# Cross-domain User Preference Learning for Cold-start Recommendation

Huiling Zhou*, Jie Liu*, Zhikang Li*, Jin Yu, Hongxia Yang

{zhule.zhl,sanshuai.lj,zhikang.lzk,kola.yu,yang.yhx}@alibaba-inc.com

DAMO Academy, Alibaba Group

China

## ABSTRACT

Cross-domain cold-start recommendation is an increasingly emerging issue for recommender systems. Existing works mainly focus on solving either cross-domain user recommendation or cold-start content recommendation. However, when a new domain evolves at its early stage, it has potential users similar to the source domain but with much fewer interactions. It is critical to learn a user's preference from the source domain and transfer it into the target domain, especially on the newly arriving contents with limited user feedback. To bridge this gap, we propose a self-trained Cross-dOmain User Preference LEarning (COUPLE) framework, targeting cold-start recommendation with various *semantic tags*, such as attributes of items or genres of videos. More specifically, we consider three levels of preferences, including user history, user content and user group to provide reliable recommendation. With user history represented by a domain-aware sequential model, a frequency encoder is applied to the underlying tags for user content preference learning. Then, a hierarchical memory tree with orthogonal node representation is proposed to further generalize user group preference across domains. The whole framework updates in a contrastive way with a First-In-First-Out (FIFO) queue to obtain more distinctive representations. Extensive experiments on two datasets demonstrate the efficiency of COUPLE in both user and content cold-start situations. By deploying an online A/B test for a week, we show that the Click-Through-Rate (CTR) of COUPLE is superior to other baselines used on Taobao APP. Now the method is serving online for the cross-domain cold micro-video recommendation.

## CCS CONCEPTS

• **Computer systems organization** → **Embedded systems**; *Redundancy*; Robotics; • **Networks** → Network reliability.

## KEYWORDS

cross-domain, cold-start recommendation, user preference learning
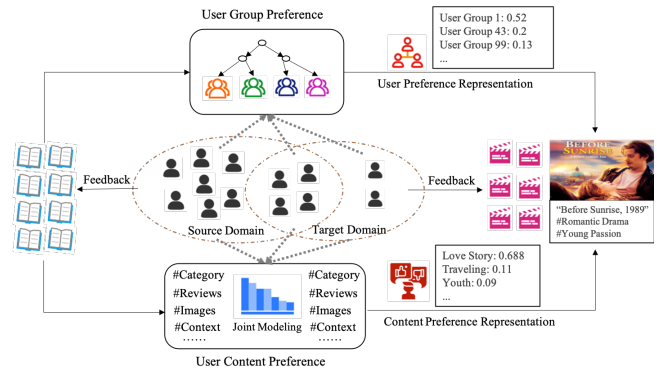


**Figure 1: Overview of the proposed cross-domain cold-start recommendation system. A user's preference is captured from different perspectives with jointly modeling behaviors in all domains.**

## 1 INTRODUCTION

Personalized recommendation system (RS) plays a vital role in an e-commerce platform, where effective strategies have been made to alleviate information overload and facilitate better user experiences. In recent years, deep learning has been widely applied in RS to overcome obstacles of conventional recommendation techniques. Great effort has been made to achieve better performances in heterogeneous multi-modal recommendation [11, 13, 38], debiased recommendation [14, 48], reliable and explainable recommendation [6, 43], etc. Among all, cross-domain recommendation [10, 20, 46] and cold-start recommendation [7, 30, 34] problems draw a lot of attention. Cross-domain recommendation systems aim to transfer knowledge available in other domains (known as the source domain) to the target domain where users have much sparser feedback or interactions. Overlapped users are often selected to learn the mapping of interest between two domains so that same pattern can be applied to those cold-start users of target domain [20, 46]. On the other hand, cold-start recommendation is often referred to

cold-start content recommendation where no or few user feedback can be found when new items arrive. It is challenging to recommend such items to users due to the unavailability of the statistical information which RS relies heavily on. Common solutions extract and analyze content information of items to match users' potential interests based on their historical behaviors [7, 13].

Previous works normally treat cold-start user and cold-start content recommendation as two separate scenarios. However, when a new domain evolves at its early stage in real-world, it may suffer both issues at the same time. For example, Taobao is the largest e-commerce platform in China where billions of products are interacted with hundreds of millions of users each day. When micro-videos were first served to the platform, they shared the same set of users as product domain but had few (or none) user interactions. How to perform reliable recommendation on such emerging new contents appears to be one of the most challenging problems for RS. A recent work [44] proposes an internal contextual attention network to deal with the cross-domain cold-start recommendation on contents with very few interactions. However, this approach has limited ability to model brand new contents which are never seen before.

In this paper, we further extend the work [44] and provide a unified Cross-dOmain User Preference LEarning (COUPLE) framework where items with limited interactions can be reliably recommended. COUPLE differs from previous cross-domain works in that we jointly model all users' behaviors across domains instead of using just the overlapped users, as shown in Figure.1. With historical feedback modeled, a user's preference across domains is obtained from content perspective and user group perspective. When a new item comes, we push it to the corresponding users with matched interest based on its content feature and the user's general preference learnt by the system. Main contributions of our proposed work are summarized as follows:

- We propose a unified user preference learning framework to solve cross-domain cold-start problem, which delicately models a user's history preference, content preference and group preference. With proper choice of a First-In-First-Out (FIFO) queue, the whole framework can be self-trained in an efficient contrastive way.
- With semantic tags extracted for item representation, a frequency encoder based on cross-layer attention is utilized. It further boosts the interactions between high-level representation with low-level content source input. To the best of our knowledge, it is the pioneering work to solve the cross-domain cold-start issue with the use of tag information.
- A hierarchical memory tree with orthogonality property is proposed to learn users' preference across domains. It can generalize user's diverse interest with top $K$ leaf-node representations based on his/her historical behaviors. This is proved to be a good way of user preference generalization in cross-domain scenario.

The rest of the paper is organized as follows. In Section 2, we review the related work. Section 3 presents the details of our proposed method. Offline experimental results and ablation study are presented in Section 4. We introduce the online deployment of the system in Section 5 and conclude our work in Section 6.

## 2 RELATED WORK

In this section, we introduce the related work on personalized recommendation system, user preference learning, and the orthogonal regularization mechanism we use in the paper.

### 2.1 Personalized Recommendation System

A standard matching (i.e. candidate generation) stage of a personalized recommender system aims to retrieve a small subset of items from the huge pool for the sophisticated ranking scheme. Previous methods typically employ collaborative filtering (CF) strategies to exploit user preferences from their explicit or implicit feedback [19, 22]. These methods achieve good performance where user-item interactions are dense but are not able to handle the cold-start situations where the interactions are limited. Content-based methods alleviate this problem by mining the content features of items or meta information about users [11–13, 16, 37]. Among them, tags are proved to be useful [12, 16], especially for cross-domain recommendation where semantically similar tags may exist in both domains and bridge the gap to some extent. Despite their success, both types of the methods treat user-item interactions in a static way and cannot capture the dynamic interests of users over time. As a result, sequential recommender systems are proposed to model the sequential nature of users' interactions [5, 32, 35]. Besides, attention mechanism has been widely adapted to recommender systems since its out-breaking success in NLP [9, 33]. It is leveraged between users and contents to capture the most representative information or provide good interpretability for recommendation [3, 27, 29].

In real-world applications, hybrid recommendation framework is usually applied which makes up for the shortcomings of using a single method and maximizes the advantages. Our work is typically within the hybrid system for cross-domain cold-start recommendation, where all strategies mentioned above are involved.

### 2.2 User Preference Learning

The core of a recommender system is to model the dynamic user preferences. Previous works usually represent a user's historical behavior with one single latent vector which may suffer from correlation loss [5]. Clustering-based approaches offer an alternative to traditional model-based methods and cluster similar users or items together [31, 40]. However, these methods rely heavily on manual selection of a user's attribute and profile, or treat different domains separately. Recently, recommendations with memory network [5, 49] are proposed to model user preference and store the long term interest. The basic idea is to maintain a key-value style memory matrix with numbers of slots. A user's embedding is fed into the matrix where similarity between the user and each key vector is computed and converted to relevance probability using softmax function. In our work, a hierarchical memory tree is utilized for user preference learning. Different from previous memory-based network, our tree structure has a much larger number of memory slots (i.e. few hundred or thousand vs. few dozens). What is more, we only pick the top memory slots each time with calculated correlation score instead of using all the memory slots.

## 2.3 Orthogonal Regularization

The orthogonality implies energy preservation and is proved to be efficient for stabilizing the distribution of activations over layers with CNNs. Orthogonal regularizers are extended to fully-connected layers and a Spectral Restricted Isometry Property (SRIP) regularizer is proposed to guarantee better convergence [2]. Later on, more investigations [4, 26] further extend SRIP regularization bounds and bring compatible results on applications like image classification and Person Re-Identification. We also apply such orthogonal regularization scheme for the user group preference learning, which will be illustrated in detail in Section 3.

## 3 THE PROPOSED METHOD

In this section, we present the problem formulation of cross-domain cold-start recommendation first. Then we explain in detail how we model user preference from three different aspects in COUPLE and train the framework efficiently in a contrastive way. The overall design of our proposed COUPLE network is shown in Figure 2.

## 3.1 Problem Formulation and Notations

*3.1.1 Problem Formulation.* Cross-domain recommendation in real-world may involve several channels [44]. Without loss of generality, we classify those channels into source domain(s) $A$ and target domain(s) $B$, where user interactions are much denser in $A$ compared to $B$. We denote $\mathcal{U}_A$, $\mathcal{U}_B$ as the sets of users and $\mathcal{I}_A$, $\mathcal{I}_B$ as the sets of interacted items with content features (semantic tags in our case) respectively. We also have a cold-start item set in target domain denoted as $\mathcal{I}'_B$ where items are with limited (or none) interactions. For the cross-domain cold-start recommendation, we have a) user overlap $\mathcal{U}_A \cap \mathcal{U}_B \neq \emptyset$ ; b) item overlap $\mathcal{I}_A \cap \mathcal{I}_B = \emptyset$ ; c) cold-start item overlap $\mathcal{I}_B \cap \mathcal{I}'_B = \emptyset$. We define two tasks with regard to the training and inference processes of COUPLE, respectively:

- **Joint recommendation for training**, i.e., make recommendation on items in $\mathcal{I}_A \cup \mathcal{I}_B$ to users $\mathcal{U}_A \cup \mathcal{U}_B$.
- **Cold-start recommendation for inference**, i.e., make recommendation on items in $\mathcal{I}'_B$ to users $\mathcal{U}_A \cup \mathcal{U}_B$.

It is noted that our work can be easily extended to recommendation on items with limited interactions as well. Notations are summarized in Table 1.

### Table 1: Notations Used in Section 3.

| Notation | Description |
|---|---|
| $\mathcal{U}, \mathcal{I}, \mathcal{T}$ | user, item and tag set |
| $d \in \mathbb{N}$ | latent vector dimension |
| $n \in \mathbb{N}$ | number of tags per item |
| $l \in \mathbb{N}$ | max sequence length |
| $e_u, e_i, e_t \in \mathbb{R}^{d \times 1}$ | user, item, tag embedding feature |
| $e_u^{(h)}, e_u^{(c)}, e_u^{(g)} \in \mathbb{R}^{d \times 1}$ | user history, content, group preference |
| $s_{i,j} \in \mathbb{R}^{d \times 1}$ | tree node vector at position $j$ in layer $i$ |
| $\mathbf{M}$ | attention map |
| $\mathbf{W}$ | weight matrix |
| $\mathcal{L}$ | loss |
| $\|.\|$ | norm |
| $[.]$ | concatenation |

## 3.2 Cross-domain User Preference Learning

As illustrated in Figure 2, we model a user's preference from three interactive stages. First, user's history preference is modeled with a domain-aware multi-head attention scheme for sequential behaviors from different domains. Then it will attend to the underlying tag-level features to encode tag importance for user content preference. A hierarchical memory tree is used to represent user's historical preference with "orthogonal" leaf-node vectors. Finally, user preferences from three aspects are weighted aggregated.

*3.2.1 User History Preference.* We model a user's historical behavior using a deep sequential model, where item-level attention scheme is applied. In standard sequential models, a sequence of item index is often given for embedding lookup [27, 35]. However, in cold-start content recommendation, items may be never seen or interacted before. To solve this issue, we use the global semantic tags to represent each item. To be more specific, we obtain a cross-domain tag set $\mathcal{T}$ from items across domains $\mathcal{I}_A \cup \mathcal{I}_B$. The tags can be extracted from hashtag made by users, item attributes or genres, or from images and texts pre-processed by vision and language models. A fair assumption can be made that the tag set $\mathcal{T}$ is able to cover most of the items in $\mathcal{I}'_B$ from the target domain.

Our input for each item is a group of tags $[t_1, t_2, \ldots, t_n]$, where $t_x \in \{1, 2, \ldots, N\}$ is the index of the tag in the whole set $\mathcal{T}$ with size $N$. A tag embedding table $\mathbf{H}_t \in \mathbb{R}^{d \times N}$ takes tag index sequence as the input and outputs the representation of the tag features $[e_{t,1}, e_{t,2}, \ldots, e_{t,n}]$, where $e_t \in \mathbb{R}^{d \times 1}$. In this way, we represent an item $e_i \in \mathbb{R}^{d \times 1}$ by aggregating the tag features by average pooling:

$$e_i = \text{AvgPool}(\{e_t\}_{k=1}^n) = \frac{1}{n} \sum_{k=1}^{n} e_{t,k}, \quad (1)$$

where $n$ is the number of tag features capped for each item. Note that the index-based embedding lookup scheme can be easily extended to using multi-modal embeddings from any pre-trained models. To explicitly involve the domain-specific knowledge, we follow the design in Bert [9] and add a domain embedding table $\mathbf{H}_d \in \mathbb{R}^{d \times Q}$ ($Q$ is the number of channels) to the current item feature embedding, for training. For inference and online serving, we only use the tag features for item embedding, as suggested by the recent recommendation work [48].

To model user history preference, we apply a multi-head attention (MHAttn) mechanism [27] to the sequence of temporally sorted items $[e_{i,1}, e_{i,2}, \ldots, e_{i,l}]$ clicked (rated) by users. Self-attention is applied to aggregate all item embeddings with adaptive weights and followed by two-layer feed-forward networks (FFN) to increase non-linearity.

$$(e_{u,1}, \ldots, e_{u,m}) = \text{MHAttn}(\{e_{i,k}\}_{k=1}^l), \quad (2)$$

where $l$ is the max sequence length and $m$ is the number of attention heads. Furthermore, we perform weighted aggregation (WA) to obtain user history preference $e_u^{(h)}$ on the output representations $\{e_u\}_{k=1}^m$ of Eq. (2) as:

$$e_u^{(h)} = \text{WA}(\{e_u\}_{k=1}^m) = \text{WA}([e_{u,1}, \ldots, e_{u,m}]), \quad (3)$$
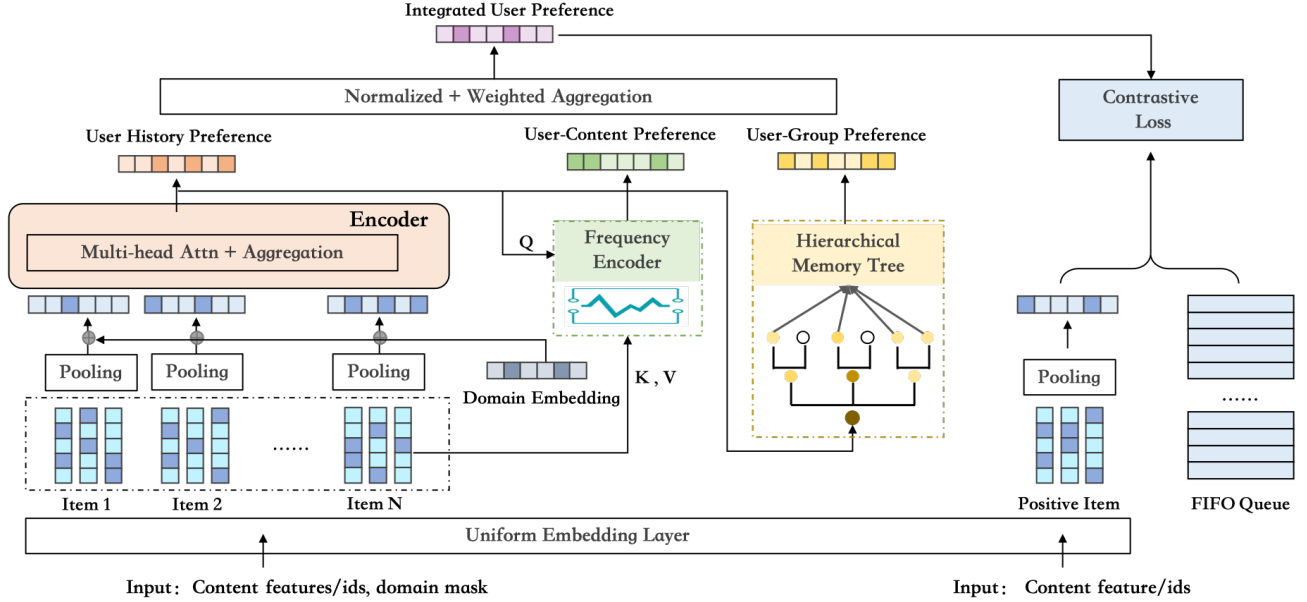
**Figure 2: Framework of our proposed method. Three levels of user preferences are modeled, including user history, user content and user group. The whole framework updates in a contrastive way with a Fisrt-In-First-Out (FIFO) queue, where a large number of negative samples from latest steps is maintained for more consistant training.**

with a parameter matrix $\mathbf{M} \in \mathbb{R}^{d \times d}$ learnt and updated, the following procedure is performed:

$$
\begin{aligned}
\boldsymbol{e}_m &= \frac{1}{m} \sum_{k=1}^{m} \boldsymbol{e}_{u,k}, \quad \boldsymbol{d}_k = \boldsymbol{e}_{u,k}^T \mathbf{M} \boldsymbol{e}_m, \\
a_k &= \frac{\exp(\boldsymbol{d}_k)}{\sum_{k'=1}^{m} \exp(\boldsymbol{d}_{k'})}, \quad \boldsymbol{e}_u^{(h)} = \sum_{k=1}^{m} a_k \boldsymbol{e}_{u,k}.
\end{aligned}
\tag{4}
$$

*3.2.2 User Content Preference.* User history feature obtained from Section 3.2.1 focuses on the item-level attention, where each tag contributes equally for an item representation. However, same tags appearing multiple times in the item sequence are supposed to contribute more to a user's interest. The similar conclusion on the importance of item frequency for next-bucket recommendation has been made recently in [21]. As a result, we propose a frequency encoder which emphasizes the appearances of the same tags. It is implemented by a cross-layer attention process.

We concatenate all the tag embeddings in the item sequence together as $\{\boldsymbol{e}_t\}_{k=1}^{n \times l} = [\boldsymbol{e}_{t,1}, \boldsymbol{e}_{t,2}, \ldots, \boldsymbol{e}_{t,n \times l}]$, and let them go through a one-layer MLP with a hyperbolic tangent function. This results in a hidden representation $\boldsymbol{h}_t$. Then the softmax function is applied between $\boldsymbol{h}_t$ and the user's history representation $\boldsymbol{e}_u^{(h)}$ to obtain the attention weights. The user content representation $\boldsymbol{e}_u^{(c)}$ is computed

as the weighted sum of attention weights and the tag features:

$$
\boldsymbol{h}_{t,i} = \tanh(\mathbf{W}_w \boldsymbol{e}_{t,i} + b_w), \quad a_i = \frac{\exp\left(\boldsymbol{e}_u^{(h)} \cdot \boldsymbol{h}_{t,i}\right)}{\sum_{i'=1}^{n \times l} \exp\left(\boldsymbol{e}_u^{(h)} \cdot \boldsymbol{h}_{t,i'}\right)},
$$

$$
\boldsymbol{e}_u^{(c)} = \sum_{i=1}^{n \times l} a_i \boldsymbol{e}_{t,i}.
\tag{5}
$$

In this way, the same tag with multiple attention weights will be augmented in the final representation of user content preference. A simple illustration of how this frequency encoder works is shown in Figure.3. From the design point of view, the user history representation is modeled based on a bottom-up attention while the user content preference is in a top-down fashion. Such cross-layer attention is proved to be critical to boost interactions between the high-level representation with the low-level content source [1, 33].

*3.2.3 User Group Preference.* Unlike previous cross-domain works where only the overlapped users are considered, we take into account all the users' behavior patterns and try to find out a uniform representation for user group preference across domains. Consequently, a hierarchical memory tree is designed to remember a user's general interest and automatically classify it into groups with closest interests. It differs from the original memory networks in that : a) we are able to have more memory slots as proved in the recent work [41] (i.e. few hundreds or thousands vs. few dozens); and b) we don't use all slots as memory networks do, instead we pick the top $K$ leaf-node slots each time with highest scores.
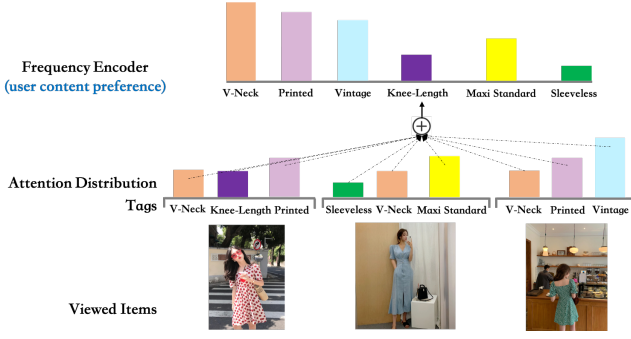
**Figure 3: Tag-level Frequency Encoder. Representation of the same tag will be aggregated and augmented after user-tag attention scheme (best viewed in color).**

**Tree Structure Design.** With user history representation from equation (4), we follow the similar procedure as [41] to use a hierarchically fully-connected tree where the root node allocates weights to children nodes and finally leaf nodes. To enforce further interpretability, the sum of the node weights at the same layer is equal to 1 and the sum of children node weights inherited from the same parent node is equal to the parent node weight. Given user history representation $e_u^{(h)}$ and a parent node at position $j$ in layer $i$ with $h$ children nodes, whose weight is denoted as $w_{i,j}$, the weights of the children nodes in layer $i + 1$ can be calculated using softmax function:

$$
\begin{aligned}
w_{i+1,k} &= w_{i,j} \cdot \text{softmax}(e_u^{(h)} \cdot s_{i+1,k}) \\
&= w_{i,j} \cdot \frac{\exp(e_u^{(h)} \cdot s_{i+1,k})}{\sum_l \exp(e_u^{(h)} \cdot s_{i+1,l}))}, \quad \forall l = 1, 2, \cdots, h
\end{aligned}
\tag{6}
$$

where $s_{i+1,k}$ is the memory slot vector at position $k$ of layer $i + 1$.

**Dead-node Solution.** During the training process of retrieving top $K$ leaf-node vectors, we met the same problem as [41] where almost the same set of nodes where visited at each step. The paper dealt with this problem by uniformly choosing $K$ random candidates each time thus solving the dead-node problem at the early stage of training. However, as the model parameters update along time, real top-$K$ node vectors are desired to be retrieved.

In our work, we use Gumbel-Softmax [25] which is an ideal sampling trick for distribution on discrete vectors (also known as reparameterization trick). The basic idea of Gumbel-Softmax is to add a Gumbel noise and temperature factor to the original softmax function so that it is able to approximate a probability distribution made up of discrete categories (which are the normalized leaf-node weights in our case).

With calculated leaf-node weights (i.e. probabilities) $w_1, w_2, \ldots, w_k$ from softmax function using Eq. (6), the Gumbel-Softmax is conducted as:

$$
\begin{aligned}
y_i &= \text{softmax}((\log(w_i) + g_i)/\tau) \\
&= \frac{\exp((\log(w_i) + g_i)/\tau)}{\sum_j \exp((\log(w_j) + g_j)/\tau)}, \quad \forall j = 1, 2, \cdots, k
\end{aligned}
\tag{7}
$$

where $g_1, g_2, \ldots, g_k$ are i.i.d samples drawn from the standard Gumbel distribution having $\mu$ and $\beta$ as 0 and 1 respectively, with PDF (Probability Density Function) of $e^{-(x+e^{-x})}$. In practice, $g_k$ can

be sampled using inverse transform sampling by drawing $u_k \sim$ uniform$(0, 1)$ and calculated as $g_k = -\log(-\log(u_k))$ [36].

One of the great property of the Gumbel-Softmax is that the output value approaches the real distribution with low temperature factor $\tau$ and tends to be uniform sampling with larger $\tau$. As a result, we start the training with a big temperature and then anneal it towards small values. In this way, the model tends to explore random nodes for update in the beginning and gradually stick to the real top $K$ selection for a better convergence.

For inference, the Top $K$ leaf-node vectors with highest weight scores are used to form the user group representation $e_u^{(g)}$:

$$
e_u^{(g)} = \sum_{k=1}^{K} w_k^{(leaf)} s_k^{(leaf)}.
\tag{8}
$$

**Orthogonality**. To make the nodes at each layer more diverse and representative, orthogonality regularizer is applied to the fully-connected layers $\mathbf{W}$ in the tree:

$$
\mathcal{L}_O = \lambda \cdot \sigma(\mathbf{W}^T \mathbf{W} - \mathbf{I}),
\tag{9}
$$

where $\lambda$ is the penalty parameter and $\sigma(\mathbf{W}) = \sup_{x \in \mathbb{R}^{1 \times n}, x \neq 0} \frac{\|\mathbf{W}x\|}{\|x\|}$ is the spectral norm of $\mathbf{W}$, i.e. the largest singular value of $\mathbf{W}$. Though computation of Eq. (9) involves expensive eigen-decomposition, it can be approximated via power iteration method [2]. Starting with a random initialized vector $v \in \mathbb{R}^d$, we iteratively perform the following procedures a small number of times (2 by default):

$$
u \leftarrow (\mathbf{W}^T \mathbf{W} - \mathbf{I})v, v \leftarrow (\mathbf{W}^T \mathbf{W} - \mathbf{I})u, \sigma(\mathbf{W}^T \mathbf{W} - \mathbf{I}) \leftarrow \frac{\|v\|}{\|u\|}.
\tag{10}
$$

Our hierarchical tree-based memory module is similar to previous work [41] in design. However, the main purpose and realization is quite different:

- The PreHash module proposed in [41] is designed to solve user embedding problem where certain anchor vectors are manually selected and kept unchanged during training. In our work, the memory tree is designed for user group classification and is fully automatic with parameters updating. All the user behavior representations will go through this module and fall into most related leaf slots.
- Orthogonality property is applied to our memory tree as it is a strong regularization for the diverse representation and ensures better convergence of our tree module while PreHash in [41] doesn't have this for model update.
- Sampling strategy during training is also different. Prehash uses the uniform sampling throughout training while the Gumbel-softmax sampling applied in COUPLE is more efficient for node representation update.

After obtaining user history representation (in Section 3.2.1), user content representation (in Section 3.2.2) and user group representation (in Section 3.2.3), we apply weighted aggregation (WA) again to get the final user representation $e_u$ using Eq. (3):

$$
e_u = \text{WA}([e_u^{(h)}, e_u^{(c)}, e_u^{(g)}]).
\tag{11}
$$

## 3.3 Contrastive Learning

Unsupervised representation learning in a contrastive way is highly successful in recently research [9, 17, 39]. Following a standard

**Table 2: Dataset statistics.**

| Datasets | #Users | #Item | #Tags | #Interact |
|----------|--------|-------|-------|-----------|
| Amazon Book | 1,672,201 | 1,677,035 | 555,685 | 11,637,960 |
| Amazon Movie | 1,080,338 | 154,451 | 73,430 | 2,561,003 |
| Taobao Product | 951,069 | 2,794,416 | 85,398 | 158,758,876 |
| Taobao m-Video | 829,930 | 969,145 | 56,291 | 47,332,179 |

contrastive definition [15], we define a contrastive loss between our user representaion $e_u$ and item representation $e_{i,0}, e_{i,1}, ..., e_{i,m}$:

$$\mathcal{L}_N = -\mathbb{E}[\log \frac{\exp(e_u \cdot e_i^+/\omega)}{\exp(e_u \cdot e_i^+/\omega) + \sum_{j=1}^{m-1} \exp(e_u \cdot e_j^-/\omega)}], \quad (12)$$

where $\omega$ is a temperature hyper-parameter. The contrasitve loss has a low value when $e_u$ similar to its positive item sample $e_i^+$ and dissimilar to all other negative item samples $\{e_j^-\}$.

In earlier study, softmax-based classifier is often utilized to classify between positive and negative samples. In [17], an efficient momentum contrastive scheme is proposed to maintain the dictionary as a queue of data samples. In this way, a rich set of negative samples can be involved for training and achieve positive results in various computer vision tasks. In our framework, a First-In-First-Out (FIFO) quque is maintained and updated with the current batch enqueued and the oldest batch dequeued for each training step. So that each batch of positive samples can be encoded with the hard negative samples over the latest steps. It makes the loss tracking more consistent and has a debiasing effect in large-scale production environment as proved in [48]. Thus, our network is trained under the loss function $\mathcal{L} = \mathcal{L}_N + \mathcal{L}_O$ consisting of a contrastive loss and a orthogonal constraint on weights of memory tree.

# 4 EXPERIMENTS

In this section, we evaluate our proposed cross-domain user preference learning (COUPLE) model for top-$K$ retrieval task under different scenarios. We introduce the experimental setup and the comparing baselines first, and present the performance comparison with other state-of-the-arts. Then the ablation study is carried out both quantitatively and qualitatively.

## 4.1 Experimental Setup

*4.1.1 Dataset Overview.* We compare the performance of our proposed framework with related methods for cross-domain cold-start recommendation on two publicly accessible datasets. Amazon review dataset[1] is one of the most widely-used public dataset [18] for e-commerce recommendations. We use subsets of Books as the source domain while Movies and TV the target domain. The other dataset is called Tao-Product-Micro-Video (TPMV), selected from the real-world online service of Taobao Recommendation. User interactions are much denser in product domain than micro-video domain, making the dataset suitable for cross-domain cold-start recommendation. The detailed statistics of the two datasets are demonstrated in Table 2.

---

[1]http://jmcauley.ucsd.edu/data/amazon/

*4.1.2 Experimental Setup.* To compose valid training and testing datasets for cross-domain cold-start recommendation, we use "click" as the implicit feedback for TPMV dataset and "review rating" for Amazon dataset. For training and evaluation, clicked items of TPMV and review ratings over 3 (ratings range from 1 to 5) of Amazon are regarded as positive samples and temporally sorted with timestamps of a user's action. For testing, we use leave-one-out strategy and make the last-position items from target domain as the ground-truth, while erasing these items from the training pool to make sure they are totally cold-start contents. To further make the situation tougher as the real circumstances, we randomly select half of the testing users, and erase all the items of the target domain (if any) from their historical sequences to make them cold-start users. For all the models compared in this paper, the detailed numbers of training and testing samples are listed in Table 3.

**Table 3: Traing and testing splits.**

| Datasets | #Train | #Test | #Cold-start User |
|----------|--------|-------|------------------|
| Amazon | 5,797,557 | 95,884 | 67,063 |
| TPMV | 5,333,879 | 62,230 | 34,934 |

*4.1.3 Evaluation metrics.* To evaluate the performance of all the methods, Hit Ratio (HR), and Normalized Discounted cumulative gain (NDCG) are used. HR measures whether the positive item retrieved within the top-$K$ and NDCG penalizes the score if positive item positioned lower in the ranking list. Higher values in these metrics indicate better recommendation performance. For each ground-truth positive item, we pair it with 100 randomly sampled negative items from the pool where users have no historical interactions with, following the prior work [19, 23]. To make the comparison more reliable, 50 of them are sampled randomly, while the other 50 items are sampled according to the popularity [23].

## 4.2 Competitors

We compare our proposed framework with six baselines with neural deep learning structures, including DSSMs, Youtube DNNs and the Sequential methods.

- **DSSMs**. The original Deep Semantic Similarity Model (DSSM) [24] serves as a strong baseline widely used in information retrieval, which can perform semantic matching between a query and a document. DeepCoNN [47] further extends DSSM to two parallel convolutional neural networks. One of them models user behaviors and the other models item properties from the review texts. It achieves a strong performance in content recommendation.
- **Youtube DNNs**. Youtube [8] is a classical deep-based matching model for building user and item embeddings collaboratively. Layers of depth on the top is proved effective to model non-linear interactions between features. The latest work [44] solves multi-channel cold-start issue based on Youtube DNN structure where SOTA result is achieved.
- **Sequential methods**. SASRec [27] encodes user's behavior based on a variant of Transformer and is used as the backbone for many sequential models including ours. And we also compared to the latest sequential method ComiRec

**Table 4: HitRate and NDCG of different methods on the two datasets, where best performance is in boldface. HP denotes hyperparameters, including $n$ the number of tags and $l$ the item length for sequential modeling.**

| Dataset | TPMV Data | | | | Amazon Data | | | |
|---|---|---|---|---|---|---|---|---|
| HP | $n = 15, l = 5$ | | | | $n = 10, l = 8$ | | | |
| Metric | HitRate@5 | HitRate@10 | NDCG@5 | NDCG@10 | HitRate@5 | HitRate@10 | NDCG@5 | NDCG@10 |
| DSSM [24] | 0.07764 | 0.12490 | 0.05397 | 0.06908 | 0.05119 | 0.09463 | 0.03093 | 0.04658 |
| DeepCoNN [47] | 0.19080 | 0.27351 | 0.13435 | 0.16062 | 0.05352 | 0.10278 | 0.03319 | 0.04800 |
| Youtube [8] | 0.19705 | 0.28187 | 0.13899 | 0.16627 | 0.06777 | 0.12398 | 0.04174 | 0.05970 |
| ICAN [44] | 0.20610 | 0.28502 | 0.14711 | 0.17248 | 0.07431 | 0.13071 | 0.04601 | 0.06405 |
| SASRec [27] | 0.20145 | 0.28491 | 0.14709 | 0.17154 | 0.06918 | 0.12410 | 0.04352 | 0.06109 |
| ComiRec [3] | 0.20171 | 0.28539 | 0.14418 | 0.17110 | 0.07036 | 0.12684 | 0.04322 | 0.06127 |
| Ours | **0.22444** | **0.30625** | **0.16433** | **0.19062** | **0.07753** | **0.13358** | **0.04920** | **0.06711** |
| Improvement | + 8.89% | + 7.44% | + 11.71% | + 10.52% | + 4.30% | + 2.19% | + 6.93% | + 4.77% |

[3] which further improves for controllable multi-interest matching based on the previous MIND [29].

For methods mentioned above, we follow their implementations with tuned parameters for better comparison. Original DSSM [24] is implemented by maximizing the conditional likelihood of clicked items given only a closest query (item in our case). For DeepCoNN [47], it models similarity between user and item representations. We feed in a sequence of items represented by tag features where convolution is applied. In terms of sequential methods, we first average tag representations for an item as our method does. Then item sequence is modeled the same way as the original papers, where position embedding is added. To implement ICAN [44], we adjust the ID input feature with tag embeddings to fit in the content cold-start scenario, while keeping all other settings unchanged.

## 4.3 Performance Comparison

The overall performance of all the methods on both datasets is summarized in Table 4. We can observe that our proposed COUPLE network, which models user preference from different perspectives, consistently yields the best performance in terms of HitRate@N and NDCG@N (N = 5, 10). We would like to note that, among three different types of baselines, Youtube DNNs and Sequential methods obtain better performance than DSSMs. It implies that sufficient interactions between user and item features are crucial for content-based recommendation. Within two-towel models, Deep-CoNN gains apparent improvement over original DSSM by modeling series of items for user presentation instead of only one item. By adding channel-wise attention between source and target domain, ICAN is able to achieve better performance over original Youtube network. The same improvement can be found on ComiRec over simple sequentially self-attention method SASRec, where multi interests of users are learnt across domains.

Furthermore, we observe that the improvement of our method over the strongest baselines is more significant on TPMV data (7.44% vs. 2.19% on HitRate@10) and the performances of all the methods on TPMV Dataset are generally better than those on Amazon Dataset. The reasons may be that: a) TPMV data have more

overlapped tags in target domain with source domain than Amazon (60% vs. 25%), so that models can learn a better correlation between domains; and b) user behaviors are more consistent on TPMV data which is collected from a range of a few days compared to Amazon with a wider range of months (over 4 years). Both attributes of TPMV favors our model where user-content attention and user-group preference learning are performed to generalize and transfer the interests across domains.

## 4.4 Ablation Study

We perform ablation study to further explore the effect of individual components of our COUPLE model on performance. As shown in Table 5, our base model with user history and domain embeddings achieves similar performance with SASRec. An apparent improvement is observed by applying user group preference and user content preference learning modules. The variant with hierarchical memory tree outstands others and raise the HitRate@10 by a large margin of 3.88%, further illustrating the effectiveness of the generalizing user preference across domains. After combining all modules together to represent the user, the performance is further improved. We also experiment with the batch negative sampling strategy compared to the contrastive learning using a FIFO queue. It can be seen that with a larger number of negative samples, model is able to distinguish data with higher efficiency. With all modules put together, our proposed model achieve a competitive performance for the real-world cold-start recommendation on Taobao platform.

## 4.5 Source-domain Data Sensitivity

One critical assumption of our work is that it helps with the cross-domain recommendation by joint modeling user behaviors from both source and target domains. In this section, We conduct experiments on two datasets the sensitivity of our model against the amount of source domain data. By manipulating the amount of training data, we gradually decrease the portion of data from the source domain. As part of the testing cases is cold-user recommendation where user behaviors only exist in the source domain, we further disassemble HR@10 and NDCG@10 for different situations.

**Table 5: Ablation Study on TPMV Dataset.**

| Variants of Our Method | Metrics | |
|---|---|---|
| | HR@10 | NDCG@10 |
| Base (User History + Domain) | 0.28411 | 0.17042 |
| Base + User-Content Preference | 0.28928 | 0.17531 |
| Base + Memory Tree | 0.29554 | 0.18392 |
| All User + Batch_neg | 0.3011 | 0.18787 |
| **All User + FIFO Queue** | **0.30625** | **0.19062** |



(a) Impact on HR of TPMV  (b) Impact on NDCG of TPMV

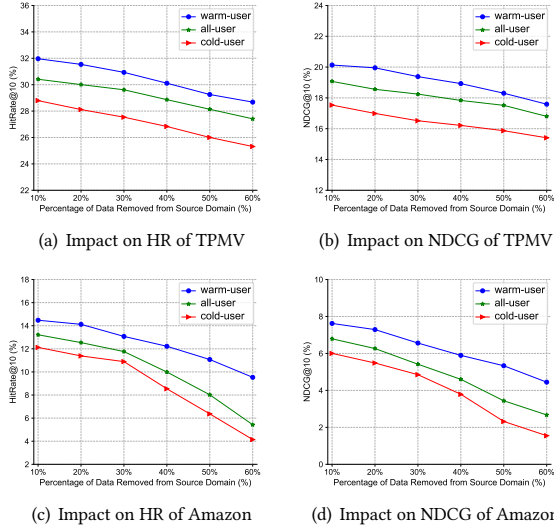(c) Impact on HR of Amazon  (d) Impact on NDCG of Amazon

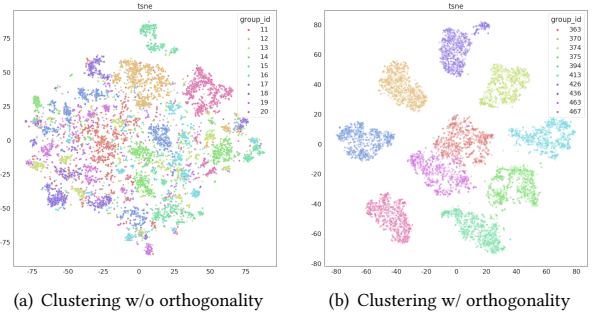**Figure 4: The impact of the source domain data.**

Figure 4 shows the impact of source domain data on the performance. We can see a clear drop on both HitRate and NDCG as the percentage of training data from the source domain decreases. Recommending cold items to cold users are more difficult as the overall performance is worse than recommending to users with interactions in target domain already. What is more, with less data from the source domain, the performance on cold-user recommendation decreases more than the warm-user situations with increasing gap between two lines. An interesting finding is that cold-user recommendation on Amazon data is more sensitive to the varying source data than that on TPMV data. We believe that it is related to the correlations between tags extracted from both domains. With less overlapping tags, it becomes more difficult to match interest between domains with progressively sparser source triggers.

### 4.6 Visualization of User Group Preference

A good orthogonal representation tree should be able to project and classify users into distinct groups. To further verify whether orthogonal regularization works, We visualize the user group representations with t-distributed stochastic neighbor embedding (t-SNE) [42]. First, we obtain the memory slot vector of leaf nodes $\{s\}_{k=1}^{h}$ via fully-connected projection of user history representation $e_u^{(h)}$. Then the leaf node group $k$ with largest weight calculated using

Eq. (6) is assigned to the user. Following t-SNE, We treat the each component of $s_k$ as an individual point and keep only the two components that have the highest confidence levels. For better illustration, we randomly sample 10 user groups with 10,000 points within each and the result is shown in Figure 5.

It is obvious that the clustering result with orthogonal regularization has higher similarity within a cluster and can separate different clusters better. When node vectors approach orthogonal, they become de-correlated so that the responses are much less redundant. By applying orthogonality constraints together with Gumbel-Softmax sampling trick, we observe a stable convergence and layer-wise distribution of the hierarchical memory tree.



(a) Clustering w/o orthogonality  (b) Clustering w/ orthogonality

**Figure 5: t-SNE of clustering samples of the memory tree.**

## 5 ONLINE DEPLOYMENT

We conduct online experiments by deploying our work to guess-you-like recommendation after purchase at Taobao platform, where brand new micro-videos with extracted tags are displayed to users with related interest. This scenario has over dozens of millions of active users and hundreds of thousands of new micro-movies uploaded each day. For comparison, all the matching methods share the same ranking stage, so that click-through-rate (CTR) can be a fair metric and is used for evaluation. We compare our work with two baselines that are already running online, one is author-based CF where new videos created by the same author will be retrieved and the other is Youtube DNN where tag features are fed into the network for similarity measurement. Users assigned to the experiment group experience COUPLE recommendation, while the other two control groups experience author-CF and Youtube DNN, respectively. We run the A/B test in the heavy-traffic scenario for a week. The CTR of the experiment group with COUPLE improves over the control group by **8.855%**. Further compared to author-CF, COUPLE has a lower CTR rate (-5.6%) but provides significant gain on exposure rate by **58.91%**. And It illustrates that with content-based recommendation scheme, more cold micro-videos get the opportunity to be shown to users which contributes to the sustainable development of recommender systems. In general COUPLE + author-CF improve CTR over Youtube + author-CF by **3.15%**. We replace Youtube DNN with COUPLE and it is now serving online for the cross-domain cold micro-video recommendation.

# 6 CONCLUSION

In this paper, we proposed a unified framework for cross-domain cold-start recommendation in real world. It learns a user's preference from three different perspectives including user history, user content and user group. By jointly modeling the user's behaviors with extracted semantic tags across domains, our work is able to perform reliable recommendation in a new domain with limit user interactions. Fed into a contrastive learning scheme, COUPLE achieve state-of-the-art results on two datasets and demonstrates superior performance compared to other baselines running at Taobao APP. It has been deployed in production at scale for cold micro-video recommendation. For the future, we plan to integrate more heterogeneous content features into our framework to further help with the cold-start recommendation. Another exciting direction is to incorporate knowledge graph for information propagation towards explainable recommendation.

# REFERENCES

[1] Peter Anderson, Xiaodong He, Chris Buehler, Damien Teney, Mark Johnson, Stephen Gould, and Lei Zhang. 2018. Bottom-up and top-down attention for image captioning and visual question answering. In *Proceedings of the IEEE conference on computer vision and pattern recognition*. 6077–6086.

[2] Nitin Bansal, Xiaohan Chen, and Zhangyang Wang. 2018. Can we gain more from orthogonality regularizations in training deep networks? *Advances in Neural Information Processing Systems* 31 (2018), 4261–4271.

[3] Yukuo Cen, Jianwei Zhang, Xu Zou, Chang Zhou, Hongxia Yang, and Jie Tang. 2020. Controllable Multi-Interest Framework for Recommendation. In *Proceedings of the 26th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining*. 2942–2951.

[4] Tianlong Chen, Shaojin Ding, Jingyi Xie, Ye Yuan, Wuyang Chen, Yang Yang, Zhou Ren, and Zhangyang Wang. 2019. Abd-net: Attentive but diverse person re-identification. In *Proceedings of the IEEE International Conference on Computer Vision*. 8351–8361.

[5] Xu Chen, Hongteng Xu, Yongfeng Zhang, Jiaxi Tang, Yixin Cao, Zheng Qin, and Hongyuan Zha. 2018. Sequential recommendation with user memory networks. In *Proceedings of the eleventh ACM international conference on web search and data mining*. 108–116.

[6] Zhongxia Chen, Xiting Wang, Xing Xie, Mehul Parsana, Akshay Soni, Xiang Ao, and Enhong Chen. 2020. Towards Explainable Conversational Recommendation. IJCAI.

[7] Deborah Cohen, Michal Aharon, Yair Koren, Oren Somekh, and Raz Nissim. 2017. Expediting exploration by attribute-to-feature mapping for cold-start recommendations. In *Proceedings of the Eleventh ACM Conference on Recommender Systems*. 184–192.

[8] Paul Covington, Jay Adams, and Emre Sargin. 2016. Deep neural networks for youtube recommendations. In *Proceedings of the 10th ACM conference on recommender systems*. 191–198.

[9] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2018. Bert: Pre-training of deep bidirectional transformers for language understanding. *arXiv preprint arXiv:1810.04805* (2018).

[10] Ali Mamdouh Elkahky, Yang Song, and Xiaodong He. 2015. A multi-view deep learning approach for cross domain user modeling in recommendation systems. In *Proceedings of the 24th International Conference on World Wide Web*. 278–288.

[11] Alessandro Epasto and Bryan Perozzi. 2019. Is a single embedding enough? learning node representations that capture multiple social contexts. In *The World Wide Web Conference*. 394–404.

[12] Fatih Gedikli and Dietmar Jannach. 2013. Improving recommendation accuracy based on item-specific tag preferences. *ACM Transactions on Intelligent Systems and Technology (TIST)* 4, 1 (2013), 1–19.

[13] Xue Geng, Hanwang Zhang, Jingwen Bian, and Tat-Seng Chua. 2015. Learning image and user features for recommendation in social networks. In *Proceedings of the IEEE International Conference on Computer Vision*. 4274–4282.

[14] Ruocheng Guo, Xiaoting Zhao, Adam Henderson, Liangjie Hong, and Huan Liu. 2020. Debiasing grid-based product search in e-commerce. In *Proceedings of the 26th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining*. 2852–2860.

[15] Raia Hadsell, Sumit Chopra, and Yann LeCun. 2006. Dimensionality reduction by learning an invariant mapping. In *2006 IEEE Computer Society Conference on Computer Vision and Pattern Recognition (CVPR'06)*, Vol. 2. IEEE, 1735–1742.

[16] Khalid Haruna, Maizatul Akmar Ismail, Suhendroyono Suhendroyono, Damiasih Damiasih, Adi Cilik Pierewan, Haruna Chiroma, and Tutut Herawan. 2017.

[17] Context-aware recommender system: A review of recent developmental process and future research direction. *Applied Sciences* 7, 12 (2017), 1211.

[17] Kaiming He, Haoqi Fan, Yuxin Wu, Saining Xie, and Ross Girshick. 2020. Momentum contrast for unsupervised visual representation learning. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 9729–9738.

[18] Ruining He and Julian McAuley. 2016. Ups and downs: Modeling the visual evolution of fashion trends with one-class collaborative filtering. In *proceedings of the 25th international conference on world wide web*. 507–517.

[19] 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.

[20] Guangneng Hu, Yu Zhang, and Qiang Yang. 2018. Conet: Collaborative cross networks for cross-domain recommendation. In *Proceedings of the 27th ACM International Conference on Information and Knowledge Management*. 667–676.

[21] Haoji Hu, Xiangnan He, Jinyang Gao, and Zhi-Li Zhang. 2020. Modeling Personalized Item Frequency Information for Next-basket Recommendation. *arXiv preprint arXiv:2006.00556* (2020).

[22] Yifan Hu, Yehuda Koren, and Chris Volinsky. 2008. Collaborative filtering for implicit feedback datasets. In *2008 Eighth IEEE International Conference on Data Mining*. Ieee, 263–272.

[23] Jin Huang, Wayne Xin Zhao, Hongjian Dou, Ji-Rong Wen, and Edward Y Chang. 2018. Improving sequential recommendation with knowledge-enhanced memory networks. In *The 41st International ACM SIGIR Conference on Research & Development in Information Retrieval*. 505–514.

[24] Po-Sen Huang, Xiaodong He, Jianfeng Gao, Li Deng, Alex Acero, and Larry Heck. 2013. Learning deep structured semantic models for web search using clickthrough data. In *Proceedings of the 22nd ACM international conference on Information & Knowledge Management*. 2333–2338.

[25] Eric Jang, Shixiang Gu, and Ben Poole. 2016. Categorical reparameterization with gumbel-softmax. *arXiv preprint arXiv:1611.01144* (2016).

[26] Kui Jia, Shuai Li, Yuxin Wen, Tongliang Liu, and Dacheng Tao. 2019. Orthogonal deep neural networks. *IEEE transactions on pattern analysis and machine intelligence* (2019).

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

[28] Diederik P Kingma and Jimmy Ba. 2014. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980* (2014).

[29] 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*. 2615–2623.

[30] Jingjing Li, Mengmeng Jing, Ke Lu, Lei Zhu, Yang Yang, and Zi Huang. 2019. From zero-shot learning to cold-start recommendation. In *Proceedings of the AAAI Conference on Artificial Intelligence*, Vol. 33. 4189–4196.

[31] U Liji, Yahui Chai, and Jianrui Chen. 2018. Improved personalized recommendation based on user attributes clustering and score matrix filling. *Computer Standards & Interfaces* 57 (2018), 59–67.

[32] Qiang Liu, Shu Wu, Diyi Wang, Zhaokang Li, and Liang Wang. 2016. Context-aware sequential recommendation. In *2016 IEEE 16th International Conference on Data Mining (ICDM)*. IEEE, 1053–1058.

[33] Jiasen Lu, Dhruv Batra, Devi Parikh, and Stefan Lee. 2019. Vilbert: Pretraining task-agnostic visiolinguistic representations for vision-and-language tasks. In *Advances in Neural Information Processing Systems*. 13–23.

[34] Yuanfu Lu, Yuan Fang, and Chuan Shi. 2020. Meta-learning on heterogeneous information networks for cold-start recommendation. In *Proceedings of the 26th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining*. 1563–1573.

[35] Jianxin Ma, Chang Zhou, Hongxia Yang, Peng Cui, Xin Wang, and Wenwu Zhu. 2020. Disentangled Self-Supervision in Sequential Recommenders. In *Proceedings of the 26th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining*. 483–491.

[36] Chris J Maddison, Andriy Mnih, and Yee Whye Teh. 2016. The concrete distribution: A continuous relaxation of discrete random variables. *arXiv preprint arXiv:1611.00712* (2016).

[37] Prem Melville, Raymond J Mooney, and Ramadass Nagarajan. 2002. Content-boosted collaborative filtering for improved recommendations. *Aaai/iaai* 23 (2002), 187–192.

[38] Aditya Pal, Chantat Eksombatchai, Yitong Zhou, Bo Zhao, Charles Rosenberg, and Jure Leskovec. 2020. PinnerSage: Multi-Modal User Embedding Framework for Recommendations at Pinterest. In *Proceedings of the 26th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining*. 2311–2320.

[39] Alec Radford, Karthik Narasimhan, Tim Salimans, and Ilya Sutskever. 2018. Improving language understanding by generative pre-training.

[40] Dimitrios Rafailidis and Fabio Crestani. 2016. Top-n recommendation via joint cross-domain user clustering and similarity learning. In *Joint European Conference on Machine Learning and Knowledge Discovery in Databases*. Springer, 426–441.

[41] Shaoyun Shi, Weizhi Ma, Min Zhang, Yongfeng Zhang, Xinxing Yu, Houzhi Shan, Yiqun Liu, and Shaoping Ma. 2020. Beyond User Embedding Matrix: Learning to Hash for Modeling Large-Scale Users in Recommendation. In *Proceedings of the 43rd International ACM SIGIR Conference on Research and Development in Information Retrieval*. 319–328.

[42] Laurens Van der Maaten and Geoffrey Hinton. 2008. Visualizing data using t-SNE. *Journal of machine learning research* 9, 11 (2008).

[43] Xiang Wang, Xiangnan He, Yixin Cao, Meng Liu, and Tat-Seng Chua. 2019. Kgat: Knowledge graph attention network for recommendation. In *Proceedings of the 25th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining*. 950–958.

[44] Ruobing Xie, Zhijie Qiu, Jun Rao, Yi Liu, Bo Zhang, and Leyu Lin. 2020. Internal and Contextual Attention Network for Cold-start Multi-channel Matching in Recommendation. (2020).

[45] Ziqi Zhang, José Iria, Christopher Brewster, and Fabio Ciravegna. 2008. A comparative evaluation of term recognition algorithms.. In *LREC*, Vol. 5.

[46] Cheng Zhao, Chenliang Li, Rong Xiao, Hongbo Deng, and Aixin Sun. 2020. CATN: Cross-Domain Recommendation for Cold-Start Users via Aspect Transfer Network. *arXiv preprint arXiv:2005.10549* (2020).

[47] Lei Zheng, Vahid Noroozi, and Philip S Yu. 2017. Joint deep modeling of users and items using reviews for recommendation. In *Proceedings of the tenth ACM international conference on web search and data mining*. 425–434.

[48] Chang Zhou, Jianxin Ma, Jianwei Zhang, Jingren Zhou, and Hongxia Yang. 2020. Contrastive Learning for Debiased Candidate Generation in Large-Scale Recommender Systems. *arXiv preprint cs.IR/2005.12964* (2020).

[49] Xiao Zhou, Cecilia Mascolo, and Zhongxiang Zhao. 2019. Topic-enhanced memory networks for personalised point-of-interest recommendation. In *Proceedings of the 25th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining*. 3018–3028.

**clicked videos**

Side handle pot, Japanese style, make tea
Anhua Black Tea, Gaoma Erxi, Black Brick Tea
Tea infuser, tea strainer, tea maker, tea bag

**retrieved video**

Loose tea, tea leaves,
silver needles

**clicked items**

Casio, black steel, bluetooth, watch
longines, men, quartz watches, watches
Casio, bluetooth, watch, male, oak

**retrieved video**

Casual, steel belt,
men's watch, Seiko

**clicked movies**

movies & tv, kids & family, genre for featured categories,
andrew jimenez, car, mandarin
movies & tv, kids & family, toy, chinese edition
movies & tv,genre for featured categories,
kung fu, panda, mandarin

**retrieved movie**

movies & tv,
teen titans,
kids & family, sam liu

**clicked books**

books, religion, spirituality, law, attraction
books, alternative medicine, violet, flame, heal,
body mind, guide practical
books, spirituality, read, akashic, record,
accessing, archive, journey

**retrieved movie**

movies & tv, heal,
spiritual messenger, god,
inspiring

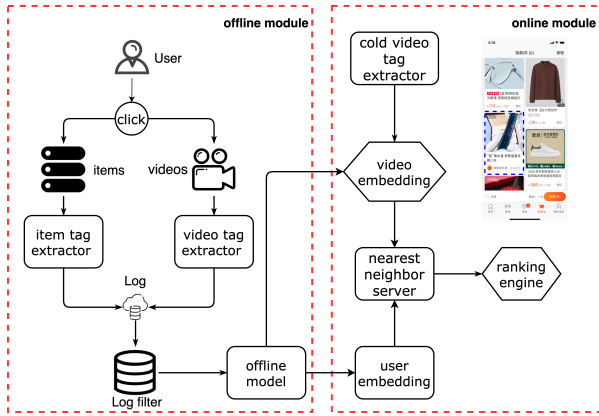**Figure 8: Retrieval cases of our model.**



**Figure 7: Online system workflow.**

# A REPRODUCTIVITY SUPPLEMENT

In this section, we provide details on both offline and online experiment settings for reproductivity purpose.

## A.1 Offline Experimental Details

*A.1.1 Tag Extraction.* The semantic tags used for representing items are extracted manually from different multi-media resources. For Amazon dataset, we extract meaningful short terms as tags from title, category, brand based on frequency and semantic annotation algorithms [45]. For TPMV dataset, richer tags can be drawn from product title, attributes, images via specialized in-house models. The tag extraction process can be referred to Figure 6.

*A.1.2 Parameter Settings of Our Method.* We implement our model with Tensorflow version 1.12 where default initialization recommended by Tensorflow is applied. We use dimension size 64 for all feature representations including tags, items and users.

**User History Representation**. To model user history representation, one block of self-attention with hidden size 64 and heads

of 4 is used. The tag numbers per item and item sequence length are capped to (10, 8) and (15, 5) for TPMV and Amazon datasets, respectively. Domain embedding is added to the item embedding to distinguish whether the items are from source or target domain.

**Memory Tree Implementation**. For the hierarchical memory tree design, we use 1-32-512 as the structure while dropout rate is set to be 0.2, the penalty parameter $\lambda$ for the SRIP regularization is set to be 0.1, and the temperature factor $\tau$ used in Gumbel-Softmax anneals using the schedule $\tau = \max(1, 20e^{-rt})$ of the global training step $t$, where $\tau$ is updated every 1000 steps and $r = 1e - 5$ are the hyperparameters used in our method.
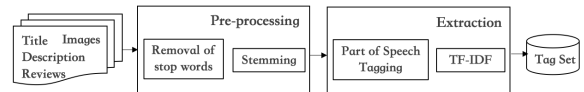


**Figure 6: Tag extraction process used in paper.**

**Contrastive Learning**. The contrastive learning is performed with a FIFO queue of size 10×256 where 256 is the mini-batch size. The temperature $\omega$ used in contrastive loss from Eq. (12) is 0.03. The Adam [28] optimizer for mini-batch gradient descent is used in our work and the learning rate is set to be 0.0001. Some extracted tags and retrieved cases of our work are visualized in Figure 8.

## A.2 Online Experimental Details

Our online system is shown in Figure 7. The offline model trained with the user behaviors is utilized for online recommendation system. A user's representation is inferred with recent interacted items (represented by tags). On the other hand, new videos will go through a tag extractor and obtain the embeddings through the model. Real-time matching will be performed via calculating the similarity score between user and video embeddings (by inner product). Videos with the highest scores are retrieved for ranking stage.