

# Adaptive Domain Scaling for Personalized Sequential Modeling in Recommenders

Zheng Chai  
chaizheng.cz@bytedance.com  
ByteDance  
Hangzhou, China

Hui Lu  
luhui.xx@bytedance.com  
ByteDance  
Hangzhou, China

Di Chen  
chend.666@bytedance.com  
ByteDance  
Beijing, China

Qin Ren  
renqin.97@bytedance.com  
ByteDance  
Beijing, China

Yuchao Zheng<sup>†</sup>  
zhengyuchao.yc@bytedance.com  
ByteDance  
Beijing, China

Xun Zhou  
zhouxun@bytedance.com  
ByteDance  
Beijing, China

## ABSTRACT

Users generally exhibit complex behavioral patterns and diverse intentions in multiple business scenarios of super applications like Douyin, presenting great challenges to current industrial multi-domain recommenders. To mitigate the discrepancies across diverse domains, researches and industrial practices generally emphasize sophisticated network structures to accommodate diverse data distributions, while neglecting the inherent understanding of user behavioral sequence from the multi-domain perspective. In this paper, we present Adaptive Domain Scaling (ADS) model, which comprehensively enhances the personalization capability in target-aware sequence modeling across multiple domains. Specifically, ADS comprises of two major modules, including personalized sequence representation generation (PSRG) and personalized candidate representation generation (PCRG). The modules contribute to the tailored multi-domain learning by dynamically learning both the user behavioral sequence item representation and the candidate target item representation under different domains, facilitating adaptive user intention understanding. Experiments are performed on both a public dataset and two billion-scaled industrial datasets, and the extensive results verify the high effectiveness and compatibility of ADS. Besides, we conduct online experiments on two influential business scenarios including Douyin Advertisement Platform and Douyin E-commerce Service Platform, both of which show substantial business improvements. Currently, ADS has been fully deployed in many recommendation services at ByteDance, serving billions of users.

## CCS CONCEPTS

• Information systems → Recommender systems.

<sup>†</sup>Corresponding Author.

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than the author(s) must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from [permissions@acm.org](mailto:permissions@acm.org).  
SIGIR '25, July 13–18, 2025, Padua, Italy

© 2025 Copyright held by the owner/author(s). Publication rights licensed to ACM.  
ACM ISBN 978-1-4503-XXXX-X/18/06  
<https://doi.org/XXXXXXX.XXXXXXX>

## KEYWORDS

Multi-Domain Learning, Sequential Modeling, Ranking, Personalized Recommender System

### ACM Reference Format:

Zheng Chai, Hui Lu, Di Chen, Qin Ren, Yuchao Zheng<sup>†</sup>, and Xun Zhou. 2025. Adaptive Domain Scaling for Personalized Sequential Modeling in Recommenders. In *Proceedings of the 48th International ACM SIGIR Conference on Research and Development in Information Retrieval (SIGIR '25)*, July 13–18, 2025, Padua, Italy. ACM, New York, NY, USA, 9 pages. <https://doi.org/XXXXXXX.XXXXXXX>

## 1 INTRODUCTION



(a) Short-Video

(b) Live-Preview

(c) Live-Slide

Figure 1: Typical business scenarios in Douyin.

With the exponential growth of digital contents and the widespread use of the internet, recommender systems have become a vital role for enhancing user experience and alleviating information overload [5, 15, 30]. In real-world applications, for improving user retention and promote business benefit, the demands of industrial recommendation are widely distributed cross multiple domains [12, 14]. For example, as shown in Fig 1, in Douyin<sup>1</sup>, one of the largest video-watching apps in the world, the major domains include *Short-videos*,

<sup>1</sup><https://www.douyin.com/>

*Live-preview*, and *Live-slide*, where users can watch short-videos, live-streams, and enjoy the e-commerce and local-life services. Besides, due to its billion-scale user volume, different user groups, like users from different countries, with different genders, highly-active or not, also contribute to different domains. As the data distributions are quite diversified across different domains, it poses a significant multi-domain modeling problem for the recommender system [31].

To this end, *de facto* common industrial practices generally build a shared-bottom with multi-heads outputs model structure, leveraging the advantages of both separate and unified mixup modeling for multiple domains [4]. To further improve this, recent approaches have put efforts in building elaborate network structures to enhance multi-domain modeling, e.g., the domain-level methods like the star topology adaptive recommender (STAR) [19], progressive layered extraction (PLE) [20], and the instance-level methods like adaptive parameter generation network (APG) [26], AdaSparse [28], etc. However, most of the existing approaches are designed for sophisticated feature interaction network structures, while the approaches to multi-domain sequential modeling attract much less attention.

Sequential modeling plays a vital role in industrial recommenders, among which the most popular and effective methods are the target-aware attention based methods, for example, deep interest network (DIN) [33], feature co-action network (CAN) [2], and multi-head attention (MHA) [21]. Despite its significance, the impacts of multi-domain discrepancies are less considered in existing target attention methods, leaving a remarkable gap for the area of multi-domain modeling. Generally speaking, the current target attention mechanism for user sequence can be formulated as a typical query-key-value modeling paradigm:  $g(Rep_{cand} \cdot Rep_{seq}) \times Rep_{seq}$ , where  $Rep_{cand}$  denotes the representation of the candidate target item whose click/convert probabilities need to be predicted,  $Rep_{seq}$  indicates the user sequence embeddings, and  $g$  calculates the attention weight between any sequence-item and target-item pairs. As discussed previously, current industrial recommenders generally follow a shared-bottom embedding paradigm, which means that 1) the embedding table of both the candidate item and user behavioral items are fully shared, with no consideration of the distinctions between items and users belonging to different domains, and 2) the candidate item serves as a shared query for different keys/values, with no consideration of the distinctions between the multiple items in user sequence which generally occurs across multiple domains. Accordingly, this poses potential challenges to the current multi-domain recommenders from two aspects:

- **Personalization of the Sequence Representations.** The multi-domain representations of identical items occurred in different users' sequences are of necessity in recommenders. For example, new users would like to watch the highly-liked videos, while some long-time users may pay attention to the video creators they follow. Thus, the same video shows different attractions to various user domains, while its embedding is a shared representation among different user sequences, which hinders the recommender to grasp the actual intention of the user.
- **Personalization of the Candidate Item.** For different users or different items in the same user's sequence, the candidate item has different influence and function due to the multi-domain impact. For example, user's shopping behaviors in Douyin Mall might be

primarily influenced by product prices, while the video creators have a more significant impact for content-preferred users in Douyin short-videos scenario. Thus an identical candidate item should be personalized across domains to accommodate different historical items in user sequence.

To overcome the limitations and fulfill the gap in multi-domain target-aware attention modeling, we propose Adaptive Domain Scaling (ADS) model, which fully mines the personalization modeling ability of the current target-attention based recommenders and provides more accurate and adaptive intention understanding ability in multi-domain tasks. Specifically, ADS comprises two modules, namely, Personalized Sequence Representation Generation (PSRG) and Personalized Candidate Representation Generation (PCRG). In PSRG, we designed a novel shared-and-private structure for learning multi-domain item representations in user behaviors, which aims to generate personalized representations for sequence items, i.e., the same item occurred in different user's sequence has different representations. In PCRG, the candidate item further enhances the personalization modeling ability via generating different target candidate representations for different sequence items. With the domain-related information as input for the generation structures, the impact of multi-domains are sufficiently injected to the sequence modeling, and thus enhances user intention understanding. Note that ADS is an efficient plug-and-play network, and can be readily integrated into current recommenders.

The contribution of this work can be summarized as follows:

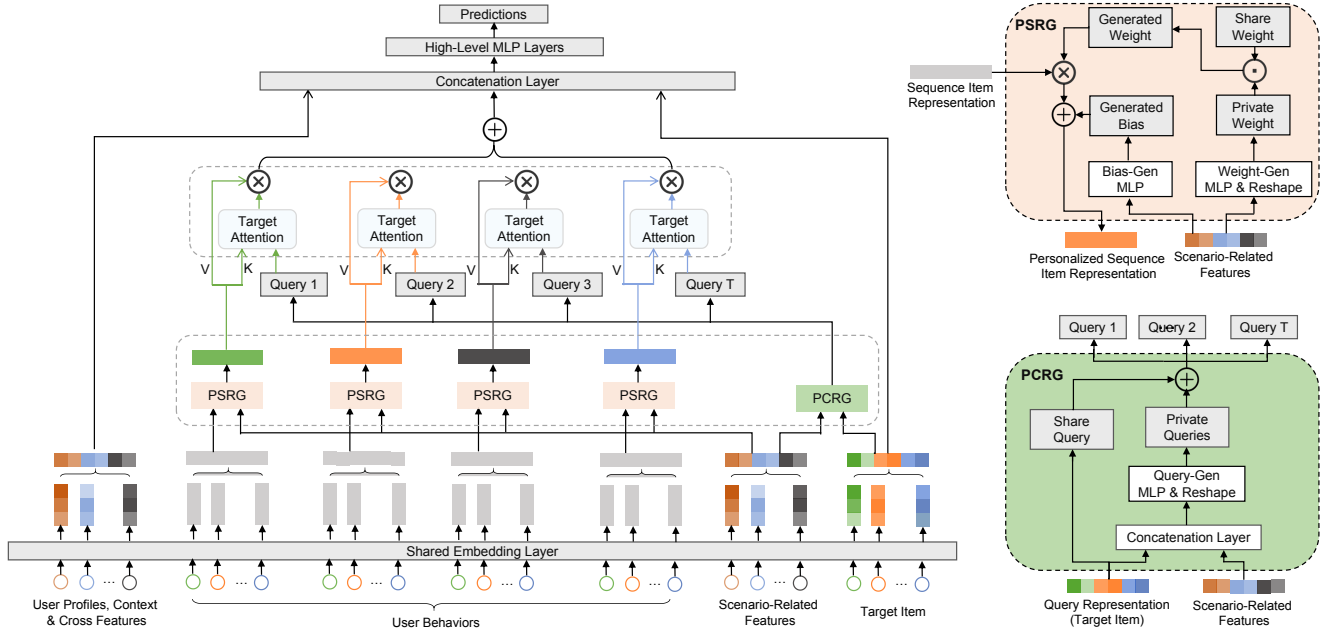
- We present Adaptive Domain Scaling (ADS) model, an effective plug-and-play personalization network structure for multi-domain user intent understanding by personalizing the target-aware attention modeling. We conduct extensive experiments on both a public dataset and two billion-scaled industrial dataset, and the results verify its superiority.
- Both personalized sequence representation generation and personalized candidate representation generation modules are developed in our framework, which captures the multi-domain characteristics from the viewpoint of users behavioral sequences and candidate target item, enhancing the multi-domain learning efficacy for current target-aware attention mechanisms.
- We deploy ADS in both the advertisement system and the e-commercial system of Douyin at ByteDance, which brings significant 1.00% and 0.79% lifts of total revenue in the Douyin ads system and the e-commercial system, respectively. Currently, ADS has been fully deployed in many recommendation systems at ByteDance, serving billions of users.

## 2 METHODOLOGY

### 2.1 Preliminaries

**2.1.1 Problem Formulation.** In this paper, we focus on the ranking modeling task in recommenders, which is a typical binary classification problem. Generally, taking the prediction of the click-through rate (CTR) as an example, the probability  $\hat{y}$  can be obtained via the following:

$$\hat{y} = f(E_U, E_I, E_O) \quad (1)$$



**Figure 2: Overview of ADS.** ADS consists of PCRGR, PSRG, and target attention module. Given scenario-related features and target item as input, PCRGR first generates multiple queries considering the co-pattern of the target item (query) and the scenario features. For PSRG, it takes scenario features as input to generate weight and bias parameters to formulate the personalized MLP, then the original sequence item embedding is passed through the generated MLP to obtain personalized representation. Both the PCRGR and PSRG share a share-and-private learning paradigm. Finally, the generated sequence is aggregated by the generated multi-queries with the target-aware attention mechanism, and the concatenation layer and high-level MLP layers are finally used to make predictions.

in which  $E(\cdot)$  indicates the embedding function, in which the raw categorical features are translated directly into embeddings, and the continuous features are first bucketed and then embedded as dense vectors.  $f$  is the MLP-based transformation function.  $U$ ,  $I$ , and  $O$  denote the user-side, target candidate item-side, and other features, respectively. User-side features generally consist of demographic feature set (for example, user locations and languages) and behavioral feature set (for example, the watch list or the shopping list of users). Item-side features include the item’s descriptive features like its category, creators, etc. Besides, the other features  $O$  generally contains contextual and user-item cross features.

## 2.2 The Proposed ADS

The structure of the proposed ADS is illustrated in Figure 2. Overall, it is composed of two major parts:

- **Personalized Sequence Representation Generation.** PSRG generates dynamic behavioral item embedding with a share-and-private learning structure, such that the same item occurred in different domains has unshared representations due to the distinction among multi-domains.
- **Personalized Candidate Representation Generation.** PCRGR captures different aspects of sequence items and generates multiple adaptive queries (i.e., candidate items) for each sequence item, such that the different query influence on diverse sequence items can be reflected.

With the adaptive queries, keys, and values generated by PCRGR and PSRG, the target-aware modeling mechanism like MHA, DIN, and CAN, can be readily integrated into this framework, facilitating the interest capture in multi-domain scenarios.

### 2.2.1 Personalized Sequence Representation Generation (PSRG).

Current large-scale industrial recommenders generally adopt the share-embedding layer to embed the raw ID and other features into dense vectors. In this manner, a specific item in the embedding table has a unified embedding, which is shared across different user sequences and neglects the impact from the difference between multiple domains.

The basic idea of PSRG is to dynamically generate a personalized layer for each item embedding in user behavioral sequence, such that the original shared representation can be diversified across multiple domains. Specifically, we use the concatenation of the domain-related features embeddings  $E_D \in \mathbb{R}^{d_D}$  as the input of the generation part of PSRG, which consists of the following two feature categories: 1) explicit-domain-indicator features, which distinguish the domain that a sample belongs to. For example, the indicator ranges from  $[0, 2]$  to indicate the three distinct business scenarios within Douyin; and 2) implicit-domain-indicator features. In recommender systems, some domains can be challenging to define explicitly. For example, whether a user is highly-active or not. Therefore, additional engineer-constructed statistical features are

necessary to be incorporated to further capture and differentiate various domains. With the two categories of features, as illustrated in Figure 2, the generation processes of the weight and bias for sequence items are designed to dynamically adjust the original item embedding.

**Sequence-Weight Gen-Net.** Denote the user sequence embedding as  $E_S \in \mathbb{R}^{T \times d_S}$ , where  $T$  and  $d_S$  represent the user sequence length and the embedding dimension of each sequence item, respectively. Based on the domain features  $E_D$ , the weight generation process consists of a *private* weight part and a *share* weight part to capture both the commonality and individuality of multi domains. For the private part, a two-layered MLP is performed to generate the private weight:

$$\mathbf{W}_{private} = \text{Sigmoid}(\text{ReLU}(E_D \mathbf{W}_1^T + \mathbf{b}_1) \mathbf{W}_2^T + \mathbf{b}_2) \quad (2)$$

where  $\mathbf{W}_1^T \in \mathbb{R}^{d_D \times d_h}$ ,  $\mathbf{W}_2^T \in \mathbb{R}^{d_h \times (d_S \times d_S)}$ ,  $\mathbf{b}_1 \in \mathbb{R}^{d_h}$ ,  $\mathbf{b}_2 \in \mathbb{R}^{(d_S \times d_S)}$ , and  $d_h$  denotes the hidden layer dimension. Note that the two-layer function instead of a single layer can not only improve the expression ability of the model, but also significantly reduce the model parameters and computational costs, as  $d_S$  is generally at the order of tens in practical cases, and  $d_h \ll (d_S \times d_S)$ .

Based on  $\mathbf{W}_{private}$ , a global weight  $\mathbf{W}_{shared} \in \mathbb{R}^{(d_S \times d_S)}$  is further defined as a learnable matrix which is shared across all the users. To enjoy the merits in learning both the commonality and individuality, the generated weight is defined as:

$$\mathbf{W}_{generated} = \eta * (\mathbf{W}_{shared} \odot \mathbf{W}_{private}) \quad (3)$$

where  $\odot$  denotes the element-wise product. As the value of  $\mathbf{W}_{private}$  ranges from  $[0, 1]$  due to introduction of *Sigmoid*, a scaling hyper-parameter  $\eta$  is further involved to enlarge the expression range of  $\mathbf{W}_{private}$ . Therefore, the calculation in Equation 3 can be viewed as a adaptive scaling method for the global parameter  $\mathbf{W}_{shared}$ .

**Sequence-Bias Gen-Net.** Similar to the above weight generation process, the bias generation can be readily obtained via the following equation:

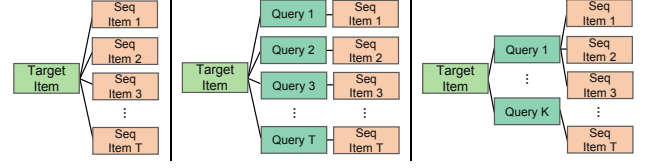
$$\mathbf{b}_{generated} = \text{ReLU}(E_D \mathbf{W}'_1^T + \mathbf{b}'_1) \mathbf{W}'_2^T + \mathbf{b}'_2 \quad (4)$$

where  $\mathbf{W}'_1^T \in \mathbb{R}^{d_D \times d_h}$ ,  $\mathbf{W}'_2^T \in \mathbb{R}^{d_h \times d_S}$ ,  $\mathbf{b}'_1 \in \mathbb{R}^{d_h}$ , and  $\mathbf{b}'_2 \in \mathbb{R}^{d_S}$ . With the generated weight and bias, the PSRG can be achieved via the following:

$$E_{S-personalized} = E_S \text{Reshape}(\mathbf{W}_{generated})^T + \mathbf{b}_{generated} \quad (5)$$

where the *Reshape* operator refers to reshape the 1-D vector-form  $\mathbf{W}_{generated}$  as a 2-D matrix-form with the shape of  $d_S \times d_S$ .

**2.2.2 Personalized Candidate Representation Generation (PCRG).** In addition to the personalized modeling for sequence, the other important part is the multi-domain modeling for the target item, which generally plays the role of query in target-aware attention. Typically, personalizing the candidate item encompasses two aspects. On one hand, similar to the sequence representation, the representation of the target item itself is also embedded through the shared embedding layer, which is not personalized across different domains. On the other hand, the candidate item plays different



**Figure 3: Comparison of traditional target-attention methods (Left) and the Multi-Query Gen-Net (Middle and Right).**

roles for different sequence items in various domains. For example, user's watchlist in Douyin mall channel reflects her shopping interests, while in short-video channel reflects her content preference. Thus, user's diverse interest should be captured via a more personalized query.

**Multi-Query Gen-Net.** To this end, we propose *Multi-Query Gen-Net*, as shown in the middle subfigure in Figure 3, which produces multiple queries corresponding to different sequence items, under the guidance of the domain-related features  $E_D$  and the original target item embedding  $E_Q \in \mathbb{R}^{d_Q}$ :

$$E_{Q-private} = \text{ReLU}((E_D \oplus E_Q) \mathbf{W}_{q1}^T + \mathbf{b}_{q1}) \mathbf{W}_{q2}^T + \mathbf{b}_{q2} \quad (6)$$

where  $\oplus$  refers to the concatenation operation,  $\mathbf{W}_{q1}^T \in \mathbb{R}^{(d_D + d_Q) \times d_h}$ ,  $\mathbf{b}_{q1} \in \mathbb{R}^{d_h}$ ,  $\mathbf{W}_{q2}^T \in \mathbb{R}^{d_h \times (T \times d_Q)}$  and  $\mathbf{b}_{q2} \in \mathbb{R}^{T \times d_Q}$ . Noted that the hidden layer dimension  $d_h \ll (T \times d_Q)$ , yielding a controllable computation cost at  $\mathcal{O}(d_h T d_Q)$ .

**Chunked-Query Generation.** For cases with long sequence with larger  $T$  at hundreds or even higher, we also devise a lightweight chunked-query generation method for improved computation efficiency. As shown in the right subfigure in Figure 3, since user's adjacent actions are prone to happen in the same domain, the raw sequence can be divided into  $G$  chunks with the adjacent items formed as a group.<sup>2</sup> Therefore, the generated  $E_{Q-private} \in \mathbb{R}^{(G \times d_Q)}$  can be further repeated to  $\mathbb{R}^{(T \times d_Q)}$  and the cost further reduces to  $\mathcal{O}(d_h G d_Q)$ .

Corresponding to the multiple private queries  $E_{Q-private} \in \mathbb{R}^{(T \times d_Q)}$ , we use the original query  $E_Q$  as the shared base, i.e.,  $E_{Q-shared} = \text{tile}(E_Q)$ , where *tile* refers to the tile operator which repeat the  $E_Q$  by  $T$  times, i.e.,  $E_{Q-shared} \in \mathbb{R}^{(T \times d_Q)}$ . Then, the final generated multiple queries can be obtained via a residual manner:

$$E_{Q-personalized} = \text{Reshape}(E_{Q-private} + E_{Q-shared}) \quad (7)$$

where *Reshape* operator reshapes the 1-D vector as a 2-D matrix-form with the shape of  $T \times d_Q$ .

**2.2.3 Target-Aware Attention & Prediction.** With the above personalized queries  $E_{Q-personalized} \in \mathbb{R}^{T \times d_Q}$  and personalized sequence items  $E_{S-personalized} \in \mathbb{R}^{T \times d_S}$ , the target attention module is performed to calculate the attention weight of each item and aggregate the sequence under the guidance of queries. Generally, the personalized queries and items can be readily integrated into many popular

<sup>2</sup>Without losing generality, here it is assumed that  $T$  is divisible by  $G$  with padding items.

**Table 1: Statistics of the two datasets.**

Dataset	#Users	#Items	#Instances	#Tasks
Taobao	101,342	500,272	24.02M	Order
Douyin-Ads	688M	270M	2.52B	CVR
Douyin-Ecom	596M	12.91M	24.65B	Click & Order

attention methods, like multi-head target attention, DIN, and CAN. Taking the multi-head target attention as an example, for each head, the candidate item and sequence items are first transformed via the following:

$$Q = E_{Q\text{-personalized}}W_Q \quad (8)$$

$$K = E_{S\text{-personalized}}W_K \quad (9)$$

$$V = E_{S\text{-personalized}}W_V \quad (10)$$

where  $W_Q \in \mathbb{R}^{d_Q \times d_A}$ ,  $W_K$  and  $W_V \in \mathbb{R}^{d_S \times d_A}$ , in which  $d_A$  refers to the dimension size in target attention. The attention weight  $z'[t]$  in the  $t$ -th query-key pair, i.e.,  $\{Q_t \in \mathbb{R}^{d_A}, K_t \in \mathbb{R}^{d_A}\}$ , can then be obtained by:

$$z'[t] = \frac{Q_t^T K_t}{\sqrt{d_A}}, \quad \text{where } 1 \leq t \leq T \quad (11)$$

following which a softmax-based operation is performed to normalize the personalized weights and aggregate the personalized sequence:

$$z = \text{softmax}(z'), \quad s = \sum_{t=1}^T (z[t] \cdot V_t) \quad (12)$$

With the sequence modeling output  $s$  and the other feature embeddings including  $E_U$ ,  $E_I$ , and  $E_O$ , a concatenation layer and several high-level MLPs are performed to merge all the information and output the prediction, and the training loss can be obtained via the binary cross-entropy function.

$$E_{all} = s \oplus E_U \oplus E_I \oplus E_O, \quad \hat{y} = \text{MLP}(E_{all}) \quad (13)$$

### 3 EXPERIMENTS

#### 3.1 Experimental Settings

**Datasets and Experimental Setup.** To sufficiently evaluate the proposed ADS, we conduct experiments on both a public dataset, i.e., Taobao Dataset<sup>3</sup> (*Taobao*), and two **billion-scale** industrial dataset from Douyin, i.e., Douyin advertising platform (*Douyin Ads*) and Douyin E-commerce platform (*Douyin Ecom*). The statistics of the three datasets are reported in Table 1.

- **Taobao.** The Taobao dataset released in [34] provides user behavior data in Taobao and is currently widely used in sequential modeling methods [3]. The former 7 days are used for training and the rest are for testing. Users with at least 200 interactions and 10 positive actions are filtered, and items with at least 10 interactions are filtered. There are 9,439 categories of items in the dataset, and we regard each category as a domain. In this

dataset, page view is considered as a negative interaction and the other actions are regarded as positive label (order).

- **Douyin Ads.** We select the Conversion Rate (CVR) prediction task in Douyin Ads, and collect a subset of online traffic logs from Dec. 14th, 2022 to Mar. 10th, 2023, 87 days and 1.73 billion samples in total. The former 77 days are used for model training and the rest 10 days are used for evaluation. In Douyin Ads platform, according to different user external actions, the dataset can be divided into three major domains including pay in live, order in live, and shopping in short-videos, denoted by domain 1, 2, and 3, respectively.
- **Douyin Ecom.** Two kinds of user shopping behaviors including click and order are selected as the prediction targets in the Ecom service in Douyin Live, the most influential scenario in E-commerce services at ByteDance. A subset of online traffic logs from Jan. 1st to Mar. 1st, 2024 is collected, including 61 days and 2.52 billion samples. The first 54 days are selected for training and the last week is for validation. The two typical scenarios in Douyin-Live, i.e., Live-Preview and Live-Slide, are involved in the dataset, denoted by domain 1 and 2, respectively.

**Comparing methods and evaluation metrics.** To comprehensively compare the proposed ADS with existing methods, we select several representative SOTA models, which can be grouped into two categories: (1) the DNN-based methods including **DNN**, **DeepFM** [9], **DCNv2** [24], **APG** [26], **AdaSparse** [28], **DFFM** [10], **MaskNet** [25]; (2) the target attention-based backbone methods including **DIN** [33], **MHA** [21], and **CAN** [2], based on which recent multi-domain embedding learning methods are added for comparison, including **FRNet** [22] and **PEPNet** [6]. All the methods are implemented based on Tensorflow[1], optimized by the cross-entropy loss. Adam [13] optimizer is adopted with an initial learning rate of 0.00002. AUC metric is used to evaluate the ranking model performance. Further, we introduce relative improvement (Imp.) [27] to measure the relative AUC gain, which is calculated as following as a random strategy yields AUC value at 0.5:

$$\text{AUC Imp.} = \left( \frac{\text{AUC}(\text{MeasuredModel}) - 0.5}{\text{AUC}(\text{BaselineModel}) - 0.5} - 1 \right) \times 100\% \quad (14)$$

#### 3.2 Experimental Results

**Overall Performance.** The comparison results of different methods on the three datasets are presented in Table 2. For Taobao dataset, only the overall results are provided due to the domain amount. For clarity, the results are reported with the comparing methods grouped into four groups, in which first group lists the DNN-based methods, and the rest-three groups list the multi-domain target attention methods based on DIN, MHA, and CAN, respectively. There are several observations from the results.

**First, it is observed that compared with the DNN-based approaches, target-attention based sequential-modeling plays a vital role in ranking models.** From the table, it is demonstrated that DIN, MHA, and CAN achieves 0.54% improvements in Taobao, 0.23%, 0.17%, and 0.08% improvements in Douyin Ads, 0.22%, 0.19%, and 0.04% improvements in click prediction, and 0.49%, 0.31%, and 0.11% improvements in order prediction tasks in Douyin Ecom, respectively, demonstrating significant improvements.

<sup>3</sup><https://tianchi.aliyun.com/dataset/dataDetail?dataId=649>

**Table 2: Comparison results of different methods on the three datasets. "D1", "D2", and "D3" are the abbreviations of "Domain 1", "Domain 2", and "Domain 3". Boldface denotes the best results in each group, and the boldface in the gray shadow area denotes that ADS significantly outperforms the second-best approach at the level of ( $p < 0.05$ ) in each group. For the two industrial datasets, note that a 0.1% Overall Imp. in Douyin Ads and 0.2% Overall Imp. in Douyin Ecom is considered to be significant improvement that can affect the performance of online A/B tests.**

Models	Taobao		Douyin Ads					Douyin E-commerce							
	Order		CVR					Click				Order			
	Overall	Overall Imp.	D1	D2	D3	Overall	Overall Imp.	D1	D2	Overall	Overall Imp.	D1	D2	Overall	Overall Imp.
DNN	0.6490	-	0.8224	0.8221	0.8717	0.8461	-	0.7949	0.8418	0.9082	-	0.8456	0.8494	0.8439	-
DeepFM	0.6503	+0.87%	0.8229	0.8233	0.8720	0.8469	+0.23%	0.7953	0.8422	0.9083	+0.02%	0.8458	0.8496	0.8442	+0.08%
DCNv2	0.6505	+1.01%	<b>0.8232</b>	<b>0.8235</b>	<b>0.8725</b>	<b>0.8470</b>	<b>+0.26%</b>	0.7950	0.8418	0.9082	+0.00%	0.8456	0.8494	0.8439	+0.00%
APG	0.6484	-0.40%	0.8225	0.8223	0.8720	0.8464	+0.08%	0.7951	0.8418	0.9082	+0.00%	0.8456	0.8494	0.8439	+0.00%
AdaSparse	<b>0.6498</b>	<b>+0.54%</b>	0.8228	0.8230	0.8721	0.8465	+0.11%	0.7951	0.8418	0.9082	+0.00%	0.8455	0.8491	0.8438	-0.03%
DFFM	0.6487	-0.20%	0.8232	0.8241	0.8722	0.8465	+0.11%	<b>0.7961</b>	<b>0.8440</b>	<b>0.9089</b>	<b>+0.17%</b>	<b>0.8469</b>	<b>0.8507</b>	<b>0.8455</b>	<b>+0.46%</b>
DIN	0.6498	+0.54%	0.8235	0.8251	0.8724	0.8469	+0.23%	0.7967	0.8443	0.9091	+0.22%	0.8472	0.8509	0.8456	+0.49%
+FRNet	0.6490	+0.00%	0.8235	0.8247	0.8722	0.8469	+0.23%	0.7967	0.8444	0.9092	+0.24%	0.8474	0.8511	0.8458	+0.55%
+PEPNet	0.6500	+0.67%	0.8235	0.8250	0.8723	0.8470	+0.26%	0.7964	0.8440	0.9090	+0.19%	0.8472	0.8507	0.8456	+0.49%
+ADS	<b>0.6507</b>	<b>+1.14%</b>	<b>0.8245</b>	<b>0.8262</b>	<b>0.8733</b>	<b>0.8477</b>	<b>+0.46%</b>	<b>0.7972</b>	<b>0.8451</b>	<b>0.9094</b>	<b>+0.29%</b>	<b>0.8477</b>	<b>0.8513</b>	<b>0.8462</b>	<b>+0.66%</b>
MHA	0.6498	+0.54%	0.8232	0.8238	0.8721	0.8467	<b>+0.17%</b>	0.7964	0.8439	0.9090	+0.19%	0.8465	0.8503	0.8450	+0.31%
+FRNet	0.6496	+0.40%	0.8233	0.8240	0.8722	0.8467	+0.17%	0.7963	0.8437	0.9089	+0.17%	0.8468	0.8506	0.8453	+0.40%
+PEPNet	0.6498	+0.54%	0.8237	0.8246	0.8725	0.8471	+0.28%	0.7964	0.8439	0.9090	+0.19%	0.8467	0.8505	0.8452	+0.37%
+ADS	<b>0.6501</b>	<b>+0.74%</b>	<b>0.8242</b>	<b>0.8256</b>	<b>0.8730</b>	<b>0.8475</b>	<b>+0.40%</b>	<b>0.7969</b>	<b>0.8448</b>	<b>0.9093</b>	<b>+0.26%</b>	<b>0.8473</b>	<b>0.8509</b>	<b>0.8458</b>	<b>+0.55%</b>
CAN	0.6498	+0.54%	0.8227	0.8238	0.8721	0.8464	+0.08%	0.7954	0.8423	0.9084	+0.04%	0.8459	0.8496	0.8443	+0.11%
+FRNet	0.6500	+0.67%	0.8234	0.8248	0.8725	0.8469	+0.23%	0.7958	0.8429	0.9086	+0.09%	0.8466	0.8503	0.8447	+0.23%
+PEPNet	0.6490	+0.00%	0.8231	0.8242	0.8723	0.8467	+0.17%	0.7958	0.8427	0.9085	+0.07%	0.8463	0.8502	0.8450	+0.31%
+ADS	<b>0.6503</b>	<b>+0.87%</b>	<b>0.8240</b>	<b>0.8256</b>	<b>0.8731</b>	<b>0.8474</b>	<b>+0.37%</b>	<b>0.7968</b>	<b>0.8446</b>	<b>0.9092</b>	<b>+0.24%</b>	<b>0.8470</b>	<b>0.8509</b>	<b>0.8456</b>	<b>+0.49%</b>

Second, the existing multi-domain methods contribute a positive effect for ranking models in general. Specifically, in DNN-based methods, it is observed that AdaSparse outperforms the baseline in Taobao and Douyin Ads. FRNet and PEPNet also show improved performance in different groups.

Finally, the proposed ADS consistently achieves the best performance in different groups with DIN, MHA, and CAN as backbone, showing its high effectiveness and compatibility. Specifically, in Taobao, the proposed ADS outperforms other SOTAs with 0.47%, 0.20%, and 0.20% with the second-best methods in the DIN-, MHA-, and CAN-based groups. In Douyin Ads, ADS beats other methods and achieves 0.20%, 0.12%, and 0.14% improvements compared with the second-best approach. In Douyin Ecom, ADS outperforms the second-best approach in the three groups by 0.05%, 0.07%, and 0.15% in click prediction task, and by 0.11%, 0.15%, and 0.18% in order prediction task, respectively. Besides, in terms of each domain of the two industrial datasets, ADS solidly outperforms the compared methods. Thus, the promising results show the superiority of the personalized target attention mechanism.

### 3.3 Ablation Study & Sensitivity Analysis

**Ablation study.** To further evaluate the performance of the two modules in ADS, i.e., PCRG and PSRG, here we perform ablation study in the challenging Douyin Ads dataset. As shown in Table 3, after ablating the PCRG module, the overall performance is dropped by 0.06%, 0.03%, and 0.23% in DIN-, MHA-, and CAN-based methods. Besides, after ablating both the PCRG and PSRG modules, the overall performance is decreased by 0.23%, 0.23%, and 0.29%, respectively. It can thus be concluded that both the personalized target item and the personalized sequences contributes positive impact to the proposed ADS, confirming the validity of these modules.

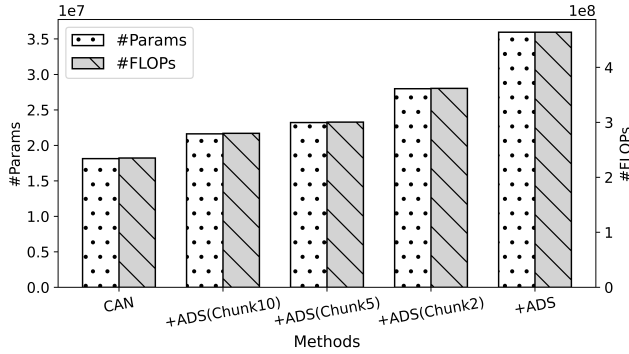
**Table 3: Ablation study of the proposed ADS.**

Groups	Methods	D1	D2	D3	Overall	Overall Imp.
DIN-based	ADS	0.8245	0.8262	0.8733	0.8477	-
	w/o PCRG	0.8242	0.8257	0.8730	0.8475	<b>-0.06%</b>
	w/o PCRG & PSRG	0.8235	0.8251	0.8724	0.8469	<b>-0.23%</b>
MHA-based	ADS	0.8242	0.8256	0.8730	0.8475	-
	w/o PCRG	0.8241	0.8252	0.8730	0.8474	<b>-0.03%</b>
	w/o PCRG & PSRG	0.8232	0.8238	0.8721	0.8467	<b>-0.23%</b>
CAN-based	ADS	0.8240	0.8256	0.8731	0.8474	-
	w/o PCRG	0.8230	0.8240	0.8723	0.8466	<b>-0.23%</b>
	w/o PCRG & PSRG	0.8227	0.8238	0.8721	0.8464	<b>-0.29%</b>

**Sensitivity Analysis to Number of Chunks in ADS.** Furthermore, to investigate the impact of number of chunks, sensitivity analysis for the ADS are conducted. Specifically, we investigate the performance patterns by varying the number of items in each chunk from [1, 2, 5, 10] in Douyin E-commerce, and the sensitivity analysis is performed from two aspects, i.e., training efficiency and model performance.

• **Training efficiency patterns.** We evaluate the model training efficiency by summarizing and comparing the model parameters and training floating point operations (FLOPs) under different chunks in ADS, and the results are illustrated in Fig 4. From the figure, it can be clearly observed that with the modeling becomes more personalized, the model parameters and training FLOPs consistently increase.





**Figure 4: Model parameter and training FLOPs patterns by varying the number of chunks in ADS, where "Chunk K" means there are K items grouped into a chunk in PCRG.**

**Table 4: Online A/B results in Douyin Ads. The boldface highlight denotes that the improvement is significant with  $p < 0.01$ .**

	Methods	CPM	ADV
Domain 1 (Short-Video)	baseline	+0.00%	+0.00%
	ADS	<b>+1.01%</b>	<b>+1.66%</b>
Domain 2 (Live)	baseline	+0.00%	+0.00%
	ADS	<b>+0.39%</b>	<b>+0.84%</b>
Overall	baseline	+0.00%	+0.00%
	ADS	<b>+0.52%</b>	<b>+1.00%</b>

- **Performance patterns.** The model performance pattern by varying the number of items in each chunk is illustrated in Fig 5. Specifically, first, in comparison with the vanilla DIN, MHA, and CAN, ADS and ADS with different chunks show obvious performance improvement in terms of both click and order prediction tasks. Besides, it can be observed that with number of items covered in the chunk decreases, the model performance continues to increase, and the most personalized model, i.e., ADS without chunk, achieves the best performance, showing that it is of significance to consider the personality characterization of the candidate item.

Overall, both the performance and the training cost increase with the personalization ability in ADS increased, while we observe that even with a small increase in training cost (ADS with Chunk 10 compared with vanilla method), the model performance still achieves promising gains, and thus practitioners can select the parameters with more flexibility according to the balance of effectiveness and efficiency.

### 3.4 Online Deployments

The ADS model is online deployed with distribution across multi GPUs with a sharding and data parallelism strategy. The low-frequency embeddings are eliminated to reduce storage. To further increase GPU throughputs, we introduce Dense Computation Asynchrony strategy. It splits the computation graph into Sparse-Forward and DenseCompute parts, which enables a pipeline effect

**Table 5: Online A/B results in Douyin Ecom. The boldface highlight denotes that the improvement is significant with  $p < 0.01$ .**

	Methods	GMV/U	Order/U	GPM
Domain 1 (Live-Preview)	baseline	+0.00%	+0.00%	+0.00%
	ADS	<b>+0.69%</b>	<b>+0.32%</b>	<b>+0.78%</b>
Domain 2 (Live-Slide)	baseline	+0.00%	+0.00%	+0.00%
	ADS	<b>+0.93%</b>	<b>+0.36%</b>	<b>+0.97%</b>
Overall	baseline	+0.00%	+0.00%	+0.00%
	ADS	<b>+0.79%</b>	<b>+0.36%</b>	<b>+0.89%</b>

and greatly improves training and inference efficiency. Benefiting from these, the offline training resource remains the same as baseline, i.e., 64 Nvidia A100s. Take the Douyin Ads as an example, the training time cost slightly increases from 41.3 to 42.8 hours (+3.6%). The online latency remains at 30ms with no significant changes.

### 3.5 Online A/B Experiments

To investigate the performance of the proposed ADS in real industrial scenarios, we perform careful online A/B testing in both the advertising system and the e-commerce system in Douyin, respectively.

- **Douyin Ads.** The online experiment on Douyin Ads is conducted from Nov. 2nd to Nov. 8th, 2023, hitting 74,079,729 users in Douyin APP. Two metrics including Cost Per Mile (CPM) and Advertiser Value (ADV) are selected for comparison. Note that the deployed scenario serves as the major traffic for advertisement in ByteDance with a strong baseline, where a 0.5% improvement in ADV or CPM is considered to be significant. The comparison results are shown in Table 4. The domains including pay in live and order in live are summarized as Domain 2 (Live). From the table, it is observed that after deploying the ADS, the overall CPM is improved by 0.52%, and the ADV is improved by 1.00%, showing the merits of adaptive sequence modeling of the proposed ADS. Besides, in the two major domains in Douyin Ads, including both live and short-video, the proposed ADS beats the baseline and achieves consistent improvements, showing the effectiveness in domain-aware sequential modeling for large-scaled industrial recommenders.
- **Douyin E-commerce.** The online experiment is performed from Jan. 23rd to Jan. 29th, 2024, on Douyin Ecom, hitting 508,926,918 users in Douyin APP. The experimental results are presented in Table 5. Three metrics are selected for comparison, i.e., the gross merchandise volume per user (GMV/U), the number of orders per user (Order/U), and GMV per Mille (GPM), which are all important commercial metrics in Douyin E-commerce. Similar to the experiment conducted on Douyin Ads, it is noted that this deployed scenario contributes the highest GMV in ByteDance with a very strong baseline, and generally, an improvement at the ratio of 0.5% in GMV is considered to be significant. As shown in Table 5, the overall GMV/U, Order/U, and GPM are uplifted by 0.79%, 0.36%, and 0.89%, respectively. Besides, consistent improvements are observed in both domains including live-preview

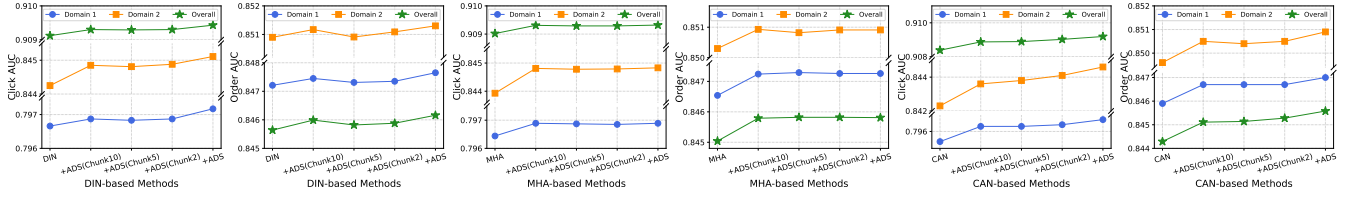


Figure 5: Performance patterns by varying the number of chunks in ADS.

and live-slide, and all the improvements are tested to be substantially statistically significant with a  $p$  – value  $< 0.01$ , verifying its efficacy.

## 4 RELATED WORKS

### 4.1 Ranking models in recommendation

The ranking model generally serves as the last stage of the industrial recommendation system funnel and plays a pivotal role in personalizing user recommendations in content and e-commerce platforms [5, 8]. Typically, they can be divided into two categories: 1) Feature interaction models which focus on modeling discrete identifiers (IDs) feature interactions and capturing the co-occurrence patterns [7, 9, 17, 23, 24, 29], and 2) Target-aware sequential models which capture user intention by evaluating the attention weight between the candidate item and the behavioral items [2, 16, 18, 21, 32, 33]. Despite the broad developments, the multi-domain problem that widely exists in large-scale recommender systems has been rarely considered in these methods.

### 4.2 Multi-Domain Recommendation

Multi-domain recommendation aims to provide accurate personalized recommendation results for users and items from multiple domains. Both coarse-grained [19, 20] and fine-grained [26, 28] multi-domain methods have been developed for sophisticated network structures, while the multi-domain learning for sequential modeling is seldom considered. Some other methods also contribute to learning personalized embeddings for multi-domain recommendation. The early idea is to use the self-gate mechanism to generate bitwise weights to reweigh the embeddings [11]. Further, the feature refinement network (FRNet) [22] takes user contexts as gate input, enabling user-level personalization on specific features. To further consider the impact of multidomains, parameter and embedding personalized network (PEPNet) [6] and domain-facilitated feature modeling (DFFM) [10] are developed, while bitwise scaling in PEPNet and linear transformation in DFFM are not flexible and sufficient enough to capture multidomains. Differently, ADS provides an approach by transforming the sequence embedding through a generated meta network, which provides a more flexible share-and-private modeling structure and is validated to be superior to existing methods. Besides, the ADS highlights that the identical target item has different influences on different sequence items. Traditional works can only obtain a novel target embedding shared for all sequence items, rather than multiple queries for different items, and this effective approach has not been explored in

recommendation, to our knowledge. Besides, the extensive experiments on both public dataset and two giant scenarios in ByteDance (billions of daily active users) have demonstrated the effectiveness compared with existing approaches.

## 5 CONCLUSION

In this paper, we propose a novel multi-domain ranking model named adaptive domain scaling (ADS) model, which builds personalized sequence representation generation (PSRG) and personalized candidate representation generation (PCRG) to separately generate personalized sequence item and candidate target item representations. Specifically, for PSRG, both a sequence-weight gen-net and sequence-bias gen-net are developed to formulate the instance-wise MLP and modify the sequence items, such that the identical item occurred in different user sequences has diverse representation. Besides, in PCRG, a multi-query gen-net is devised to generate multiple queries for sequence items, such that the identical candidate item extracts diverse intentions of users from the sequences from a multi-domain perspective. Offline experiments on both a public and two industrial datasets validate its significant and consistent improvements over existing SOTA methods, and extensive online experiments on two influential scenarios in ByteDance demonstrate its effectiveness.

ADS is currently fully deployed and achieves substantial enhancements in dozens of recommendation services at ByteDance, and now serves billions of users each day. We believe that it has offered a solid solution and propelled the advancements in multi-domain ranking from a sequential-modeling perspective, and more efficient and effective adaptive modeling structures will be explored in the future.

## REFERENCES

- [1] Martin Abadi, Paul Barham, Jianmin Chen, Zhifeng Chen, Andy Davis, Jeffrey Dean, Matthieu Devin, Sanjay Ghemawat, Geoffrey Irving, Michael Isard, et al. 2016. TensorFlow: a system for large-scale machine learning. In *12th USENIX symposium on operating systems design and implementation (OSDI 16)*. 265–283.
- [2] Weijie Bian, Kailun Wu, Lejian Ren, Qi Pi, Yujing Zhang, Can Xiao, Xiang-Rong Sheng, Yong-Nan Zhu, Zhangming Chan, Na Mou, et al. 2022. CAN: feature co-action network for click-through rate prediction. In *Proceedings of the fifteenth ACM international conference on web search and data mining*. 57–65.
- [3] Yue Cao, Xiaojiang Zhou, Jiaqi Feng, Peihao Huang, Yao Xiao, Dayao Chen, and Sheng Chen. 2022. Sampling is all you need on modeling long-term user behaviors for CTR prediction. In *Proceedings of the 31st ACM International Conference on Information & Knowledge Management*. 2974–2983.
- [4] R Caruana. 1993. Multitask learning: A knowledge-based source of inductive bias1. In *Proceedings of the Tenth International Conference on Machine Learning*. Citeseer, 41–48.
- [5] Zheng Chai, Zhihong Chen, Chenliang Li, Rong Xiao, Houyi Li, Jiawei Wu, Jingxu Chen, and Haihong Tang. 2022. User-aware multi-interest learning for candidate matching in recommenders. In *Proceedings of the 45th International ACM SIGIR Conference on Research and Development in Information Retrieval*. 1326–1335.



- [6] Jianxin Chang, Chenbin Zhang, Yiqun Hui, Dewei Leng, Yanan Niu, Yang Song, and Kun Gai. 2023. Pepnet: Parameter and embedding personalized network for infusing with personalized prior information. In *Proceedings of the 29th ACM SIGKDD Conference on Knowledge Discovery and Data Mining*. 3795–3804.
- [7] Heng-Tze Cheng, Levent Koc, Jeremiah Harmsen, Tal Shaked, Tushar Chandra, Hrishu Aradhye, Glen Anderson, Greg Corrado, Wei Chai, Mustafa Ispir, et al. 2016. Wide & deep learning for recommender systems. In *Proceedings of the 1st workshop on deep learning for recommender systems*. 7–10.
- [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] Huifeng Guo, Ruiming Tang, Yunming Ye, Zhenguo Li, and Xiuqiang He. 2017. DeepFM: a factorization-machine based neural network for CTR prediction. In *Proceedings of the 26th International Joint Conference on Artificial Intelligence*. 1725–1731.
- [10] Wei Guo, Chenxu Zhu, Fan Yan, Bo Chen, Weiwen Liu, Huifeng Guo, Hongkun Zheng, Yong Liu, and Ruiming Tang. 2023. DFFM: Domain Facilitated Feature Modeling for CTR Prediction. In *Proceedings of the 32nd ACM International Conference on Information and Knowledge Management*. 4602–4608.
- [11] Tongwen Huang, Qingyun She, Zhiqiang Wang, and Junlin Zhang. 2020. GateNet: gating-enhanced deep network for click-through rate prediction. *arXiv preprint arXiv:2007.03519* (2020).
- [12] Yuchen Jiang, Qi Li, Han Zhu, Jinbei Yu, Jin Li, Zirui Xu, Huihui Dong, and Bo Zheng. 2022. Adaptive domain interest network for multi-domain recommendation. In *Proceedings of the 31st ACM International Conference on Information & Knowledge Management*. 3212–3221.
- [13] Diederik P Kingma and Jimmy Ba. 2015. Adam: A method for stochastic optimization. In *Proceedings of ICLR*.
- [14] Chenglin Li, Yuanzhen Xie, Chenyun Yu, Bo Hu, Zang Li, Guoqiang Shu, Xiaohu Qie, and Di Niu. 2023. One for all, all for one: Learning and transferring user embeddings for cross-domain recommendation. In *Proceedings of the Sixteenth ACM International Conference on Web Search and Data Mining*. 366–374.
- [15] Hui Lu, Zheng Chai, Yuchao Zheng, Zhe Chen, Deping Xie, Peng Xu, Xun Zhou, and Di Wu. 2025. Large Memory Network for Recommendation. In *Proceedings of the ACM Web Conference*. <https://doi.org/10.1145/3701716.3715514>
- [16] Qi Pi, Guorui Zhou, Yujing Zhang, Zhe Wang, Lejian Ren, Ying Fan, Xiaoqiang Zhu, and Kun Gai. 2020. Search-based user interest modeling with lifelong sequential behavior data for click-through rate prediction. In *Proceedings of the 29th ACM International Conference on Information & Knowledge Management*. 2685–2692.
- [17] Steffen Rendle. 2010. Factorization machines. In *2010 IEEE International conference on data mining*. IEEE, 995–1000.
- [18] Qijie Shen, Hong Wen, Wanjie Tao, Jing Zhang, Fuyu Lv, Zulong Chen, and Zhao Li. 2022. Deep interest highlight network for click-through rate prediction in trigger-induced recommendation. In *Proceedings of the ACM Web Conference 2022*. 422–430.
- [19] Xiang-Rong Sheng, Liqin Zhao, Guorui Zhou, Xinyao Ding, Binding Dai, Qiang Luo, Siran Yang, Jingshan Lv, Chi Zhang, Hongbo Deng, et al. 2021. One model to serve all: Star topology adaptive recommender for multi-domain ctr prediction. In *Proceedings of the 30th ACM International Conference on Information & Knowledge Management*. 4104–4113.
- [20] Hongyan Tang, Junling Liu, Ming Zhao, and Xudong Gong. 2020. Progressive layered extraction (ple): A novel multi-task learning (mtl) model for personalized recommendations. In *Proceedings of the 14th ACM Conference on Recommender Systems*. 269–278.
- [21] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. 2017. Attention is all you need. *Advances in neural information processing systems* 30 (2017).
- [22] Fangye Wang, Yingxu Wang, Dongsheng Li, Hansu Gu, Tun Lu, Peng Zhang, and Ning Gu. 2022. Enhancing CTR prediction with context-aware feature representation learning. In *Proceedings of the 45th International ACM SIGIR Conference on Research and Development in Information Retrieval*. 343–352.
- [23] Ruoxi Wang, Bin Fu, Gang Fu, and Mingliang Wang. 2017. Deep & cross network for ad click predictions. In *Proceedings of the ADKDD'17*. 1–7.
- [24] Ruoxi Wang, Rakesh Shivanna, Derek Cheng, Sagar Jain, Dong Lin, Lichan Hong, and Ed Chi. 2021. DCN v2: Improved deep & cross network and practical lessons for web-scale learning to rank systems. In *Proceedings of the web conference 2021*. 1785–1797.
- [25] Zhiqiang Wang, Qingyun She, and Junlin Zhang. 2021. Masknet: Introducing feature-wise multiplication to CTR ranking models by instance-guided mask. *arXiv preprint arXiv:2102.07619* (2021).
- [26] Bencheng Yan, Pengjie Wang, Kai Zhang, Feng Li, Hongbo Deng, Jian Xu, and Bo Zheng. 2022. Apg: Adaptive parameter generation network for click-through rate prediction. *Advances in Neural Information Processing Systems* 35 (2022), 24740–24752.
- [27] Ling Yan, Wu-Jun Li, Gui-Rong Xue, and Dingyi Han. 2014. Coupled group lasso for web-scale ctr prediction in display advertising. In *International conference on machine learning*. PMLR, 802–810.
- [28] Xuanhua Yang, Xiaoyu Peng, Penghui Wei, Shaoguo Liu, Liang Wang, and Bo Zheng. 2022. Adaspase: Learning adaptively sparse structures for multi-domain click-through rate prediction. In *Proceedings of the 31st ACM International Conference on Information & Knowledge Management*. 4635–4639.
- [29] Buyun Zhang, Liang Luo, Xi Liu, Jay Li, Zeliang Chen, Weilin Zhang, Xiaohan Wei, Yuchen Hao, Michael Tsang, Wenjun Wang, et al. 2022. DHEN: A deep and hierarchical ensemble network for large-scale click-through rate prediction. In *Proceedings of the DLP-KDD*.
- [30] Shuai Zhang, Lina Yao, Aixin Sun, and Yi Tay. 2019. Deep learning based recommender system: A survey and new perspectives. *ACM computing surveys (CSUR)* 52, 1 (2019), 1–38.
- [31] Yiqian Zhang, Yinfu Feng, Wen-Ji Zhou, Yunan Ye, Min Tan, Rong Xiao, Haihong Tang, Jiajun Ding, and Jun Yu. 2024. Multi-Domain Deep Learning from a Multi-View Perspective for Cross-Border E-commerce Search. In *Proceedings of the AAAI Conference on Artificial Intelligence*, Vol. 38. 9387–9395.
- [32] Guorui Zhou, Na Mou, Ying Fan, Qi Pi, Weijie Bian, Chang Zhou, Xiaoqiang Zhu, and Kun Gai. 2019. Deep interest evolution network for click-through rate prediction. In *Proceedings of the AAAI conference on artificial intelligence*, Vol. 33. 5941–5948.
- [33] 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.
- [34] Han Zhu, Xiang Li, Pengye Zhang, Guozheng Li, Jie He, Han Li, and Kun Gai. 2018. Learning tree-based deep model for recommender systems. In *Proceedings of the 24th ACM SIGKDD international conference on knowledge discovery & data mining*. 1079–1088.