# Multi-trends Enhanced Dynamic Micro-video Recommendation

Yujie Lu
UC Santa Barbara
United States
yujielu@zju.edu.cn

Yingxuan Huang
The University of Hong Kong
China
eloise@connect.hku.hk

Shengyu Zhang*
Zhejiang University
China
sy_zhang@zju.edu.cn

Wei Han
Singapore University of Technology
and Design
China
henryhan88888@gmail.com

Hui Chen
Singapore University of Technology
and Design
China
hui_chen@mymail.sutd.edu.sg

Fei Wu*
Zhejiang University
China
wufei@zju.edu.cn

Zhou Zhao*
Zhejiang University
China
zhaozhou@zju.edu.cn

## ABSTRACT

The explosively generated micro-videos on content sharing platforms call for recommender systems to permit personalized micro-video discovery with ease. Recent advances in micro-video recommendation have achieved remarkable performance in mining users' current preference based on historical behaviors. However, most of them neglect the dynamic and time-evolving nature of users' preference, and the prediction on future micro-videos with historically mined preference may deteriorate the effectiveness of recommender systems. In this paper, we propose to explicitly model dynamic multi-trends of users' current preference and make predictions based on both the history and future potential trends. We devise the DMR framework, which comprises: 1) the implicit user network module which identifies sequence fragments from other users with similar interests and extracts the sequence fragments that are chronologically behind the identified fragments; 2) the multi-trend routing module which assigns each extracted sequence fragment into a trend group and update the corresponding trend vector; 3) the history-future trend prediction module jointly uses the history preference vectors and future trend vectors to yield the final click-through-rate. We validate the effectiveness of the proposed framework over multiple state-of-the-art micro-video recommenders on two publicly available real-world datasets. Relatively extensive analysis further demonstrate the superiority of modeling dynamic multi-trend for micro-video recommendation.

*Corresponding Authors.

## CCS CONCEPTS

• **Information systems** → **Recommender systems**.

## KEYWORDS

Micro-video Recommendation; Dynamic User Modeling; Future-aware; Multi-trend Routing

## 1 INTRODUCTION

In recent years, the amount of searchable micro-videos has increased dramatically and exacerbated the need for recommender systems that can effectively mine users' preference and identify potentially interested micro-videos in a personalized manner. Due to the powerful representation learning capacity, the rapid development of deep learning techniques has nourished the research field of recommendation [17, 24, 33, 41, 42, 57, 58, 62, 65, 67, 68, 70, 73, 74]. Such a development also gives rise to diverse models for video recommendation, which can be roughly categorized to collaborative filtering [2, 29], content-based filtering [11, 16, 44, 48, 77], and hybrid ones [5, 6, 72].

Compared with professional video recommendation, micro-video recommendation poses many unique challenges. For example, micro-videos typically lack of meta-data (e.g., genre, director, actor/actress, which are commonly available in professional videos), leading to semantic gap in representation [9]. Moreover, users might be interested in multiple topics of videos simultaneously, i.e., diverse interests, and yield interests to different extends (e.g., like, follow, click), i.e., multi-level interests [39]. Recent years have witnessed much progress to confront the above challenges in this vein. THACIL [9] employs temporal block splitting and hierarchical multi-head attention to model diverse interests across blocks. ALPINE [39] models users' dynamic interests by constructing temporal behavior graph

and devising the temporal graph-based LSTM. MTIN [30] considers personalized importance decay over time and diverse interests using item-level temporal mask and group routing mechanism, individually. In spite of the great advances of these works, we argue that solely modeling the historical behaviors deteriorates the capacity of user modeling capturing *diverse* and *dynamic* users' interests. For example, MTIN [30] assigns historically interacted items to one of six interest groups and accordingly updates the six interest vectors. Since users' interests are by nature dynamic, the interests learned from the logged data might be out-of-date or at least limited to the history, falling short to recommend fresh items and hurting the recommendation diversity. Therefore, capturing dynamic interest trends based on (but not limited to) historical items can be an indispensable function for high-quality recommender systems.

Towards this end, we devise the multi-trends framework for dynamic micro-video recommendation, abbreviated as DMR. We start from the perspective that trends refer to the possible future directions of the current interest implied by the logged interactions. Since we have no access to items interacted in the future, DMR encapsulates an implicit user network construction module that first identifies sequence fragments that yield similar interests as the current sequence from similar users. Then, we constructs possible trending sequences by extracting the sequence fragments that are chronologically behind the identified ones. We note that some trending sequences may share similar interests and representing each sequence as an individual interest may introduce unnecessary noises and computation costs. Towards this end, inspired by [30, 38], we devise a multi-trend routing module that transforms multiple trending sequences to fewer number of multiple trend interest vectors. However, extracting trending sequences and mapping them to trend vectors for each testing inference might hurt the serving efficiency of industrial deployment. Thus, multi-trend routing module constructs a fixed-length trend memory for each user and read-writes the memory during training. For memory read-writing, we propose to assign trending sequences to memory slots in a soft way and power the process with attention mechanisms. During inference, we directly take the off-the-shelf history/trending vectors without extracting or transforming trending sequences, and thus addressing the efficiency issue. Predictions are performed with the history-trend joint prediction module.

To this end, DMR framework makes predictions based on both the history interests implied by the historical behaviors as well as multi-trends implied in similar users, which helps to capture even more diverse and dynamic interests compared with existing micro-video recommenders. We validate the effectiveness of DMR on micro-video recommendation benchmarks. The substantial improvement over state-of-the-art comparison methods and in-depth model analysis demonstrate the superiority of modeling multi-trend for micro-video recommendation. Overall, this paper has the following contributions:

- We propose to capture even more diverse and dynamic interests beyond the historical behaviors by modeling the possible interest trends for micro-video recommendation.
- We devise the novel DMR framework that encapsulates the implicit user network construction module, which extracts

trending sequences from similar users, the multi-trend routing module, which performs dynamic trending memory read-write and improves the inference efficiency at the inference stage, and the history-trend joint prediction module.
- We conduct extensive experiments on micro-video recommendation benchmarks, of which the results show DMR framework achieves high-quality recommendation with improvement on both accuracy and diversity.

## 2 RELATED WORK

### 2.1 Video Recommendation

The methods for recommendation can be generally classified into two categories. Early algebraic approaches adopted collaborative filtering [15, 27, 35, 54] or model-based methods [13, 34, 52, 63] to estimate user-item correlations and make predictions about users' future interests. Collaborative filtering (CF) assumes that users sharing the same opinion on one issue tend to have more similar opinions on other issues [37], and thus it makes predictions specific to each user through information gleaned from other users [60]. Due to the extreme high computational complexity and data sparsity in traditional CF [13, 47], model-based methods alleviate this overhead by mapping user-item interactions into matrix entries, then apply factorization to the characteristic matrix to build nonlinear models that estimate correlations (i.e. preferences) between every pair of user and item and employ Hidden Markov models (HMM) to capture temporal trends of preferences [52].

Recently, as the major advances in deep learning techniques, a wealth of research has sprung up on incorporating them into recommender systems. Most of work reformulated traditional estimation problem as learning task based on deep neural networks [18, 25, 56]. In the field of video recommendation, representative work focuses on content-based learning [10, 12, 64, 66], in which features of videos are extracted into embedding vectors and then matched with user representations that indicate individual preference. To name a few, Chen et al. [7] tackled the item- and component-level implicit feedback issue in multimedia recommendation by learning independent video and users characteristics in a unified hierarchical attention network and then reckoning pair-wise scores as a measure of user preference. Although these works improve the accuracy of user modeling, they lack a clear partition of history and future for the given dataset and hence may encounter prediction bias due to mixing the two parts together. Our work adopts a multi-step time partition and similarity matching approach to alleviate this issue.

### 2.2 User Behavior Modeling

Modeling latent user interest from historical behaviors is commonly used in recommender systems. In the past two decades, a variety of approaches have been proposed, ranging from Markov chains [22, 23, 46, 52, 55] and traditional collaborative filtering [15, 36, 53] to deep representation learning [50, 76]. The approaches based on Markov decision processes implicitly track user state dynamics to predict future behaviors. For example, Rendle et al. [52] captured long-term user interest via personalized transition graphs over underlying Markov chains. He and McAuley [23] integrated

similarity-based methods with Markov chains smoothly in personalized sequential recommender systems. Besides, temporal collaborative filtering is proposed to deal with the drifting user preferences. Koren [36] offered a paradigm that tracks time changing behaviors throughout the life span of the data.

With the development of deep learning, more and more researchers adopted deep neural networks (DNN) to model the user dynamics in recommender systems. Particularly, Hidasi et al. [28] applied recurrent neural networks (RNN) to model the whole session and introduced a new ranking loss function to make recommendations more accurate. Tang and Wang [59] utilized convolutional filters to embed a sequence of recent items into an "image" in the time and latent spaces as well as learn sequential patterns as local features of the image. Wu et al. [71] considered session sequences as graph structured data and used graph neural networks (GNN) to capture complex transitions of items. Recently, self-attention mechanism [61] has been widely employed in recommender systems[32]. For instance, Wu et al. [69] proposed a Contextualized Temporal Attention Mechanism to weigh historical actions' influence on not only what action it is, but also when and how the action took place.

However, previous work does not consider the influence of future information when modeling user behaviors in history sequences. In this work, we constructed a user-item heterogeneous graph to capture future interactions of each user with items.

## 3 METHODOLOGY

In this section, we first formulate the micro-video recommendation problem, and then introduce the proposed framework in detail. As illustrated in Figure 1, our proposed DMR framework for dynamic micro-video recommendation mainly comprises of three modules:1) Pearson Correlation Coefficient enhanced implicit user network module; 2) A history-future multi-trend joint routing module; 3) A multi-level time-aware attention module.

### 3.1 Problem Formulation

In a typical micro-video recommendation scenario, we have a set of users and a set of micro-videos, which can be denoted as $U = \{u_1, u_2, u_3, ..., u_{|U|}\}$ and $V = \{v_1, v_2, v_3, ..., v_{|V|}\}$ respectively. Let $I_u = \{x_1^u, x_2^u, ..., I_{|I_u|}^u\}$ represent the sequence of interacted micro-videos $x \in I_u$ of user $u \in U$, which is sorted in a chronological order according to the timestamp of each interaction, and $x_t^u$ denote the micro-video that the user $u$ has interacted with at timestamp $t$. The interaction sequence $I_u$ is split into $I_+$ and $I_-$ which represent the micro-videos clicked by the user and the ones not clicked respectively. Given the user's historical micro-video interaction behaviors, the investigated goal of the micro-video recommendation task in this paper is to predict the probability that the new candidate micro-video will be clicked by user $u$. Notations are summarized in Table 1.

Specifically, each instance is represented by a tuple $(I_u, A_i)$, where $I_u$ denotes the set of items interacted by user $u$, $A_i$ the features of target item $i$ including the information of interaction timestamp and micro-video embeddings. Through implicit user network module, we extract relative future sequence of user $u$ based on $I_u$ and their similar users' historical interaction $I_{u'}, u' \in U$. The detail will be illustrated in Section 3.3.

**Table 1: Notations.**

| Notation | Description |
|----------|-------------|
| u | a user |
| v | a micro-video |
| x | an interaction |
| d | the dimension of user/micro-video embeddings |
| t | the number of trends |
| U | the set of users |
| V | the set of micro-videos |
| I | the set of interactions |
| T | the trends set |

To model diverse user preferences dynamically, DMR learns a function $f$ for mapping history trend set $T_u^h$ and future trend set $T_u^f$ into user representations, which can be formulated as:

$$\overrightarrow{e_u} = f(T_u^h, T_u^f) \tag{1}$$

where $\overrightarrow{e_u} \in \mathbb{R}^{d \times 1}$ denotes the representation vector of user $u$, $d$ the dimension. Besides, the representation vector of target micro-video $i$ is obtained by an embedding function $g$ as:

$$\overrightarrow{e_i} = g(A_i) \tag{2}$$

where $\overrightarrow{e_i} \in \mathbb{R}^{d \times 1}$ denotes the representation vector of target micro-video $i$.

Based on the learned user representation vector and micro-video representation vector, the probability of candidate micro-video is calculated using the likelihood function $P$ as:

$$p(i|U, V, X) = P(\overrightarrow{e_u}, \overrightarrow{e_i}) \tag{3}$$

where $\overrightarrow{e_i}$ is the embedding of target item $i$ from set of micro-videos $V$. Our framework outputs the click probabilities of the candidate micro-video to rank the personalized recommendation list. Then the system provide precise and diversified recommendation for each user, which entails potential preference of the specific user as they are most likely to interact with the recommended micro-videos.

The objective function for training our model is described in Section 3.6 We use the Adam optimizer to train our method.

### 3.2 Overview

The overall structure of our proposed framework DMR is illustrated in Figure 1, which is composed of an implicit user network module, a multi-trend routing module, a multi-level time-aware attention module and a prediction layer. As the relative future sequence for current user is actually the history sequence for the neighbors, the multi-trend routing algorithm is applied on both the future and history sequences using shared parameters in parallel. The framework takes the user historical interactions set $X$ as input. We use $X_{1,N-K}^u$ and $X_{N-K+1,N}^u$ to represent training and testing data of interactions sequence of user $u$ respectively. $N$ and $K$ denotes the selected total length of interaction sequence of each user $u$ and the length of training sequence respectively. For micro-videos from the set of $X_{1,N-K}^u$, embeddings are presented as $\overrightarrow{e}_{X_{1,N-K}^u}$.
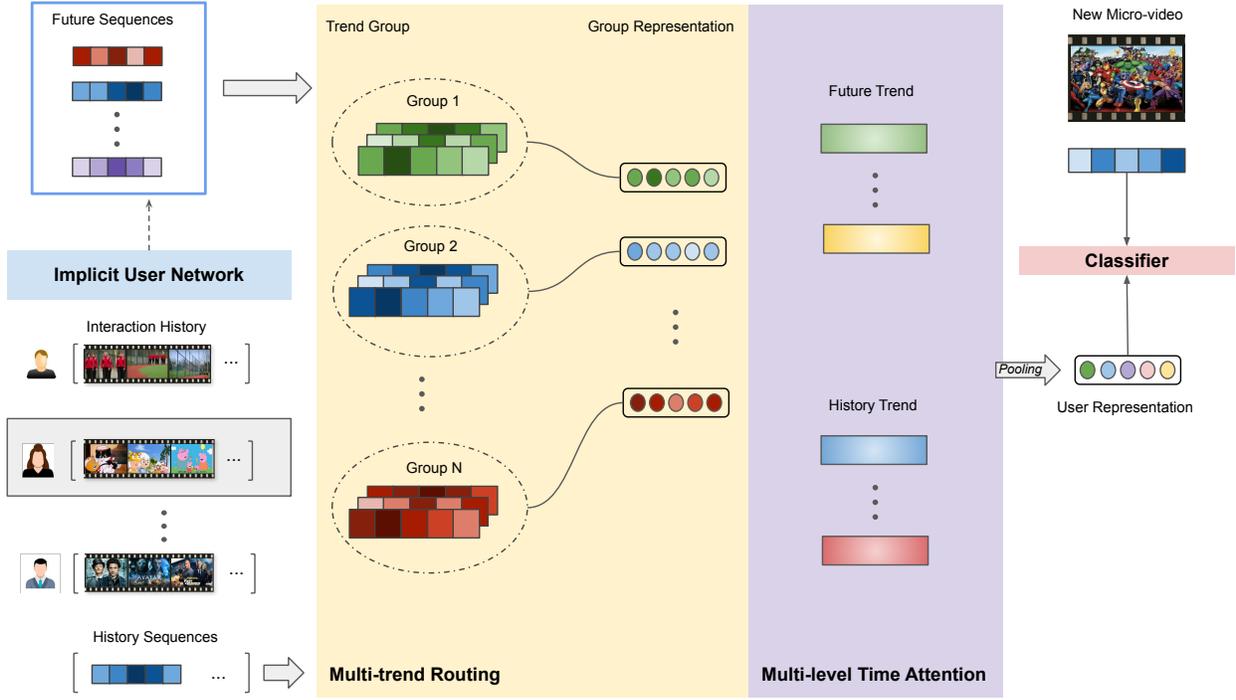
**Figure 1: Network Architecture of DMR. DMR is composed of an implicit user network module, a multi-trend routing module, a multi-level time attention layer and a prediction layer. Based on the users' historical interactions, we build a implicit user network to construct future sequences. The multi-trend module are applied on the current user's history sequences and future sequences in parallel to get representation of each trend group. The multi-level time attention mechanism are applied before the pooling layer to generate the history trend representation and future trend representation, which is further concatenated as dynamic user preference representation. Finally, the user representation and the candidate micro-video embedding are utilized for prediction in the classifier.**

The implicit user network module constructs neighbors set for each user by selecting the users that have similar micro-video preference as indicated in their past behaviors, and then extract the relative future sequences from each neighbor. The query items can be selected from the user historical interaction $I_u$, for simplicity, we solely choose the last one in the list, which can be both efficient and effective as demonstrated in the empirical analysis. The relative future behaviors are defined as the interacted items following the query item in the chronological order, aiming at representing dynamic preference of the user. The intuition in behind is that the user tend to have similar preference trend as users with similar historical behaviors, and that the user can have diverse and dynamic trends of preferences.

The multi-trend routing module is developed to obtain the neighbor centroids according to diverse motivation behind specific interactions with the micro-videos. Then we learn future-aware diverse trends based on history and future sequence jointly. Furthermore, the future sequence evolved user representation acquired by time-aware attention layer is concatenated with the historical behavior

evolved user representation to generate the dynamic user preference representation vector. Finally we compute the user's preferences over different micro-videos from the pool by the prediction decoder. Each part will be elaborated in the following sections.

### 3.3 Implicit User Network

As shown in Figure 2 ,the implicit user network is constructed based on user-item heterogeneous graph, which contains both the user nodes and item nodes. An edge in the graph represents the interaction between the user and the item. The weight of the edge indicates the temporal weight of each interacted item in a chronological order. The query items are selected in a multi-hop manner. The user nodes connected to the selected query items are considered as the candidate neighbor nodes of the current user.

Inspired by some works [19, 20], which extract social relationships in absence of explicit social networks [45] , we construct the user network from user-item correlation implicitly.

Specifically, we compare the similarity among users via collaborative filtering implicitly based on the historical interactions with micro-videos. As the Pearson Correlation Coefficient(PCC) is a
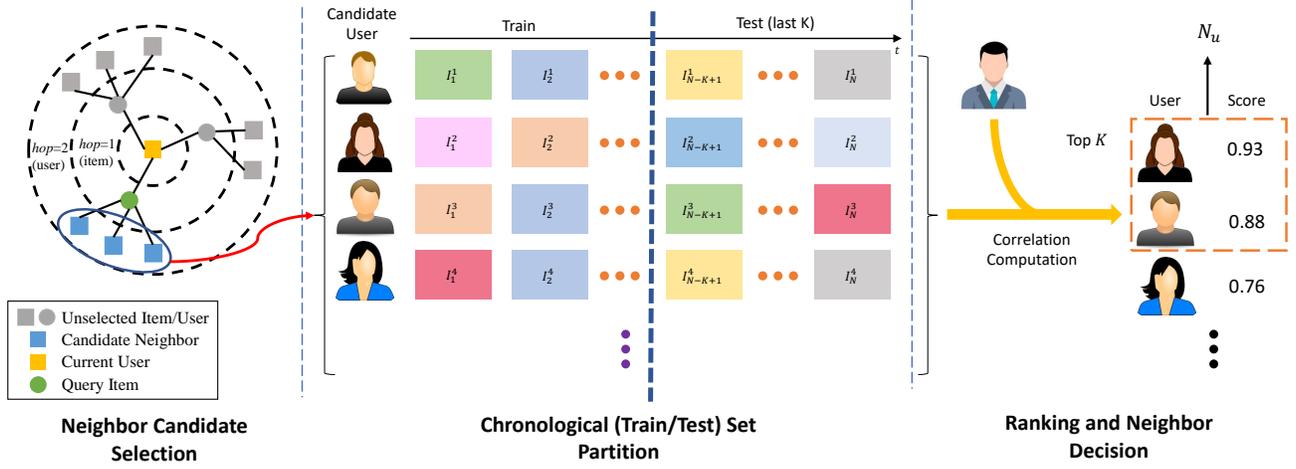
Figure 2: Architecture of the implicit user network module. The leftmost part stands for the neighbor candidate selection process based on user-item graph with the interactions by edge, user and micro-videos by node. User behaviors of the selected neighbors are then split into train and test set to compare their similarity to the current user. The relative future sequence of the most similar users are utilized to generate the future sequence as the input of multi-trend routing module, which output the future trend representation.

---

**Algorithm 1** Implicit User Network Construction

**Input:**
 The set of users $U$;
 User's historical interaction sequence $I_u$;
 User's query items sequence $K_u$ and upper bound k;
 User's candidate neighbors $G_u$ and upper bound g;
 Similarity threshold $\tau$ for neighbor selection;

**Output:**
 The extracted neighbor set of user $N_u, u \in U$;
1: **for** each $u \in U$ **do**
2:    $N_u \leftarrow \emptyset$
3:    **for** each $i \in Inverse(I_u)$ **do**
4:      **if** $|K_u| < k$ **then**
5:        $K_u \leftarrow INSERT(i)$
6:      **end if**
7:    **end for**
8: **end for**
9: **for** each $u \in U$ **do**
10:    **for** each $n \in U$ **do**
11:      $s_{un} = USER\_SIMILARITY(u, n)$
12:      **if** $n \neq u \land |G_u| < g \land s_{un} > \tau$ **then**
13:        $|G_u| \leftarrow INSERT(n)$
14:      **end if**
15:    **end for**
16:    $N_u \leftarrow TOP\_SIMILARITY(G_u)$;
17: **end for**
18: **return** $N_u$

---

widely used similarity measure, we adopt Pearson Correlation Coefficient [4] to compute a linear correlation between the user and each candidate neighbor as:

$$s_{ij} = \frac{\sum\limits_{k \in I(i) \cap I(j)} (r_{ik} - \overline{r}_i) \cdot (r_{jk} - \overline{r}_j)}{\sqrt{\sum\limits_{k \in I(i) \cap I(j)} (r_{ik} - \overline{r}_i)^2} \cdot \sqrt{\sum\limits_{k \in I(i) \cap I(j)} (r_{jk} - \overline{r}_j)^2}} \quad (4)$$

where $I(i)$ is a set of micro-videos user $i$ interacted with, $r_{ik}$ and $\overline{r}_i$ represents the level (click or not click) of interaction of user $i$ over micro-video $k$ and the average level of action of user $i$. The user similarity $s_i$ is ranging from $[-1, 1]$, and the similarity between users $i$ and $j$ is proportional to the value according to this definition. Following [43], we employ a mapping function $f(x) = (x + 1)/2$ to bound the range of PCC similarities into $[0, 1]$.

In the case of users with only one common micro-video in history, PCC similarity gets 1 when the users' preferences over the common micro-video are similar and −1 when not, which encourages diversity of neighbors while damaging the fairness of similarity calculation. To tackle this issue, we only kept less than 20% of such neighbor nodes to seek the balance.

In addition to the PCC method, we also design a filter with simple schema to extract similar users. For each user, if the historical interactions $I_u$ is split into two pieces, $I^u_{1:t_1}$ for training data, and $\hat{I}^u_{t_1:t_2}$ for testing data, the item $\hat{I}^u_k$ is defined as the last $k$ micro-videos, $k$ could be any value less than or equal to $|I^u|$, while in practice $k = 1$ can achieve good enough performance with simplicity. We extracted a list of neighbors $N = \{n_1, n_2, ..., n_{|N|}\}$ according to the query item. The detail of this process is described in Algorithm 1.

Furthermore, we constructed the future sequence of user $u$ as:

$$F_u = \{n_f, n_f \in I^n, TI(n_f) \geq TI(I^u_{|I_u|-k})\} \quad (5)$$

where Timestamp is denoted as $TI$ and the query item is denoted as $I^u_{|I_u|-k}$. $I_n$ represents the interaction set of neighbor $n$

## 3.4 Multi-trend Routing

To capture the trend information lies in both history sequence and future sequence, we devised a multi-trend routing module into a two-stage manner to generate trend represent parallelly. Specifically, we group each micro-video from both the user's historical sequence and extracted relative future sequence into diverse trends in the first stage. The micro-videos that are grouped into the same trend are considered to be similar according to users' interactions over them and their own basic features. In the second stage, the micro-videos from historical sequence and relative future sequence are utilized to generate the representation of history and future trend group in parallel.

Based on the positive historical interaction sequence $I_+$ of user $u$, we represent each micro-video $x$ in $I_+$ as an embedding vector $\overrightarrow{x} \in \mathbb{R}^d$, where $d$ is the embedding size. And we initialize positive history trend group as $T^h_u \in \mathbb{R}^{s \times d}$ for user $u$, where $s$ denotes the number of trend groups indicated from historical sequence and $d$ denotes the embedding dimension of each history trend. Specifically, each trend embedding is represented as $\overrightarrow{t} \in \mathbb{R}^d$.

Similarly, based on the extracted future sequence $F_+$ from the implicit user network. The positive future trend group is denoted as $T^f_u \in \mathbb{R}^{s \times d}$ for user $u$, where $s$ denotes the number of trend groups indicated from future sequence and $d$ denotes the embedding dimension of each future trend.

In order to fine-tune the representation of each trend, we apply attention mechanism over each micro-video and the initialized trend group. Given the micro-video embedding $\overrightarrow{x} \in \mathbb{R}^d$ and the trend embedding $\overrightarrow{t} \in \mathbb{R}^d$, we calculate the weight between the micro-video and the trend based on a co-attention memory matrix. The micro-video from the history sequence and the future sequence are put into history trend and future trend separately. As the history sequence and future sequence is processed separately, our module is capable of capturing timeliness of trends which indicates evolved user interest.

## 3.5 Multi-level Time Attention Mechanism

As for the item-level, we use the weighted sum of historical micro-video features to obtain the current micro-video representation. Finally, we get the representation of each trend by attention mechanism on each micro-video in the trend group. As for the trend-level, we utilize the time-aware attention to activate the weight of diverse trends to capture the timeliness of each trend. Specifically, the attention function takes the interaction time of item $i$, the interaction time of trends and trend embeddings as the query, key and value respectively. We compute the final representation of trend representation future sequence of user $u$ as:

$$HF_u = Attention(\overrightarrow{TI_i}, \overrightarrow{TI_{tr}}, \overrightarrow{t_u}) = \overrightarrow{t_u} softmax(pow(\overrightarrow{TI_i}, \overrightarrow{TI_{tr}})) \quad (6)$$

where Attention denotes the attention function, $TI_i$ represents the interaction time of micro-video $i$, $TI_{tr}$ represents the average

interaction time of micro-videos related to the trend group, $\overrightarrow{t_u}$ represents the embedding of the specific trend group.

The trend group generated from the user's historical sequence and future sequence are then eventually updated by adding the corresponding trend group in $T^h_u$ and $T^f_u$ with the aggregation of history trend and future trend representation respectively.

## 3.6 Prediction

After computing the trend embeddings from activated trends through time-aware attention layer, we apply sumpooling to both history and future trend representations.

$$e^h_u = sumpooling(T^{h_1}_u, ..., T^{h_s}_u), e^f_u = sumpooling(T^{f_1}_u, ..., T^{f_s}_u) \quad (7)$$

And then we concatenate the history trend representation vector $e^h_u$ and future trend representation vector $e^f_u$ to form a user preference embedding $\overrightarrow{e_u}$ as:

$$\overrightarrow{e_u} = e^h_u \frown e^f_u \quad (8)$$

Given a training sample $u, i$ with the user preference embedding $\overrightarrow{e_u}$ and micro-video embedding $\overrightarrow{e_i}$ as well as the micro-video set $V$, we can predict the possibility of the user interacting with the micro-video as

$$p(i|U, V, I) = \frac{exp(\overrightarrow{e_u}^T \overrightarrow{e_i})}{\sum_{v \in V} exp(\overrightarrow{e_u}^T \overrightarrow{e_v})} \quad (9)$$

In the same way, we calculate the prediction score $P(x|H_-)$ based on the negative interaction sequence, which aims to maximize the distance between the new micro-video embedding and user's negative trend embeddings.

The final recommendation probability $\hat{p}_{ij}$ is represented by the linear combination of $p(x|H_+)$ and $p(x|H_-)$. And the objective function of our model is as follows:

$$\mathbb{L} = -\sum_{i \in \mathbb{U}} \left( \sum_{i \in H_+} \log \sigma(\hat{p}_{ui}) + \sum_{i \in H_-} log(1 - \sigma(\hat{p}_{ui})) \right) \quad (10)$$

where $\hat{p}_{ui}$ denotes the prediction score of micro-video $i$ for user $u$, $\sigma$ represents the sigmoid activation function.

# 4 EXPERIMENTS

## 4.1 Dataset

MicroVideo-1.7M and KuaiShou were used as micro-video benchmark datasets in our experiments. Micro-video data and user-video interaction information can be found in each of these datasets. Each micro-video is represented by its features in these two datasets, and each interaction record includes the userID, micro-video ID, visited timestamp, and whether the user clicked the video. The two datasets' statistics are shown in Table 2.

- `MicroVideo-1.7M`[8]: This dataset comes from real data of micro-video sharing service in China which contains 1.7 million micro-videos.
- `KuaiShou`: This dataset is released by the Kuaishou Competition in China MM 2018 conference.

| Dataset | users | items | interactions | train | test |
|---------|-------|-------|--------------|-------|------|
| MicroVideo-1.7M | 10,986 | 1,704,880 | 12,737,619 | 8,970,310 | 3,767,309 |
| KuaiShou | 10,000 | 3,239,534 | 13,661,383 | 10,931,092 | 2,730,291 |

## 4.2 Implementation Details

We used TensorFlow on four Tesla P40 GPUs to train our model with Adam optimizer. The following are the hyper-parameters: The micro-video embedding is 512-dimensional vectors, while the user embedding is 128-dimensional vectors. The batch size was set to 32, the optimizer was Adam, the learning rate was set to 0.001, and the regularization factor was set to 0.0001.

To find the user's similar neighbors, we used the Pearson Correlation Coefficient (PCC) described earlier. In the ablation analysis, we set neighbor numbers as 5, 20, and 50. As for the future sequences, we cut off each neighbor's at most 100 interacted micro-videos after the current user's query items.

## 4.3 Evaluation Metrics

To compare the performance of different models,we use **Precision@N**, **Recall@N**, **F1-score@N** and **AUC**, where N is set to 50 as metrics for evaluation.

- Precision: Number of correctly predicted positive observations divided by the total number of predicted positive observations.

$$Precision@N = \frac{1}{|U|} \sum_{u \in U} \frac{|\hat{I}_{u,N} \cap I_r|}{|I_r|} \quad (11)$$

where $\hat{I}_{u,N}$ denotes the set of top-N recommended micro-videos for user u and $I_r$ is the total recommendation list for user u.

- Recall: Number of corrected recommended micro-videos divided by the total number of all recommended micro-videos.

$$Recall@N = \frac{1}{|U|} \sum_{u \in U} \frac{|\hat{I}_{u,N} \cap I_u|}{|I_u|} \quad (12)$$

where $\hat{I}_{u,N}$ denotes the set of top-N recommended micro-videos for user u and $I_u$ is the set of testing micro-videos for user u.

- F1-score: F1 Score is the weighted average of Precision and Recall. It's used to balance between Presicion and Recall.

$$F1 - score = 2 * \frac{Precision * Recall}{Precision + Recall} \quad (13)$$

- AUC: AUC (Area Under the ROC Curve) is used in classification analysis to determine the quality of classifiers.

## 4.4 Competitors

To validate the effectiveness of our proposed DMR framework, we conducted experiments on two publicly available real-world datasets. The comparision to other state-of-the-art micro-video recommenders are summarized in Table 3.

- BPR[51]: Trained on pairwise items, the Bayesian personalized ranking(BPR) maximize the difference between positive and negative items of each user in Bayesian approach.
- LSTM[75]: Long short-term memory(LSTM) is a sequence model. Hidden states of each unit are aggregated to form user interest representation.
- CNN: The convolutional neural network (CNN) can be utilized to generate user interest representations based on the interaction sequence. The max pooling layer and MLP layers are used for user interest extraction and prediction.
- NCF[26]: As a collaborative filtering based model, NCF learns user embedding and item embedding with a shallow network and a deep network, which is able to learn an arbitrary function from data.
- ATRank[76]: ATRank is an attention-based behavior modeling framework, which can model with heterogeneous user behaviors using only the attention model. It utilizes self-attention in multiple semantic spaces to capture behaviors interactions. The model is capable of predicting all types of user actions in a multi-task manner, which shows effectiveness over the highly optimized individual models.
- THACIL[8]: THACIL achieved the click-through prediction for micro-videos by modeling user's historical behaviors. The proposed recommendation algorithm characterizes both short-term and long-term correlation within user behaviors. It also profiles user interests at both coarse and fine granularities.
- ALPINE[40]: To intelligently route micro videos to target users, ALPINE proposed an LSTM model based on a temporal graph, which is encoded by user's historical interaction sequence. The model captures the complex and diverse interests of users via a multi-level interest modeling layer. Moreover, the model achieves better performance by utilizing true negative samles, which indicates uninterested information.
- MTIN[31]: This model is a multi-scale time-aware user interest modeling framework, which learns user interests from fine-grained interest groups. MTIN incorporates the interest group routing unit to generate user interest groups based on the interaction sequence and leverages fine-grained interest groups via item-level and group-level interest extraction unit. The distilled user interest representation is used to predict the click probabilities of micro-video candidates.

## 4.5 Results

The model performance on the two datasets is summarized in Table 3. We run experiments to dissect the effectiveness of our recommendation model. We compare the performance of DMR with several commonly used and state-of-the-art models: BPR, LSTM, CNN, NCF, ATRank, THACIL, ALPINE and MTIN. All these models are running on the two datasets introduced above: MicroVideo-1.7M and KuaiShou-Dataset. According to the results shown in Table 3, our model DMR achieve better performance on precision over KuaiShou dataset and performs better in terms of AUC, Recall and F1-score over MicroVideo-1.7M dataset.

Table 4 compares the result of different neighbor number setting of 5, 20 and 50. Considering more neighbors could result in more

**Table 3: Overall Performance Comparision. The model performance of our model and several state-of-the-art baselines on two public datasets: MicroVideo-1.7M and KuaiShou-Dataset. The best results are highlighted in bold.**

| Model | MicroVideo-1.7M | | | | KuaiShou-Dataset | | | |
|---|---|---|---|---|---|---|---|---|
| | AUC@50 | Precision@50 | Recall@50 | F1-score@50 | AUC@50 | Precision@50 | Recall@50 | F1-score@50 |
| BPR | 0.583 | 0.241 | 0.181 | 0.206 | 0.595 | 0.290 | 0.387 | 0.331 |
| LSTM | 0.641 | 0.277 | 0.205 | 0.236 | 0.731 | 0.316 | 0.420 | 0.360 |
| CNN | 0.650 | 0.287 | 0.214 | 0.245 | 0.719 | 0.312 | 0.413 | 0.356 |
| NCF | 0.672 | 0.316 | 0.225 | 0.262 | 0.724 | 0.320 | 0.420 | 0.364 |
| ATRank | 0.660 | 0.297 | 0.221 | 0.253 | 0.722 | 0.322 | 0.426 | 0.367 |
| THACIL | 0.684 | **0.324** | 0.234 | 0.269 | 0.727 | 0.325 | 0.429 | 0.369 |
| ALPINE | 0.713 | 0.300 | 0.460 | 0.362 | 0.739 | 0.331 | 0.436 | 0.376 |
| MTIN | 0.729 | 0.317 | 0.476 | 0.381 | **0.752** | 0.341 | **0.449** | **0.388** |
| DMR | **0.731** | 0.323 | **0.478** | **0.385** | 0.742 | **0.343** | 0.442 | 0.386 |

**Table 4: Effect analysis of Neighbors. The model performance with different Neighbor Number setting on two datasets: MicroVideo-1.7M and KuaiShou-Dataset. The metrics are @50. Here we set Neighbor Number to 5, 20, 50.**

| Model | MicroVideo-1.7M | | | | KuaiShou-Dataset | | | |
|---|---|---|---|---|---|---|---|---|
| | AUC@50 | Precision@50 | Recall@50 | F1-score@50 | AUC@50 | Precision@50 | Recall@50 | F1-score@50 |
| DMR-N5 | 0.689 | 0.319 | 0.425 | 0.364 | 0.674 | 0.333 | 0.439 | 0.378 |
| DMR-N20 | **0.731** | **0.323** | **0.478** | **0.385** | **0.742** | **0.343** | **0.442** | **0.386** |
| DMR-N50 | 0.668 | 0.280 | 0.282 | 0.281 | 0.652 | 0.329 | 0.404 | 0.362 |

diversity, but too many neighbors would dilute interest trends' embedding. Our model achieves improvements on neighbor number equals 20 over 5. Besides, it shows reduction if setting neighbor number from 20 to 50. This means the number of neighbors also play a crucial part in model performance.

The computational complexity of sequence layer modeling user and neighbors is $O(knd^2)$, where $k$ denotes the number of extracted neighbors, $n$ denotes the average sequence length and $d$ denotes the dimension of item's representation. Capsule layer's computational complexity depends on kernel size and number of trends. Average time complexity of capsule layer scales $O(nTr^2)$, where $r$ denotes kernel size of capsule layer and $T$ denotes the number of trends. For large-scale applications, our proposed model could reduce computational complexity by two measures: (1)encode neighbors with a momentum encoder[21].(2)adopt a light-weight Capsule network.

## 4.6 Recommendation Diversity

Aside from achieving high recommendation accuracy, diversity is also essential for the user experience. With little information of historical interactions between the users and the micro-videos, recommendation systems learned to assist users in selecting micro-videos that would be of interest to them. Recommender systems keep track of how users interacted with the micro-videos they've chosen.

Many research works [1, 3, 14, 49] have been undertaken to propose novel diversification algorithms. Our proposed module

can learn the diverse trends of user preference and provide recommendation with diversity. We define the individual diversity as below:

$$Diversity@N = \frac{\sum_{j=1}^{N} \sum_{k=j+1}^{N} \delta(CATE(\hat{i}_{u,j}) \neq CATE(\hat{i}_{u,k}))}{N \times (N-1)/2} \quad (14)$$

where $CATE$ represents the category of the item. $\hat{i}_u$ denotes item recommended for user $u$, $j$ and $k$ represents the order of the recommended items. $\delta(\cdot)$ is an indicator function.

Table 5 presents comparisons with THACIL and MTIN over the recommendation diversity metric on Micro-video dataset, which provides category infromation of micro-videos. We adopt the setting of six historical trend and six future trend evolved from 5 neighbors for our model. From the table, our module DMR achieve the optimum diversity metric indicating the recommendation it provide can effectively take neighbors' interests into account.

## 5 CONCLUSION AND FUTURE WORK

In this work, we propose to capture even more diverse and dynamic interests beyond those implied by the historical behaviors for microvideo recommendation. We refer to the future interest directions as trends and devise the DMR framework. DMR employ an implicit user network module to extract future sequence fragments from similar users. A mutli-trend routing module assigns these future sequences to different trend groups and updates the corresponding

**Table 5: Model Recommendation Diversity Comparision on Micro-video Dataset.**

| MicroVideo-1.7M | THACIL | MTIN | DMR |
|---|---|---|---|
| Diversity@10 | 1.9112 | 1.9940 | **1.9948** |
| Diversity@50 | 1.9104 | 1.9948 | **1.9956** |
| Diversity@100 | 1.9436 | 1.9950 | **1.9954** |

trending memory slot in a dynamic read-write manner. Final predictions are made based on both future evolved trends and history evolved trends with a history-future trends joint prediction module.

This work represents one of the initial attempts to explicitly capture possible interest trends for a given historical behavior sequence, especially for ranking models and micro-video recommendation. We believe that such an idea can be inspirational to future works on learning recommender systems of high diversity. In the future, though the implicit user network module does not affect serving efficiency, we would like to explore whether more efficient and effective solutions exist to boost the training since introducing information from other users might also introduce inevitable noises. Moreover, we plan to extend the multi-trend capturing idea to more applications in recommender systems and address application-specific challenges.

## REFERENCES

[1] G. Adomavicius and Y. Kwon. 2012. Improving Aggregate Recommendation Diversity Using Ranking-Based Techniques. *IEEE Transactions on Knowledge and Data Engineering* 24, 5 (2012), 896–911.

[2] Shumeet Baluja, Rohan Seth, D. Sivakumar, Yushi Jing, Jay Yagnik, Shankar Kumar, Deepak Ravichandran, and Mohamed Aly. 2008. Video suggestion and discovery for youtube: taking random walks through the view graph.. In *Proceedings of the 17th International Conference on World Wide Web, WWW 2008, Beijing, China, April 21-25, 2008*.

[3] Rubi Boim, Tova Milo, and Slava Novgorodov. 2011. Diversification and Refinement in Collaborative Filtering Recommender. In *Proceedings of the 20th ACM International Conference on Information and Knowledge Management* (Glasgow, Scotland, UK) *(CIKM '11)*. Association for Computing Machinery, New York, NY, USA, 739–744. https://doi.org/10.1145/2063576.2063684

[4] John S. Breese, David Heckerman, and Carl Kadie. 2013. Empirical Analysis of Predictive Algorithms for Collaborative Filtering. arXiv:1301.7363 [cs.IR]

[5] Bisheng Chen, Jingdong Wang, Qinghua Huang, and Tao Mei. 2012. Personalized video recommendation through tripartite graph propagation.. In *Proceedings of the 20th ACM Multimedia Conference, MM '12, Nara, Japan, October 29 - November 02, 2012*.

[6] Jingyuan Chen, Xuemeng Song, Liqiang Nie, Xiang Wang, Hanwang Zhang, and Tat-Seng Chua. 2016. Micro Tells Macro: Predicting the Popularity of Micro-Videos via a Transductive Model.. In *Proceedings of the 2016 ACM Conference on Multimedia Conference, MM 2016, Amsterdam, The Netherlands, October 15-19, 2016*.

[7] Jingyuan Chen, Hanwang Zhang, Xiangnan He, Liqiang Nie, Wei Liu, and Tat-Seng Chua. 2017. Attentive collaborative filtering: Multimedia recommendation with item-and component-level attention. In *Proceedings of the 40th International ACM SIGIR conference on Research and Development in Information Retrieval*. 335–344.

[8] Xusong Chen, Dong Liu, Zhengjun Zha, W. Zhou, Zhiwei Xiong, and Y. Li. 2018. Temporal Hierarchical Attention at Category- and Item-Level for Micro-Video Click-Through Prediction. *Proceedings of the 26th ACM international conference on Multimedia* (2018).

[9] Xusong Chen, Dong Liu, Zheng-Jun Zha, Wengang Zhou, Zhiwei Xiong, and Yan Li. 2018. Temporal Hierarchical Attention at Category- and Item-Level for Micro-Video Click-Through Prediction.. In *2018 ACM Multimedia Conference on Multimedia Conference, MM 2018, Seoul, Republic of Korea, October 22-26, 2018*.

[10] Paul Covington, Jay Adams, and Emre Sargin. 2016. Deep Neural Networks for YouTube Recommendations. In *Proceedings of the 10th ACM Conference on Recommender Systems*. New York, NY, USA.

[11] Peng Cui, Zhiyu Wang, and Zhou Su. 2014. What Videos Are Similar with You?: Learning a Common Attributed Representation for Video Recommendation.. In *Proceedings of the ACM International Conference on Multimedia, MM '14, Orlando, FL, USA, November 03 - 07, 2014*.

[12] Yashar Deldjoo, Mehdi Elahi, Paolo Cremonesi, Franca Garzotto, Pietro Piazzolla, and Massimo Quadrana. 2016. Content-based video recommendation system based on stylistic visual features. *Journal on Data Semantics* 5, 2 (2016), 99–113.

[13] Mukund Deshpande and George Karypis. 2004. Item-based top-n recommendation algorithms. *ACM Transactions on Information Systems (TOIS)* 22, 1 (2004), 143–177.

[14] Tommaso Di Noia, Vito Claudio Ostuni, Jessica Rosati, Paolo Tomeo, and Eugenio Di Sciascio. 2014. An Analysis of Users' Propensity toward Diversity in Recommendations. In *Proceedings of the 8th ACM Conference on Recommender Systems* (Foster City, Silicon Valley, California, USA) *(RecSys '14)*. Association for Computing Machinery, New York, NY, USA, 285–288. https://doi.org/10.1145/2645710.2645774

[15] Yi Ding and Xue Li. 2005. Time weight collaborative filtering. In *Proceedings of the 14th ACM international conference on Information and knowledge management*. 485–492.

[16] Jianfeng Dong, Xirong Li, Chaoxi Xu, Gang Yang, and Xun Wang. 2018. Feature Re-Learning with Data Augmentation for Content-based Video Recommendation.. In *2018 ACM Multimedia Conference on Multimedia Conference, MM 2018, Seoul, Republic of Korea, October 22-26, 2018*.

[17] Xiaoyu Du, Xiang Wang, Xiangnan He, Zechao Li, Jinhui Tang, and Tat-Seng Chua. 2020. How to Learn Item Representation for Cold-Start Multimedia Recommendation?. In *MM '20: The 28th ACM International Conference on Multimedia, Virtual Event / Seattle, WA, USA, October 12-16, 2020*.

[18] Wenqi Fan, Yao Ma, Qing Li, Yuan He, Eric Zhao, Jiliang Tang, and Dawei Yin. 2019. Graph Neural Networks for Social Recommendation. In *The World Wide Web Conference* (San Francisco, CA, USA) *(WWW '19)*. Association for Computing Machinery, New York, NY, USA, 417–426. https://doi.org/10.1145/3308558.3313488

[19] Cricia Felicio, Klérisson Paixão, Guilherme Alves, Sandra Amo, and Philippe Preux. 2016. Exploiting Social Information in Pairwise Preference Recommender System. *Journal of Information and Data Management* 7 (08 2016), 99.

[20] J. Guo, Y. Zhu, A. Li, Q. Wang, and W. Han. 2016. A Social Influence Approach for Group User Modeling in Group Recommendation Systems. *IEEE Intelligent Systems* 31, 5 (2016), 40–48.

[21] Kaiming He, Haoqi Fan, Yuxin Wu, Saining Xie, and Ross B. Girshick. 2019. Momentum Contrast for Unsupervised Visual Representation Learning. *CoRR* abs/1911.05722 (2019). arXiv:1911.05722 http://arxiv.org/abs/1911.05722

[22] Ruining He, Chen Fang, Zhaowen Wang, and Julian J. McAuley. 2016. Vista: A Visually, Socially, and Temporally-aware Model for Artistic Recommendation. *CoRR* abs/1607.04373 (2016). arXiv:1607.04373 http://arxiv.org/abs/1607.04373

[23] Ruining He and Julian McAuley. 2016. Fusing similarity models with markov chains for sparse sequential recommendation. In *2016 IEEE 16th International Conference on Data Mining (ICDM)*. IEEE, 191–200.

[24] Xiangnan He, Kuan Deng, Xiang Wang, Yan Li, Yong-Dong Zhang, and Meng Wang. 2020. LightGCN: Simplifying and Powering Graph Convolution Network for Recommendation.. In *Proceedings of the 43rd International ACM SIGIR conference on research and development in Information Retrieval, SIGIR 2020, Virtual Event, China, July 25-30, 2020*.

[25] Xiangnan He, Zhankui He, Jingkuan Song, Zhenguang Liu, Yu-Gang Jiang, and Tat-Seng Chua. 2018. Nais: Neural attentive item similarity model for recommendation. *IEEE Transactions on Knowledge and Data Engineering* 30, 12 (2018), 2354–2366.

[26] Xiangnan He, Lizi Liao, Hanwang Zhang, Liqiang Nie, Xia Hu, and Tat-Seng Chua. 2017. Neural Collaborative Filtering. *CoRR* abs/1708.05031 (2017). arXiv:1708.05031 http://arxiv.org/abs/1708.05031

[27] Xiangnan He, Lizi Liao, Hanwang Zhang, Liqiang Nie, Xia Hu, and Tat-Seng Chua. 2017. Neural collaborative filtering. In *Proceedings of the 26th international conference on world wide web*. 173–182.

[28] Balázs Hidasi, Alexandros Karatzoglou, Linas Baltrunas, and Domonkos Tikk. 2015. Session-based Recommendations with Recurrent Neural Networks. arXiv:1511.06939 [cs.LG]

[29] Yanxiang Huang, Bin Cui, Jie Jiang, Kunqian Hong, Wenyu Zhang, and Yiran Xie. 2016. Real-time Video Recommendation Exploration.. In *Proceedings of the 2016 International Conference on Management of Data, SIGMOD Conference 2016, San Francisco, CA, USA, June 26 - July 01, 2016*.

[30] Hao Jiang, Wenjie Wang, Yinwei Wei, Zan Gao, Yinglong Wang, and Liqiang Nie. 2020. What Aspect Do You Like: Multi-scale Time-aware User Interest Modeling for Micro-video Recommendation.. In *MM '20: The 28th ACM International Conference on Multimedia, Virtual Event / Seattle, WA, USA, October 12-16, 2020*.

[31] Hao Jiang, Wenjie Wang, Yinwei Wei, Zan Gao, Yinglong Wang, and Liqiang Nie. 2020. What Aspect Do You Like: Multi-Scale Time-Aware User Interest Modeling for Micro-Video Recommendation. In *Proceedings of the 28th ACM International Conference on Multimedia* (Seattle, WA, USA) *(MM '20)*. Association for Computing Machinery, New York, NY, USA, 3487–3495. https://doi.org/10.

1145/3394171.3413653

[32] Wang-Cheng Kang and Julian McAuley. 2018. Self-attentive sequential recommendation. In *2018 IEEE International Conference on Data Mining (ICDM)*. IEEE, 197–206.

[33] Wang-Cheng Kang and Julian J. McAuley. 2018. Self-Attentive Sequential Recommendation.. In *IEEE International Conference on Data Mining, ICDM 2018, Singapore, November 17-20, 2018*.

[34] Younghoon Kim and Kyuseok Shim. 2014. TWILITE: A recommendation system for Twitter using a probabilistic model based on latent Dirichlet allocation. *Information Systems* 42 (2014), 59–77.

[35] Joseph A Konstan, Bradley N Miller, David Maltz, Jonathan L Herlocker, Lee R Gordon, and John Riedl. 1997. Grouplens: Applying collaborative filtering to usenet news. *Commun. ACM* 40, 3 (1997), 77–87.

[36] Yehuda Koren. 2009. Collaborative filtering with temporal dynamics. In *Proceedings of the 15th ACM SIGKDD international conference on Knowledge discovery and data mining*. 447–456.

[37] Yehuda Koren and Robert Bell. 2015. Advances in collaborative filtering. *Recommender systems handbook* (2015), 77–118.

[38] Chao Li, Zhiyuan Liu, Mengmeng Wu, Yuchi Xu, Huan Zhao, Pipei Huang, Guoliang Kang, Qiwei Chen, Wei Li, and Dik Lun Lee. 2019. Multi-Interest Network with Dynamic Routing for Recommendation at Tmall.. In *Proceedings of the 28th ACM International Conference on Information and Knowledge Management, CIKM 2019, Beijing, China, November 3-7, 2019*.

[39] Yongqi Li, Meng Liu, Jianhua Yin, Chaoran Cui, Xin-Shun Xu, and Liqiang Nie. 2019. Routing Micro-videos via A Temporal Graph-guided Recommendation System.. In *Proceedings of the 27th ACM International Conference on Multimedia, MM 2019, Nice, France, October 21-25, 2019*.

[40] Yongqi Li, Meng Liu, Jianhua Yin, Chaoran Cui, Xin-Shun Xu, and Liqiang Nie. 2019. Routing Micro-Videos via A Temporal Graph-Guided Recommendation System. In *Proceedings of the 27th ACM International Conference on Multimedia* (Nice, France) *(MM '19)*. Association for Computing Machinery, New York, NY, USA, 1464–1472. https://doi.org/10.1145/3343031.3350950

[41] Zhaopeng Li, Qianqian Xu, Yangbangyan Jiang, Xiaochun Cao, and Qingming Huang. 2020. Quaternion-Based Knowledge Graph Network for Recommendation.. In *MM '20: The 28th ACM International Conference on Multimedia, Virtual Event / Seattle, WA, USA, October 12-16, 2020*.

[42] Yujie Lu, Shengyu Zhang, Yingxuan Huang, Luyao Wang, Xinyao Yu, Zhou Zhao, and Fei Wu. 2020. Future-Aware Diverse Trends Framework for Recommendation. *CoRR* (2020).

[43] Hao Ma. 2013. An experimental study on implicit social recommendation. 73–82. https://doi.org/10.1145/2484028.2484059

[44] Tao Mei, Bo Yang, Xian-Sheng Hua, and Shipeng Li. 2011. Contextual Video Recommendation by Multimodal Relevance and User Feedback. *ACM Trans. Inf. Syst.* (2011).

[45] Subhabrata Mukherjee and Stephan Guennemann. 2019. GhostLink: Latent Network Inference for Influence-aware Recommendation. arXiv:1905.05955 [cs.SI]

[46] James R Norris and James Robert Norris. 1998. *Markov chains*. Number 2. Cambridge university press.

[47] Manos Papagelis, Dimitris Plexousakis, and Themistoklis Kutsuras. 2005. Alleviating the sparsity problem of collaborative filtering using trust inferences. In *International conference on trust management*. Springer, 224–239.

[48] Jonghun Park. 2010. An online video recommendation framework using view based tag cloud aggregation. *IEEE Multimedia, 2010* (2010).

[49] Wichian Premchaiswadi, Pitaya Poompuang, Nipat Jongswat, and Nucharee Premchaiswadi. 2013. Enhancing Diversity-Accuracy Technique on User-Based Top-N Recommendation Algorithms. In *Proceedings of the 2013 IEEE 37th Annual Computer Software and Applications Conference Workshops (COMPSACW '13)*. IEEE Computer Society, USA, 403–408. https://doi.org/10.1109/COMPSACW.2013.68

[50] Yanru Qu, Han Cai, Kan Ren, Weinan Zhang, Yong Yu, Ying Wen, and Jun Wang. 2016. Product-based neural networks for user response prediction. In *2016 IEEE 16th International Conference on Data Mining (ICDM)*. IEEE, 1149–1154.

[51] Steffen Rendle, Christoph Freudenthaler, Zeno Gantner, and Lars Schmidt-Thieme. 2012. BPR: Bayesian Personalized Ranking from Implicit Feedback. *CoRR* abs/1205.2618 (2012). arXiv:1205.2618 http://arxiv.org/abs/1205.2618

[52] 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.

[53] Ruslan Salakhutdinov and Andriy Mnih. 2008. Probabilistic matrix factorization. *Advances in Neural Information Processing Systems* 20 (2008), 1257–1264.

[54] Badrul Sarwar, George Karypis, Joseph Konstan, and John Riedl. 2001. Item-Based Collaborative Filtering Recommendation Algorithms. In *Proceedings of the 10th International Conference on World Wide Web* (Hong Kong, Hong Kong) *(WWW '01)*. Association for Computing Machinery, New York, NY, USA, 285–295. https://doi.org/10.1145/371920.372071

[55] Guy Shani, David Heckerman, Ronen I Brafman, and Craig Boutilier. 2005. An MDP-based recommender system. *Journal of Machine Learning Research* 6, 9 (2005).

[56] Elena Smirnova and Flavian Vasile. 2017. Contextual sequence modeling for recommendation with recurrent neural networks. In *Proceedings of the 2nd workshop on deep learning for recommender systems*. 2–9.

[57] 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, CIKM 2019, Beijing, China, November 3-7, 2019*.

[58] Jiaxi 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, WSDM 2018, Marina Del Rey, CA, USA, February 5-9, 2018*.

[59] Jiaxi 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.

[60] Loren Terveen and Will Hill. 2001. Beyond recommender systems: Helping people help each other. *HCI in the New Millennium* 1, 2001 (2001), 487–509.

[61] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Lukasz Kaiser, and Illia Polosukhin. 2017. Attention is All you Need. In *NIPS*.

[62] Jizhe Wang, Pipei Huang, Huan Zhao, Zhibo Zhang, Binqiang Zhao, and Dik Lun Lee. 2018. Billion-scale Commodity Embedding for E-commerce Recommendation in Alibaba.. In *Proceedings of the 24th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining, KDD 2018, London, UK, August 19-23, 2018*.

[63] Pengfei Wang, Jiafeng Guo, Yanyan Lan, Jun Xu, Shengxian Wan, and Xueqi Cheng. 2015. Learning hierarchical representation model for nextbasket recommendation. In *Proceedings of the 38th International ACM SIGIR conference on Research and Development in Information Retrieval*. 403–412.

[64] Peng Wang, Yunsheng Jiang, Chunxu Xu, and Xiaohui Xie. 2019. Overview of Content-Based Click-Through Rate Prediction Challenge for Video Recommendation. In *Proceedings of the 27th ACM International Conference on Multimedia*. 2593–2596.

[65] Xiang Wang, Xiangnan He, Meng Wang, Fuli Feng, and Tat-Seng Chua. 2019. Neural Graph Collaborative Filtering.. In *Proceedings of the 42nd International ACM SIGIR Conference on Research and Development in Information Retrieval, SIGIR 2019, Paris, France, July 21-25, 2019*.

[66] Yinwei Wei, Xiang Wang, Weili Guan, Liqiang Nie, Zhouchen Lin, and Baoquan Chen. 2019. Neural multimodal cooperative learning toward micro-video understanding. *IEEE Transactions on Image Processing* 29 (2019), 1–14.

[67] Yinwei Wei, Xiang Wang, Liqiang Nie, Xiangnan He, and Tat-Seng Chua. 2020. Graph-Refined Convolutional Network for Multimedia Recommendation with Implicit Feedback.. In *MM '20: The 28th ACM International Conference on Multimedia, Virtual Event / Seattle, WA, USA, October 12-16, 2020*.

[68] Yinwei Wei, Xiang Wang, Liqiang Nie, Xiangnan He, Richang Hong, and Tat-Seng Chua. 2019. MMGCN: Multi-modal Graph Convolution Network for Personalized Recommendation of Micro-video.. In *Proceedings of the 27th ACM International Conference on Multimedia, MM 2019, Nice, France, October 21-25, 2019*.

[69] Jibang Wu, Renqin Cai, and Hongning Wang. 2020. DéJà vu: A Contextualized Temporal Attention Mechanism for Sequential Recommendation *(WWW '20)*. Association for Computing Machinery, New York, NY, USA, 11 pages. https://doi.org/10.1145/3366423.3380285

[70] Shu Wu, Yuyuan Tang, Yanqiao Zhu, Liang Wang, Xing Xie, and Tieniu Tan. 2019. 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*.

[71] Shu Wu, Yuyuan Tang, Yanqiao Zhu, Liang Wang, Xing Xie, and Tieniu Tan. 2019. Session-based recommendation with graph neural networks. In *Proceedings of the AAAI Conference on Artificial Intelligence*, Vol. 33. 346–353.

[72] Ming Yan, Jitao Sang, and Changsheng Xu. 2015. Unified YouTube Video Recommendation via Cross-network Collaboration.. In *Proceedings of the 5th ACM on International Conference on Multimedia Retrieval, Shanghai, China, June 23-26, 2015*.

[73] Xuewen Yang, Dongliang Xie, Xin Wang, Jiangbo Yuan, Wanying Ding, and Pengyun Yan. 2020. Learning Tuple Compatibility for Conditional Outfit Recommendation.. In *MM '20: The 28th ACM International Conference on Multimedia, Virtual Event / Seattle, WA, USA, October 12-16, 2020*.

[74] Xuzheng Yu, Tian Gan, Yinwei Wei, Zhiyong Cheng, and Liqiang Nie. 2020. Personalized Item Recommendation for Second-hand Trading Platform.. In *MM '20: The 28th ACM International Conference on Multimedia, Virtual Event / Seattle, WA, USA, October 12-16, 2020*.

[75] Yuyu Zhang, Hanjun Dai, Chang Xu, Jun Feng, Taifeng Wang, Jiang Bian, Bin Wang, and Tie-Yan Liu. 2014. Sequential Click Prediction for Sponsored Search with Recurrent Neural Networks. *CoRR* abs/1404.5772 (2014). arXiv:1404.5772 http://arxiv.org/abs/1404.5772

[76] Guorui Zhou, Xiaoqiang Zhu, Chenru Song, Ying Fan, Han Zhu, Xiao Ma, Yanghui Yan, Junqi Jin, Han Li, and Kun Gai. 2018. Deep interest network for click-through

rate prediction. In *Proceedings of the 24th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining.* 1059–1068.

[77] Xiangmin Zhou, Lei Chen, Yanchun Zhang, Longbing Cao, Guangyan Huang, and Chen Wang. 2015. Online Video Recommendation in Sharing Community.. In *Proceedings of the 2015 ACM SIGMOD International Conference on Management of Data, Melbourne, Victoria, Australia, May 31 - June 4, 2015.*