

| Dataset-Base | ML-1M | | | | | QB-Video | | | | | Brightkite | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| Model | HR@5↑ | NDCG@5↑ | HR@10↑ | NDCG@10↑ | Avg. Impv.↑ | HR@5↑ | NDCG@5↑ | HR@10↑ | NDCG@10↑ | Avg. Impv.↑ | HR@5↑ | NDCG@5↑ | HR@10↑ | NDCG@10↑ | Avg. Impv.↑ |
| EyeMLP | 0.12562 | 0.08582 | 0.18097 | 0.10360 | - | 0.02426 | 0.01491 | 0.04274 | 0.02079 | - | 0.04025 | 0.02106 | 0.06493 | 0.02901 | - |
| SqrMLP | 0.00530 | 0.00316 | 0.01210 | 0.00529 | -95.08% | 0.00546 | 0.00274 | 0.01402 | 0.00545 | -75.02% | 0.00665 | 0.00360 | 0.01120 | 0.00508 | -82.91% |
| TriMLP_G | 0.14418 | 0.09651 | 0.21064 | 0.11792 | +14.36% | 0.04127 | 0.02479 | 0.07397 | 0.03527 | +69.77% | 0.04498 | 0.02688 | 0.05828 | 0.03118 | +9.16% |
| TriMLP_L | 0.13656 | 0.09207 | 0.20865 | 0.11520 | +10.62% | 0.03670 | 0.02255 | 0.06673 | 0.03215 | +53.32% | 0.04760 | 0.02872 | 0.05775 | 0.03194 | +13.42% |
| TriMLP | **0.16390** | **0.11196** | **0.23434** | **0.13454** | **+30.07%** | **0.04284** | **0.02550** | **0.07445** | **0.03563** | **+73.30%** | **0.04848** | **0.03016** | **0.05863** | **0.03335** | **+17.23%** |

| Dataset-Large | Yelp | | | | | Gowalla | | | | | ML-10M | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| Model | HR@5↑ | NDCG@5↑ | HR@10↑ | NDCG@10↑ | Avg. Impv.↑ | HR@5↑ | NDCG@5↑ | HR@10↑ | NDCG@10↑ | Avg. Impv.↑ | HR@5↑ | NDCG@5↑ | HR@10↑ | NDCG@10↑ | Avg. Impv.↑ |
| EyeMLP | 0.01116 | 0.00687 | 0.01934 | 0.00948 | - | 0.04749 | 0.02822 | 0.08632 | 0.04303 | - | 0.10652 | 0.07313 | 0.15292 | 0.08804 | - |
| SqrMLP | 0.00210 | 0.00112 | 0.00469 | 0.00194 | -80.04% | 0.00487 | 0.00240 | 0.00909 | 0.00376 | -90.49% | 0.08097 | 0.05532 | 0.11866 | 0.06746 | -23.53% |
| TriMLP_G | 0.01491 | 0.00903 | **0.02812** | 0.01325 | +37.55% | 0.05629 | 0.03489 | 0.09726 | 0.04801 | +16.60% | 0.12317 | 0.08480 | 0.18006 | 0.10312 | +16.62% |
| TriMLP_L | 0.01241 | 0.00724 | 0.02053 | 0.01019 | +7.56% | 0.05809 | 0.03624 | 0.09822 | 0.04982 | +20.08% | 0.12609 | 0.08699 | 0.18347 | 0.10539 | +19.25% |
| TriMLP | **0.01554** | **0.00963** | 0.02784 | **0.01357** | **+41.63%** | **0.06175** | **0.03968** | **0.09997** | **0.05194** | **+26.79%** | **0.13900** | **0.09636** | **0.20273** | **0.11685** | **+31.89%** |

- **SqrMLP** replaces Triangular Mixer with two standard MLPs where cross-token interactions are no-limited.
- **TriMLP_G** only reserves the global mixing in Triangular Mixer to model the long-range dependencies.
- **TriMLP_L** only reserves the local mixing in Triangular Mixer to capture the local patterns.

We consider EyeMLP as baseline, and measure the effectiveness of different variants with the corresponding Average Improvement (Avg. Impv.) against baseline. According to Table IV, we have the following findings:

*Finding 1: Fully-connection profoundly impairs performance.* Compared to EyeMLP, SqrMLP erodes performances on all validated 12 datasets averagely by 77.99%. It unveils that the incompatibility between the fully-connection structure and the auto-regressive training fashion. The resulting information leakage is a serious and non-negligible issue which originally motivates this paper. We also explore the feasibility of adopting the auto-encoding manner to ingratiate the bidirectional particularity of MLP in section V-K.

*Finding 2: Triangular design does matter.* Both TriMLP_G and TriMLP_L remarkably boost the performance against EyeMLP, where the leading margins achieve up to 27.55% and 26.05%, singly. It demonstrates that our triangular design sufficiently evokes the sequential modeling potential in MLP under the auto-regressive training mode.

*Finding 3: Two mixing layers complement each other.* Triangular Mixer constantly attains superior performance than solely employing either global or local mixing. It shows that jointly utilizing these two mixing layers is productive to the fine-grained modeling of long and short-term preferences. We compare the performance of serial-connected mixing layers with other internal structures in Triangular Mixer in Section V-H1, and visualize how is the mutual influence among these two mixing branches in Section V-I2.

### H. Micro-Design of Triangular Mixer (RQ 3)

This section analyzes the intrinsic properties of Triangular Mixer by decomposing it into various internal structures and operation components.

*1) Sensitivity w.r.t. Different Internal Structures:* Recall that the vanilla Triangular Mixer is denoted as $\mathbf{Mix\_G} \rightarrow \mathbf{Mix\_L}$, we consider the following 3 variants:

- $\mathbf{Mix\_L} \rightarrow \mathbf{Mix\_G}$ follows the serial-connection. It first encodes the local patterns by local mixing and then models the long-range dependency with global mixing.
- $\mathbf{Mix\_G} + \mathbf{Mix\_L}$ employs the parallel structure that combines the results of independent global and local mixing branches with element-wise addition.
- $\mathbf{Mix\_G} \parallel \mathbf{Mix\_L}$ employs the parallel structure that concatenates the results of independent global and local mixing branches and merges them with linear layer.

As depicted in Figure 5, the serial combinations generally perform better than the parallel ones on most datasets, and $\text{Mix\_G} \rightarrow \text{Mix\_L}$ achieves more stable scores. Note that $\text{Mix\_L} \rightarrow \text{Mix\_G}$ has extremely poor performances on Tiny datasets. It might be caused by the split of historical sequences, where the inappropriate truncation leads to the non-uniform distribution of time intervals among user-item interactions. Thus, the short and long-term preferences might vary greatly, and the encoded local patterns would mislead the global one.

*2) Sensitivity w.r.t. Various Operation Components:* Recall that the vanilla Triangular Mixer drops the positional information and Feed-Forward Network, and utilizes the 1-0 initialization and Softmax normalization, we consider the following 4 alternatives:

- **w. PE** injects the absolute order information into the sequence representation by employing the positional embedding [6] after the Embedding layer.
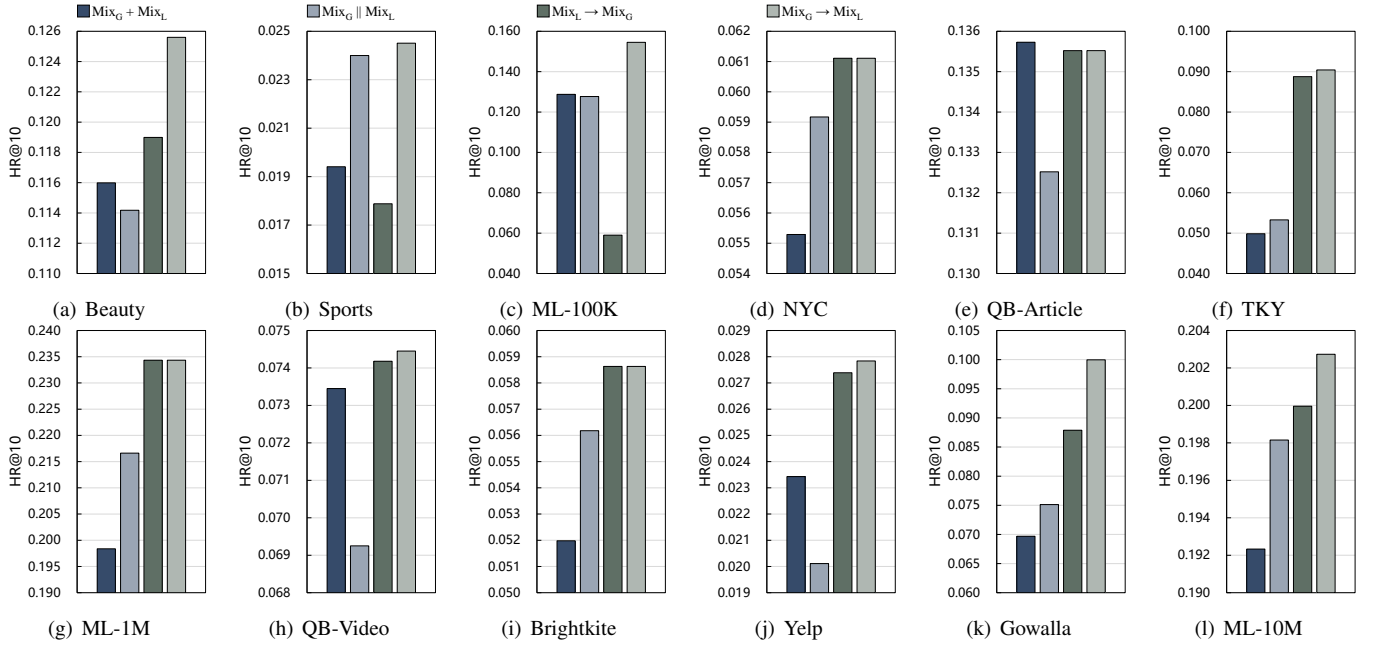
Fig. 5. Recommendation performance comparison (HR@10) of Triangular Mixer with different internal structures.

- **w. FFN** adds the Feed-Forward Network [6] after Triangular Mixer, along with the pre-layer-normalization [49], [50] and residual connection [51].
- **w.o. 1-0 Init** initializes the mixing kernels with the default kaiming uniform distribution [52] in PyTorch.
- **w.o. Softmax** removes the Softmax operation conducted on the mixing kernels.

According to the experimental results listed in Table V, we find the following properties of Triangular Mixer:

*Property 1: Positional embedding has conflicts with Triangular Mixer.* Since the mixing layers bust the symmetry of mixing kernels and explicitly endow cross-token interactions in chronological order, the extra positional information becomes redundant and dramatically damages the performances.

*Property 2: FFN brings limited profits in certain scenario.* Through all validated 12 datasets, adding FFN slightly works on ML-1M. Since FFN significantly increases the parameter-scale, it provokes all-MLP architectures to trap in the data-hungry issue on Tiny, Small and Base datasets.

*Property 3: 1-0 initialization proves to be helpful.* Compared to the uniform distributed initialization, our method make all the tokens contribute equally to the targets during the early training stage. It is conducive to avoid the local optimal, especially when the static parameters in MLP are more likely to be troubled with the under-fitting issue.

*Property 4: Softmax prominently promotes performances.* As reported in [10], [11], the weights might be irregular, messy and disorganized in standard MLPs. Conducting Softmax on the global and local mixing kernels, i.e., transforms the learnable parameter into probabilities, is instrumental in evolving the weights towards exhuming the relations among tokens.

*I. Sensitivity w.r.t. Hyper-parameter Setting*

We mainly verify the influence brought by setting different session number $s$ (or session length $l$) in local mixing, which decides the short-term preference modeling.

*1) Influence on Performance:* Recall that the historical sequence of length $n$ should be divided into $s$ non-overlapped sessions of length $l$ that $n = l \times s$, we set $(s, l) = \{(1, n), (2, n/2), ...(n, 1)\}$. Note that the local mixing works in the global manner with $(s, l) = (1, n)$, and degrades as the identical mapping with $(s, l) = (n, 1)$.

As shown in Figure 6, we find that independently modeling the short-term preferences in shorter sessions improves the performance on most datasets. On QB-Article, Brightkite and Yelp, we observe that modeling the dependency from the global perspective ($s = 1$) is more suitable, while Gowalla and TKY reveal the contrast situation ($s = n$). The possible reason also lies in the split of historical sequences (as explained in Section V-H1). Generally, TriMLP achieves the best performance when setting $s = 2$ on NYC, QB-Video, ML-10M, $s = 8$ on ML-1M, $s = 16$ on Beauty, Sports and $s = 32$ on ML-100K.

*2) Influence on Reception Field:* We visualize the weights of global and local kernels on ML-1M ($n = 128$) to explore how these two mixing layers complement each other. The corresponding 8 heat maps with different session numbers $s$ (or session lengths $l$) are plotted in Figure 8. Accordingly, we observe the following characteristics:

*Characteristic 1: Local kernels sustain more attention on the tokens around the current time step.* We observe that the weights in all local kernels share the similar distribution, that the elements nearing the diagonal have greater absolute values than others, i.e., more active. When setting $s = n = 128$

TABLE V

SENSITIVITY W.R.T. VARIOUS OPERATION COMPONENTS IN GLOBAL AND LOCAL MIXING. THE BEST PERFORMANCE IS BOLDFACED. THE VANILLA IMPLEMENTATION IS MARKED WITH PURPLE SHADING.

| Dataset-Tiny | | Beauty | | | | Sports | | | | ML-100K | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| Variant | Alternative | HR@5↑ | NDCG@5↑ | HR@10↑ | NDCG@10↑ | HR@5↑ | NDCG@5↑ | HR@10↑ | NDCG@10↑ | HR@5↑ | NDCG@5↑ | HR@10↑ | NDCG@10↑ |
| Layer | w. PE | 0.00060 | 0.00060 | 0.00601 | 0.00227 | 0.00102 | 0.00077 | 0.00664 | 0.00257 | 0.02253 | 0.01492 | 0.04614 | 0.02244 |
| | w. FFN | 0.09195 | 0.06939 | 0.12260 | 0.07910 | 0.01634 | 0.01057 | 0.01839 | 0.01118 | 0.06009 | 0.03815 | 0.11695 | 0.05597 |
| Normalization | w.o. 1-0 Init. | 0.08534 | 0.06288 | 0.11599 | 0.07282 | 0.01532 | 0.01022 | 0.02145 | 0.01211 | 0.08369 | 0.05541 | 0.13627 | 0.07189 |
| | w.o. Softmax | 0.00781 | 0.00444 | 0.01262 | 0.00603 | 0.00000 | 0.00000 | 0.00102 | 0.00030 | 0.01395 | 0.00726 | 0.03004 | 0.01245 |
| **Triangular Mixer** | | **0.09615** | **0.07070** | **0.12560** | **0.08003** | **0.01839** | **0.01258** | **0.02451** | **0.01442** | **0.08691** | **0.05848** | **0.15451** | **0.07988** |
| Dataset-Small | | NYC | | | | QB-Article | | | | TKY | | | |
| Variant | Alternative | HR@5↑ | NDCG@5↑ | HR@10↑ | NDCG@10↑ | HR@5↑ | NDCG@5↑ | HR@10↑ | NDCG@10↑ | HR@5↑ | NDCG@5↑ | HR@10↑ | NDCG@10↑ |
| Layer | w. PE | 0.02619 | 0.01469 | 0.05723 | 0.02445 | 0.03961 | 0.02387 | 0.08563 | 0.03858 | 0.02955 | 0.01819 | 0.04985 | 0.02464 |
| | w. FFN | 0.03977 | 0.02226 | 0.05529 | 0.02750 | 0.06193 | 0.03899 | 0.12096 | 0.05734 | **0.05382** | 0.03088 | **0.09263** | 0.04309 |
| Normalization | w.o. 1-0 Init. | 0.04268 | 0.02529 | 0.05723 | 0.02990 | 0.07193 | 0.04433 | **0.13766** | 0.06523 | 0.04014 | 0.02498 | 0.08690 | 0.04003 |
| | w.o. Softmax | 0.03298 | 0.02056 | 0.04850 | 0.02548 | 0.01905 | 0.00989 | 0.04624 | 0.01854 | 0.02955 | 0.01865 | 0.04940 | 0.02501 |
| **Triangular Mixer** | | **0.04365** | **0.02574** | **0.06111** | **0.03121** | **0.07621** | **0.04751** | 0.13552 | **0.06639** | 0.05293 | **0.03175** | 0.09043 | **0.04384** |
| Dataset-Base | | ML-1M | | | | QB-Video | | | | Brightkite | | | |
| Variant | Alternative | HR@5↑ | NDCG@5↑ | HR@10↑ | NDCG@10↑ | HR@5↑ | NDCG@5↑ | HR@10↑ | NDCG@10↑ | HR@5↑ | NDCG@5↑ | HR@10↑ | NDCG@10↑ |
| Layer | w. PE | 0.00829 | 0.00503 | 0.02121 | 0.00917 | 0.01885 | 0.01157 | 0.03171 | 0.01569 | 0.04550 | 0.02475 | 0.05880 | 0.02909 |
| | w. FFN | 0.16042 | 0.10813 | **0.24130** | 0.13408 | 0.03381 | 0.02029 | 0.06284 | 0.02955 | 0.03903 | 0.02114 | **0.06860** | 0.03066 |
| Normalization | w.o. 1-0 Init. | 0.16341 | 0.11194 | 0.23616 | **0.13538** | 0.04048 | 0.02427 | 0.07418 | 0.03505 | 0.04830 | 0.03009 | 0.05915 | 0.03349 |
| | w.o. Softmax | 0.00795 | 0.00394 | 0.01972 | 0.00770 | 0.00635 | 0.00357 | 0.01449 | 0.00616 | 0.04235 | **0.03522** | 0.05250 | **0.03848** |
| **Triangular Mixer** | | **0.16390** | **0.11196** | 0.23434 | 0.13454 | **0.04284** | **0.02550** | **0.07445** | **0.03563** | **0.04848** | 0.03016 | 0.05863 | 0.03335 |
| Dataset-Large | | Yelp | | | | Gowalla | | | | ML-10M | | | |
| Variant | Alternative | HR@5↑ | NDCG@5↑ | HR@10↑ | NDCG@10↑ | HR@5↑ | NDCG@5↑ | HR@10↑ | NDCG@10↑ | HR@5↑ | NDCG@5↑ | HR@10↑ | NDCG@10↑ |
| Layer | w. PE | 0.00221 | 0.00133 | 0.00370 | 0.00182 | 0.00786 | 0.00466 | 0.01372 | 0.00651 | 0.01553 | 0.00904 | 0.02970 | 0.00136 |
| | w. FFN | 0.01213 | 0.00765 | 0.02301 | 0.01112 | 0.06113 | 0.03838 | **0.10013** | 0.05090 | **0.14109** | **0.09749** | **0.20442** | **0.11786** |
| Normalization | w.o. 1-0 Init. | 0.01512 | 0.00915 | 0.02779 | 0.01319 | 0.06051 | 0.03879 | 0.09902 | 0.05115 | 0.13782 | 0.09557 | 0.20112 | 0.11593 |
| | w.o. Softmax | 0.00120 | 0.00076 | 0.00207 | 0.00104 | 0.00663 | 0.00393 | 0.01267 | 0.00588 | 0.00638 | 0.00346 | 0.01491 | 0.00618 |
| **Triangular Mixer** | | **0.01554** | **0.00963** | **0.02784** | **0.01357** | **0.06175** | **0.03968** | 0.09997 | **0.05194** | 0.13900 | 0.09636 | 0.20273 | 0.11685 |

(Figure 8 (h)), the local kernel completely degrades into the identity matrix where all remaining diagonal elements are 1.
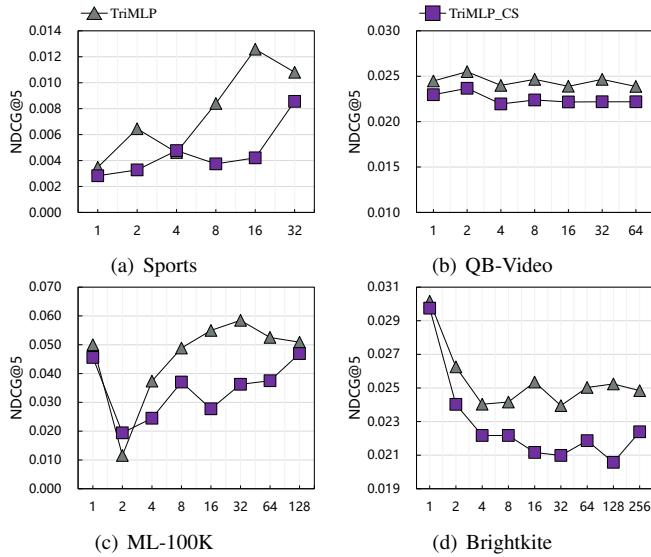


Fig. 7. NDCG@5 comparison between w./w.o. cross-session interactions. The axes of all sub-figures stand for the variable session number $s$.

*Characteristic 2: Shorter sessions encourage the global kernel attaching more importance to previous tokens.* Compared to Figure 8 (a), the upper-right elements in global kernels have greater absolute values (as Figure 8 (d)-(g)), which are responsible for the long-term user-item interactions. It indicates that the shorter sessions call for the larger reception field in global kernels to support the sufficient sequential dependency modeling. Unfortunately, such pattern seems to lose efficacy with $s = \{2, 4, 128\}$ (as Figure 8 (b, c, h)), i.e., global kernels are no longer sharp to long-range tokens, which accord to the inferior performances in Figure 6 (g).

*Characteristic 3: Suitable session settings produce superior performances.* As the experimental results in Figure 6 (g), TriMLP achieves comparable and preferred scores when setting $s = \{8, 32\}$ on ML-1M. Combined with the corresponding heatmaps in Figure 8 (d, f), besides the adequate short-term patterns offered by the local kernels, both of the global kernels in these two cases are more perceptive to the previous tokens. It proves the effectiveness of the serial structure in Triangular Mixer, that the global mixing layer and local mixing layer mutually assit each other indeed to realize the fine-grained modeling of sequential dependency.

### J. Cross-session Communications in Local Mixing

Recall another variant of local mixing that endows the cross-session communications (Eq. 7), denoted as TriMLP_CS, we verify the performance on 4 datasets Sports, QB-Video, ML-100K and Brightkite, where the maximum sequence length $n$ ranges from 32 to 256. TriMLP and TriMLP_CS share the same local mixing layer when $s = \{1, n\}$.

Figure 7 reports the NDCG@5 scores correlated with different session number $s$. It shows that connection sessions decreases the performance. The possible reasons are twofold. On the one hand, cross-session connections bring more previous information to the current overlapped token after
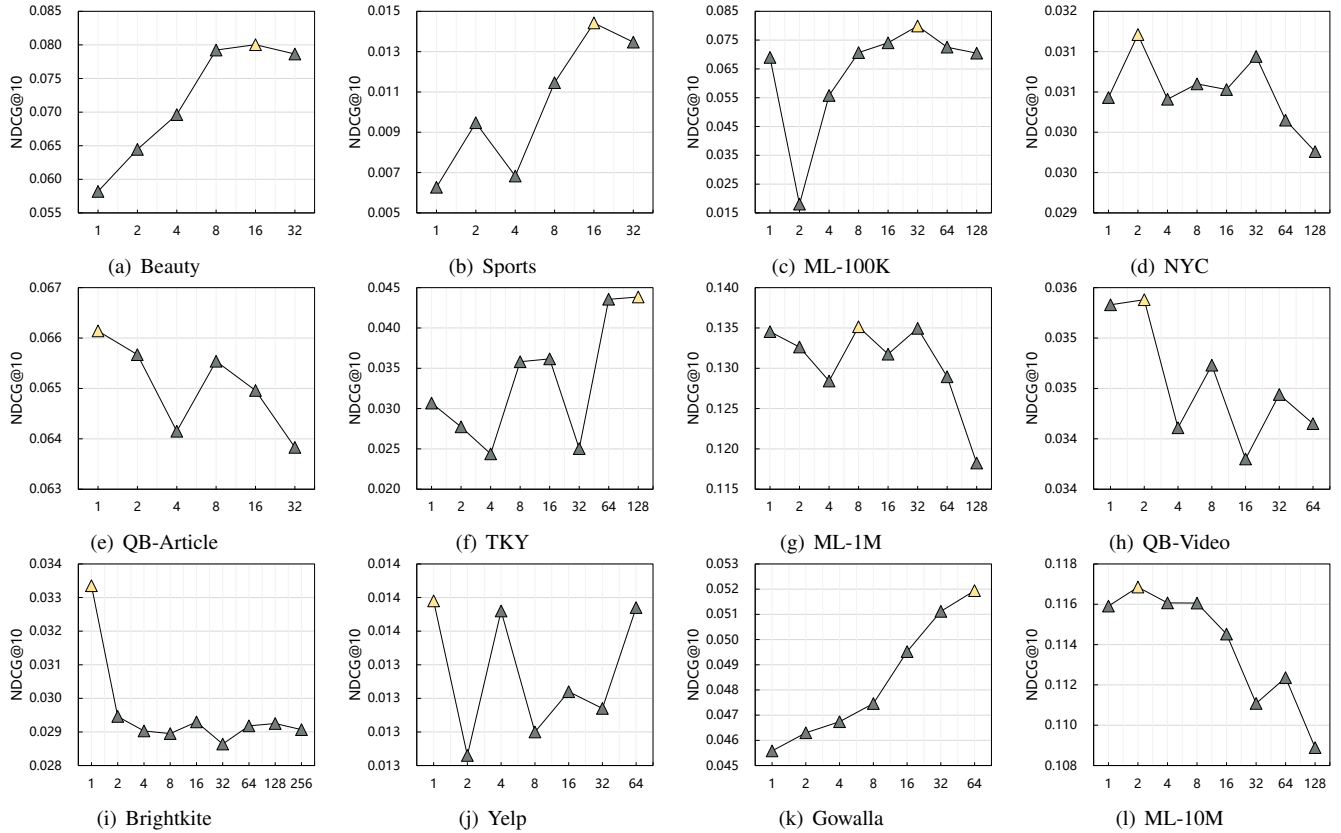
Fig. 6. Recommendation performance comparison (NDCG@10) w.r.t. different session number $s$. The horizontal axes of all sub-figures are the variable $s$.

multiple iterations, and degrade the local mixing as the global scheme, which disables the short-term preference modeling. On the other hand, since such connections are build solely upon the overlapped tokens, the corresponding neurons in the mixing kernel might accumulate more errors especially when the static weights are agnostic to the input sequence, which leads to the biased local pattern.

### K. Auto-regressive V.S. Auto-encoding

Auto-encoding is the other popular training fashion, represented by [8], [23], that utilizes the past and future tokens to predict the current one. To verify whether or not the auto-encoding is compatible with MLP, we derive another variant BiMLP which stacks 2 MLP layers as the basic architecture, and train it on 4 datasets of different scales `Beauty`, `NYC`, `ML-1M` and `Yelp` under the auto-encoding manner. Specifically, we randomly mask the tokens in historical sequences with setting mask ratio $r = \{0.1, 0.2, 0.3, 0.4, 0.5, 0.6, 0.7, 0.8, 0.9\}$. Among these masked tokens, 80% are padded with `mask token`, 10% are replaced with other tokens and 10% reserve the original ones.

Figure 9 reports the HR@5 scores correlated with different mask ratio. We find that the recommendation performance of BiMLP stands far behind TriMLP. The possible reason lies in that the auto-encoding training mode is more inclined to suffer from the data-hungry issue. It is more compatible with dense datasets for learning better mask representations, while se-

quential recommendation datasets are always extremely sparse (sparsity is usually larger than 99%). Moreover, this observation is also in line with [13], [29], [53], [54], that unidirectional models offer better results than bidirectional ones. Thus, we exploit the triangular design and put forward TriMLP under the unidirectional auto-regressive training scheme.
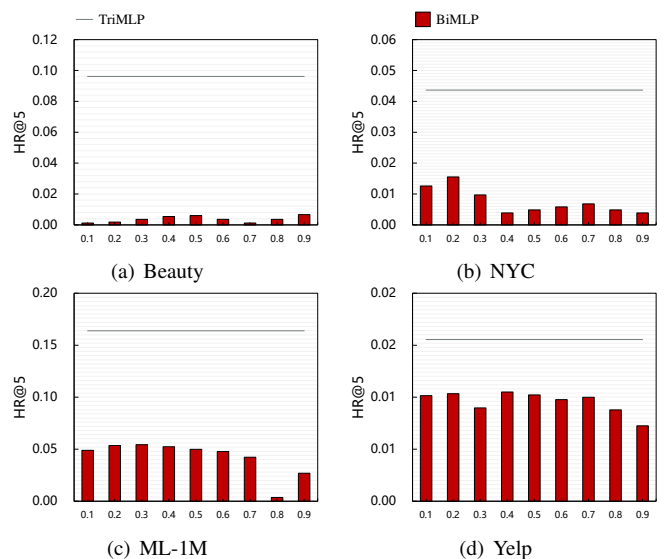


Fig. 9. HR@5 comparison between auto-regressive and auto-encoding. The axes of all sub-figures denote the variable mask ratio $r$.
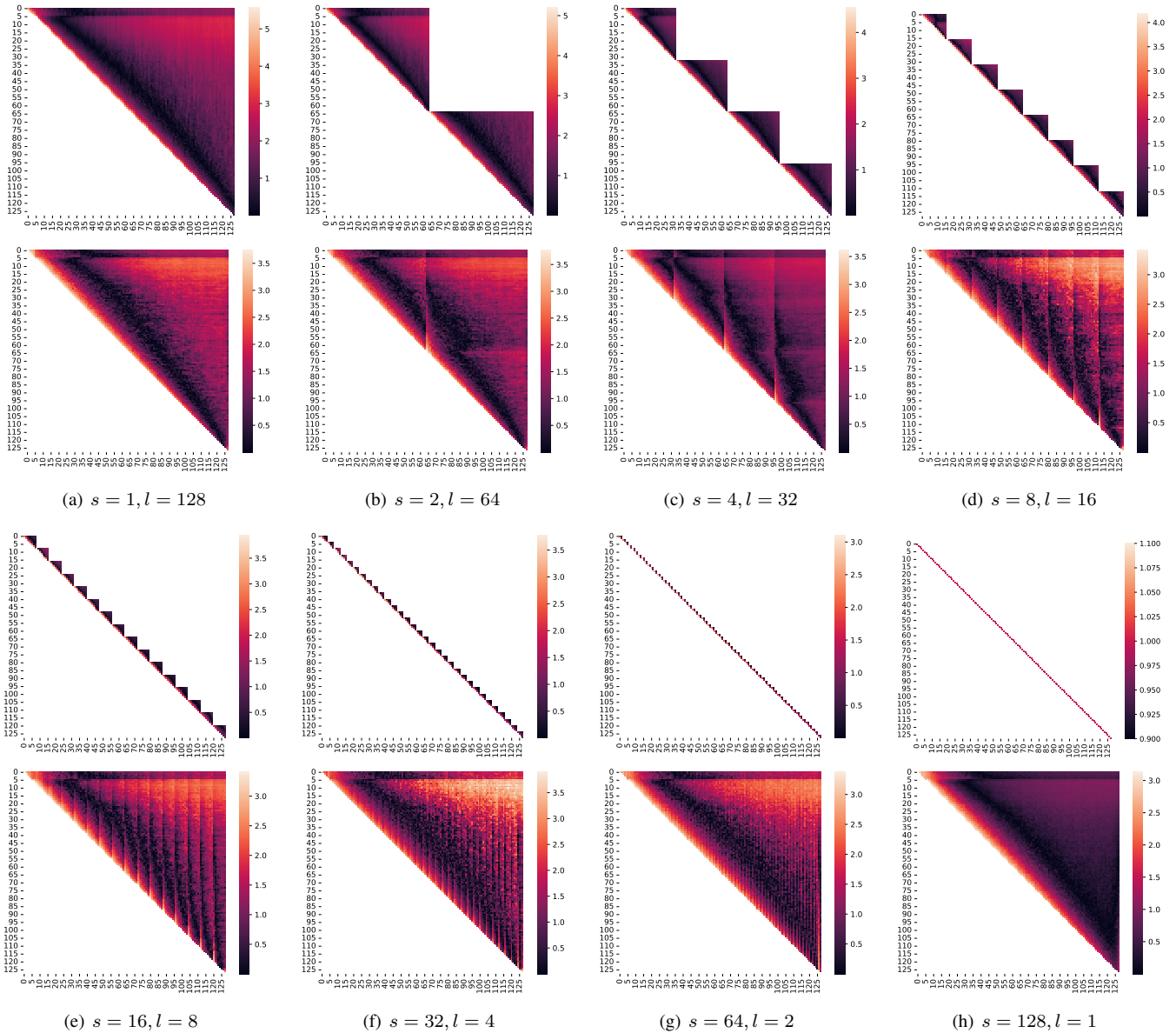
Fig. 8. Visualizing the weights of global and local mixing kernels with different session number $s$ on ML-1M. The global and local kernel are separately plotted in the upper and lower part of each sub-figure. Black indicates that the weight is 0, and the brighter, the greater the weight's absolute value.

## VI. CONCLUSION

In this paper, we make our exploration to study the capacity of MLP in sequential recommendation. We present the MLP-like sequential recommender TriMLP with a novel Triangular Mixer. Credited to the chronological cross-token communication and the serial mixing structure in Triangular Mixer, TriMLP successfully realizes the fine-grained modeling of sequential dependency. The experimental results on 12 datasets demonstrate that TriMLP attains stable, competitive and even better performance than several state-of-the-art baselines under the essential auto-regressive training mode with prominent less inference time, which well performs the "Revenge of MLP in Sequential Recommendation".

In the future, we will further improve TriMLP by introducing auxiliary information like temporal factors and item attributes, data augmentation and pre-training techniques. Moreover, it is intriguing to deliberate how to decouple the strong correlation between the sequence length and MLP shape, which will enable MLP to flexibly handle sequences of different lengths.

## REFERENCES

[1] H. Zhang, E. Yuan, W. Guo, Z. He, J. Qin, H. Guo, B. Chen, X. Li, and R. Tang, "Disentangling past-future modeling in sequential recommendation via dual networks," in *Proceedings of the 31st ACM International Conference on Information & Knowledge Management, Atlanta, GA, USA, October 17-21, 2022*, M. A. Hasan and L. Xiong, Eds. ACM, 2022, pp. 2549–2558. [Online]. Available: https://doi.org/10.1145/3511808.3557289

[2] B. Hidasi, A. Karatzoglou, L. Baltrunas, and D. Tikk, "Session-based recommendations with recurrent neural networks," in *4th International Conference on Learning Representations, ICLR 2016, San Juan, Puerto Rico, May 2-4, 2016, Conference Track Proceedings*,

Y. Bengio and Y. LeCun, Eds., 2016. [Online]. Available: http://arxiv.org/abs/1511.06939

[3] D. Yang, B. Fankhauser, P. Rosso, and P. Cudré-Mauroux, "Location prediction over sparse user mobility traces using rnns: Flashback in hidden states!" in *Proceedings of the Twenty-Ninth International Joint Conference on Artificial Intelligence, IJCAI 2020*, C. Bessiere, Ed. ijcai.org, 2020, pp. 2184–2190. [Online]. Available: https://doi.org/10.24963/ijcai.2020/302

[4] F. Yuan, A. Karatzoglou, I. Arapakis, J. M. Jose, and X. He, "A simple convolutional generative network for next item recommendation," in *Proceedings of the Twelfth ACM International Conference on Web Search and Data Mining, WSDM 2019, Melbourne, VIC, Australia, February 11-15, 2019*, J. S. Culpepper, A. Moffat, P. N. Bennett, and K. Lerman, Eds. ACM, 2019, pp. 582–590. [Online]. Available: https://doi.org/10.1145/3289600.3290975

[5] J. Tang and K. Wang, "Personalized top-n sequential recommendation via convolutional sequence embedding," in *Proceedings of the Eleventh ACM International Conference on Web Search and Data Mining, WSDM 2018, Marina Del Rey, CA, USA, February 5-9, 2018*, Y. Chang, C. Zhai, Y. Liu, and Y. Maarek, Eds. ACM, 2018, pp. 565–573. [Online]. Available: https://doi.org/10.1145/3159652.3159656

[6] A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. N. Gomez, L. Kaiser, and I. Polosukhin, "Attention is all you need," in *Advances in Neural Information Processing Systems 30: Annual Conference on Neural Information Processing Systems 2017, December 4-9, 2017, Long Beach, CA, USA*, I. Guyon, U. von Luxburg, S. Bengio, H. M. Wallach, R. Fergus, S. V. N. Vishwanathan, and R. Garnett, Eds., 2017, pp. 5998–6008. [Online]. Available: https://proceedings.neurips.cc/paper/2017/hash/3f5ee243547dee91fbd053c1c4a845aa-Abstract.html

[7] W. Kang and J. J. McAuley, "Self-attentive sequential recommendation," in *IEEE International Conference on Data Mining, ICDM 2018, Singapore, November 17-20, 2018*. IEEE Computer Society, 2018, pp. 197–206. [Online]. Available: https://doi.org/10.1109/ICDM.2018.00035

[8] F. Sun, J. Liu, J. Wu, C. Pei, X. Lin, W. Ou, and P. Jiang, "Bert4rec: Sequential recommendation with bidirectional encoder representations from transformer," in *Proceedings of the 28th ACM International Conference on Information and Knowledge Management, CIKM 2019, Beijing, China, November 3-7, 2019*, W. Zhu, D. Tao, X. Cheng, P. Cui, E. A. Rundensteiner, D. Carmel, Q. He, and J. X. Yu, Eds. ACM, 2019, pp. 1441–1450. [Online]. Available: https://doi.org/10.1145/3357384.3357895

[9] E. Wang, Y. Jiang, Y. Xu, L. Wang, and Y. Yang, "Spatial-temporal interval aware sequential POI recommendation," in *38th IEEE International Conference on Data Engineering, ICDE 2022, Kuala Lumpur, Malaysia, May 9-12, 2022*. IEEE, 2022, pp. 2086–2098. [Online]. Available: https://doi.org/10.1109/ICDE53745.2022.00202

[10] I. O. Tolstikhin, N. Houlsby, A. Kolesnikov, L. Beyer, X. Zhai, T. Unterthiner, J. Yung, A. Steiner, D. Keysers, J. Uszkoreit, M. Lucic, and A. Dosovitskiy, "Mlp-mixer: An all-mlp architecture for vision," in *Advances in Neural Information Processing Systems 34: Annual Conference on Neural Information Processing Systems 2021, NeurIPS 2021, December 6-14, 2021, virtual*, M. Ranzato, A. Beygelzimer, Y. N. Dauphin, P. Liang, and J. W. Vaughan, Eds., 2021, pp. 24261–24272. [Online]. Available: https://proceedings.neurips.cc/paper/2021/hash/cba0a4ee5ccd02fda0fe3f9a3e7b89fe-Abstract.html

[11] H. Touvron, P. Bojanowski, M. Caron, M. Cord, A. El-Nouby, E. Grave, A. Joulin, G. Synnaeve, J. Verbeek, and H. Jégou, "Resmlp: Feedforward networks for image classification with data-efficient training," *CoRR*, vol. abs/2105.03404, 2021. [Online]. Available: https://arxiv.org/abs/2105.03404

[12] R. Liu, Y. Li, L. Tao, D. Liang, and H. Zheng, "Are we ready for a new paradigm shift? A survey on visual deep MLP," *Patterns*, vol. 3, no. 7, p. 100520, 2022. [Online]. Available: https://doi.org/10.1016/j.patter.2022.100520

[13] G. Yuan, F. Yuan, Y. Li, B. Kong, S. Li, L. Chen, M. Yang, C. YU, B. Hu, Z. Li, Y. Xu, and X. Qie, "Tenrec: A large-scale multipurpose benchmark dataset for recommender systems," in *Thirty-sixth Conference on Neural Information Processing Systems Datasets and Benchmarks Track*, 2022. [Online]. Available: https://openreview.net/forum?id=PfuW84q25y9

[14] Q. Hou, Z. Jiang, L. Yuan, M. Cheng, S. Yan, and J. Feng, "Vision permutator: A permutable mlp-like architecture for visual recognition," *CoRR*, vol. abs/2106.12368, 2021. [Online]. Available: https://arxiv.org/abs/2106.12368

[15] H. Liu, Z. Dai, D. R. So, and Q. V. Le, "Pay attention to mlps," in *Advances in Neural Information Processing Systems 34: Annual Conference on Neural Information Processing Systems 2021, NeurIPS 2021, December 6-14, 2021, virtual*, M. Ranzato, A. Beygelzimer, Y. N. Dauphin, P. Liang, and J. W. Vaughan, Eds., 2021, pp. 9204–9215. [Online]. Available: https://proceedings.neurips.cc/paper/2021/hash/4cc05b35c2f937c5bd9e7d41d3686fff-Abstract.html

[16] Y. Tang, K. Han, J. Guo, C. Xu, Y. Li, C. Xu, and Y. Wang, "An image patch is a wave: Phase-aware vision MLP," in *IEEE/CVF Conference on Computer Vision and Pattern Recognition, CVPR 2022, New Orleans, LA, USA, June 18-24, 2022*. IEEE, 2022, pp. 10925–10934. [Online]. Available: https://doi.org/10.1109/CVPR52688.2022.01066

[17] C. Tang, Y. Zhao, G. Wang, C. Luo, W. Xie, and W. Zeng, "Sparse MLP for image recognition: Is self-attention really necessary?" in *Thirty-Sixth AAAI Conference on Artificial Intelligence, AAAI 2022, Thirty-Fourth Conference on Innovative Applications of Artificial Intelligence, IAAI 2022, The Twelveth Symposium on Educational Advances in Artificial Intelligence, EAAI 2022 Virtual Event, February 22 - March 1, 2022*. AAAI Press, 2022, pp. 2344–2351. [Online]. Available: https://ojs.aaai.org/index.php/AAAI/article/view/20133

[18] D. Lian, Z. Yu, X. Sun, and S. Gao, "AS-MLP: an axial shifted MLP architecture for vision," in *The Tenth International Conference on Learning Representations, ICLR 2022, Virtual Event, April 25-29, 2022*. OpenReview.net, 2022. [Online]. Available: https://openreview.net/forum?id=fvLLcIYmXb

[19] J. Guo, Y. Tang, K. Han, X. Chen, H. Wu, C. Xu, C. Xu, and Y. Wang, "Hire-mlp: Vision MLP via hierarchical rearrangement," in *IEEE/CVF Conference on Computer Vision and Pattern Recognition, CVPR 2022, New Orleans, LA, USA, June 18-24, 2022*. IEEE, 2022, pp. 816–826. [Online]. Available: https://doi.org/10.1109/CVPR52688.2022.00090

[20] G. Wei, Z. Zhang, Z. Lan, Y. Lu, and Z. Chen, "Activemlp: An mlp-like architecture with active token mixer," *CoRR*, vol. abs/2203.06108, 2022. [Online]. Available: https://doi.org/10.48550/arXiv.2203.06108

[21] T. Yu, X. Li, Y. Cai, M. Sun, and P. Li, "$S^2$-mlp: Spatial-shift MLP architecture for vision," in *IEEE/CVF Winter Conference on Applications of Computer Vision, WACV 2022, Waikoloa, HI, USA, January 3-8, 2022*. IEEE, 2022, pp. 3615–3624. [Online]. Available: https://doi.org/10.1109/WACV51458.2022.00367

[22] ——, "$S^2$-mlpv2: Improved spatial-shift MLP architecture for vision," *CoRR*, vol. abs/2108.01072, 2021. [Online]. Available: https://arxiv.org/abs/2108.01072

[23] J. Devlin, M. Chang, K. Lee, and K. Toutanova, "BERT: pre-training of deep bidirectional transformers for language understanding," in *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, NAACL-HLT 2019, Minneapolis, MN, USA, June 2-7, 2019, Volume 1 (Long and Short Papers)*, J. Burstein, C. Doran, and T. Solorio, Eds. Association for Computational Linguistics, 2019, pp. 4171–4186. [Online]. Available: https://doi.org/10.18653/v1/n19-1423

[24] X. Ding, H. Chen, X. Zhang, J. Han, and G. Ding, "Repmlpnet: Hierarchical vision MLP with re-parameterized locality," in *IEEE/CVF Conference on Computer Vision and Pattern Recognition, CVPR 2022, New Orleans, LA, USA, June 18-24, 2022*. IEEE, 2022, pp. 568–577. [Online]. Available: https://doi.org/10.1109/CVPR52688.2022.00066

[25] X. He, L. Liao, H. Zhang, L. Nie, X. Hu, and T. Chua, "Neural collaborative filtering," in *Proceedings of the 26th International Conference on World Wide Web, WWW 2017, Perth, Australia, April 3-7, 2017*, R. Barrett, R. Cummings, E. Agichtein, and E. Gabrilovich, Eds. ACM, 2017, pp. 173–182. [Online]. Available: https://doi.org/10.1145/3038912.3052569

[26] H. Guo, R. Tang, Y. Ye, Z. Li, and X. He, "Deepfm: A factorization-machine based neural network for CTR prediction," in *Proceedings of the Twenty-Sixth International Joint Conference on Artificial Intelligence, IJCAI 2017, Melbourne, Australia, August 19-25, 2017*, C. Sierra, Ed. ijcai.org, 2017, pp. 1725–1731. [Online]. Available: https://doi.org/10.24963/ijcai.2017/239

[27] G. Zhou, X. Zhu, C. Song, Y. Fan, H. Zhu, X. Ma, Y. Yan, J. Jin, H. Li, and K. Gai, "Deep interest network for click-through rate prediction," in *Proceedings of the 24th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining, KDD 2018, London, UK, August 19-23, 2018*, Y. Guo and F. Farooq, Eds. ACM, 2018, pp. 1059–1068. [Online]. Available: https://doi.org/10.1145/3219819.3219823

[28] G. Zhou, N. Mou, Y. Fan, Q. Pi, W. Bian, C. Zhou, X. Zhu, and K. Gai, "Deep interest evolution network for click-through rate prediction,"

in *The Thirty-Third AAAI Conference on Artificial Intelligence, AAAI 2019, The Thirty-First Innovative Applications of Artificial Intelligence Conference, IAAI 2019, The Ninth AAAI Symposium on Educational Advances in Artificial Intelligence, EAAI 2019, Honolulu, Hawaii, USA, January 27 - February 1, 2019.* AAAI Press, 2019, pp. 5941–5948. [Online]. Available: https://doi.org/10.1609/aaai.v33i01.33015941

[29] K. Zhou, H. Yu, W. X. Zhao, and J. Wen, "Filter-enhanced MLP is all you need for sequential recommendation," in *WWW '22: The ACM Web Conference 2022, Virtual Event, Lyon, France, April 25 - 29, 2022*, F. Laforest, R. Troncy, E. Simperl, D. Agarwal, A. Gionis, I. Herman, and L. Médini, Eds. ACM, 2022, pp. 2388–2399. [Online]. Available: https://doi.org/10.1145/3485447.3512111

[30] M. Li, X. Zhao, C. Lyu, M. Zhao, R. Wu, and R. Guo, "Mlp4rec: A pure MLP architecture for sequential recommendations," in *Proceedings of the Thirty-First International Joint Conference on Artificial Intelligence, IJCAI 2022, Vienna, Austria, 23-29 July 2022*, L. D. Raedt, Ed. ijcai.org, 2022, pp. 2138–2144. [Online]. Available: https://doi.org/10.24963/ijcai.2022/297

[31] S. Wu, Y. Tang, Y. Zhu, L. Wang, X. Xie, and T. Tan, "Session-based recommendation with graph neural networks," in *The Thirty-Third AAAI Conference on Artificial Intelligence, AAAI 2019, The Thirty-First Innovative Applications of Artificial Intelligence Conference, IAAI 2019, The Ninth AAAI Symposium on Educational Advances in Artificial Intelligence, EAAI 2019, Honolulu, Hawaii, USA, January 27 - February 1, 2019.* AAAI Press, 2019, pp. 346–353. [Online]. Available: https://doi.org/10.1609/aaai.v33i01.3301346

[32] Y. Liu, S. Yang, Y. Xu, C. Miao, M. Wu, and J. Zhang, "Contextualized graph attention network for recommendation with item knowledge graph," *IEEE Trans. Knowl. Data Eng.*, vol. 35, no. 1, pp. 181–195, 2023. [Online]. Available: https://doi.org/10.1109/TKDE.2021.3082948

[33] C. Li, X. Niu, X. Luo, Z. Chen, and C. Quan, "A review-driven neural model for sequential recommendation," in *Proceedings of the Twenty-Eighth International Joint Conference on Artificial Intelligence, IJCAI 2019, Macao, China, August 10-16, 2019*, S. Kraus, Ed. ijcai.org, 2019, pp. 2866–2872. [Online]. Available: https://doi.org/10.24963/ijcai.2019/397

[34] Z. Wang, J. Zhang, H. Xu, X. Chen, Y. Zhang, W. X. Zhao, and J. Wen, "Counterfactual data-augmented sequential recommendation," in *SIGIR '21: The 44th International ACM SIGIR Conference on Research and Development in Information Retrieval, Virtual Event, Canada, July 11-15, 2021*, F. Diaz, C. Shah, T. Suel, P. Castells, R. Jones, and T. Sakai, Eds. ACM, 2021, pp. 347–356. [Online]. Available: https://doi.org/10.1145/3404835.3462855

[35] S. Zhang, D. Yao, Z. Zhao, T. Chua, and F. Wu, "Causerec: Counterfactual user sequence synthesis for sequential recommendation," in *SIGIR '21: The 44th International ACM SIGIR Conference on Research and Development in Information Retrieval, Virtual Event, Canada, July 11-15, 2021*, F. Diaz, C. Shah, T. Suel, P. Castells, R. Jones, and T. Sakai, Eds. ACM, 2021, pp. 367–377. [Online]. Available: https://doi.org/10.1145/3404835.3462908

[36] S. Bian, W. X. Zhao, K. Zhou, J. Cai, Y. He, C. Yin, and J. Wen, "Contrastive curriculum learning for sequential user behavior modeling via data augmentation," in *CIKM '21: The 30th ACM International Conference on Information and Knowledge Management, Virtual Event, Queensland, Australia, November 1 - 5, 2021*, G. Demartini, G. Zuccon, J. S. Culpepper, Z. Huang, and H. Tong, Eds. ACM, 2021, pp. 3737–3746. [Online]. Available: https://doi.org/10.1145/3459637.3481905

[37] S. Bian, W. X. Zhao, K. Zhou, X. Chen, J. Cai, Y. He, X. Luo, and J. Wen, "A novel macro-micro fusion network for user representation learning on mobile apps," in *WWW '21: The Web Conference 2021, Virtual Event / Ljubljana, Slovenia, April 19-23, 2021*, J. Leskovec, M. Grobelnik, M. Najork, J. Tang, and L. Zia, Eds. ACM / IW3C2, 2021, pp. 3199–3209. [Online]. Available: https://doi.org/10.1145/3442381.3450109

[38] K. Zhou, H. Wang, W. X. Zhao, Y. Zhu, S. Wang, F. Zhang, Z. Wang, and J. Wen, "S3-rec: Self-supervised learning for sequential recommendation with mutual information maximization," in *CIKM '20: The 29th ACM International Conference on Information and Knowledge Management, Virtual Event, Ireland, October 19-23, 2020*, M. d'Aquin, S. Dietze, C. Hauff, E. Curry, and P. Cudré-Mauroux, Eds. ACM, 2020, pp. 1893–1902. [Online]. Available: https://doi.org/10.1145/3340531.3411954

[39] Y. Tatsunami and M. Taki, "Raftmlp: Do mlp-based models dream

[40] Z. Wang, W. Jiang, Y. Zhu, L. Yuan, Y. Song, and W. Liu, "Dynamixer: A vision MLP architecture with dynamic mixing," in *International Conference on Machine Learning, ICML 2022, 17-23 July 2022, Baltimore, Maryland, USA*, ser. Proceedings of Machine Learning Research, K. Chaudhuri, S. Jegelka, L. Song, C. Szepesvári, G. Niu, and S. Sabato, Eds., vol. 162. PMLR, 2022, pp. 22 691–22 701. [Online]. Available: https://proceedings.mlr.press/v162/wang22i.html

[41] S. Chen, E. Xie, C. Ge, R. Chen, D. Liang, and P. Luo, "Cyclemlp: A mlp-like architecture for dense prediction," in *The Tenth International Conference on Learning Representations, ICLR 2022, Virtual Event, April 25-29, 2022.* OpenReview.net, 2022. [Online]. Available: https://openreview.net/forum?id=NMEceG4v69Y

[42] H. Zheng, P. He, W. Chen, and M. Zhou, "Mixing and shifting: Exploiting global and local dependencies in vision mlps," *CoRR*, vol. abs/2202.06510, 2022. [Online]. Available: https://arxiv.org/abs/2202.06510

[43] L. Melas-Kyriazi, "Do you even need attention? A stack of feed-forward layers does surprisingly well on imagenet," *CoRR*, vol. abs/2105.02723, 2021. [Online]. Available: https://arxiv.org/abs/2105.02723

[44] J. J. McAuley and J. Leskovec, "Hidden factors and hidden topics: understanding rating dimensions with review text," in *Seventh ACM Conference on Recommender Systems, RecSys '13, Hong Kong, China, October 12-16, 2013*, Q. Yang, I. King, Q. Li, P. Pu, and G. Karypis, Eds. ACM, 2013, pp. 165–172. [Online]. Available: https://doi.org/10.1145/2507157.2507163

[45] F. M. Harper and J. A. Konstan, "The movielens datasets: History and context," *ACM Trans. Interact. Intell. Syst.*, vol. 5, no. 4, pp. 19:1–19:19, 2016. [Online]. Available: https://doi.org/10.1145/2827872

[46] M. Weimer, A. Karatzoglou, Q. V. Le, and A. J. Smola, "COFI RANK - maximum margin matrix factorization for collaborative ranking," in *Advances in Neural Information Processing Systems 20, Proceedings of the Twenty-First Annual Conference on Neural Information Processing Systems, Vancouver, British Columbia, Canada, December 3-6, 2007*, J. C. Platt, D. Koller, Y. Singer, and S. T. Roweis, Eds. Curran Associates, Inc., 2007, pp. 1593–1600. [Online]. Available: https://proceedings.neurips.cc/paper/2007/hash/f76a89f0cb91bc419542ce9fa43902dc-Abstract.html

[47] W. Krichene and S. Rendle, "On sampled metrics for item recommendation," *Commun. ACM*, vol. 65, no. 7, pp. 75–83, 2022. [Online]. Available: https://doi.org/10.1145/3535335

[48] D. P. Kingma and J. Ba, "Adam: A method for stochastic optimization," in *3rd International Conference on Learning Representations, ICLR 2015, San Diego, CA, USA, May 7-9, 2015, Conference Track Proceedings*, Y. Bengio and Y. LeCun, Eds., 2015. [Online]. Available: http://arxiv.org/abs/1412.6980

[49] L. J. Ba, J. R. Kiros, and G. E. Hinton, "Layer normalization," *CoRR*, vol. abs/1607.06450, 2016. [Online]. Available: http://arxiv.org/abs/1607.06450

[50] Q. Wang, B. Li, T. Xiao, J. Zhu, C. Li, D. F. Wong, and L. S. Chao, "Learning deep transformer models for machine translation," in *Proceedings of the 57th Conference of the Association for Computational Linguistics, ACL 2019, Florence, Italy, July 28- August 2, 2019, Volume 1: Long Papers*, A. Korhonen, D. R. Traum, and L. Màrquez, Eds. Association for Computational Linguistics, 2019, pp. 1810–1822. [Online]. Available: https://doi.org/10.18653/v1/p19-1176

[51] K. He, X. Zhang, S. Ren, and J. Sun, "Deep residual learning for image recognition," in *2016 IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2016, Las Vegas, NV, USA, June 27-30, 2016.* IEEE Computer Society, 2016, pp. 770–778. [Online]. Available: https://doi.org/10.1109/CVPR.2016.90

[52] ——, "Delving deep into rectifiers: Surpassing human-level performance on imagenet classification," in *2015 IEEE International Conference on Computer Vision, ICCV 2015, Santiago, Chile, December 7-13, 2015.* IEEE Computer Society, 2015, pp. 1026–1034. [Online]. Available: https://doi.org/10.1109/ICCV.2015.123

[53] A. Dallmann, D. Zoller, and A. Hotho, "A case study on sampling strategies for evaluating neural sequential item recommendation models," in *RecSys '21: Fifteenth ACM Conference on Recommender Systems, Amsterdam, The Netherlands, 27 September 2021 - 1 October 2021*, H. J. C. Pampín, M. A. Larson, M. C. Willemsen, J. A. Konstan, J. J. McAuley, J. Garcia-Gathright, B. Huurnink, and

of winning over computer vision?" *CoRR*, vol. abs/2108.04384, 2021. [Online]. Available: https://arxiv.org/abs/2108.04384

E. Oldridge, Eds. ACM, 2021, pp. 505–514. [Online]. Available: https://doi.org/10.1145/3460231.3475943

[54] C. Liu, X. Liu, R. Zheng, L. Zhang, X. Liang, J. Li, L. Wu, M. Zhang, and L. Lin, "C$^2$-rec: An effective consistency constraint for sequential recommendation," *CoRR*, vol. abs/2112.06668, 2021. [Online]. Available: https://arxiv.org/abs/2112.06668