# Deep Pattern Network for Click-Through Rate Prediction

Hengyu Zhang
zhang-hy21@mails.tsinghua.edu.cn
Tsinghua Shenzhen International
Graduate School, Tsinghua University
Shenzhen, China

Junwei Pan
jonaspan@tencent.com
Tencent
Shenzhen, China

Dapeng Liu
rocliu@tencent.com
Tencent
Shenzhen, China

Jie Jiang
zeus@tencent.com
Tencent
Shenzhen, China

Xiu Li*
li.xiu@sz.tsinghua.edu.cn
Tsinghua Shenzhen International
Graduate School, Tsinghua University
Shenzhen, China

## ABSTRACT

Click-through rate (CTR) prediction tasks play a pivotal role in real-world applications, particularly in recommendation systems and online advertising. A significant research branch in this domain focuses on user behavior modeling. Current research predominantly centers on modeling co-occurrence relationships between the target item and items previously interacted with by users in their historical data. However, this focus neglects the intricate modeling of user behavior patterns. In reality, the abundance of user interaction records encompasses diverse behavior patterns, indicative of a spectrum of habitual paradigms. These patterns harbor substantial potential to significantly enhance CTR prediction performance. To harness the informational potential within user behavior patterns, we extend Target Attention (TA) to Target Pattern Attention (TPA) to model pattern-level dependencies. Furthermore, three critical challenges demand attention: the inclusion of unrelated items within behavior patterns, data sparsity in behavior patterns, and computational complexity arising from numerous patterns. To address these challenges, we introduce the Deep Pattern Network (DPN), designed to comprehensively leverage information from user behavior patterns. DPN efficiently retrieves target-related user behavior patterns using a target-aware attention mechanism. Additionally, it contributes to refining user behavior patterns through a pre-training paradigm based on self-supervised learning while promoting dependency learning within sparse patterns. Our comprehensive experiments, conducted across three public datasets, substantiate the superior performance and broad compatibility of DPN.

## CCS CONCEPTS

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

## KEYWORDS

User Behavior Pattern, Click-Through Rate Prediction, Recommendation System
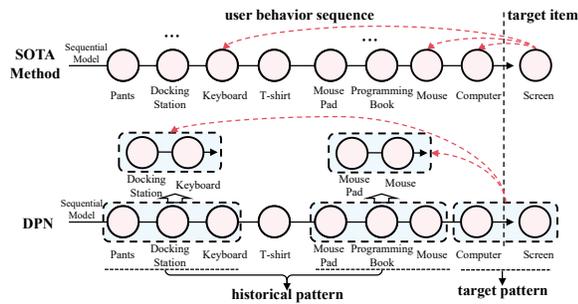
---

*Corresponding author.

## 1 INTRODUCTION

Click-Through Rate (CTR) prediction holds paramount significance in industrial scenarios such as online advertising and recommender systems, where the goal is to predict the probability of the specific user clicking on a target item. Effective CTR prediction not only serves to maximize revenue for advertisers but also enhances the user experience by delivering more pertinent content. The pursuit of better prediction models has led to the widespread adoption of deep learning techniques in the field of click-through rate prediction.

Deep-learning-based CTR prediction methods have taken the forefront in current research and achieved remarkable success [47], which aim to capture intricate feature interactions using neural networks. User behavior modeling, which seeks to extract latent interest from the extensive user history behavior records via various techniques including pooling, RNN-based [20, 49], CNN-based [39], Graph-based [27, 28, 41], and attention [6, 21, 37, 46] model, represent one of the most crucial research branches [25, 49–51].
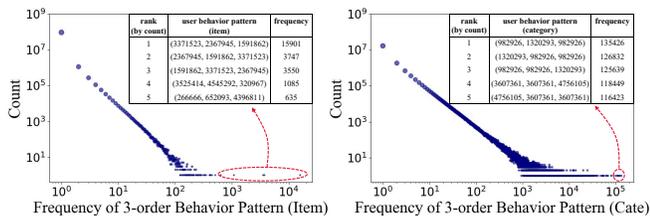
DIN [50] stands out as a notable milestone in user behavior modeling. User interest is diverse and varies for different target items. DIN handles this issue by introducing a target-attention mechanism that enables the extraction of specific interests for a given target item. Due to its superiority, DIN-based methods have become the mainstream branch in CTR prediction based on user behavior modeling. DIEN [49] further improves the model performance by incorporating GRU [9] networks for dynamic evolutionary modeling. Recent research endeavors have extended this progress in terms of modeling long sequences of user behavior [5, 29, 30] and co-occurrence behavior representation modeling [2]. In summary, the state-of-the-art (SOTA) methods typically apply a sequential model to handle the user behavior sequence and use target attention to aggregate the item representations, as shown in Figure 1.

According to psychological studies [16, 22], the entire user personality is linked to a variety of behavior patterns. Specifically, a *behavior pattern* can be defined as a *subsequence of two or more actions that occur in a prescribed arrangement*, such as purchasing

**Figure 1: Illustration of the dependency modeling that the existing SOTA methods and DPN focus on. The red dotted line indicates Target Attention or Target Pattern Attention.**

a *mouse pad* and then buying a *mouse* as shown in Figure 1. The user behavior sequence is composed of diverse behavior patterns driven by varying user interest. However, the existing SOTA methods uniformly model user behavior sequences containing diverse interests via a unified sequential model, overlooking the intricate modeling of the varied behavior patterns within. These methods solely focus on item-level dependency, neglecting the critical aspect of dependency modeling among patterns. Furthermore, many fixed behavior patterns frequently appear in recommendation scenarios as shown in Figure 2, demonstrating they embed meaningful dependency information. Hence, exploring how to effectively leverage diverse behavior patterns in CTR prediction tasks is an intriguing research question.



**Figure 2: Statistical results for 3-order user behavior patterns in the Taobao dataset (log-scale plot). Each scatter point with coordinates $(x, y)$ represents that the number of 3-order user behavior patterns with occurrence frequency $x$ in the Taobao dataset is $y$. The figure shows the Top-5 high-frequency user behavior patterns for items and categories, respectively.**

To fully leverage the information from behavior patterns, we extend Target Attention (TA) into **T**arget **P**attern **A**ttention (**TPA**) to perform *pattern-level* interest aggregation. TPA models the dependency between the target behavior pattern and the historical patterns, where the target pattern consists of the user's recent behaviors and the interaction to be predicted with the target item. However, due to the multitude and complexity of behavior patterns contained within user behavior sequences, there are several challenges making it a non-trivial problem:

- **C1:** **Unrelated items mixed into behavior patterns.** In some cases, behavior patterns may not be continuous segments, which

means some irrelevant items misclicked or driven by other interest may mix in. For example in Figure 1, the user interacts with the mouse pad, programming book, and mouse successively, but the programming book may not have a strong relationship with the pattern (mouse pad, mouse). So effective refinement of behavior patterns is necessary, which guarantees patterns to be more meaningful, i.e., with stronger intra-dependencies within the items in behavior patterns, and reduces the risks of introducing noise.

- **C2:** **Data Sparsity of Behavior Patterns.** Behavior patterns are subsequences of multiple items interacted with by the user, implying high-order item dependencies. Some behavior patterns exhibit high sparsity, leading to difficulty in dependency learning within them.

- **C3:** **Computational complexity arising from numerous patterns.** User historical behavior records contain a plethora of complex user behavior patterns driven by different interest. Modeling all the behavior patterns incurs unacceptable computation costs, and the behavior patterns irrelevant to the target item may introduce noises. Consequently, efficiently retrieving target-related behavior patterns from the abundance of user historical behaviors is a crucial issue to address.

To address the challenges above, we introduce a novel **D**eep **P**attern **N**etwork (**DPN**) for the Click-Through Rate prediction task. In contrast to existing user behavior modeling methods, DPN distinguishes itself by not only efficiently discovering effective behavior patterns but also introducing pattern-level dependency modeling, as illustrated in Figure 1.

Firstly, for the efficiency of behavior pattern modeling (**C3**), we introduce a Target-aware Pattern Retrieval Module (TPRM) to identify the Top-K target-related user behavior patterns. Secondly, we ingeniously integrate the two major desideratas of effectively refining behavior patterns (**C1**) and addressing pattern sparsity (**C2**). To achieve this, we design an innovative self-supervised pattern refinement module, thereby avoiding excessive time and storage overhead associated with complex structural designs. SPRM is pre-trained by a self-supervised denoising task, specifically removing noise introduced by random augmentation of behavior patterns, enhancing the process of pattern refinement and optimizing pattern representation learning. Thirdly, DPN extends the concept of target attention to pattern level, proposing Target Pattern Attention to model the dependencies between patterns.

From a model expansion standpoint, our DPN can be regarded as **extend** the width of Query, Key, and Value within the attention mechanism, providing a new perspective that deviates significantly from previous efforts concentrated on making attention module in recommendation longer [3, 5, 30] or deeper [45]. This idea shares some similarities with works in CV [12, 13] and NLP [11, 33]. For instance, the concept of widening query, key, and value from a single item to the pattern (consecutive items) resembles extending convolution windows from $1 \times 1$ to larger ones, enhancing the model's perception of the local context. To some extent, our work can also be aligned with Vision Transformers (ViT), where items can be regarded as pixels and patterns as patches. In the realm of NLP, some works like ZEN [11] also explore the expansion of models from uni-gram to N-gram to enhance transformer representations.

Our contributions can be summarized as follows:

- We highlight the significance of behavior patterns and summarize the challenges of efficiently exploiting them.
- To address these challenges, we propose DPN, which includes target-aware retrieval, self-supervised refinement, and pattern-level interest aggregation to fully leverage behavior patterns.
- We conduct comprehensive experiments on three public real-world datasets to demonstrate the superior effectiveness and broad compatibility of our proposed DPN. Furthermore, further in-depth analysis of how user behavior patterns are utilized is conducted.

## 2 RELATED WORKS

In early research work, click-through rate prediction models mainly focused on improving CTR prediction performance by means of feature engineering [4, 26, 32]. In recent years, deep learning techniques have developed rapidly and achieved impressive results in areas such as computer vision [18, 35] and natural language processing [38, 40]. Therefore, researchers are currently working on capturing feature interactions automatically through deep neural networks. Approaches like DeepFM [17] and its derivatives xDeepFM [24] employ neural networks inspired by factorization machines to autonomously learn feature interactions, eliminating the need for manual feature engineering. PNN [31] introduces a product layer capable of capturing feature interactions by utilizing either inner or outer products as factorization functions. Lastly, ONN [44] devises operation-aware embedding layers to enhance feature interaction through enriched embeddings.

A significant research branch in CTR prediction focuses on user behavior modeling, giving birth to a series of representative works [2, 25, 49, 50]. The CTR prediction model based on user behavior modeling is dedicated to capturing the implicit interest representations of users from the rich history of their behavior records. DIN [50] is a pioneer work in this field of research. It proposes a target attention mechanism, which treats target items as queries and user history as keys and values, to aggregate the user's specific interests for a given target item from the user's history of interactions. DIEN [49], an improved version of DIN, captures dynamic interest representations of users by modeling the evolution of their interests through GRU [9]. AutoAttention [48] automated field selection in the target attention mechanism. MIMN [29] employs a memory network [15] to consolidate historical user behavior data, effectively tackling the challenge of modeling users' long-term interests. Subsequent research, exemplified by SIM [30] and ETA [5], concentrates on the efficient retrieval-based modeling of long-term interests. CAN [2] disentangles the process of modeling feature interactions from the initial feature modeling, by directing features into a compact multi-layer perceptron (mini-MLP) generated by other features. TIN [51] incorporates target-aware temporal encoding into the target attention mechanism to capture both semantic and temporal correlation.

## 3 METHOD

The overall architecture of DPN is depicted in Figure 3. The inputs to DPN encompass the target item to be predicted in the CTR prediction task, the user's historical behavior sequence, and other relevant features. The goal of DPN is to predict the probability of a user interacting with the target item.

### 3.1 Base Model

*3.1.1 Embedding Layer.* For a given set of users $\mathcal{U}$, the historical behavior records of user $u \in \mathcal{U}$ can be formulated as a chronological sequence $S^u = \{b_i^u, b_2^u, \cdots, b_{N_u}^u\}$, where $b_k^u$ denotes the $k$-th behavior record of user $u$ and $N_u$ denotes the length of user interaction sequence. Notably, the superscript $u$ will be omitted for simplicity in the following presentation. To enable parallel computation in tensor form, the behavior sequence should be truncated or padded to a fixed length $N$, resulting in a fix-length behavior sequence $S = \{b_1, b_2, \cdots, b_N\}$.

Typically, the behavior record $b_k$ consists of several features, such as item ID $i_k$ and category ID $c_k$ for our experiments. So for each interaction $b_k$, we can map it to a dense embedding $\mathbf{e}_k$ by concatenating the embeddings of corresponding features obtained via LookUp tables, which can be formulated as follows:

$$\mathbf{e}_k = \left[ \text{LookUP}(i_k, \mathbf{E}^{\mathcal{I}}); \text{LookUP}(c_k, \mathbf{E}^C) \right], \quad (1)$$

where $[\cdot; \cdot]$ means concatenateing the embedding vectors, and $\text{LookUp}(a, \mathbf{E})$ denotes the operation of retrieving the $a$-th embedding vector in embedding matrix $\mathbf{E}$. $\mathbf{E}^{\mathcal{I}} \in \mathbb{R}^{|\mathcal{I}| \times d}$, $\mathbf{E}^C \in \mathbb{R}^{|C| \times d}$ are embedding matrixes of items and categories, respectively, where $d$ denotes the embedding dimensionality and $|\mathcal{I}|, |C|$ means the volume of item space and category space.

Thus, we can get the sequence of user behavior embeddings:

$$\mathbf{E}_S = \{\mathbf{e}_1, \mathbf{e}_2, \cdots, \mathbf{e}_N\}. \quad (2)$$

*3.1.2 User Interest Extraction Module.* For a fair comparison, we adopt the typical user-behavior-based CTR prediction model DIN [50] as the user interest extraction backbone, following [5, 30]. DIN proposes the target attention mechanism to capture the specific user interest representation $\mathbf{v}_T$ for the target item $i_T$ of the target category $c_T$.

The target attention mechanism can be formulated as:

$$\mathbf{v}_T = \sum_{k=1}^{N} a(\mathbf{e}_T, \mathbf{e}_k) \cdot \mathbf{e}_k, \quad (3)$$
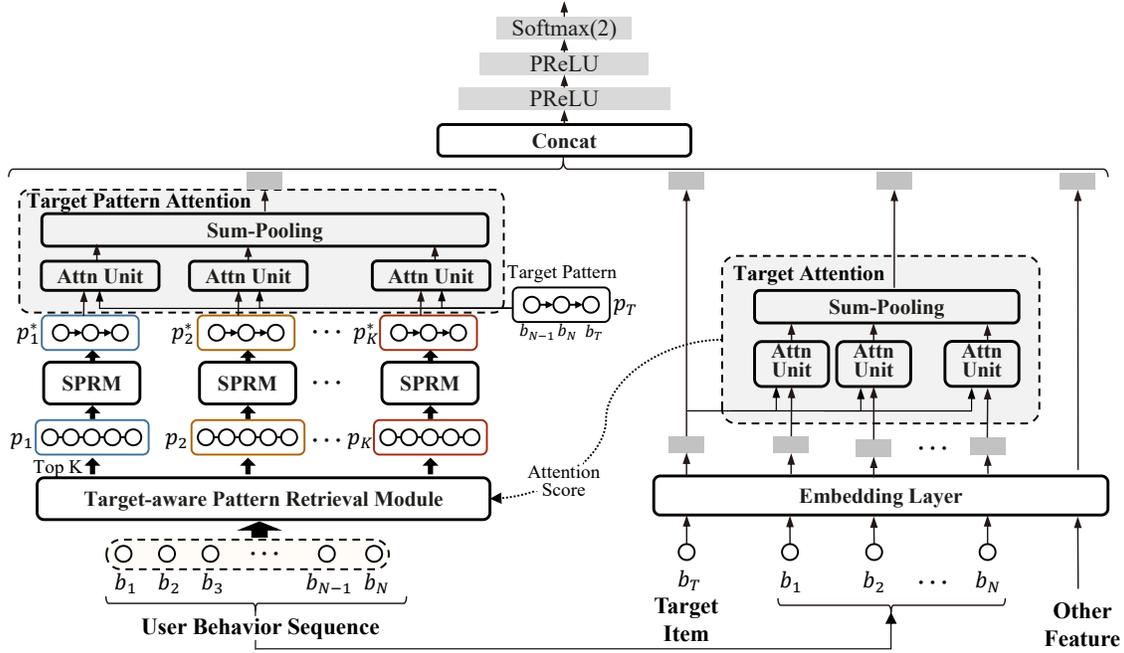
where $a(\cdot, \cdot)$ means the attention function, $\mathbf{e}_T$ denotes the embedding of the target item $\mathbf{i}_T$ and the target category $\mathbf{c}_T$, which is to be predicted whether to be clicked by the user. The attention function used in DIN [50] is $a(\mathbf{q}, \mathbf{k}) = \text{softmax}(\text{MLP}([\mathbf{q}; \mathbf{k}; \mathbf{q} - \mathbf{k}; \mathbf{q} \circ \mathbf{k}]))$, where $\circ$ denotes the Hadamard product, i.e., element-wise multiplication.

### 3.2 Target-aware Pattern Retrieval Module (TPRM)

Target-aware Pattern Retrieval Module aims to search for the user behavior pattern relative to the target item and category.

We retrieve the positions with Top-K attention scores in the User Interest Extraction Module to obtain similar behavior records with the target item $i_T$ and category $c_T$, which can be presented as follows:

$$idx_1, idx_2, \cdots, idx_K = \underset{K}{\arg\text{topk}} \ a(\mathbf{e}_T, \mathbf{e}_1), \cdots, a(\mathbf{e}_T, \mathbf{e}_N), \quad (4)$$

Figure 3: Overall architecture of the proposed DPN. Target-ware Pattern Retrieval Model (TPRM) searches the Top-K target-related pattern from user behavior sequence according to target attention scores. SPRM adopts a pre-trained refinement network to perform fine-grained denoising of the behavior patterns, whose detailed illustration is shown in Figure 4. TPA models the dependency between the target behavior pattern and refined historical behavior patterns.

where argtopk retrieves the corresponding indexes of the Top-K
$K$
largest attention scores. The retrieval method is not the core of our work, it can be replaced by other approaches such as locality-sensitive hashing for ultra-long sequence modeling.
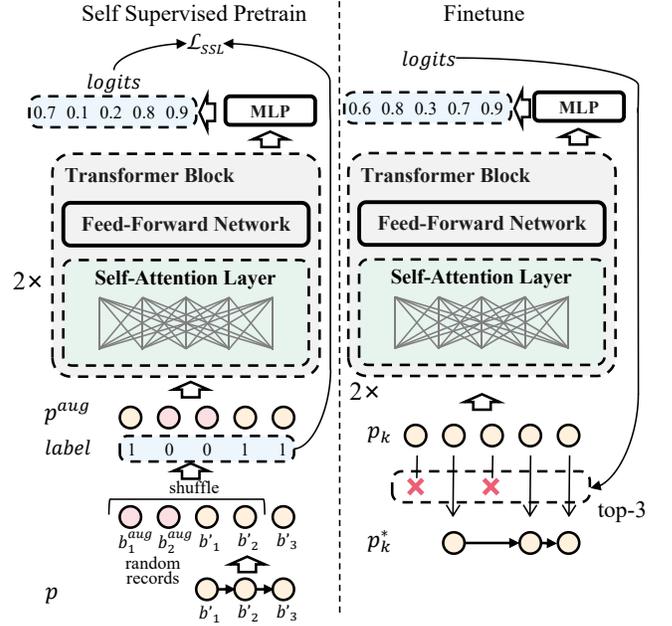
The $k$-th user behavior pattern $p_k$ of length $l$ consists of the $idx_k$-th interaction and the $l-1$ interactions preceding it, which can be denoted as:

$$p_k = \{b_{idx_k-l+1}, \cdots, b_{idx_k-1}, b_{idx_k}\}, \quad k = 1, 2, \cdots, K. \quad (5)$$

As a result, we obtain $K$ user behavior patterns. Since this is a target-aware retrieval process, these retrieved behavior patterns reflect the behavior paradigms that users may engage in before clicking on the target item and category.

### 3.3 Self-supervised Pattern Refinement Module (SPRM)

User behavior patterns are not always continuous segments, which means some irrelevant items misclicked or driven by other interest may mix in. For example, if the user interacts with the iPhone, T-shirt, and Airpods successively, only the transition from iPhone to Airpods is a meaningful pattern. Therefore, we propose a self-supervised pattern refinement network to extract true user behavior patterns from the retrieved raw user behavior patterns, as presented in Figure 4. Firstly, the refinement network is pretrained by a self-supervised denoising task. Then DPN applies the pre-trained refinement network to refine the retrieved pattern to be more meaningful, which also ensures that subsequent pattern-level dependency modeling is more meaningful.



Figure 4: Illustration of Self-supervised Pattern Refinement Module (SPRM). On the left, random data augmentation and denoising objectives are applied to pretrain refinement network based on self-supervised learning. On the right, the refinement network is utilized to generate a logit vector for pattern refinement in CTR prediction tasks.

**Data Augmentation.**

We take sequential interaction records in the dataset, i.e., one user behavior pattern, as the raw input $p = \{b'_1, b'_2, \cdots, b'_s\}$. We then perform data augmentation by randomly mixing in interaction records with some sampled interaction records into the original behavior pattern, resulting in an augmented pattern $p^{aug} = \{b''_1, b''_2, \cdots, b''_l\}$. Notably, the last record of the augmented behavior pattern sequence and the original behavior pattern sequence should be consistent, which indicates to the refinement network what the decision goal of the behavior pattern is.

**Refinement Network.**

We use a two-layer Transformer network as an encoder to extract behavior pattern representations from the augmented sequence of user behavior patterns. For a specific augmented behavior pattern $p^{aug} = \{b''_1, b''_2, \cdots, b''_l\}$, the encoder can be formalized as follows:

$$\mathbf{X}^{(0)} = \left[\mathbf{e}''_1 + \mathbf{e}^{pos}_1, \mathbf{e}''_2 + \mathbf{e}^{pos}_2, \cdots, \mathbf{e}''_l + \mathbf{e}^{pos}_l\right], \qquad (6)$$

$$\mathbf{H}^{(m)} = \text{LayerNorm}\left(\mathbf{X}^{(m-1)} + \text{SA}\left(\mathbf{X}^{(m-1)}\right)\right),$$
$$\mathbf{X}^{(m)} = \text{LayerNorm}\left(\mathbf{H}^{(m)} + \text{FFN}\left(\mathbf{H}^{(m)}\right)\right), \qquad (7)$$

where $\mathbf{e}''_j$ and $\mathbf{e}^{pos}_j$ are the embedding vectors of $j$-th record $b''_j$ in augmented behavior pattern $p^{aug}$ and the $j$-th positional embedding, respectively. SA and FFN denote the Self-Attention Layer and Feed-Forward Network, LayerNorm means layer normalization. $\mathbf{X}^{(m)}, \mathbf{H}^{(m)}$ means the hidden representation and the intermediate representation at layer $m$, and $\mathbf{X}^{(0)}$ is the input of transformer.

Thus, we can get the representation vector $\mathbf{X}^{(M)}$ of the augmented behavior pattern $p^{aug}$, where the depth $M$ of the Transformer network is set to 2. Then, we encode the representation into a logits vector $\mathbf{o}$ of length $l$ by MLP to predict which interaction records in the augmented behavior pattern $p^{aug}$ correspond to the raw behavior pattern $p$, which can be formulated as:

$$\mathbf{o} = \text{MLP}\left(\mathbf{X}^{(M))}\right) \in \mathbb{R}^l, \qquad (8)$$

where the positions corresponding to the Top-$s$ largest values in logits vector $\mathbf{o}$ indicate the predicted interaction records belonging in raw pattern $p$. In the CTR prediction task, SPRM selects Top-$s$ interaction records according to the logits vector $\mathbf{o}$ as the refined behavior pattern, as shown in the right part of Figure 4.

**Loss Function.**

Cross-entropy loss is used for pre-training based on self-supervised learning to guide the model in the refinement of behavior patterns. The loss function can be represented as:

$$\mathcal{L}_{SSL} = -\sum_{j=1}^{l} \mathbb{I}_p(b''_j) * \log(o_j), b''_j \in p^{aug}, \qquad (9)$$

where $\mathbb{I}_p(b''_j)$ denotes the indicator function to point out whether $b''_j$ belongs to behavior pattern $p$. If so, the function output is 1, and vice versa 0.

### 3.4 Target Pattern Attention (TPA)

User behavior patterns reflect the implicit psychological decision paradigm of users. Modeling the dependency between the historical

behavior patterns and the current target behavior pattern can effectively help the model determine whether the clicks on the target items and categories are in line with the user's habitual paradigm, thus improving the model's recommendation performance.

DPN extends the well-known target attention to Target Pattern Attention to achieve inter-pattern dependency modeling. After obtaining the refined target-related user behavior patterns $p_k^* = \{b_1^*, b_2^*, \cdots, b_s^*\}$, DPN applies the attention units to capture the dependency between the refined pattern $p_k$ and the target behavior pattern $p_T = \{b_{N-s+2}, \cdots, b_N, b_T\}$ resulting in pattern-level interest representation $\mathbf{v}_p$ that are used in the final CTR prediction task:

$$\mathbf{v}_p = \sum_{k=1}^{K} a(\mathbf{E}_{p_T}, \mathbf{E}_{p_k^*}) \cdot \mathbf{E}_{p_k^*}, \qquad (10)$$

where $\mathbf{E}_{p_k^*}$ and $\mathbf{E}_{p_T}$ denote the representations of $p_k^*$ and $p_T$ generated by Transformer encoders, respectively.

### 3.5 Training Objective

The task of click-through rate prediction can be structured as a binary classification problem, specifically, determining whether the target instance will be clicked or not. Consequently, the final output of the DPN model can be expressed as:

$$(\underset{\substack{\uparrow \\ \text{click}}}{\hat{y}}, \underset{\substack{\uparrow \\ \text{unclick}}}{1-\hat{y}}) = \text{softmax}(\text{MLP}(\left[\mathbf{e}_T; \sum_{j=1}^{N} \mathbf{e}_j; \mathbf{v}_T; \mathbf{v}_p \circ \mathbf{E}_{p_T}\right]). \quad (11)$$

Here, $\hat{y}$ and $1 - \hat{y}$ represent the predicted logit values for click and unclick, respectively. MLP denotes a compact neural network with three layers and a 2-dimensional output.

The training objective of the CTR prediction task can be defined as follows:

$$\min_{\Theta} \mathcal{L}_{CTR} = -\frac{1}{B} \sum_{\mathcal{B}} (y \cdot \log(\hat{y}) + (1 - y) \cdot \log(1 - \hat{y})). \qquad (12)$$

In this equation, $y$ is the label indicating whether the user will interact with the target item. $\Theta$ denotes the set of trainable parameters in DPN, and $\mathcal{B}$ represents a batch of instances with a batch size of $B$.

### 3.6 Model Discussion

*3.6.1 Time Complexity Analysis.* DPN mainly contains the Base Model and three core components. The time complexity of the Base Model, which defaults to DIN, is $O(BNd)$. TPRM is proposed to retrieve Top-K target-related behavior patterns from user behavior sequences, so its time complexity is $O(BN \log K)$ if Top-K retrieval is performed using heap structure. A two-layer Transformer is used in SPRM to encode the raw user behavior patterns. The time complexity of self-attention layer and feed-forward network are $O(BKl^2d)$ and $O(BKld^2)$, respectively, so the complexity amount of SPRM is $O(BKld(l+d))$. TPA involves only attention units between historical behavior patterns and target behavior patterns, which has a time complexity of $O(BKd)$.

**Table 1: Statistics of evaluation datasets.**

| Dataset | #User | #Item | #Interactions |
|---------|-------|-------|---------------|
| Tmall | 424,170 | 1,090,390 | 54,925,330 |
| Taobao | 987,994 | 4,162,024 | 100,150,807 |
| Alipay | 963,923 | 2,353,207 | 44,528,127 |

## 4 EXPERIMENTS

### 4.1 Experiment Settings

*4.1.1 Dataset Descriptions.* To comprehensively assess the performance of our proposed DPN, we conducted experiments on three public recommendation datasets. These datasets are as follows:

- **Tmall Dataset**[1]. This dataset comprises user behavior records on Tmall.com, spanning from May 2015 to November 2015. To construct the user behavior sequences for our experiments, we arrange the interaction records from the Tmall dataset in chronological order based on their action time.
- **Taobao Dataset**[2]. This dataset, provided by Alibaba Group, captures diverse user shopping activities on Taobao.com, including clicks, adding items to the cart, and making purchases. In our experiments, we treat all these user behaviors as interactions.
- **Alipay Dataset**[3]. Alipay dataset encompasses a comprehensive record of user activities both online and on-site during the period spanning from July 1st, 2015, to November 30th, 2015. This extensive dataset draws its data from Tmall.com, Taobao.com, and the Alipay app, providing valuable insights into user behaviors and interactions with merchants.

*4.1.2 Evaluation Metrics.* In our research trials, we assess the effectiveness of the models by employing the Area Under the Curve (AUC) metric, which is the most frequently used metric in Click-Through Rate (CTR) prediction tasks. AUC gauges the model's ability to prioritize positive samples over randomly selected negative samples. A greater AUC value signifies improved performance in predicting CTR.

For a prediction model with random parameters, the model's AUC value is approximately 0.5. Therefore, the worst-case scenario for AUC is not 0, as with other metrics, but rather 0.5. This is why its relative gain is somewhat different. The definition of relative improvement of AUC is as follows [43, 50]:

$$RelaImpr. = \left( \frac{\text{AUC} - 0.5}{\text{AUC}_{base} - 0.5} - 1 \right) \times 100\%, \quad (13)$$

*4.1.3 Comparision Baselines.* In order to thoroughly assess the effectiveness of our DPN proposal, we conduct a comprehensive comparison with eleven state-of-the-art CTR prediction baseline models, including **DNN** [10], **NCF** [19], **Wide&Deep** [8] , **PNN** [31], **RUM** [7], **ONN** [44], **DIN** [50], **GRU4Rec** [20], **DIEN** [49], **MIMN** [29], **SIM** [30], **ETA** [5], **CAN** [2].

*4.1.4 Implement Details.* All the baseline techniques and our proposed DPN are implemented using the TensorFlow framework [1]. The source code for DPN will be available soon. To ensure a fair

---

[1]https://tianchi.aliyun.com/dataset/dataDetail?dataId=42
[2]https://tianchi.aliyun.com/dataset/dataDetail?dataId=649
[3]https://tianchi.aliyun.com/dataset/dataDetail?dataId=53

comparison, we establish the following parameter settings: a maximum sequence length of $N = 100$ and an embedding dimensionality of $d = 16$ for items and categories. The number of items retrieved by the retrieval-enhanced methods (SIM, ETA) is aligned with the number of items contained in utilized patterns in DPN, i.e., $K \times s$. We employ the Adam optimizer [23] with a learning rate of 0.001 for training all methods. The batch size is set to 256 for both the training dataset and the testing dataset.

During the pretraining process of SPRM, we aligned the lengths of the behavior patterns before and after augmentation, denoted as $l$ and $s$ respectively, with the settings of the CTR prediction task. We employed the Adam optimizer with a learning rate of 0.001 for the pretraining phase. The batch size was set to 8196, and the entire pretraining process comprised 10 epochs.

Regarding the hyperparameters for DPN, we have configured the number of retrieved behavior patterns in TPRM, $K$, to be 5, and the length of the retrieved behavior patterns, $l$, to be 5. The length of refined behavior patterns in SPRM, $s$, is set to 3.

### 4.2 Overall Performances

**Table 2: Overall performance of various methods on three publicly available recommendation datasets, the most effective method is highlighted in bold. An asterisk (\*) is used to indicate the statistical significance (with p < 0.05) when comparing DPN to the top-performing baseline results.**

| Method | AUC | | |
|--------|-------|--------|--------|
| | Tmall | Taobao | Alipay |
| DNN(Sum-Pooling) | 0.9268 | 0.8731 | 0.9037 |
| NCF | 0.9256 | 0.8762 | 0.9046 |
| Wide&Deep | 0.9362 | 0.8751 | 0.9040 |
| PNN | 0.9408 | 0.8839 | 0.9058 |
| RUM | 0.9386 | 0.9046 | 0.9232 |
| GRU4REC | 0.9392 | 0.9056 | 0.9191 |
| DIN | 0.9307 | 0.8931 | 0.9185 |
| DIEN | 0.9459 | 0.9058 | 0.9251 |
| ONN | 0.9459 | 0.9062 | 0.9122 |
| MIMN | 0.9458 | 0.9162 | 0.9280 |
| SIM | 0.9494 | 0.9281 | 0.9205 |
| ETA | 0.9466 | 0.9259 | 0.9202 |
| CAN | 0.9504 | 0.9327 | 0.9311 |
| DPN (Ours) | **0.9573**\* | **0.9431**\* | **0.9438**\* |

The comprehensive performance evaluation of our novel DPN is presented in Table 2. Based on the findings displayed in the table, we can draw the following conclusions:

- DPN outperforms other baselines on all three public real-world datasets, showcasing its remarkable effectiveness and superiority. In contrast to the leading baseline approach, CAN, DPN exhibited relative performance improvements of 1.53%, 2.40%, and 2.95% on these datasets according to Equation 13. Additionally, when compared to the original DIN backbone, DPN displayed even more substantial relative improvements of 6.18%, 12.72%, and 6.05% according to Equation 13. To delve further into DPN's compatibility with mainstream CTR prediction models, we present

compatibility analysis experiments in Section 4.4. The superior performance of DPN is attributed to the efficient extraction and sufficient utilization of user behavior patterns, and the follow-up in-depth analysis further demonstrates the significance of user behavior pattern information.

- Enhanced modeling of user behavior associations can lead to improved model performance. DIEN and MIMN introduce improvements by modeling the dynamic evolution of user behavior sequences through GRU and Neural Turing Machines. They captured positional relationships between behavior sequences, resulting in enhanced performance. CAN independently model the co-occurrence relationships between the target item and historical interactions on the basis of DIEN, achieving a significant performance boost. These methods have all achieved performance gains through more comprehensive modeling of behavior associations. However, their improvements have remained focused on aggregating behavior representations and modeling co-occurrence relationships between behaviors, while overlooking higher-order behavior pattern information and the pattern-level dependencies. In contrast, DPN efficiently retrieves and fully utilizes target-related user behavior patterns, leading to a remarkable performance improvement

## 4.3 Ablation Study

As outlined in Section 3, the DPN model comprises three primary components, specifically the TPRM, SPRM, and TPA. To assess the effectiveness of these three elements, we conduct an ablation analysis by eliminating each component from DPN, yielding the subsequent three variants:

- **DPN without TPRM:** Removing TPRM from DPN means no longer retrieving target-related user behavior patterns, i.e., every three sequential user interaction records are considered as behavior patterns to be utilized.
- **DPN without SPRM:** The SPRM in DPN is removed, i.e., the retrieved user behavior patterns are not denoised and are directly utilized.
- **DPN without TPA:** Remove the TPA in DPN, i.e., instead of dependency modeling between the target behavior pattern and the refined historical behavior patterns via attention mechanism, the historical patterns and the target pattern are directly sumpooled and taken as the feature.

**Table 3: Ablation analysis results on three public real-world datasets for three key components in DPN.**

| Method | AUC | | |
| --- | --- | --- | --- |
| | Tmall | Taobao | Alipay |
| Base Model | 0.9307 | 0.8931 | 0.9185 |
| DPN w/o TPRM | 0.9507 | 0.9308 | 0.9371 |
| DPN w/o SPRM | 0.9409 | 0.9320 | 0.9364 |
| DPN w/o TPA | 0.9456 | 0.9312 | 0.9347 |
| DPN | 0.9573 | 0.9431 | 0.9438 |

The comparison results on three public real-world datasets are shown in Table 3. We can draw the following observations from the experiment results:

- Evidently, the removal of any individual component within our proposed DPN unequivocally leads to a reduction in performance, thereby substantiating the indispensability and excellence of each constituent element. The retrieval, refinement, and dependency modeling of user behavior patterns are all crucial processes for effectively utilizing behavior patterns.
- The roles of different modules vary across distinct datasets. The Taobao dataset comprises over a hundred million interaction records, harboring a rich diversity of user interests. Consequently, it is crucial to effectively retrieve target-relevant user behavior patterns within the Taobao dataset, as behavior patterns reflective of unrelated user interests introduce noise into the click-through rate (CTR) prediction task. Conversely, for the Tmall dataset, the refinement of behavior patterns holds paramount importance, possibly due to the greater variety in which these behavior patterns manifest within the interaction records.

## 4.4 Compatibility Analysis

DPN is not just a particular model; it's a general framework that can seamlessly work with several well-known CTR prediction models. In order to gauge the compatibility of the DPN framework, we utilize prominent CTR models as the base model within DPN. We then assess the performance of these resultant models in comparison to the original backbone model. The selected mainstream CTR models for this analysis encompass DNN, DIN, and DIEN, all of which serve as widely employed backbones in various related studies. The findings from our compatibility analysis experiments are presented in Table 4.

**Table 4: Compatibility analysis results on three public datasets with various mainstream CTR prediction methods as the Base Model in DPN.**

| Method | AUC | | |
| --- | --- | --- | --- |
| | Tmall | Taobao | Alipay |
| DNN | 0.9268 | 0.8731 | 0.9037 |
| DPN(DNN) | 0.9542 | 0.9309 | 0.9326 |
| DIN | 0.9307 | 0.8931 | 0.9185 |
| DPN(DIN) | 0.9573 | 0.9431 | 0.9438 |
| DIEN | 0.9459 | 0.9058 | 0.9251 |
| DPN(DIEN) | 0.9601 | 0.9430 | 0.9387 |

According to the experimental results presented in Table 4, DPN consistently demonstrates significant performance improvements across all backbone models. The enhancement observed in performance can be attributed to the efficient utilization of rich user behavior pattern information by DPN. The Base Model aggregates user interest representations from historical user interaction records, undertaking the modeling of co-occurrence relationships between the target item and historical interaction records. Building upon this foundation, DPN models user behavior patterns and their interrelationships, thereby achieving significant performance gains. The aforementioned results underscore the remarkable superiority and broad applicability of DPN, while also providing evidence for the crucial research significance of user behavior pattern mining in CTR prediction tasks.

## 4.5 Hyperparameter Analysis

The key hyperparameters in DPN include: 1) the number of target-related user behavior patterns retrieved in TPRM, i.e., $K$; 2) the length of refined behavior patterns in SPRM, i.e., $s$. To investigate the impact of these hyperparameters on DPN, we perform a comparison analysis with various settings of the key hyperparameter on public recommendation datasets. When exploring the interested hyperparameters $K$ and $s$, we keep all other hyperparameters constant.
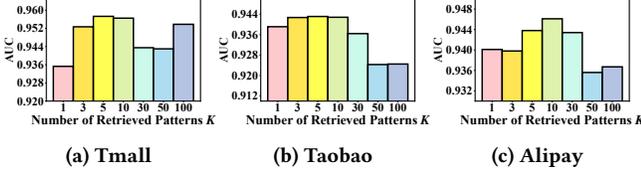


**(a) Tmall**  **(b) Taobao**  **(c) Alipay**

**Figure 5: Performances with different numbers of retrieved patterns $K$.**
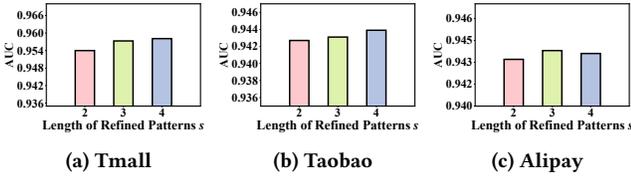


**(a) Tmall**  **(b) Taobao**  **(c) Alipay**

**Figure 6: Performances with different lengths of refined patterns $s$.**

The comparison results with different hyperparameter settings are shown in Figure 5 and 6. We can draw the following conclusions:

- **Number of retrieved patterns $K$.** In the case of DPN, optimal performance is achieved when the value of $K$ is moderate across all datasets. This is attributed to the fact that when $K$ is too small, the retrieved user behavior patterns are too limited to adequately reflect the habitual paradigm of users regarding the target item. Conversely, when $K$ is excessively large, the retrieved user behavior patterns become contaminated with numerous target-irrelevant noise patterns, consequently leading to a degradation in the predictive performance of the model.

- **Length of refined patterns $s$.** Generally, higher-order behavior patterns lead to greater performance improvements. This phenomenon can be attributed to two key factors. On one hand, higher-order behavior patterns to some extent reflect information from lower-order patterns. On the other hand, higher-order patterns possess unique pattern characteristics, enabling them to achieve superior performance compared to lower-order patterns. However, higher-order patterns do not always yield better results. For instance, when comparing a 4-order DPN to a 3-order DPN on the Alipay dataset, there was a slight decrease in performance. This could be attributed to the dataset's lack of high-order behavior patterns. Therefore, an increase in the length of behavior patterns did not introduce valuable pattern characteristics and potentially introduced unwanted noise.

## 4.6 Case Study

To validate that DPN operates as designed, we employ a case study approach to elucidate the internal structure of DPN. Figure 7 illustrates the process of handling a real-world instance in DPN, where TPRM retrieves target-relevant behavior patterns, SPRM further refines the retrieved patterns and TPA aggregates the refined behavior pattern representations while capturing the dependency among patterns.
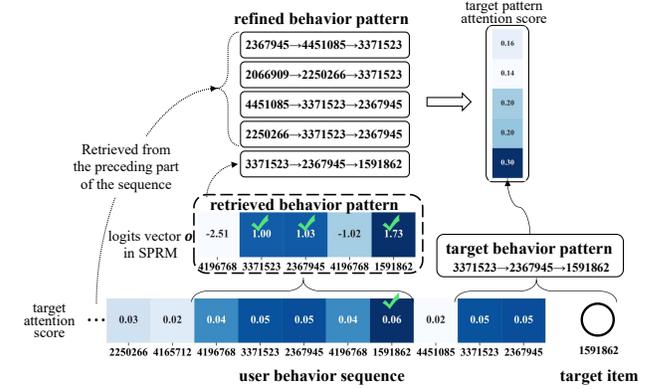


**Figure 7: Illustration of Case Study on DPN. This figure visualizes the process of DPN handling a real-world user behavior sequence. For simplification, only the last 10 behavior records are shown and the preceding part is omitted. The numbers under the heatmap are the item IDs.**
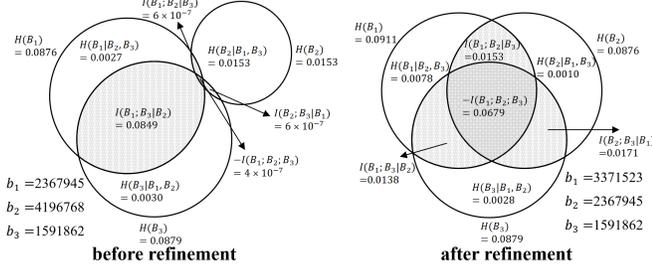
In this case, the target behavior pattern is *3371523 → 2367945 → 1591862*, which is exactly the most high-frequency behavior pattern of Taobao shown in Figure 2. Firstly, TPRM retrieves the Top-5 target-related behavior patterns, and the last retrieved one is *(4196768, 3371523, 2367945, 4196768, 1591862)*. Subsequently, the SPRM further refines this pattern, presenting the refined behavior pattern *3371523 → 2367945 → 1591862*, and it is identical to the target pattern. Finally, the TPA models the dependency between the target pattern and refined patterns via the attention mechanism.

From the insights in the case study, DPN seems to handle behavior patterns rightly. To further analyze the effectiveness of DPN, we design experiments to discuss the ability of DPNs to refine behavior patterns and model dependencies between behavior patterns by means of a more global visualization in the following Subsections 4.7 and 4.8. We did not delve extensively into pattern retrieval due to space constraints, as it primarily aims to avoid the computational complexity introduced by behavior pattern modeling, which is not the core focus of our work. The core of our work lies in learning pattern information and modeling dependencies between patterns.

## 4.7 Validation of Pattern Refinement

The actions of the user are driven by diverse interest. So the behavior patterns driven by different interest may be mutually permeated, and not be perfectly continuous segments. Moreover, misclick is also not a rare event, which introduces noise into behavior patterns. To capture more meaningful patterns and mitigate the risk of introducing noise, we design SPRM to extract $s$-length behavior

patterns from the retrieved behavior patterns. Meaningful patterns refer to which have stronger intra-dependencies, which typically can be measured by the concept of mutual information (MI). To validate the effectiveness of refinement, we define Intra-Pattern Mutual Information (IntraPMI) to evaluate the intra-dependencies within the behavior pattern, referencing [42] which extends MI to the multi-variables case. IntraPMI measures all the possible components of correlations among the behaviors within the behavior pattern, so a pattern with higher IntraPMI is a more meaningful pattern.



**Figure 8: Venn diagram of information measurement among items in behavior patterns before/after refinement (using the same case in Figure 7 as an example). The dotted area corresponds to the IntraPMI, reflecting the dependencies among the items. It is obvious that the intra-pattern dependency is more closely knit after refinement.**

*Definition 4.1.* **Intra-Pattern Mutual Information.** The intra-pattern mutual information of a given behavior pattern $p = \{b_1, b_2, b_3\}$ is defined as:
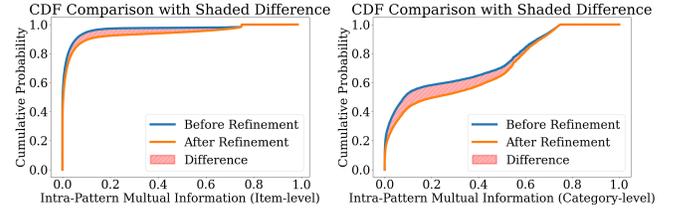
$$IntraPMI(p) = H(B_1) + H(B_2) + H(B_3) - H(B1; B2; B3),$$
$$B_i = \{0, 1 | \mathbb{I}(b_i \text{ interacted by the user})\} \quad (14)$$

where $B_i$ is the binary variable indicating whether the user interacts with $b_i$, $\mathbb{I}(\cdot)$ is the indicator function that outputs 1 when the condition is valid and 0 vice versa. $H(B_i)$ means the information entropy of variable $B_i$, and $H(B1; B2; B3)$ is the joint entropy of variables $B_1, B_2, B_3$.

To illustrate what IntraPMI means, a Venn diagram of information measurement among items within behavior patterns before/after refinement is shown in Figure 8. IntraPMI corresponds to the dotted area (the overlap of circles) in Figure 8, which measures the amount of all the correlations among the behaviors within the pattern. According to Figure 8, the IntraPMI of the pattern after refinement is significantly larger, i.e., the circles corresponding to items are more closely knit, showing more stronger intra-pattern dependencies.

To conduct a more general analysis instead of a case study, we visualize the Cumulative Distribution Function (CDF) of the IntraPMI of patterns before/after refinement on the Taobao dataset as shown in Figure 9. Notably, we normalize the metric to range [0, 1] following [36]. The CDF curves corresponding to the refined patterns are located below the ones before refinement, indicating that our proposed SPRM can make the mutual information of the patterns larger, i.e. more meaningful.



**Figure 9: Cumulative Distribution Function (CDF) of the Intra-Pattern Mutual Information of patterns before/after refinement on both item and category level. The difference in CDFs is indicated by the red diagonal texture.**

## 4.8 What does DPN capture?

To investigate which dependencies among patterns are valuable for the CTR prediction task and what models can capture, we use the concept of conditional mutual information which is widely used in feature selection [14, 34] to reflect the correlation between features and labels for visualization.

*4.8.1 Metrics.* For a given historical behavior pattern $p_i$, whether $p_i$ occurs in the user behavior sequence can be described as a binary variable $\mathcal{P}_H^{(p_i)} = \{0, 1 | \mathbb{I}(p_i \text{ occurs in the user behavior records})\}$, where $\mathbb{I}(\cdot)$ denotes the indicator function to discriminate the condition is true or false. Similarly, $\mathcal{P}_T^{(p_t)} = \{0, 1 | \mathbb{I}(p_t \text{ is the target pattern})\}$ indicates whether the behavior pattern $p_t$ is the target pattern.
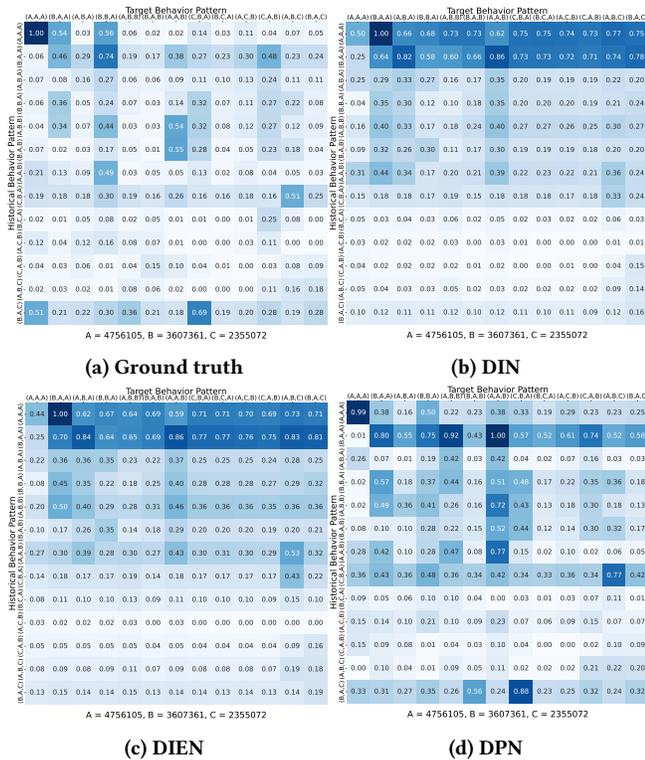
The conditional mutual information $I(\mathcal{P}_H^{(p_i)}; \mathcal{Y} | \mathcal{P}_T^{(p_t)})$ measures the dependency relationship between $\mathcal{P}_H^{(p_i)}$ and the labels $\mathcal{Y}$ given $\mathcal{P}_T^{(p_t)}$, reflecting the value of historical pattern $p_i$ for CTR prediction under a given target pattern $p_t$. It is the ground truth that demonstrates which dependencies among patterns are valuable for the CTR prediction.

The conditional mutual information $I(\mathcal{P}_H^{(p_i)}; \hat{\mathcal{Y}} | \mathcal{P}_T^{(p_t)})$ measures the dependency relationship between $\mathcal{P}_H^{(p_i)}$ and the click probability $\hat{\mathcal{Y}}$ predicted by the model given $\mathcal{P}_T^{(p_t)}$. $I(\mathcal{P}_H^{(p_i)}; \hat{\mathcal{Y}} | \mathcal{P}_T^{(p_t)})$ on test dataset reflects the effect of pattern-level dependencies on CTR prediction learned by the model.

*4.8.2 Analysis.* We select the Top-3 frequent categories 4756105, 3607361, and 2355072 (denoted as A, B, C for simplification) in the test dataset of Taobao dataset for analysis. There are many possible patterns that can be formed by A, B, and C. We visualize $I(\mathcal{P}_H^{(p_i)}; \mathcal{Y} | \mathcal{P}_T^{(p_t)})$ for some of them as shown in the Figure 10 (a). It is obvious that some dependency between the historical patterns and the target pattern is meaningful (with a higher value in the matrix) to CTR prediction tasks, such as (A, A, A) & (A, A, A), (B, A, A) & (B, B, A) and (B, A, C) & (C, B, A).

Furthermore, we visualize $I(\mathcal{P}_H^{(p_i)}; \hat{\mathcal{Y}} | \mathcal{P}_T^{(p_t)})$ for the DPN and the most respective behavior-based CTR prediction models, i.e., DIN and DIEN on the test dataset of Taobao dataset as shown in Figure 10 (b-d). According to the visualization results, we can draw the following conclusions:

- DPN can capture the valuable pattern-level dependency well. (a) and (d) in Figure 10 exhibit a noticeable correlation. On the other

**Figure 10: Visualization of the conditional mutual information.**

hand, DIN and DIEN are completely unable to do so, as (a) is significantly different from (b)/(c).

- DIN and DIEN are not aware of the target pattern. For example, regarding target patterns with the target category A in Figure 10 (ending with A, such as (B, A, A), (A, B, A), (B, B, A), (C, B, A), (B, C, A)), their corresponding columns in the conditional mutual information matrix exhibit strong cosine similarity (>0.9) pairwise.

- DIN and DIEN are not sensitive to the historical patterns. For the patterns consisting of the same items, such as (B, A, A), (A, B, A), (A, A, B), their corresponding rows in the matrix are obviously similar. DIN can not handle the sequential dependency, so it can not capture the pattern information. DIEN handles the user behavior sequence via a unified sequential model, modeling the global evolution of user interest while neglecting the local behavior patterns.

In summary, DPN not only learns the pattern information but also captures the valuable dependencies among behavior patterns, thus achieving significant performance gain for the CTR prediction.

## 5 CONCLUSIONS

In this article, we highlight the significant importance of diverse user behavior patterns hidden within massive historical interaction records and the neglect of this idea by existing CTR prediction methods. To fully leverage user behavior pattern information, we propose the Deep Pattern Network (DPN) for click-through rate

prediction. DPN not only learns the information of behavior patterns but also captures the valuable dependencies among behavior patterns. Comprehensive experiments conducted on three different datasets thoroughly demonstrate the outstanding superiority and broad compatibility of our proposed DPN. Further analysis elucidates how the rich behavior patterns enhance the performance of DPN.

## REFERENCES

[1] Martín Abadi, Paul Barham, Jianmin Chen, Zhifeng Chen, Andy Davis, Jeffrey Dean, Matthieu Devin, Sanjay Ghemawat, Geoffrey Irving, Michael Isard, Manjunath Kudlur, Josh Levenberg, Rajat Monga, Sherry Moore, Derek Gordon Murray, Benoit Steiner, Paul A. Tucker, Vijay Vasudevan, Pete Warden, Martin Wicke, Yuan Yu, and Xiaoqiang Zheng. 2016. TensorFlow: A System for Large-Scale Machine Learning. In *12th USENIX Symposium on Operating Systems Design and Implementation, OSDI 2016, Savannah, GA, USA, November 2-4, 2016*, Kimberly Keeton and Timothy Roscoe (Eds.). USENIX Association, 265–283. https://www.usenix.org/conference/osdi16/technical-sessions/presentation/abadi

[2] Weijie Bian, Kailun Wu, Lejian Ren, Qi Pi, Yujing Zhang, Can Xiao, Xiang-Rong Sheng, Yong-Nan Zhu, Zhangming Chan, Na Mou, Xinchen Luo, Shiming Xiang, Guorui Zhou, Xiaoqiang Zhu, and Hongbo Deng. 2022. CAN: Feature Co-Action Network for Click-Through Rate Prediction. In *WSDM '22: The Fifteenth ACM International Conference on Web Search and Data Mining, Virtual Event / Tempe, AZ, USA, February 21 - 25, 2022*, K. Selcuk Candan, Huan Liu, Leman Akoglu, Xin Luna Dong, and Jiliang Tang (Eds.). ACM, 57–65. https://doi.org/10.1145/3488560.3498435

[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, Atlanta, GA, USA, October 17-21, 2022*, Mohammad Al Hasan and Li Xiong (Eds.). ACM, 2974–2983. https://doi.org/10.1145/3511808.3557082

[4] Yin-Wen Chang, Cho-Jui Hsieh, Kai-Wei Chang, Michael Ringgaard, and Chih-Jen Lin. 2010. Training and Testing Low-degree Polynomial Data Mappings via Linear SVM. *J. Mach. Learn. Res.* 11 (2010), 1471–1490. https://doi.org/10.5555/1756006.1859899

[5] Qiwei Chen, Changhua Pei, Shanshan Lv, Chao Li, Junfeng Ge, and Wenwu Ou. 2021. End-to-End User Behavior Retrieval in Click-Through RatePrediction Model. *CoRR* abs/2108.04468 (2021). arXiv:2108.04468 https://arxiv.org/abs/2108.04468

[6] Qiwei Chen, Huan Zhao, Wei Li, Pipei Huang, and Wenwu Ou. 2019. Behavior Sequence Transformer for E-commerce Recommendation in Alibaba. *CoRR* abs/1905.06874 (2019). arXiv:1905.06874 http://arxiv.org/abs/1905.06874

[7] 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, WSDM 2018, Marina Del Rey, CA, USA, February 5-9, 2018*, Yi Chang, Chengxiang Zhai, Yan Liu, and Yoelle Maarek (Eds.). ACM, 108–116. https://doi.org/10.1145/3159652.3159668

[8] Heng-Tze Cheng, Levent Koc, Jeremiah Harmsen, Tal Shaked, Tushar Chandra, Hrishi Aradhye, Glen Anderson, Greg Corrado, Wei Chai, Mustafa Ispir, Rohan Anil, Zakaria Haque, Lichan Hong, Vihan Jain, Xiaobing Liu, and Hemal Shah. 2016. Wide & Deep Learning for Recommender Systems. In *Proceedings of the 1st Workshop on Deep Learning for Recommender Systems, DLRS@RecSys 2016, Boston, MA, USA, September 15, 2016*, Alexandros Karatzoglou, Balázs Hidasi, Domonkos Tikk, Oren Sar Shalom, Haggai Roitman, Bracha Shapira, and Lior Rokach (Eds.). ACM, 7–10. https://doi.org/10.1145/2988450.2988454

[9] Junyoung Chung, Çaglar Gülçehre, KyungHyun Cho, and Yoshua Bengio. 2014. Empirical Evaluation of Gated Recurrent Neural Networks on Sequence Modeling. *CoRR* abs/1412.3555 (2014). arXiv:1412.3555 http://arxiv.org/abs/1412.3555

[10] Paul Covington, Jay Adams, and Emre Sargin. 2016. Deep Neural Networks for YouTube Recommendations. In *Proceedings of the 10th ACM Conference on Recommender Systems, Boston, MA, USA, September 15-19, 2016*, Shilad Sen, Werner Geyer, Jill Freyne, and Pablo Castells (Eds.). ACM, 191–198. https://doi.org/10.1145/2959100.2959190

[11] Shizhe Diao, Jiaxin Bai, Yan Song, Tong Zhang, and Yonggang Wang. 2020. ZEN: Pre-training Chinese Text Encoder Enhanced by N-gram Representations. In *Findings of the Association for Computational Linguistics: EMNLP 2020, Online Event, 16-20 November 2020 (Findings of ACL, Vol. EMNLP 2020)*. Association for Computational Linguistics, 4729–4740. https://doi.org/10.18653/V1/2020.FINDINGS-EMNLP.425

[12] Xiaohan Ding, Xiangyu Zhang, Jungong Han, and Guiguang Ding. 2022. Scaling Up Your Kernels to 31×31: Revisiting Large Kernel Design in CNNs. In *IEEE/CVF Conference on Computer Vision and Pattern Recognition, CVPR 2022, New Orleans, LA, USA, June 18-24, 2022*. IEEE, 11953–11965. https://doi.org/10.1109/CVPR52688.

2022.01166

[13] Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, Jakob Uszkoreit, and Neil Houlsby. 2021. An Image is Worth 16x16 Words: Transformers for Image Recognition at Scale. In *9th International Conference on Learning Representations, ICLR 2021, Virtual Event, Austria, May 3-7, 2021*. OpenReview.net. https://openreview.net/forum?id=YicbFdNTTy

[14] François Fleuret. 2004. Fast Binary Feature Selection with Conditional Mutual Information. *J. Mach. Learn. Res.* 5 (2004), 1531–1555. http://jmlr.org/papers/volume5/fleuret04a/fleuret04a.pdf

[15] Alex Graves, Greg Wayne, and Ivo Danihelka. 2014. Neural Turing Machines. *CoRR* abs/1410.5401 (2014). arXiv:1410.5401 http://arxiv.org/abs/1410.5401

[16] Daniel W Greeno, Montrose S Sommers, and Jerome B Kernan. 1973. Personality and implicit behavior patterns. *Journal of Marketing Research* 10, 1 (1973), 63–69.

[17] 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 Twenty-Sixth International Joint Conference on Artificial Intelligence, IJCAI 2017, Melbourne, Australia, August 19-25, 2017*, Carles Sierra (Ed.). ijcai.org, 1725–1731. https://doi.org/10.24963/ijcai.2017/239

[18] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. 2016. Deep Residual Learning for Image Recognition. In *2016 IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2016, Las Vegas, NV, USA, June 27-30, 2016*. IEEE Computer Society, 770–778. https://doi.org/10.1109/CVPR.2016.90

[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, WWW 2017, Perth, Australia, April 3-7, 2017*, Rick Barrett, Rick Cummings, Eugene Agichtein, and Evgeniy Gabrilovich (Eds.). ACM, 173–182. https://doi.org/10.1145/3038912.3052569

[20] Balázs Hidasi, Alexandros Karatzoglou, Linas Baltrunas, and Domonkos Tikk. 2016. Session-based Recommendations with Recurrent Neural Networks. In *4th International Conference on Learning Representations, ICLR 2016, San Juan, Puerto Rico, May 2-4, 2016, Conference Track Proceedings*, Yoshua Bengio and Yann LeCun (Eds.). http://arxiv.org/abs/1511.06939

[21] Wang-Cheng Kang and Julian J. McAuley. 2018. Self-Attentive Sequential Recommendation. In *IEEE International Conference on Data Mining, ICDM 2018, Singapore, November 17-20, 2018*. IEEE Computer Society, 197–206. https://doi.org/10.1109/ICDM.2018.00035

[22] Otto F Kernberg. 2016. What is personality? *Journal of personality disorders* 30, 2 (2016), 145–156.

[23] Diederik P. Kingma and Jimmy Ba. 2015. Adam: A Method for Stochastic Optimization. In *3rd International Conference on Learning Representations, ICLR 2015, San Diego, CA, USA, May 7-9, 2015, Conference Track Proceedings*, Yoshua Bengio and Yann LeCun (Eds.). http://arxiv.org/abs/1412.6980

[24] Jianxun Lian, Xiaohuan Zhou, Fuzheng Zhang, Zhongxia Chen, Xing Xie, and Guangzhong Sun. 2018. xDeepFM: Combining Explicit and Implicit Feature Interactions for Recommender Systems. In *Proceedings of the 24th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining, KDD 2018, London, UK, August 19-23, 2018*, Yike Guo and Faisal Farooq (Eds.). ACM, 1754–1763. https://doi.org/10.1145/3219819.3220023

[25] Weiwen Liu, Wei Guo, Yong Liu, Ruiming Tang, and Hao Wang. 2023. User Behavior Modeling with Deep Learning for Recommendation: Recent Advances. In *Proceedings of the 17th ACM Conference on Recommender Systems, RecSys 2023, Singapore, Singapore, September 18-22, 2023*, Jie Zhang, Li Chen, Shlomo Berkovsky, Min Zhang, Tommaso Di Noia, Justin Basilico, Luiz Pizzato, and Yang Song (Eds.). ACM, 1286–1287. https://doi.org/10.1145/3604915.3609496

[26] H. Brendan McMahan, Gary Holt, David Sculley, Michael Young, Dietmar Ebner, Julian Grady, Lan Nie, Todd Phillips, Eugene Davydov, Daniel Golovin, Sharat Chikkerur, Dan Liu, Martin Wattenberg, Arnar Mar Hrafnkelsson, Tom Boulos, and Jeremy Kubica. 2013. Ad click prediction: a view from the trenches. In *The 19th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining, KDD 2013, Chicago, IL, USA, August 11-14, 2013*, Inderjit S. Dhillon, Yehuda Koren, Rayid Ghani, Ted E. Senator, Paul Bradley, Rajesh Parekh, Jingrui He, Robert L. Grossman, and Ramasamy Uthurusamy (Eds.). ACM, 1222–1230. https://doi.org/10.1145/2487575.2488200

[27] Chang Meng, Hengyu Zhang, Wei Guo, Huifeng Guo, Haotian Liu, Yingxue Zhang, Hongkun Zheng, Ruiming Tang, Xiu Li, and Rui Zhang. 2023. Hierarchical Projection Enhanced Multi-behavior Recommendation. In *Proceedings of the 29th ACM SIGKDD Conference on Knowledge Discovery and Data Mining, KDD 2023, Long Beach, CA, USA, August 6-10, 2023*. ACM, 4649–4660. https://doi.org/10.1145/3580305.3599838

[28] Erxue Min, Yu Rong, Tingyang Xu, Yatao Bian, Da Luo, Kangyi Lin, Junzhou Huang, Sophia Ananiadou, and Peilin Zhao. 2022. Neighbour Interaction based Click-Through Rate Prediction via Graph-masked Transformer. In *SIGIR '22: The 45th International ACM SIGIR Conference on Research and Development in Information Retrieval, Madrid, Spain, July 11 - 15, 2022*. ACM, 353–362. https://doi.org/10.1145/3477495.3532031

[29] Qi Pi, Weijie Bian, Guorui Zhou, Xiaoqiang Zhu, and Kun Gai. 2019. Practice on Long Sequential User Behavior Modeling for Click-Through Rate Prediction.

In *Proceedings of the 25th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining, KDD 2019, Anchorage, AK, USA, August 4-8, 2019*, Ankur Teredesai, Vipin Kumar, Ying Li, Rómer Rosales, Evimaria Terzi, and George Karypis (Eds.). ACM, 2671–2679. https://doi.org/10.1145/3292500.3330666

[30] 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 *CIKM '20: The 29th ACM International Conference on Information and Knowledge Management, Virtual Event, Ireland, October 19-23, 2020*, Mathieu d'Aquin, Stefan Dietze, Claudia Hauff, Edward Curry, and Philippe Cudré-Mauroux (Eds.). ACM, 2685–2692. https://doi.org/10.1145/3340531.3412744

[31] Yanru Qu, Han Cai, Kan Ren, Weinan Zhang, Yong Yu, Ying Wen, and Jun Wang. 2016. Product-Based Neural Networks for User Response Prediction. In *IEEE 16th International Conference on Data Mining, ICDM 2016, December 12-15, 2016, Barcelona, Spain*, Francesco Bonchi, Josep Domingo-Ferrer, Ricardo Baeza-Yates, Zhi-Hua Zhou, and Xindong Wu (Eds.). IEEE Computer Society, 1149–1154. https://doi.org/10.1109/ICDM.2016.0151

[32] Steffen Rendle. 2010. Factorization Machines. In *ICDM 2010, The 10th IEEE International Conference on Data Mining, Sydney, Australia, 14-17 December 2010*, Geoffrey I. Webb, Bing Liu, Chengqi Zhang, Dimitrios Gunopulos, and Xindong Wu (Eds.). IEEE Computer Society, 995–1000. https://doi.org/10.1109/ICDM.2010.127

[33] Aurko Roy, Rohan Anil, Guangda Lai, Benjamin Lee, Jeffrey Zhao, Shuyuan Zhang, Shibo Wang, Ye Zhang, Shen Wu, Rigel Swavely, Tao Yu, Phuong Dao, Christopher Fifty, Zhifeng Chen, and Yonghui Wu. 2022. N-Grammer: Augmenting Transformers with latent n-grams. *CoRR* abs/2207.06366 (2022). https://doi.org/10.48550/ARXIV.2207.06366 arXiv:2207.06366

[34] Alexander Shishkin, Anastasia A. Bezzubtseva, Alexey Drutsa, Ilia Shishkov, Ekaterina Gladkikh, Gleb Gusev, and Pavel Serdyukov. 2016. Efficient High-Order Interaction-Aware Feature Selection Based on Conditional Mutual Information. In *Advances in Neural Information Processing Systems 29: Annual Conference on Neural Information Processing Systems 2016, December 5-10, 2016, Barcelona, Spain*, Daniel D. Lee, Masashi Sugiyama, Ulrike von Luxburg, Isabelle Guyon, and Roman Garnett (Eds.). 4637–4645. https://proceedings.neurips.cc/paper/2016/hash/d5e2fbef30a4eb668a203060ec8e5eef-Abstract.html

[35] Karen Simonyan and Andrew Zisserman. 2015. Very Deep Convolutional Networks for Large-Scale Image Recognition. In *3rd International Conference on Learning Representations, ICLR 2015, San Diego, CA, USA, May 7-9, 2015, Conference Track Proceedings*, Yoshua Bengio and Yann LeCun (Eds.). http://arxiv.org/abs/1409.1556

[36] Gustavo Sosa-Cabrera, Miguel García-Torres, Santiago Gómez-Guerrero, Christian E. Schaerer, and Federico Divina. 2019. A multivariate approach to the symmetrical uncertainty measure: Application to feature selection problem. *Information Sciences* 494 (2019), 1–20. https://doi.org/10.1016/j.ins.2019.04.046

[37] Fei Sun, Jun Liu, Jian Wu, Changhua Pei, Xiao Lin, Wenwu Ou, and Peng Jiang. 2019. BERT4Rec: Sequential Recommendation with Bidirectional Encoder Representations from Transformer. In *Proceedings of the 28th ACM International Conference on Information and Knowledge Management, CIKM 2019, Beijing, China, November 3-7, 2019*, Wenwu Zhu, Dacheng Tao, Xueqi Cheng, Peng Cui, Elke A. Rundensteiner, David Carmel, Qi He, and Jeffrey Xu Yu (Eds.). ACM, 1441–1450. https://doi.org/10.1145/3357384.3357895

[38] Ilya Sutskever, Oriol Vinyals, and Quoc V. Le. 2014. Sequence to Sequence Learning with Neural Networks. In *Advances in Neural Information Processing Systems 27: Annual Conference on Neural Information Processing Systems 2014, December 8-13 2014, Montreal, Quebec, Canada*, Zoubin Ghahramani, Max Welling, Corinna Cortes, Neil D. Lawrence, and Kilian Q. Weinberger (Eds.). 3104–3112. https://proceedings.neurips.cc/paper/2014/hash/a14ac55a4f27472c5d894ec1c3c743d2-Abstract.html

[39] Jiaxi Tang and Ke Wang. 2018. Personalized Top-N Sequential Recommendation via Convolutional Sequence Embedding. In *Proceedings of the Eleventh ACM International Conference on Web Search and Data Mining, WSDM 2018, Marina Del Rey, CA, USA, February 5-9, 2018*, Yi Chang, Chengxiang Zhai, Yan Liu, and Yoelle Maarek (Eds.). ACM, 565–573. https://doi.org/10.1145/3159652.3159656

[40] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N. Gomez, Lukasz Kaiser, and Illia Polosukhin. 2017. Attention is All you Need. In *Advances in Neural Information Processing Systems 30: Annual Conference on Neural Information Processing Systems 2017, December 4-9, 2017, Long Beach, CA, USA*, Isabelle Guyon, Ulrike von Luxburg, Samy Bengio, Hanna M. Wallach, Rob Fergus, S. V. N. Vishwanathan, and Roman Garnett (Eds.). 5998–6008. https://proceedings.neurips.cc/paper/2017/hash/3f5ee243547dee91fbd053c1c4a845aa-Abstract.html

[41] Wen Wang, Wei Zhang, Shukai Liu, Qi Liu, Bo Zhang, Leyu Lin, and Hongyuan Zha. 2020. Beyond Clicks: Modeling Multi-Relational Item Graph for Session-Based Target Behavior Prediction. In *WWW '20: The Web Conference 2020, Taipei, Taiwan, April 20-24, 2020*. ACM / IW3C2, 3056–3062. https://doi.org/10.1145/3366423.3380077

[42] Satosi Watanabe. 1960. Information theoretical analysis of multivariate correlation. *IBM Journal of research and development* 4, 1 (1960), 66–82.

[43] Ling Yan, Wu-Jun Li, Gui-Rong Xue, and Dingyi Han. 2014. Coupled Group Lasso for Web-Scale CTR Prediction in Display Advertising. In *Proceedings of the 31th International Conference on Machine Learning, ICML 2014, Beijing, China, 21-26 June 2014 (JMLR Workshop and Conference Proceedings, Vol. 32)*. JMLR.org, 802–810. http://proceedings.mlr.press/v32/yan14.html

[44] Yi Yang, Baile Xu, Shaofeng Shen, Furao Shen, and Jian Zhao. 2020. Operation-aware Neural Networks for user response prediction. *Neural Networks* 121 (2020), 161–168. https://doi.org/10.1016/j.neunet.2019.09.020

[45] Jiaqi Zhai, Lucy Liao, Xing Liu, Yueming Wang, Rui Li, Xuan Cao, Leon Gao, Zhaojie Gong, Fangda Gu, Michael He, Yinghai Lu, and Yu Shi. 2024. Actions Speak Louder than Words: Trillion-Parameter Sequential Transducers for Generative Recommendations. *CoRR* abs/2402.17152 (2024). https://doi.org/10.48550/ARXIV.2402.17152 arXiv:2402.17152

[46] Hengyu Zhang, Enming Yuan, Wei Guo, Zhicheng He, Jiarui Qin, Huifeng Guo, Bo Chen, Xiu Li, and Ruiming Tang. 2022. Disentangling Past-Future Modeling in Sequential Recommendation via Dual Networks. In *Proceedings of the 31st ACM International Conference on Information & Knowledge Management, Atlanta, GA, USA, October 17-21, 2022*. ACM, 2549–2558. https://doi.org/10.1145/3511808.3557289

[47] Weinan Zhang, Jiarui Qin, Wei Guo, Ruiming Tang, and Xiuqiang He. 2021. Deep Learning for Click-Through Rate Estimation. In *Proceedings of the Thirtieth International Joint Conference on Artificial Intelligence, IJCAI 2021, Virtual Event /*

*Montreal, Canada, 19-27 August 2021*, Zhi-Hua Zhou (Ed.). ijcai.org, 4695–4703. https://doi.org/10.24963/ijcai.2021/636

[48] Zuowu Zheng, Xiaofeng Gao, Junwei Pan, Qi Luo, Guihai Chen, Dapeng Liu, and Jie Jiang. 2022. Autoattention: automatic field pair selection for attention in user behavior modeling. In *2022 IEEE International Conference on Data Mining (ICDM)*. IEEE, 803–812.

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

[50] Guorui Zhou, Xiaoqiang Zhu, Chengru 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, KDD 2018, London, UK, August 19-23, 2018*, Yike Guo and Faisal Farooq (Eds.). ACM, 1059–1068. https://doi.org/10.1145/3219819.3219823

[51] Haolin Zhou, Junwei Pan, Xinyi Zhou, Xihua Chen, Jie Jiang, Xiaofeng Gao, and Guihai Chen. 2023. Temporal Interest Network for Click-Through Rate Prediction. *CoRR* abs/2308.08487 (2023).