

CETN: Contrast-enhanced Through Network for CTR Prediction

HONGHAO LI, Anhui University, China

LEI SANG, Anhui University, China

YI ZHANG, Anhui University, China

XUYUN ZHANG, Macquarie University, Australia

YIWEN ZHANG, Anhui University, China

Click-through rate (CTR) Prediction is a crucial task in personalized information retrievals, such as industrial recommender systems, online advertising, and web search. Most existing CTR Prediction models utilize explicit feature interactions to overcome the performance bottleneck of implicit feature interactions. Hence, deep CTR models based on parallel structures (e.g., DCN, FinalMLP, xDeepFM) have been proposed to obtain joint information from different semantic spaces. However, these parallel subcomponents lack effective supervision and communication signals, making it challenging to efficiently capture valuable multi-views feature interaction information in different semantic spaces. To address this issue, we propose a simple yet effective novel CTR model: Contrast-enhanced Through Network for CTR (CETN), so as to balance the diversity and homogeneity of feature interaction information. Specifically, CETN employs product-based feature interactions and the augmentation (perturbation) concept from contrastive learning to segment different semantic spaces, each with distinct activation functions. This improves diversity in the feature interaction information captured by the model. Additionally, we introduce self-supervised signals and through connection within each semantic space to ensure the homogeneity of the captured feature interaction information. The experiments and research conducted on four real datasets demonstrate that our model consistently outperforms twenty baseline models in terms of AUC and Logloss.

CCS Concepts: • **Information systems** → **Information retrieval**; **Recommender systems**; • **Networks**;

Additional Key Words and Phrases: Contrastive Learning, Feature Interaction, Neural Network, Recommender Systems, CTR Prediction

ACM Reference Format:

Honghao Li, Lei Sang, Yi Zhang, Xuyun Zhang, and Yiwen Zhang. xxx. CETN: Contrast-enhanced Through Network for CTR Prediction. *ACM Trans. Inf. Syst.* 0, 0, Article 0 (xxx), 31 pages. <https://doi.org/XXXXXXX.XXXXXXX>

1 INTRODUCTION

Accurately predicting user responses to items (e.g., products, movies, and advertisements) under certain contexts (e.g., websites, and apps) plays a pivotal role in commercial personalized information

This work was supported by the National Natural Science Foundation of China [No.62272001, No.62206002]. The corresponding author of this paper is Yiwen Zhang.

Authors' addresses: Honghao Li, salmon1802li@gmail.com, Anhui University, 111 Jiulong Rd, Hefei, Anhui Province, China, 230601; Lei Sang, Anhui University, 111 Jiulong Rd, Hefei, Anhui Province, China, 230601, sanglei@ahu.edu.cn; Yi Zhang, Anhui University, 111 Jiulong Rd, Hefei, Anhui Province, China, 230601, zhangyi.ahu@gmail.com; Xuyun Zhang, Macquarie University, Balaclava Rd, Macquarie Park NSW, Australia, 2109, xuyun.zhang@mq.edu.au; Yiwen Zhang, Anhui University, 111 Jiulong Rd, Hefei, Anhui Province, China, 230601, zhangyiwen@ahu.edu.cn.

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.

© xxx Copyright held by the owner/author(s). Publication rights licensed to ACM.

ACM 1046-8188/xxx/0-ART0

<https://doi.org/XXXXXXX.XXXXXXX>

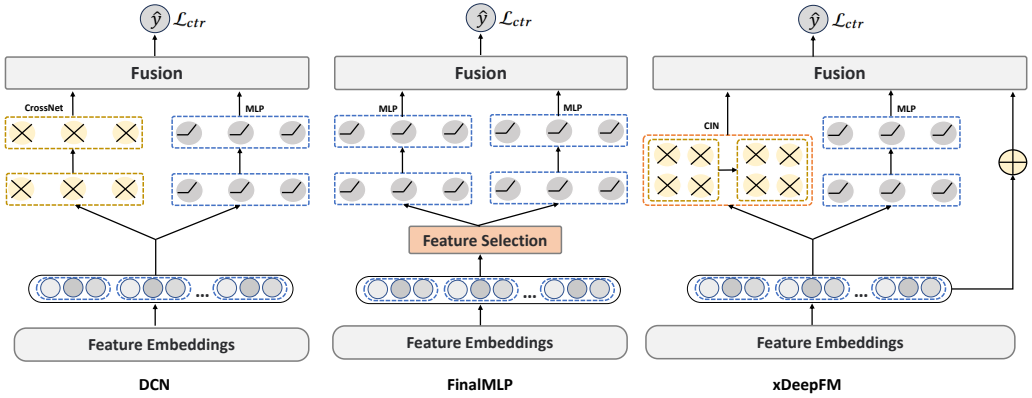


Fig. 1. The architecture of three strong baseline models with parallel structure: DCN, FinalMLP, and xDeepFM.

retrieval (IR) scenarios [5, 10, 34, 54, 68]. The user’s probability of clicking on an item serves as an intuitive evaluation metric, widely applied across various downstream scenarios, such as online advertising [8, 12, 60], recommender systems [18, 42, 56], and web search [1, 11]. Its primary objective is to assess the probability of users clicking on recommended items within a system. On the one hand, the revenue of the vast majority of commercial IR systems is closely tied to user click-through rates (CTR), making the accuracy of CTR prediction directly impact the profitability of these systems. On the other hand, user satisfaction is closely linked to the performance metrics of recommender system. A well-performing CTR prediction model aids in swiftly discerning user interests, thereby enhancing the user experience.

Capturing effective feature interactions is one of the crucial strategies to enhance the performance of CTR models. In the initial period, researchers employ explicit product-based feature interaction to build models [22, 40–42], aiming to mitigate data sparsity issues. However, due to model complexity constraints, they could only capture low-order feature interactions in practical applications [28, 36, 42, 53]. With the recent proliferation of deep learning, deep learning-based CTR models have emerged, where the multi-layer perceptron (MLP) is widely employed [14, 34, 40, 53–55]. MLP implicitly captures high-order feature interaction information. Nevertheless, some existing work has pointed out that the expressive power of a single MLP is inefficient for capturing feature interaction information, and it may even struggle to learn simple inner-product operations [43, 45]. To overcome the performance bottleneck of MLPs in capturing feature interaction information, researchers have attempted to combine explicit product-based feature interactions with the ability of MLPs to implicitly capture high-order feature interactions, leading to significant performance improvements [70]. Depending on the integration approach, this can be divided into two structures, namely the parallel structure and the stacked structure [3, 14, 34, 40]. The stacked structure attempts to first perform explicit product-based feature interactions on the features and then feed them as input to a MLP [40, 54, 55, 64]. In contrast, the parallel structure jointly trains explicit and implicit components of feature interactions in a parallel manner, utilizing a fusion layer to obtain joint information from different semantic spaces [14, 30, 34, 54, 55]. However, many existing CTR models suffer from three main issues, even if they use the two structures mentioned above, still result in a narrow semantic space available for capturing information, and are unable to efficiently capture diversiform feature interaction information in different semantic spaces, as detailed follow.

Lackluster segmentation in the semantic space. Most CTR models with parallel structures share embedding layers among their parallel subcomponents and directly use concatenated embeddings as inputs for the model [14, 30, 53, 54]. This results in the information from different semantic spaces remaining similar, and relying solely on the information capture capacity of a single subcomponent in an attempt to improve diversity in the obtained information is undoubtedly inefficient. For models that employ parallel subcomponents with the same structure [34, 52, 55], the problem of semantic space segmentation becomes even more pronounced. Without capturing feature interaction information in different semantic spaces differently, it often leads to sub-optimal performance and fails to guarantee the diversity of feature interaction information.

Feeble supervisory signals in multi-semantic space. CTR models based on parallel structures often require a fusion layer to aggregate the feature interaction information captured by various subcomponents to obtain the final prediction [30, 34, 53, 55]. However, prior to information aggregation, there is a lack of communication and effective supervisory signals between the individual subcomponents. This hinders the model from capturing high-quality, non-redundant information and diminishes the diversity of the information captured.

Excessive noise in a multi-semantic space. While we pursue diversity within semantic spaces, it inevitably leads to the challenge of "being too different". Empirically, we have observed that if different parallel subcomponents exclusively capture entirely distinct information, it results in the model aggregating a substantial amount of noise signals in the fusion layer, thereby diminishing model performance. Hence, ensuring homogeneity in the information captured by individual subcomponents is crucial.

To explain the three challenges we proposed in more detail, as illustrated in Figure 1, we present three strong baseline models as examples. In the case of DCN [53], its two subcomponents share an embedding layer. It explicitly captures low-order feature interactions using CrossNet and implicitly captures high-order feature interactions through an MLP. In the fusion layer, a simple summation is used to aggregate information from two semantic spaces. However, this straightforward information capture approach fails to address the three challenges we have identified. FinalMLP [34] begins by segmenting two semantic spaces using a feature selection layer, ensuring diversity in the information contained within these spaces. It then employs two MLPs to implicitly capture feature interaction information within each semantic space. Finally, an explicit product operation is introduced in the fusion layer. This model resolves the issue of semantic space segmentation, ensuring information diversity in these spaces and achieving outstanding performance. However, it lacks effective supervisory signals within each semantic space, making it difficult to ensure that MLP can capture diverse information in semantic space through a self-supervised manner. It also lacks differentiation between primary and auxiliary semantic spaces, failing to guarantee the homogeneity of captured information across the model. Similarly, xDeepFM [30] divides data into three semantic spaces and models feature interactions in a more complex manner to ensure diversity in the captured feature interaction information. However, it faces challenges similar to DCN, which cannot guarantee the diversity and homogeneity of the captured information.

To address the widespread issues among parallel subcomponents, we introduce a novel CTR prediction model in this paper, named **Contrast-enhanced Through Network (CETN)**. Specifically, to improve diversity in the captured feature interaction information, we employ the semantic space using the product & perturbation paradigm, then utilize multiple Key-Value Blocks with different activation functions as parallel subcomponents of the model. Furthermore, to address issues stemming from diversity, we propose a Through Network to ensure homogeneity in the information captured by individual subcomponents. Before the fusion layer, we introduce Denominator-only InfoNCE (Do-InfoNCE) and cosine loss to self-supervised learning processes within each semantic space. Extensive experiments conducted on four real-world datasets demonstrate that this simple

yet efficient model consistently outperforms twenty baseline models in terms of both AUC and Logloss.

In summary, the primary contributions of this work can be summarized as follows:

- By analyzing the existing CTR models with parallel structures, we summarize three challenges common to these models, i.e., lackluster segmentation, feeble supervisory signals, and excessive noise. To address these three challenges, we introduce the complementary principles of diversity and homogeneity, leading to the proposal of a new model, CETN.
- To enhance the diversity of information captured by the model, we further segment the semantic space, utilize Key-Value Blocks with different activation functions as subcomponents, and introduce Do-InfoNCE to provide auxiliary signals for capturing feature interactions.
- To ensure the homogeneity of the information captured by the subcomponents, we introduced element-wise Through Connection and cosine loss as communication bridges between multiple semantic spaces.
- We conduct fair and extensive experiments on four real public benchmark datasets to evaluate the performance of our proposed model. Based on the experimental results, we demonstrate that CETN consistently outperforms twenty state-of-the-art baseline models.

2 PRELIMINARIES

2.1 Problem Definition

The click-through rate prediction task aims to enhance the profitability and user satisfaction of commercial recommendation systems by predicting the likelihood of a current user clicking on items recommended by the system. Therefore, about the CTR prediction tasks based on feature interaction, we can define it formally as follows.

DEFINITION 1 (CTR Prediction Based on Feature Interactions). For a given user u and item v , three groups of features are extracted from them:

- *User profiles* (x_p): age, gender, occupation, etc.
- *Item attributes* (x_a): brand, price, category, etc.
- *Context* (x_c): timestamp, device, position, etc.

As shown in Table 1, in general, x_p , x_a , and x_c are multi-field categorical data and are represented using one-hot encoding, and then we utilize an embedding layer to transform them into low-dimensional dense vectors: $\mathbf{e}_i = E_i x_i$, where $E_i \in \mathbb{R}^{d \times s_i}$ and s_i separately indicate the embedding matrix and the vocabulary size for the i -th field, d represents the embedding dimension. Analogously, for numerical data (e.g., age, price), we typically start by discretizing them into categorical data using bucketing method¹ [51, 70] and then transform them into embedding vectors using embedding methods for categorical features. Subsequently, we can obtain the result representation of the embedding layer: $\mathbf{E} = [\mathbf{e}_1, \mathbf{e}_2, \dots, \mathbf{e}_f]$, where f denotes the number of fields.

Table 1. An example of multi-field categorical data.

Click	User (x_p)			Item (x_a)			Context (x_c)		
	age	gender	occupation	brand	price	category	timestamp	device	position
1	30	female	engineer	Armani	200	clothing	2022/8/19	Huawei	London
0	20	male	student	Vuitton	1000	clothing	2022/11/3	Apple	New York
1	40	male	engineer	Gucci	600	clothing	2022/3/16	Samsung	Hong Kong
0	20	female	student	Dior	2000	cosmetic	2022/5/26	Lenovo	Tokyo
Number	10	2	500	1000	5000	1000	365	1000	1000

¹<https://www.csie.ntu.edu.tw/~r01922136/kaggle-2014-criteo.pdf>

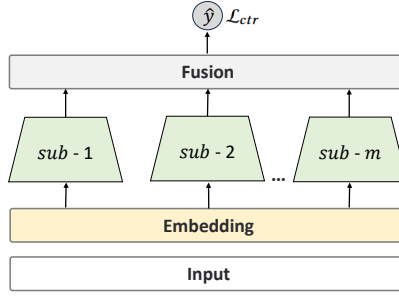


Fig. 2. The primary backbone structures of common CTR prediction models

Variable $y \in \{0, 1\}$ is an associated label for user click behavior:

$$y = \begin{cases} 1, & u \text{ has clicked } v, \\ 0, & \text{otherwise.} \end{cases} \quad (1)$$

Let $o^{(1)}$ denotes E (namely, the first-order feature interaction representation). If we need to obtain higher-order feature interactions we can use the following form:

$$o^{(i+1)} = g(o^{(1)}, o^{(i)}), \quad (2)$$

where $g(\cdot)$ denotes an interaction function and $o^{(i)}$ denotes the i -th order feature interactions. Therefore, for a simple CTR model, the common framework is formulated as:

$$\hat{y} = \text{Model}(u, v, \{o^{(1)}, o^{(2)}, \dots, o^{(n)}\}; \theta), \quad (3)$$

where Model is the CTR model with parameters θ , and o^n represents feature interactions of appropriate order.

DEFINITION 2 (Semantic Space in CTR). Consider an enhanced embedding as an additional auxiliary semantic space, and the original embedding as the original semantic space. For any semantic space, we can define it as follows:

$$\text{space}_i = \text{segment}(\mathbf{E}), \quad (4)$$

where the segment represents various enhancement or no operations, such as gating mechanisms, masking mechanisms, noise, and so on.

DEFINITION 3 (Capturing Feature Interactions from Multi-Semantic Spaces). Learning feature interactions in only one semantic space can result in a limited range of available feature interaction information. Therefore, we need to expand the semantic space further, as shown in Figure 2. We further refine the basic CTR model, enabling its formalization as follows:

$$\text{sub}_i = \text{subcomponent}_i(u, v, \{o^{(1)}, o^{(2)}, \dots, o^{(s)}\}; \theta_i), \quad (5)$$

where sub_i denotes the feature interaction information implicit in the i -th semantic space, and o^s, θ_i denote the appropriate feature interaction order and parameters in the current semantic space, respectively. After that, we utilize the fusion layer to further aggregate information from each semantic space, yielding the ultimate prediction.

$$\hat{y} = \text{fusion}(\text{sub}_1, \text{sub}_2, \dots, \text{sub}_m), \quad (6)$$

where m represents the number of suitable semantic spaces, and fusion denotes the aggregate function.

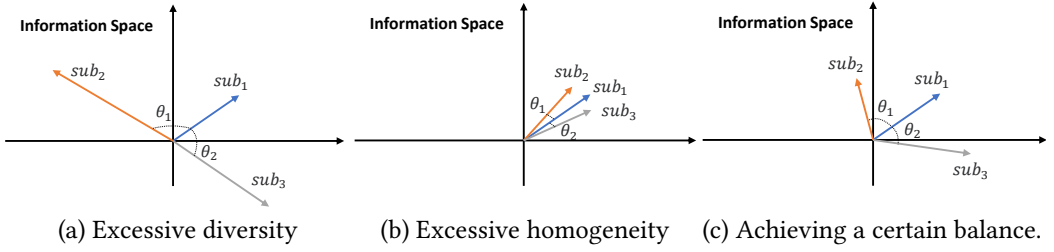


Fig. 3. An illustration of the diversity and homogeneity in \mathbb{R}^2 .

DEFINITION 4 (Diversity and Homogeneity). The concepts of diversity and homogeneity have garnered attention in sociology [2, 31]. Sociologists have delved deeply into these two concepts, seeking to comprehend different aspects and dimensions within social structures. Diversity underscores differences among individuals or groups, emphasizing the uniqueness of each member. Conversely, homogeneity emphasizes commonality and similarity, highlighting shared features and commonalities among social members. The dynamic balance between these complementary attributes constitutes diverse and interconnected social systems, aiding sociologists in gaining a better understanding of the complexity of society [26, 39].

Motivated by the insights from the aforementioned studies, we introduce these two concepts to assist the model in capturing better feature interaction information. Specifically, we can incorporate additional self-supervisory signals in the fusion layer to balance both the diversity and homogeneity of information across various semantic spaces. To illustrate the benefits of this approach more intuitively, we can illustrate these two concepts in \mathbb{R}^2 . In Figure 3, if the model captures feature interaction information that is too dissimilar across different semantic spaces, it results in the scenario depicted in Figure 3 (a). Distinct subspaces acquire significantly different information, leading to excessively large angles between θ_1 and θ_2 . Conversely, as shown in Figure 3 (b), if the model captures overly similar information, redundant information is generated, resulting in suboptimal performance. Figure 3 (c) represents a balanced scenario. Assuming sub_1 is the semantic space with the highest information quality, we can use self-supervisory signals to compel sub_2 and sub_3 to converge towards sub_1 (i.e., reducing the angles θ_1 and θ_2), while preserving their individual information to some extent.

2.2 Contrastive Learning

Contrastive learning (CL) is one of the mainstream methods in self-supervised learning [57, 61–63], which first **augments** the input samples to obtain representations from multiple perspectives, then tries to encourage consistency between pairs of positive samples (**alignment**) and minimize the agreement between pairs of negative samples (**uniformity**), ultimately enhancing the model's performance. InfoNCE [35] loss plays a pivotal role in CL, in which the CL loss is:

$$\mathcal{L}_{cl} = \sum_{i \in \mathcal{B}} -\log \frac{\exp(\text{sim}(\mathbf{x}'_i, \mathbf{x}''_i)/\tau)}{\sum_{j \in \mathcal{B}} \exp(\text{sim}(\mathbf{x}'_i, \mathbf{x}''_j)/\tau)}, \quad (7)$$

where \mathcal{B} is the batch size, \mathbf{x}' , \mathbf{x}'' are input instance representations learned from two different augmentations, τ represents the temperature coefficient, $\text{sim}(\cdot, \cdot)$ is employed to evaluate the mutual information score between them. However, as pointed out in some previous works [15, 62, 63], the augmentation and alignment operations on positive samples are not always effective. Therefore, we

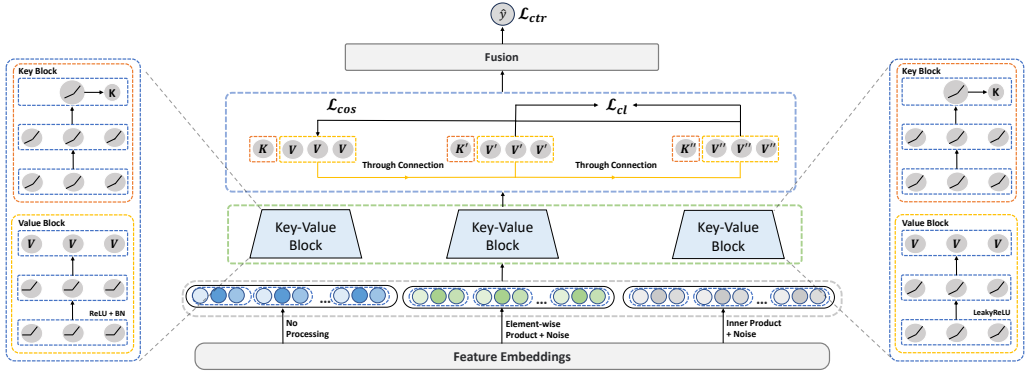


Fig. 4. The architecture of CETN

redefine a contrastive loss more suitable for multi-semantic space learning, which will be presented in Section 3.2.

3 CONTRAST-ENHANCED THROUGH NETWORK (CETN)

In this section, we will introduce the Contrast-enhanced Through Network (CETN) model. We will approach it from three perspectives and describe the architecture of the CETN model in detail.

3.1 How to Simply and Efficiently Capture Feature Interactions from Multi-Semantic Spaces

Most existing parallel-structured CTR models attempt to capture feature interaction information from multi-semantic spaces using **subcomponents**, such as DCN [53], DCNv2 [54], AutoInt [47], xDeepFM [30], and EDCN [3]. From their model architectures, we can learn that if we want to efficiently capture feature interaction information from multi-semantic spaces, we need multiple subcomponents to model feature interactions in parallel. However, as we pointed out in Section 1, relying only on the ability of subcomponents to capture information is undoubtedly inefficient. Therefore, in order to better capture the feature interaction information in a multi-semantic space, we need to further **segment** the information implicit in the semantic space. What can be further considered is that we also need a suitable fusion layer to aggregate information from multi-semantic spaces. Some existing work performs simple summing [5, 8, 14] or more complex operations [3, 28, 34] (concat, product, linear transform) as a fusion layer for the outputs of each subcomponent, but does not take into account the information weights of the semantic spaces. Therefore, we need a spatial-level attention mechanism to learn the appropriate **weights** for each semantic space.

In summary, we identify three key elements for extracting feature interaction information from multiple semantic spaces: subcomponent, segmentation, and weight. So, how do we construct a simple and effective CTR model from these elements?

Firstly, considering the efficiency and flexibility of MLP, which performs well in parallel, we utilize it as the subcomponent of the model in each semantic space, and capture information in each semantic space:

$$\begin{aligned}
 \mathbf{v} &= \text{MLP}_v(\mathbf{E}), \\
 \mathbf{v}' &= \text{MLP}'_v(\mathbf{SEP}), \\
 \mathbf{v}'' &= \text{MLP}''_v(\mathbf{SIP}),
 \end{aligned} \tag{8}$$

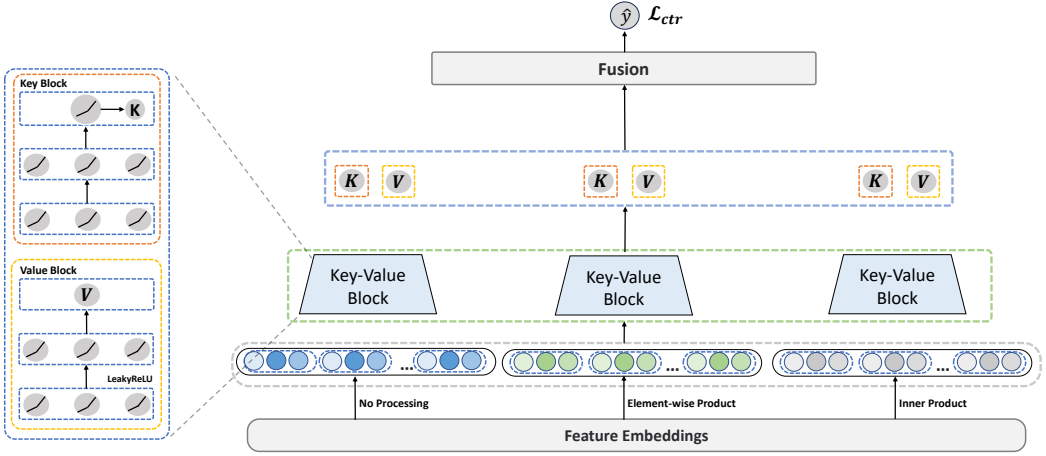


Fig. 5. The architecture of simMHN

Table 2. Performance comparison of simMHN with three strong baseline models.

Models	Avazu		Criteo		Movielens	
	Logloss↓	AUC(%)↑	Logloss↓	AUC(%)↑	Logloss↓	AUC(%)↑
DNN (base)	0.372142	79.2717	0.438233	81.3728	0.208941	96.9276
DCN	0.372353	79.3142	0.438091	81.4103	0.204643	97.0140
FinalMLP	0.372084	<u>79.3177</u>	0.437631	81.4472	0.196641	97.2373
xDeepFM	<u>0.371944</u>	79.3121	0.437820	81.4291	0.206501	96.9769
simMHN	0.371091	0.794810	<u>0.437681</u>	<u>81.4355</u>	<u>0.199238</u>	<u>97.0165</u>

where MLP_v mentioned here is a customizable multi-layer perceptron that uses LeakyReLU as an activation function for its hidden layer and no activation function for its output layer, S denotes the augmented embeddings, v is a scalar and represents the final real-values of the feature interaction information in the current semantic space.

Secondly, the product operation is highly effective for modeling feature interactions [24, 41, 69]. Different product operations often yield distinct feature interaction information. Hence, we conduct pairwise product operations on the feature embeddings, represented as segmented E , to yield S .

$$\begin{aligned} S_{EP} &= \parallel_{i=1}^n \parallel_{j=i}^n \mathcal{F}_{EP}(\mathbf{e}_i, \mathbf{e}_j), \\ S_{IP} &= \parallel_{i=1}^n \parallel_{j=i}^n \mathcal{F}_{IP}(\mathbf{e}_i, \mathbf{e}_j), \end{aligned} \quad (9)$$

where \parallel denotes the concat operation, $\mathcal{F}_{(\cdot)}$ indicates a product operation, and $\forall \mathbf{e}_i, \mathbf{e}_j \in E$. EP denotes element-wise product, while IP represents vector inner product.

Thirdly, attention mechanisms have been widely employed in CTR [10, 28, 47, 58, 68], but these attention mechanisms are only feature-level and do not work well to weight the information in the semantic space. Therefore, we introduce a spatial-level attention mechanism to better aggregate information from various semantic spaces within the fusion layer, in which the spatial-level attention mechanism is:

$$\mathbf{K}_i = \text{MLP}_k^i(\text{space}_i), \quad (10)$$

where MLP_k and MLP_v are nearly identical, with the difference being the use of LeakyReLU activation in the output layer instead of no activation function. \mathbf{K} is a scalar that represents

the weight of information implied by the current semantic space (e.g., \mathbf{E} and \mathbf{S}). To facilitate the discussion, we collectively refer to the subcomponent consisting of MLP_k and MLP_v as the *Key-Value Block*.

In the end, we obtain the final prediction result by performing a straightforward weighted summation pooling as the fusion layer.

$$\hat{y} = \sum_{i=1}^H \mathbf{K}_i \mathbf{V}_i, \quad (11)$$

where H represents the number of semantic spaces. We call this model, which simply and efficiently captures feature interaction information from multi-semantic spaces, the **Simple Multi-Head Network (simMHN)**. This model is the predecessor of CETN, whose performance and architecture are shown in Table 2 and Figure 5.

It can be observed that this simple model achieves state-of-the-art performance, which obtains sub-optimal results and even outperforms the strongest baseline model FinalMLP [34].

Next, we attempt to enhance the model's performance by focusing on the diversity and homogeneity of information, without significantly altering the complexity of the simMHN model.

3.2 How to Improve the Diversity of Information Captured

Perturbation operations in contrastive learning can further differentiate information within semantic spaces [56, 63]. Therefore, we employ product & perturbation as the better segmentation approach for semantic spaces and distinguish between primary and auxiliary semantic spaces (the semantic space processed by segmentation is called the auxiliary space, and vice versa). In this way we can get an augmented \mathbf{E}' :

$$\begin{aligned} \mathbf{E}' &= \mathbf{E} + \Delta', \\ \Delta' &= \bar{\Delta} \odot \text{sign}(\mathbf{E}), \\ \bar{\Delta} &\in \mathbb{R}^d \sim U(0, 1), \end{aligned} \quad (12)$$

where Δ' represents a noise signal, $\bar{\Delta}$ is a sample drawn from a uniform distribution between 0 and 1. Afterwards, we further differentiate information within the semantic space using the product operation:

$$\begin{aligned} \mathbf{S}'_{EP} &= \|\|_{i=1}^n \|\|_{j=i}^n \mathcal{F}_{EP}(\mathbf{e}'_i, \mathbf{e}'_j), \\ \mathbf{S}'_{IP} &= \|\|_{i=1}^n \|\|_{j=i}^n \mathcal{F}_{IP}(\mathbf{e}'_i, \mathbf{e}'_j), \end{aligned} \quad (13)$$

where $\forall \mathbf{e}'_i, \mathbf{e}'_j \in \mathbf{E}'$. After segmentation, we improve the diversity of information within the semantic space. However, it becomes apparent that there is no inherent supervisory signal between multiple subcomponents to ensure their ability to capture distinct feature interactions. Relying solely on a self-adaptive joint learning strategy among subcomponents often leads to suboptimal performance. Therefore, we introduce auxiliary supervisory signals for the real-valued semantic space.

Given that we imitate the augmentation concept from contrastive learning during segmentation, we employ a contrastive loss to supervise the information-capturing behavior of subcomponents within the auxiliary semantic space. However, it's worth noting that the alignment strategy in contrastive learning is not always effective, as indicated in previous studies [62, 63]. A similar observation holds for CTR tasks based on user historical behavior sequences [15]. As a result, we modify the contrastive loss, abandoning alignment while retaining uniformity. More specifically, we consider the information in both auxiliary spaces to be negative samples even if they are obtained from the same \mathbf{S} , and the modified loss is formulated to minimize the agreement with these samples,

as follows:

$$\mathcal{L}_{cl} = \sum_{i \in \mathcal{B}} -\log \frac{\exp(1/\tau)}{\sum_{j \in \mathcal{B}} \exp(\text{sim}(\mathbf{V}'_i, \mathbf{V}''_j)/\tau)}, \quad (14)$$

where \mathbf{V}' , $\mathbf{V}'' \in \mathbb{R}^{d_v}$ represent the real-values (obtained from different *Key-Value Blocks*) within the auxiliary semantic space after segmentation, d_v represents the real-values vector dimension. We refer to this modified loss as **Denominator-only InfoNCE** (Do-InfoNCE).

For subcomponents, we seek to further enhance their inherent ability to capture information from different semantic spaces. Some studies suggest that choosing appropriate activation functions can significantly impact the performance of neural networks in various application scenarios [9, 23, 59]. For instance, when there is excessive noise in the current semantic space, we can consider using ReLU as the activation function to filter out irrelevant information. Therefore, utilizing suitable activation functions in different semantic spaces can aid the model in effectively capturing diverse feature interactions. We employ different activation functions within MLP_v belonging to different semantic spaces.

3.3 How to Ensure Homogeneity of Information

After segmenting the semantic space and enhancing the ability of individual subcomponents to capture diverse information, while the model can better capture a variety of information, the increase in the number of feature spaces can lead to an issue of excessive noise. Therefore, we further introduce the concept of homogeneity, which complements diversity. Specifically, we aim for the information captured by various semantic spaces to be as distinct as possible (diversity). However, we also seek this information to be fundamentally similar (homogeneity), without being entirely dissimilar. To achieve this goal, we draw inspiration from residual networks [16] and introduce a **Through Network** to ensure the homogeneity of the captured information across various semantic spaces. Formally, we define this **Through Connection** as:

$$\begin{aligned} \mathbf{V} &= \text{MLP}_v(\mathbf{E}), \\ \mathbf{V}' &= \text{MLP}'_v(\mathbf{S}'_{EP}) + \mathbf{V}, \\ \mathbf{V}'' &= \text{MLP}''_v(\mathbf{S}'_{IP}) + \mathbf{V}, \end{aligned} \quad (15)$$

where $\mathbf{V} \in \mathbb{R}^{d_v}$ denotes the real-values information in the main semantic space, \mathbf{S}' represents the augmented embeddings in the auxiliary semantic space. By constructing the model in this manner, we establish a framework that is able to connect subcomponents in multiple semantic spaces, ensuring their homogeneity in capturing information while mitigating overfitting. Additionally, we also simply introduce the cosine similarity as an auxiliary loss, which takes the following form:

$$\begin{aligned} \mathcal{L}'_{cos} &= \sum_{i \in \mathcal{B}} 1 - \text{sim}(\mathbf{V}_i, \mathbf{V}'_i), \\ \text{sim}(\mathbf{V}, \mathbf{V}') &= \frac{\mathbf{V}^\top \mathbf{V}'}{\|\mathbf{V}\| \|\mathbf{V}'\|}. \end{aligned} \quad (16)$$

The \mathcal{L}'_{cos} encourages homogeneity between \mathbf{V} and \mathbf{V}' , other similarity functions can also be utilized here.

Algorithm 1: The overall training process of CETN

Require: feature embeddings \mathbf{E} ;

- 1: Initialize all parameters.
- 2: **if** use perturbing **then**
- 3: Get \mathbf{E}' according to Equation (12);
- 4: **else**
- 5: Get $\mathbf{E}' = \mathbf{E}$;
- 6: **end**
- 7: Construct \mathbf{S}'_{EP} and \mathbf{S}'_{IP} according to Equation (13);
- 8: **while** CETN not converge **do**
- 9: Calculate spatial information weights \mathbf{K} , \mathbf{K}' , and \mathbf{K}'' according to Equation (10);
- 10: Calculate real-values of each semantic space \mathbf{V} , \mathbf{V}' , and \mathbf{V}'' according to Equation (15);
- 11: Calculate \mathcal{L}_{cl} according to Equation (14);
- 12: Calculate \mathcal{L}_{cos} according to Equation (16);
- 13: Fusion information from various semantic spaces according to Equation (17);
- 14: Update the model parameters using Equation (19);
- 15: **end while**

3.4 Fusion Layer and Multi-task Training

After capturing information concurrently from H semantic spaces, the real-values (\mathbf{V}_i) and weights (\mathbf{K}_i) from each semantic space are aggregated by the fusion layer:

$$\hat{y} = \sum_{i=1}^H \mathbf{K}_i (\mathbf{W}^\top \mathbf{V}_i + \mathbf{b}), \quad (17)$$

where $\mathbf{W} \in \mathbb{R}^{d_v}$ and \mathbf{b} are the weight and bias parameters. At this point, we have obtained a new CTR model, **Contrast-enhanced Through Network (CETN)**, whose architecture is illustrated in Figure 4.

As we are targeting the click-through rate, which is a binary classification task, we have employed the widely-used logloss [28, 30, 47, 70] as the loss function for our model:

$$\mathcal{L}_{ctr} = -\frac{1}{N} \sum_{i=1}^N (y_i \log(\hat{y}_i) + (1 - y_i) \log(1 - \hat{y}_i)), \quad (18)$$

where N is the total number of training samples, i represents the sample index, y and \hat{y} represent the true label and the predicted result of CETN. Next, we can obtain the total loss, denoted as \mathcal{L}_{total} :

$$\mathcal{L}_{total} = \mathcal{L}_{ctr} + \alpha \cdot \mathcal{L}_{cl} + (\beta' \cdot \mathcal{L}'_{cos} + \beta'' \cdot \mathcal{L}''_{cos}), \quad (19)$$

where α , β' , β'' represent hyperparameters that control the balance between the loss functions.

3.5 Training Procedure

In the above sections, we engaged in a detailed discussion on how to simply yet effectively capture feature interactions across various semantic spaces, as well as how to balance the diversity and homogeneity of the captured information. To facilitate a deeper understanding of our proposed CETN model, we provide a comprehensive depiction of its training process in Algorithm 1.

Initially, we decide whether to perturb the embedding \mathbf{E} for augmentation (lines 2-6), followed by conducting product-based augmentation (line 7) to distinguish the feature interaction information

Table 3. Comparison of Analytical Time Complexity
 $N \gg s > |W_\Psi| \approx |W_{gate}| > d_1 \approx d_2 \approx |h_1| > f \approx d \approx f_s \approx M \approx U$

Model	Embedding	Implicit interaction	Explicit interaction	Objective Function
DNN	$O(2dfs)$	$O(W_\Psi)$	-	$O(N)$
DCN	$O(2dfs)$	$O(W_\Psi)$	$O(4dfL)$	$O(N)$
FinalMLP	$O(2dfs + 2dsf_s)$	$O(2 W_\Psi)$	$O(2d_1d_2 + 2 W_{gate} + 2df_s)$	$O(N)$
xDeepFM	$O(2dfs)$	$O(W_\Psi)$	$O(dfM(1 + ML))$	$O(N)$
AFN+	$O(4dfs)$	$O(W_\Psi)$	$O(2df(1 + U))$	$O(N)$
CLACTR	$O(2dfs)$	$O(3 W_\Psi)$	-	$O(N + N h_1 + \frac{dfN(1+N)}{2} + \frac{dfN^2(f-1)}{2})$
CETN	$O(2dfs)$	$O(6 W_\Psi)$	$O(df(1 + f))$	$O(N + 2Nd_o)$

implied in different semantic spaces. At this point, we have made the necessary segmentation of the semantic spaces. Subsequently, we employ the *Key-Value Block* to capture feature interaction information within each semantic space (lines 9-10) and calculate the self-supervised signal (lines 11-12). Finally, we aggregate the information from all semantic spaces in the fusion layer and update the parameters (lines 13-14).

3.6 Model Analysis

3.6.1 Model Size. To effectively capture feature interaction information across various semantic spaces, CETN employs the *Key-Value Block* and product & perturbation. For ease of discussion, we regard W_Ψ as the set of weights in the corresponding MLP and ignore the embedding parameters. In the *Key-Value Block*, it can be simply viewed as six parallel MLPs, so its corresponding space complexity is $O(6|W_\Psi|)$. For the product & perturbation operation, we further divide additional auxiliary semantic spaces, hence its corresponding space complexity is $O(df + \frac{df(f+1)}{2} + \frac{f(f+1)}{2})$. As the space complexity represented by the inner product $\frac{f(f+1)}{2}$ is of constant level, therefore, the space complexity of CETN is $O(6|W_\Psi| + df + \frac{df(f+1)}{2})$. For more detailed information on the parameter size, refer to Tables 5 and 6.

3.6.2 Time Complexity. In a manner similar to the calculation of space complexity, the inference time of CETN is primarily due to the *Key-Value Block* and product & perturbation. Therefore, the time complexity during inference is $O(6|W_\Psi| + 2dfs + \frac{df(f+1)}{2} + df(f+1))$. It's worth mentioning that due to the relationship between the Hadamard product and the inner product (the latter being the sum of the former), in practical operations, we can optimize it as follows: $O(6|W_\Psi| + 2dfs + df(f+1))$.

Furthermore, for a more detailed comparison, we present the time complexities in the training of DNN, DCN, FinalMLP, xDeepFM, AFN+, and our proposed CETN model in Table 3. We let L , U , and M represent the number of explicit interaction layers, logarithmic neurons, and feature maps, respectively. d_1 , d_2 , W_{gate} , and f_s represent the output dimension of the two MLPs in FinalMLP, the number of parameters in the gate unit, and the feature fields to be filtered, respectively. (h_1, h_2) represent the outputs of the two feature interaction (FI) encoders in CL4CTR.

It can be concluded that CETN, compared to other models that can also capture feature interaction information in multiple semantic spaces, has a comparable time complexity. It is worth noting that, due to the use of contrastive loss in CETN, this leads to a longer training time. However, it does not affect the corresponding inference speed in actual applications. In the case of implicit interactions, although CETN requires six MLPs, the practical time complexity does not increase proportionally. This is because MLPs are parallel-friendly and simple yet effective, which ensures the time complexity remains manageable.

3.6.3 Comparison with AFN+. The Adaptive Feature Network (AFN) [6] is a solid mathematical model that cleverly applies logarithmic operation rules for adaptive order explicit feature interaction. Its enhanced version, AFN+, employs DNN and uses a Logarithmic Transformation (LT) Layer and MLP as parallel components to capture feature interaction information in two semantic spaces. Unlike previous works [14, 53], it neither shares an embedding layer nor uses gating units to segment the semantic space, instead simply maintaining separate embedding matrices for the two parallel components. However, maintaining two sets of embedding matrices is inefficient due to the majority of parameters and computational costs being taken up by the embedding operation (the performance of AFN+ will be demonstrated in Section 4). The LT layer in AFN+ uses logarithmic operations to simulate the Hadamard product, which can be formulated as follows:

$$\exp\left(\sum_{i=1}^m w_{ij} \ln \mathbf{e}_i\right) = \mathbf{e}_1^{w_{1j}} \odot \mathbf{e}_2^{w_{2j}} \odot \dots \odot \mathbf{e}_m^{w_{mj}}, \quad (20)$$

where w_{ij} is the coefficient of the j -th neuron in the i -th feature field and \odot denotes Hadamard product. It is not difficult to observe that when the model learns the correct value for w_{ij} , it can adaptively learn feature interactions of any order. However, since w_{ij} is usually non-zero, the interaction learned by AFN is essentially an exponentially weighted, fixed full-order feature interaction. Interestingly, it is pointed out in some existing works that high-order feature interactions do not bring the expected performance improvement [32, 69], leading to subpar performance of AFN+. In CETN, we only use the Hadamard product to model second-order feature interactions, serving as an enhanced auxiliary semantic space, thus achieving better performance in both model complexity and performance metrics.

3.6.4 Comparison with CL4CTR. CL4CTR [52] is the first to introduce the contrastive learning paradigm into feature interaction-based CTR prediction models. It uses Euclidean distance loss for feature alignment to make feature representations in the same feature field as similar as possible, and employs cosine loss for field uniformity to maximize the difference between feature representations in different feature fields. Additionally, it uses a perturbation operation to divide into two new feature interaction auxiliary spaces and compares the encoded feature interaction information through a feature interaction (FI) encoder. Moreover, the CL4CTR still uses Euclidean distance loss as a contrastive loss to minimize the encoding differences between the two FI encoders, as shown in Equation (21), without introducing the widely acclaimed InfoNCE loss.

$$\mathcal{L}_{CL4CTR} = \frac{1}{B} \sum_{i=1}^B \|h_{i,1} - h_{i,2}\|_2^2, \quad (21)$$

The fundamental difference between CETN and CL4CTR lies in their approach to handling information captured in the auxiliary semantic spaces. Where CL4CTR aims to make the captured information increasingly similar, CETN uses Do-InfoNCE to increase the diversity (i.e., dissimilarity) of captured information, while ensuring homogeneity (i.e., similarity) with the information in the original semantic space. Furthermore, as shown in Table 3, the time complexity of the contrastive loss in CL4CTR is very high, making it challenging to realistically apply in real-world production scenarios.

4 EXPERIMENTS

In this section, we provide a detailed account of our experimental setup and substantiate the superiority of CETN over other state-of-the-art models through a fair and extensive series of experiments. Subsequently, we conduct ablation experiments to investigate the impact of our

configured hyperparameters on model performance and assess the rationale behind the presence of various modules.

Table 4. Dataset statistics

Dataset	#Instances	#Fields	#Features	#Split
Avazu	40M	24	8.37M	8:1:1
Criteo	46M	39	5.55M	8:1:1
MovieLens	2,006,859	3	90,445	7:2:1
Frappe	288,609	10	5,382	7:2:1

4.1 Experimental Settings

To ensure a fair comparison, we closely followed the settings of the [34, 70] work and selected the same CTR benchmark datasets originating from real production environments. Table 4 below provides detailed information about these datasets.

- **Avazu**²: This dataset contains 10 days of data on user clicks to ads while using mobile devices, as well as 15 explicit and 9 anonymous feature fields.
- **Criteo**³: It is the well-known CTR benchmark dataset, which contains a 7-day stream of real data from Criteo, covering 39 anonymous feature fields.
- **MovieLens**⁴: It consists of users' tagging records on movies. The datasets have been widely used in various research on recommender systems.
- **Frappe**⁵: It contains app usage logs from users under different contexts (e.g., daytime, location). The target value indicates whether the user has used the app under the context.

Data preprocessing: We follow the approach outlined in [70]. For the Avazu dataset, we transform the timestamp field it contains into three new feature fields: hour, weekday, and weekend. For the Criteo dataset, we discretize the numerical feature fields by rounding down each numeric value x to $\lfloor \log^2(x) \rfloor$ if $x > 2$, and $x = 1$ otherwise. For all datasets' categorical features, infrequent features ($\text{min_count}=2$) are uniformly replaced with a default "OOV" token.

4.1.1 Evaluation Metrics. In order to compare the performance, we utilize two commonly used metrics in CTR models: AUC and logloss [14, 28, 40, 47].

- **AUC:** AUC stands for Area Under the ROC Curve. It measures the probability that a positive instance will be ranked higher than a randomly chosen negative one. A higher AUC indicates a better performance.
- **Logloss:** logloss is the calculation result of Equation (18). A lower logloss suggests a better capacity for fitting the data.

It's worth noting that even a slight improvement in AUC is meaningful in the context of CTR prediction tasks. Typically, CTR researchers consider improvements at the **0.001-level (0.1%)** to be statistically significant, as often highlighted in existing literature [3, 5, 14, 52, 53, 70]. Following previous work [46, 52, 55, 68], we further use Relaimpr to measure the relative improvement of

²<https://www.kaggle.com/c/avazu-ctr-prediction>

³<https://www.kaggle.com/c/criteo-display-ad-challenge>

⁴<https://grouplens.org/datasets/movielens/>

⁵<http://baltrunas.info/research-menu/frappe>

AUC and Logloss, as defined:

$$\begin{aligned} \text{RelaImpr}_{AUC} &= \left(\frac{\text{AUC}(\text{target model}) - 0.5}{\text{AUC}(\text{base model}) - 0.5} - 1 \right) \times 100\%, \\ \text{RelaImpr}_{\text{Logloss}} &= \left(\frac{\text{Logloss}(\text{base model}) - \text{Logloss}(\text{target model})}{\text{Logloss}(\text{base model})} \right) \times 100\%, \end{aligned} \quad (22)$$

4.1.2 Baselines. We compared CETN with some classical state-of-the-art (SOTA) models. Given that deep CTR models often perform better, for models that have both non-DNN and DNN versions, we tend to choose the latter. The list of models we have chosen in chronological order of publication is as follows:

- **LR** [44]: Logistic regression (LR) is a simple baseline model for CTR prediction.
- **FM** [42]: This model employs factorization techniques to address the challenge of learning on sparse datasets.
- **DNN** [7]: This approach utilizes a feedforward neural network that takes a straightforward concatenation of feature embeddings as input.
- **IPNN** [40]: The model is an inner product-based feedforward neural network.
- **Wide & Deep** [5]: It encompasses logistic regression (wide network) and the integration of a feedforward neural network (deep network).
- **DeepFM** [14]: This model combines FM and feedforward neural networks in parallel by sharing embedding layers.
- **NFM** [17]: This method vertically combines FM and feed-forward neural networks through the Bi-interaction layer.
- **AFM** [58]: The model incorporates an attention mechanism to discern the importance of different feature interactions.
- **DCN** [53]: The model introduces the CrossNet that can explicitly model feature interactions, and integrate feedforward neural networks in parallel.
- **xDeepFM** [30]: Similar to DCN, this model introduces the Compressed Interaction Network (CIN), enhancing feature interaction from bit-wise to vector-wise.
- **FiGNN** [28]: This model pioneers the use of graph neural networks to model feature interactions.
- **AutoInt+** [47]: This model is the first to learn higher-order feature interactions using a multi-headed attention mechanism.
- **AFN+** [6]: The model utilizes a logarithmic transformation layer to learn adaptive-order feature interactions.
- **DCNv2** [54]: Expanding upon DCN, this model enhances the projection matrix's dimensionality and introduces a mixture of low-rank expert systems to optimize model inference speed.
- **EDCN** [3]: This model introduces both a bridge module and a regulation module, thereby enhancing the performance of the DCN model.
- **MaskNet** [55]: This model utilizes the MaskBlock as its foundational structure. The introduced Instance-Guided Mask in the MaskBlock assists the model in acquiring high-quality information more effectively.
- **GraphFM** [29]: This model employs graph structure learning to address FM's limitations in selecting and learning appropriate higher-order feature interactions.
- **FinalMLP** [34]: The model demonstrates the efficacy of the two-stream MLP model for implicit feature interaction learning.

Table 5. Performance comparison of different models for CTR prediction. We highlight the top-5 best results in each dataset. "+": Integrating the original model with DNN networks.

Year	Models	Avazu			Criteo		
		Logloss↓	AUC↑	#Params	Logloss↓	AUC↑	#Params
2007	LR [44]	0.381727	0.777286	3.7M	0.456573	0.793558	5.5M
2010	FM [42]	0.376240	0.786085	78.7M	0.444295	0.807856	116.5M
2016	DNN [7]	0.372142	0.792717	78.5M	0.438233	0.813728	115.7M
2016	IPNN [40]	0.371156(2)	0.794330(2)	75.6M	0.438217	0.813945	114.6M
2016	Wide & Deep [5]	0.372015	0.792982	82.2M	0.438110	0.813814	120.4M
2017	DeepFM [14]	0.371890(5)	0.793116	82.4M	0.437676(3)	0.814200(5)	117.6M
2017	NFM [17]	0.373843	0.789997	79.2M	0.446362	0.805424	118.6M
2017	AFM [58]	0.378901	0.782119	78.7M	0.444468	0.807100	116.5M
2017	DCN [53]	0.372353	0.793142	76.5M	0.438091	0.814103	132.5M
2018	xDeepFM [30]	0.371944	0.793121	79.4M	0.437820	0.814291(4)	120.5M
2019	FiGNN [28]	0.374346	0.789236	75.0M	0.438860	0.813340	111.1M
2019	AutoInt+ [47]	0.372085	0.792639	77.5M	0.438919	0.812950	170.5M
2020	AFN+ [6]	0.372007	0.793513(5)	177.9M	0.439244	0.813098	226.3M
2021	DCNv2 [54]	0.372111	0.793050	77.2M	0.438392	0.813951	114.9M
2021	EDCN [3]	0.371627(4)	0.793874(3)	75.7M	0.437742(4)	0.814292(3)	112.8M
2021	MaskNet [55]	0.371623(3)	0.793677(4)	77.8M	0.439455	0.812424	116.8M
2022	GraphFM [29]	0.372646	0.791912	75.0M	0.441261	0.810482	111.0M
2023	FinalMLP [34]	0.372084	0.793177	80.0M	0.437631(2)	0.814472(2)	116.2M
2023	CL4CTR [52]	0.373158	0.791745	76.2M	0.437794(5)	0.814159	112.5M
2023	EulerNet [49]	0.376032	0.792391	75.6M	0.442632	0.811003	113.4M
ours	CETN	0.370402(1)	0.796238(1)	85.1M	0.437319(1)	0.814804(1)	140.1M

- **CL4CTR** [52]: This model pioneers the integration of contrastive learning into CTR prediction based on feature interaction by introducing the concepts of feature alignment and field uniformity.
- **EulerNet** [49]: This model employs Euler’s formula to explicitly model feature interactions, integrating it with linear layers to adaptively learn feature interactions of the arbitrary order.

4.1.3 Implementation Details. We implemented all models using Pytorch [38] and refer to existing works [20, 70]. We employ the Adam optimizer [25] to optimize all models, with a default learning rate set to 0.001. For the sake of fair comparison, we set the embedding dimension to 20, and the batch size to 10,000 for all models. The hyperparameters of the baseline model were configured based on the *optimal values* provided in [20, 70] and their original paper. To prevent overfitting, we employed early stopping with a patience value of 2. To facilitate reproducible research, we have open-sourced the code and running logs of CETN⁶.

4.2 Overall Comparison

The performance of CETN and the baseline model is shown in Table 5 and Table 6. It can be observed that deep Click-Through Rate (CTR) models, with DeepFM [14] as their representative, consistently outperform shallow models, typified by FM [42]. This underscores the effectiveness of combining implicit high-order feature interactions with explicit shallow feature interactions, resulting in improved click-through rate predictions. Concurrently, it emphasizes the necessity of leveraging explicit feature interactions to overcome the performance bottlenecks associated with MLP.

⁶<https://github.com/salmon1802/CETN>

Table 6. Performance comparison of different models for CTR prediction. We highlight the top-5 best results in each dataset. "+": Integrating the original model with DNN networks.

Year	Models	MovieLens			Frappe		
		Logloss↓	AUC↑	#Params	Logloss↓	AUC↑	#Params
2007	LR [44]	0.345537	0.931499	88597	0.364208	0.931528	5384
2010	FM [42]	0.276611	0.942978	1.8M	0.193952	0.969701	0.1M
2016	DNN [7]	0.208941	0.969276	2.1M	0.169329	0.980640	0.5M
2016	IPNN [40]	0.206157	0.970354(3)	2.1M	0.151626(5)	0.984333(4)	0.5M
2016	Wide & Deep [5]	0.207965	0.969820	2.2M	0.152807	0.984018	0.5M
2017	DeepFM [14]	0.205239(4)	0.969749	2.2M	0.153563	0.983501	0.5M
2017	NFM [17]	0.302908	0.939400	2.1M	0.160586	0.980831	0.4M
2017	AFM [58]	0.277981	0.942794	1.8M	0.242329	0.955657	0.1M
2017	DCN [53]	0.204643	0.970140(5)	2.1M	0.155340	0.983995	0.5M
2018	xDeepFM [30]	0.206501	0.969769	2.3M	0.155676	0.983409	0.5M
2019	FiGNN [28]	0.256296	0.952781	1.8M	0.226627	0.964828	0.2M
2019	AutoInt+ [47]	0.203297(3)	0.970191(4)	2.2M	0.150433(2)	0.984076(5)	0.7M
2020	AFN+ [6]	0.205358(5)	0.969492	8.0M	0.156007	0.980949	3.8M
2021	DCNv2 [54]	0.206726	0.969712	2.1M	0.150948(4)	0.984368(3)	0.6M
2021	EDCN [3]	0.228215	0.968716	1.7M	0.161859	0.983283	0.2M
2021	MaskNet [55]	0.242475	0.968499	2.8M	0.189036	0.982834	1.6M
2022	GraphFM [29]	0.222897	0.964445	1.7M	0.288119	0.939295	0.1M
2023	FinalMLP [34]	0.196641(2)	0.972373(2)	2.2M	0.150553(3)	0.984854(2)	0.6M
2023	CL4CTR [52]	0.206971	0.969982	2.4M	0.152105	0.983712	1.0M
2023	EulerNet [49]	0.213534	0.965486	1.7M	0.153704	0.980542	0.2M
ours	CETN	0.185652(1)	0.973957(1)	1.9M	0.150283(1)	0.985710(1)	1.6M

Table 7. Relative improvement of AUC and Logloss with CETN. Δ Logloss and Δ AUC denote the average performance improvement.

Model	Avazu _{Relaimpr}		Criteo _{Relaimpr}		MovieLens _{Relaimpr}		Frappe _{Relaimpr}		Δ Logloss ↓	Δ AUC ↑
	Logloss↓	AUC(%)↑	Logloss↓	AUC(%)↑	Logloss↓	AUC(%)↑	Logloss↓	AUC(%)↑		
LR [44]	2.97%	6.83%	4.22%	7.24%	46.27%	9.84%	58.74%	12.56%	28.05%	8.44%
FM [42]	1.55%	3.55%	1.57%	2.26%	32.88%	6.99%	22.52%	3.41%	14.63%	4.95%
DNN [7]	0.47%	1.20%	0.21%	0.34%	11.15%	1.00%	11.25%	1.05%	5.77%	0.89%
IPNN [40]	0.20%	0.65%	0.20%	0.27%	9.95%	0.77%	0.89%	0.28%	2.81%	0.61%
Wide & Deep [5]	0.43%	1.11%	0.18%	0.32%	10.73%	0.88%	1.65%	0.35%	3.25%	0.80%
DeepFM [14]	0.40%	1.07%	0.08%	0.19%	9.54%	0.90%	2.14%	0.46%	3.04%	0.76%
NFM [17]	0.92%	2.15%	2.03%	3.07%	38.71%	7.86%	6.42%	1.01%	12.02%	5.24%
AFM [58]	2.24%	5.00%	1.61%	2.51%	33.21%	37.04%	37.98%	6.60%	18.76%	5.40%
DCN [53]	0.52%	1.06%	0.18%	0.22%	9.28%	0.81%	3.26%	0.35%	3.31%	0.73%
xDeepFM [30]	0.41%	1.06%	0.11%	0.16%	10.10%	0.89%	3.46%	0.48%	3.52%	0.75%
FiGNN [28]	1.05%	2.42%	0.35%	0.47%	27.56%	4.68%	33.69%	4.49%	15.66%	3.06%
AutoInt+ [47]	0.45%	1.23%	0.36%	0.59%	8.68%	0.80%	0.10%	0.34%	2.40%	0.86%
AFN+ [6]	0.43%	0.93%	0.44%	0.54%	9.60%	0.95%	3.67%	0.99%	3.53%	0.84%
DCNv2 [54]	0.46%	1.09%	0.24%	0.27%	10.19%	0.90%	0.44%	0.28%	2.83%	0.79%
EDCN [3]	0.33%	0.80%	0.10%	0.16%	18.65%	1.12%	7.15%	0.50%	6.56%	0.80%
MaskNet [55]	0.33%	0.87%	0.49%	0.76%	23.43%	1.16%	20.50%	0.60%	11.19%	0.99%
GraphFM [29]	0.60%	1.48%	0.89%	1.39%	16.71%	2.05%	47.84%	10.57%	16.51%	1.74%
FinalMLP [34]	0.45%	1.04%	0.07%	0.11%	5.59%	0.34%	0.18%	0.18%	1.57%	0.46%
CL4CTR [52]	0.74%	1.54%	0.11%	0.21%	10.30%	0.85%	1.20%	0.41%	3.09%	0.86%
EulerNet [49]	1.50%	1.32%	1.20%	1.22%	13.06%	1.82%	2.23%	1.08%	4.50%	1.54%

Focusing on the models that achieved top-5 performance in the experiments, it can be observed that all of these models are based on parallel or stacked structures. Interestingly, over time, when all models are configured with their respective optimal parameters, the performance improvements of

these SOTA models are relatively modest. Notably, the FinalMLP [34], consistently delivers strong performance across multiple datasets, further highlighting the importance of subcomponents and segments.

Taking an overall view, the CETN consistently maintains the highest AUC and Logloss performance, even when all models are configured with their respective optimal parameters. As is shown in Table 7, in comparison to the best baseline model, CETN exhibits improvements of 0.65% (compared to IPNN), 0.11% (compared to FinalMLP), 0.34% (compared to FinalMLP), and 0.18% (compared to FinalMLP) in AUC on the four datasets, while achieving corresponding reductions in Logloss by 0.2% (compared to IPNN), 0.07% (compared to FinalMLP), 5.59% (compared to FinalMLP), and 0.1% (compared to AutoInt+). Compared to the xDeepFM model, which also utilizes three semantic spaces, CETN achieves AUC average improvements of 0.75% in AUC and 3.52% in Logloss. This performance enhancement can be attributed to its ability to capture diversity within the feature space while preserving information homogeneity. It is noteworthy that, in contrast to intricate explicit feature interaction networks, the proposed CETN model simply integrates contrastive learning into simMHN, effectively improving its ability to capture feature interaction information.

Compared to the CL4CTR model, which also employs the contrastive learning paradigm, CETN demonstrates an average relative improvement of 3.09% in Logloss and 0.86% in AUC across the four datasets. Interestingly, in CL4CTR, the contrastive loss does not adopt the popular InfoNCE paradigm used in graph neural network-based recommender. Instead, it utilizes Euclidean distance loss as the contrastive loss, resulting in the maximization of similarity in the enhanced auxiliary semantic space embeddings. This stands in contrast to our Do-InfoNCE approach and is similar to \mathcal{L}_{cos} . This suggests that CL4CTR overlooks the connection between information in the auxiliary semantic space and the primary semantic space, emphasizing the necessity of homogeneity and diversity as unsupervised guiding signals for model learning.

4.3 Ablation Study

In this section, we conducted extensive ablation studies to assess the contributions of individual modules of CETN to its overall performance. Several variants were designed to validate the effectiveness of the various CETN modules:

- **CETN (-A):** To further enhance the diversity of information captured by the model, we employ different activation functions in multiple semantic spaces. To assess the necessity of this approach, we only specify the use of ReLU activation functions in the Multi-Layer Perceptron (MLP) for modeling feature interactions within multiple semantic spaces. This will reduce the diversity of information captured by the model.
- **CETN (-CL):** Contrastive loss serves to self-supervise and augment the diversity of information captured by the model. To determine whether it aids in improving model performance, we experiment by removing it from the model.
- **CETN (-COS):** The cosine loss primarily reinforces the homogeneity of the model from a supervised signal perspective. To ascertain its necessity, we exclude it from the model.
- **CETN (-K):** It serves to fine-tune the contribution weights of various subcomponents to the final prediction, further refining the model's predictive results. To establish the necessity of the MLP_k (spatial-level attention) within the *Key-Value Block*, we eliminate it from the model.
- **CETN (-P):** To verify the effectiveness of our proposed product & perturbation method, we replace it with the original embeddings \mathbf{E} without any additional processing.

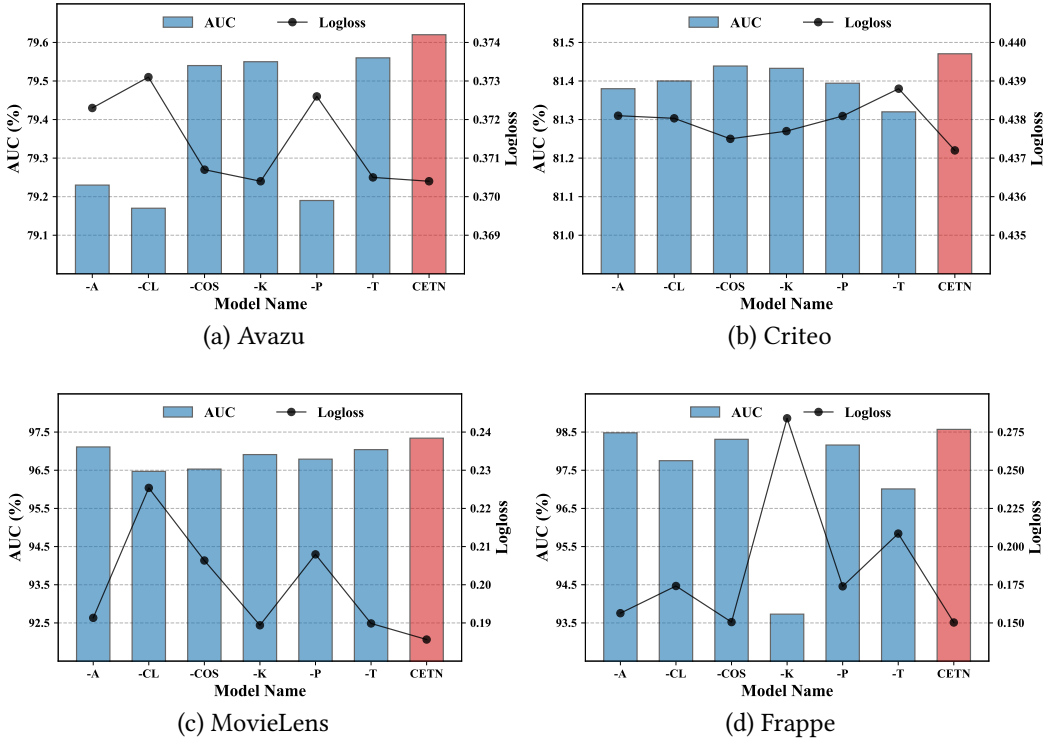


Fig. 6. Ablation study of CETN on Avazu (a), Criteo (b), MovieLens (c), and Frappe (d) datasets.

- **CETN (-T):** The Through Connection ensures the model maintains homogeneity in the information captured across multiple semantic spaces. To evaluate its usefulness, we remove it from the model.

Figure 6 illustrates the performance of CETN and its six variants. We can observe that CETN outperforms all the ablation models, providing compelling evidence for the necessity of each component in CETN. To delve deeper into this, let’s break down the performance across different datasets.

Specifically, on the Avazu dataset, we notice that the performance drop is most pronounced for CETN (-A, -CL, -P). This indicates that the Avazu dataset benefits significantly from diverse feature interaction information to boost model performance. On the other hand, the lower performance loss of CETN (-COS, -T) suggests that the Avazu dataset is more in need of diverse feature interaction information to enhance model performance.

On the Criteo dataset, we observe the most significant performance drop in CETN (-T). This highlights that an increase in the number of feature fields often introduces more noise signals. Therefore, it becomes crucial to ensure information homogeneity while enhancing information diversity.

On the MovieLens dataset, We empirically think that due to its limited number of feature fields and relatively fewer instances, it’s susceptible to overfitting. Consequently, CETN (-CL, -COS), both complementing each other, play a significant role in constraining the scope of information captured by the model, allowing it to capture high-quality interaction information, thereby preventing

overfitting to some extent. The correctness of this hypothesis can be intuitively observed in Figure 6 (c), where the model's performance experiences a noticeable decline when CETN removes \mathcal{L}_{cl} , \mathcal{L}_{cos} .

On the Frappe dataset, the performance of CETN (-K) exhibits a cliff-like decline, providing evidence for the effectiveness of the spatial-level attention mechanism. It is noteworthy that the Frappe dataset has only 0.7% and 0.6% of the data volume compared to the Avazu and Criteo datasets, respectively. This effectively underscores the challenge of models to adaptively assess the importance of information in various semantic spaces when dealing with relatively fewer data. Consequently, this leads to a poorer predictive performance of the model.

Table 8. Single performance of each semantic space. RelImpr denotes the relative improvements compared with the simMHN. The underscore indicates the performance in the optimal semantic space.

Model	Avazu		Criteo		MovieLens		Frappe	
	Logloss	AUC	Logloss	AUC	Logloss	AUC	Logloss	AUC
S_{EP}	<u>0.370488</u>	<u>0.795332</u>	0.441558	0.811001	0.236711	0.964678	0.166236	<u>0.982632</u>
S_{IP}	0.373471	0.790807	0.439198	0.813027	0.227524	0.960124	<u>0.161219</u>	<u>0.977373</u>
E	0.372142	0.792717	<u>0.438233</u>	<u>0.813728</u>	<u>0.208941</u>	<u>0.969276</u>	0.169329	0.980640
simMHN	0.371091	0.794810	0.437681	0.814355	0.199238	0.970165	0.180469	0.983528
CETN	0.370402	0.796238	0.437319	0.814804	0.185652	0.973957	0.150283	0.985710
RelImpr	0.18%	0.48%	0.08%	0.14%	6.82%	0.81%	16.72%	0.45%

4.4 Which Semantic Space is More Useful

For the three semantic Spaces we define and their subcomponents, we conduct experiments to explore their individual contributions to the model's final predictions. The performance of each separate subcomponent on the four datasets is presented in Table 8. On the Avazu dataset, S_{EP} achieved outstanding performance, even surpassing IPNN (the best baseline model). However, its performance declined under the simple fusion of simMHN, demonstrating the drawbacks of this simple fusion approach. On the Criteo and MovieLens datasets, the original embeddings E showed better results, and simMHN further improved performance through simple fusion. On the Frappe dataset, both S_{EP} and S_{IP} respectively achieved the best AUC and Logloss, but in simMHN, although the AUC performance improved further, Logloss increased, indicating that simMHN might struggle to effectively predict true click probabilities.

To further enhance the performance of the simMHN model without significantly increasing its complexity, we introduced diversity- and homogeneity-guided self-supervised signals. Additionally, we incorporated skip connections and various activation functions to balance both homogeneity and diversity. Consequently, CETN achieved improvements in Logloss by 0.18%, 0.08%, 6.82%, and 16.72% on the four datasets, and in AUC by 0.48%, 0.14%, 0.81%, and 0.45%, respectively. This demonstrates the effectiveness of ensuring both homogeneity and diversity in the information captured by the model.

4.5 Hyper-parameter Analysis

4.5.1 Impact of the Weights in the \mathcal{L}_{cos} . We conduct a further investigation into the impact of different weight parameter combinations on the model's performance. For the sake of discussion, we confine the ranges of β' and β'' to $[0.1 \sim 0.4, 0.5 \sim 0.8]$ while keeping other settings fixed. The results are presented as a heat map in Fig. 7, and it is evident that CETN achieves its optimal performance when $\beta' = 0.3$ and $\beta'' = 0.2$ on the Criteo dataset. In the heatmap of MovieLens, multiple chunked regions in the model's performance are observed. The performance continues

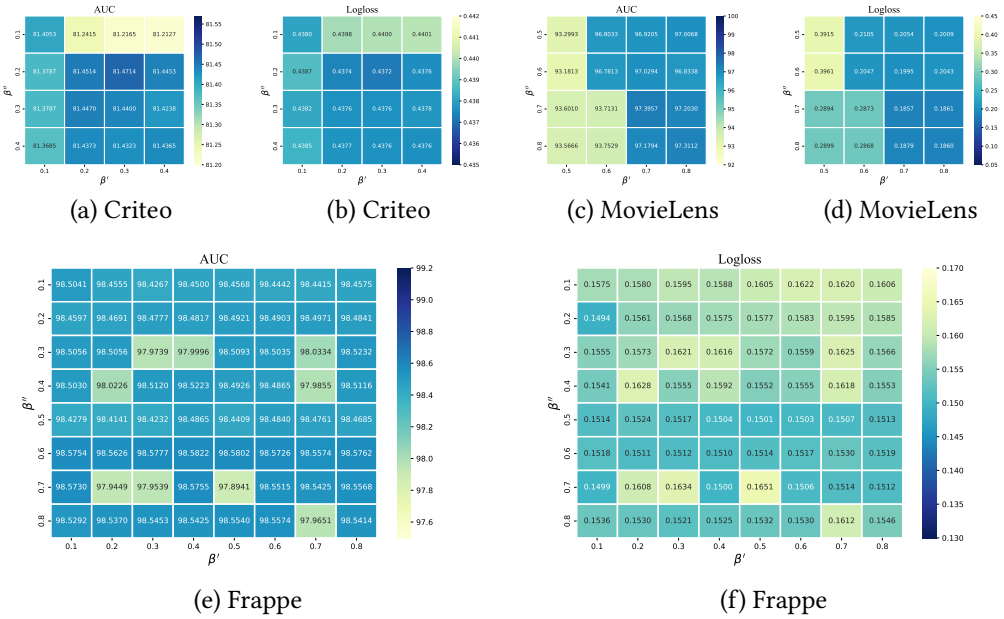


Fig. 7. Performance comparison of different weights of cosine loss on Criteo (a,b), MovieLens (c,d), and Frappe (e,f) datasets.

to decrease when β' and β'' are set to 0.5 or 0.6, but reaches optimal performance at 0.7 and 0.8. This phenomenon precisely validates the results presented in Table 8. Due to the relatively lower effective information content in the auxiliary semantic space, setting a larger \mathcal{L}_{cos} becomes necessary to prevent the model from capturing noise and ensure homogeneity of information. On the Frappe dataset, it is evident that the model performs well when β'' is 0.6 ~ 0.8 and achieves optimal results when β' equals 0.1, aligning with our previous reporting in Table 8. Notably, the model attains the lowest Logloss when $\beta'=0.1$ and $\beta''=0.2$. This is reasonable, as on the Frappe dataset, even though S_{IP} exhibits poor AUC performance, which achieves the lowest Logloss in multiple semantic spaces.

4.5.2 Impact of α in the \mathcal{L}_{cl} . We make modifications to the contrastive loss weight on the Avazu and Criteo datasets with a step size of 0.1, and in the smaller datasets MovieLens and Frappe, the modification step size is set to 0.05. The results are depicted in Figure 8. On the Criteo dataset, it's observed that the model achieves better performance when the contrastive loss weight is set to 0.2 or 0.3. Subsequently, as the weight increases, the model's performance starts to decline. Notably, On the MovieLens dataset, due to its limited three feature fields, the model is more sensitive to hyperparameter variations. Therefore, there is a substantial performance change when the weight is increased from 0.15 to 0.2. Compared to the MovieLens dataset, the performance of different values for the parameter α on the Avazu dataset is more stable. However, there is still a trend of declining performance with the increase in weight values. It's worth mentioning that compared to the simMHN model, when $\alpha=0.1$, the model's AUC performance improves by 0.48%. This demonstrates the double-edged nature of diversity. On one hand, it can help the model capture additional information, and on the other hand, it can easily introduce noise signals. On the Frappe dataset, the model's performance peaks when $\alpha=0.1$, after which the model's performance shows a

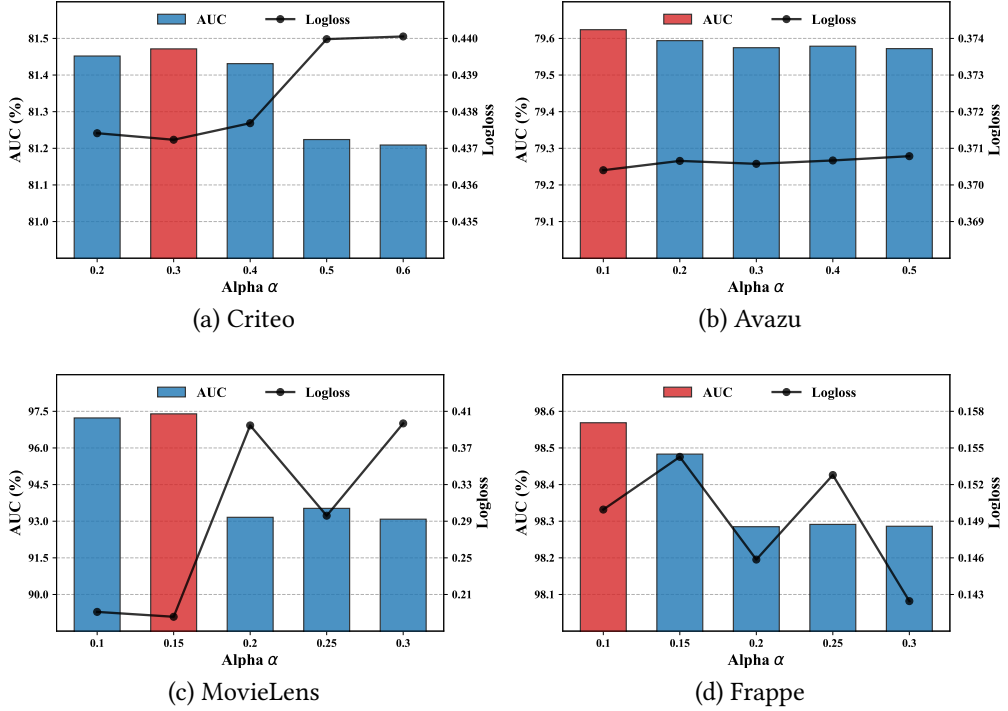


Fig. 8. Performance comparison of different weights of contrast loss on Criteo (a), Avazu (b), MovieLens (c), and Frappe (d) datasets.

precipitous decline at 0.2, consistent with the trend observed on the MovieLens dataset. Looking at the overall Figure 8, we find that when the CETN model achieves optimal performance in various datasets, the value of α is between 0.15 and 0.3. This suggests that although diversity of information is necessary, it still needs to be managed moderately. Otherwise, it will bring unnecessary noise signals to the model in various semantic spaces.

4.5.3 Impact of τ . In many models based on InfoNCE loss [56, 62, 63], the temperature coefficient is commonly set to 0.2 by default. However, due to the sensitivity of click-through rate prediction tasks to performance metrics, we further explored its impact on the model's performance. We keep other parameters fixed and change the values for τ with a step size of 0.1 in the four datasets, as shown in Figure 9. As can be observed, with an increase in the temperature coefficient, the model's performance gradually improves and then gradually decreases. However, MovieLens is an exception. To achieve good performance, the model requires an extreme temperature coefficient. This suggests that at this point, the model needs to rely on a stronger force to differentiate the semantic information between different input instances. In summary, we recommend fine-tuning the temperature coefficient of the contrastive loss within the range of 0.1 to 1.0 for optimal performance.

4.6 Denominator-only InfoNCE vs InfoNCE

4.6.1 Role of Denominator-only InfoNCE. To make InfoNCE more suitable for feature interaction-based CTR tasks, we make modifications to it. By simplifying and deriving the formula for InfoNCE

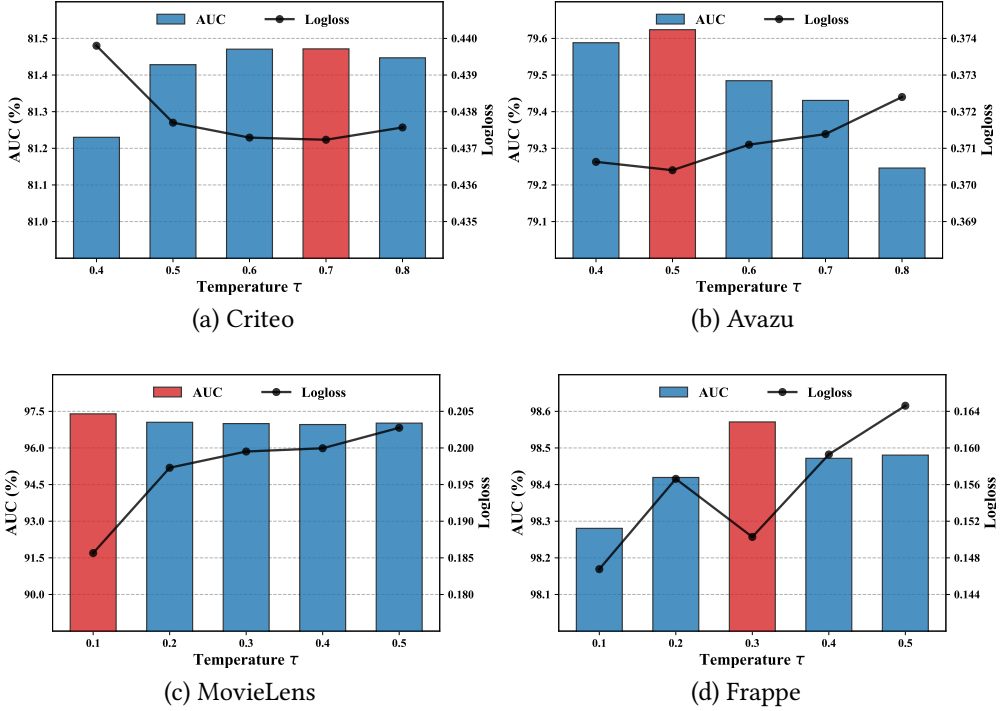


Fig. 9. Performance comparison of different temperature coefficients of CETN on Criteo (a), Avazu (b), MovieLens (c), and Frappe (d) datasets.

(Eq. 7), we can transform it into the following form:

$$\begin{aligned} & \sum_{i \in \mathcal{B}} -\log \frac{\exp(\text{sim}(\mathbf{V}'_i, \mathbf{V}''_i)/\tau)}{\sum_{j \in \mathcal{B}} \exp(\text{sim}(\mathbf{V}'_i, \mathbf{V}''_j)/\tau)}, \\ \Rightarrow & \sum_{i \in \mathcal{B}} -\log \frac{\exp(1/\tau)}{\sum_{j \in \mathcal{B}} \exp(\text{sim}(\mathbf{V}'_i, \mathbf{V}''_j)/\tau)}, \end{aligned} \quad (23)$$

$$\Rightarrow \sum_{i \in \mathcal{B}} \log \left(\exp(\text{sim}(\mathbf{V}'_i, \mathbf{V}''_i)/\tau) + \sum_{j \in \mathcal{B}/\{i\}} \exp(\text{sim}(\mathbf{V}'_i, \mathbf{V}''_j)/\tau) \right), \quad (24)$$

in Equation (23), we discarded alignment, retaining only uniformity, thereby obtaining Denominator-only InfoNCE. In Equation (24), we disregarded $1/\tau$, it becomes apparent that when the model optimizes this loss, it ensures that the model acquires dissimilar information across different semantic spaces, even if \mathbf{V}_i and \mathbf{V}_j originate from the same input.

In some studies [56, 62, 63], this uniformity is interpreted as minimizing similarity. This perspective is especially prevalent in the context of graph neural networks, where the aim is to disperse nodes away from dense clusters in the representation space, leading to a more uniformly distributed representation. In fact, looking at this from the perspective of information capture, the reason for the improvement in model performance due to this uniform distribution can be simply attributed

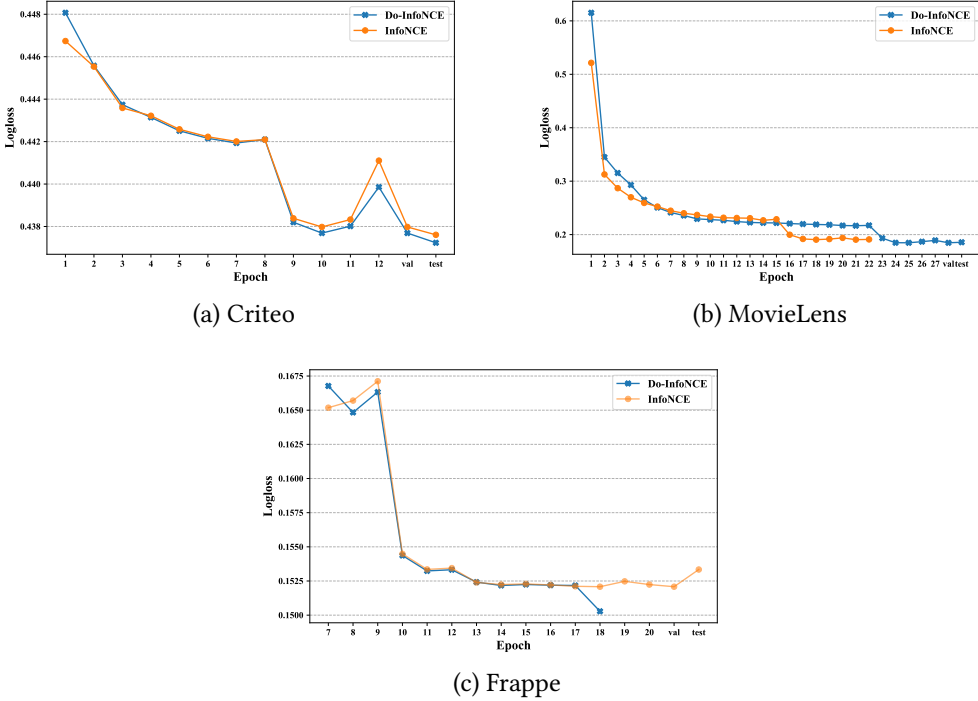


Fig. 10. Optimization Process of CETN on Criteo (a), MovieLens (b), and Frappe (c) datasets.

to the model capturing more diverse information. In other words, we ensure the diversity of the captured information through the InfoNCE loss function.

4.6.2 Optimization Process. To further investigate the impact of Do-InfoNCE and InfoNCE on the model's loss optimization process, we visualize the model's training process, and the results are presented in Figure 10. On the Criteo dataset, while InfoNCE initially yields better results in the first epoch, after lowering the learning rate with the Adam optimizer, the learning process of the model exhibits an issue which is insufficient capacity to capture information. Therefore, the optimization performance of InfoNCE is less favorable. On the MovieLens dataset, with the early-stop patience set uniformly to two, InfoNCE even prematurely terminates the training process, so it obtains sub-optimal results. In the smaller Frappe dataset, Do-InfoNCE demonstrates superior performance. It stops training at the 16th epoch and achieves better results than InfoNCE in the test set (at the 18th epoch). In contrast, InfoNCE not only slows down the model's convergence but also leads to serious overfitting issues. Looking at the overall picture, our proposed Do-InfoNCE achieves better performance, which can aid the model in finding local optimum more effectively.

4.6.3 Visualization of information within semantic spaces. To further explore the impact of self-supervisory signals on capturing feature interaction information within various semantic spaces, we randomly sampled 1,000 instances and visualized the information captured by the model, as shown in Figure 11. In Figure 11 (a) and (b), the representation distribution is more dispersed due to the influence of \mathcal{L}_{cl} , indicating that the model has captured richer and more diverse information. On the other hand, in the case of TriDNN, the feature interaction information captured by the

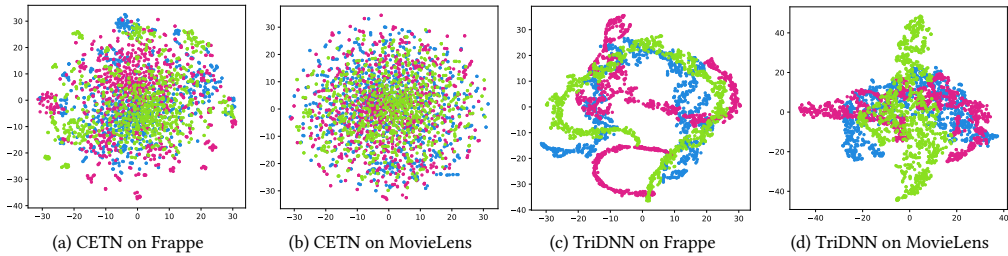


Fig. 11. The visualization of information within semantic spaces on Frappe and MovieLens datasets. The three colors in the plot represent three semantic spaces. TriDNN denotes the use of three independent MLPs as subcomponents without employing self-supervisory signals.

model across the three semantic spaces tends to be more similar. Consequently, in Figure 11 (c) and (d), the representation distribution is more concentrated, suggesting that the model's acquired information is narrower, thus reducing its effectiveness.

4.7 Through Network vs Residual Network

4.7.1 The Widespread Phenomenon of Shallow Networks in Recommender Systems. In the field of computer vision, neural networks often tend to develop in a deeper direction [16, 48], but in the field of recommender systems, shallow networks are often sufficient to achieve the expected task objectives. For example, in graph neural network-based recommender systems [18, 63, 66], the number of layers in the graph neural network is often set to 3. Similarly, in the click-through rate prediction tasks based on feature interaction [34, 70], the number of layers in the MLP is often also set to 3. The primary reason researchers set it this way is due to the issues of over-smoothing and degradation. This leads to a rapid descent of the neural network as the number of layers increases, eventually causing the model to collapse. Moreover, this severe data sparsity issue can even lead to a unique phenomenon in the field of recommender systems known as the one-epoch phenomenon [67]. Therefore, in CTR prediction tasks, neural networks often tend to widen rather than deepen.

4.7.2 Reproduction of the Collapse Phenomenon in CTR prediction. To verify the degradation phenomenon of deep neural networks in the field of recommender systems, we conducted experiments on the MovieLens dataset using MLP with different numbers of layers. We also visualized the training process of the model, as shown in Figure 12. In order to intuitively observe the impact of the number of neural network layers on model training, we stopped using methods to prevent model overfitting, such as dropout and batch norm, and only retained L_2 normalization to prevent the model from experiencing the one-epoch phenomenon. From the experimental results, it can be observed that as the number of MLP layers gradually increased from 3 to 8 and then to 16, the performance of the model does not show a significant improvement, but it does consume more computational resources. Therefore, further deepening the depth of the neural network is not a good choice. Next, when we increased the number of layers in the MLP to 20, the model experienced a collapse, and the AUC performance was close to 0.5. Even if we further deepened the depth of the MLP, the situation did not improve.

4.7.3 Further Theoretical Analysis. In previous research represented by residual networks, they use vertical, element-wise addition skip connections to pass information from earlier layers to later layers in the neural network, preventing the model from experiencing the degradation phenomenon.

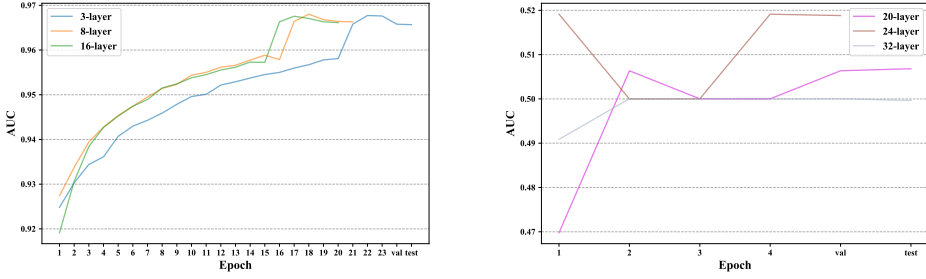


Fig. 12. Optimization process of plain DNN with different number of layers on MovieLens datasets.

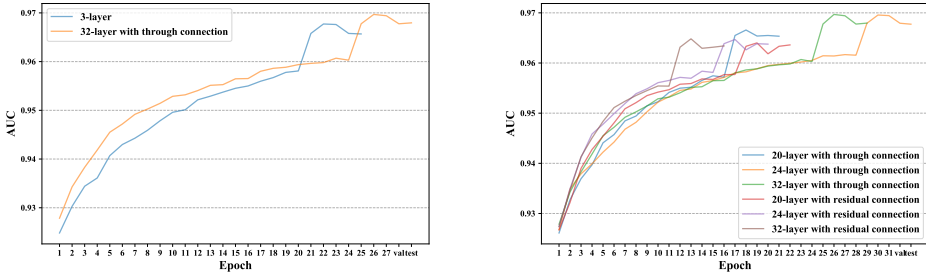


Fig. 13. Optimization Process on MovieLens datasets. Left: comparison between a 3-layer plain DNN and a 32-layer DNN using through connection. Right: comparison between the residual network and the through network.

This residual connection can be simply defined as:

$$y_{\text{deep}} = \mathcal{F}(\mathbf{x}) + \mathbf{x}, \quad (25)$$

in this formula, \mathbf{x} represents the input of a certain layer in the network, \mathcal{F} is the transformation applied to \mathbf{x} by the layer, and y is the output. The key idea here is that the output y is the sum of the input \mathbf{x} and its residual (i.e., the difference between \mathbf{x} and y). Next, by horizontally generalizing it, we arrive at the fundamental definition of through connections:

$$\begin{aligned} y_{\text{shallow}} &= \mathcal{F}_{\text{shallow}}(\mathbf{x}), \\ y_{\text{deep}} &= \mathcal{F}_{\text{deep}}(\mathbf{x}) + y_{\text{shallow}}, \end{aligned} \quad (26)$$

where we first obtain the output y_{shallow} through a shallow network, and then add the output of a deep network to the output of the shallow network element-wise to get y_{deep} . In fact, the residual connections can be regarded as a special case of the through connections. We can further redefine the residual connection to the following form:

$$\begin{aligned} y_{\text{shallow}} &= \mathcal{F}_{\text{shallow}}(\mathbf{x}), \\ y_{\text{deep}} &= \mathcal{F}_{\text{deep}}(y_{\text{shallow}}) + y_{\text{shallow}}, \end{aligned} \quad (27)$$

From Equation (27), we can see that the difference between residual connections and through connections is simply that the input to $\mathcal{F}_{\text{deep}}$ has changed from raw \mathbf{x} to y_{shallow} .

4.7.4 Validation of Effectiveness. By changing the skip connections in the residual network from vertical to horizontal, we can generalize to the Through Network, which can help the model overcome the model collapse phenomenon brought about by the deepening of network layers. To demonstrate this hypothesis, we combined a 32-layer DNN with a 3-layer DNN using through connections. The experimental results are shown in Figure 13 (left). What we can observe is that after using through connections, the 32-layer DNN not only avoids the model collapse phenomenon but also further improves the performance of the 3-layer DNN.

To compare the performance differences between through networks and residual networks in CTR prediction tasks, we conduct experiments using the layer settings where model collapse occurred. The results are shown in Figure 13. One observation is that the performance of the model using through connections exceeds that of residual connections. The residual connections perform best at 24 layers, while through connections perform best at 32 layers. This implies that through connections can better utilize deep neural networks to fit the target function on sparse datasets. Another point worth mentioning is that in our subsequent experiments, as the depth of the network continued to increase, the performance of the DNN model using the through connections continued to improve. This suggests that the through connections not only prevent the model from collapsing in deeper settings but also contribute to better performance as depth increases.

5 RELATED WORK

5.1 CTR Models Based on Feature Interactions

Most existing CTR models based on feature interaction predominantly follow the Embedding & Cross paradigm. They begin by encoding categorical and continuous features through embedding techniques. Subsequently, they employ a variety of complex cross-operations to augment the first-order feature data, achieving the goal of feature interaction, and ultimately enhancing model performance. MLP has played an indelible role in implicit feature interaction (Cross), greatly improving the benchmark performance of deep CTR models. However, many researchers have highlighted the inefficiency of MLP in learning product-based feature interactions (inner product, outer product, or Hadamard product) [43, 45]. Therefore, researchers have endeavored to employ additional feature data augmentation techniques to overcome the performance bottleneck of MLPs. Enhanced MLP models can primarily be categorized into two types: those enhanced based on stack structures and those enhanced based on parallel structures.

NFM [17], PNN [40], MaskNet [55], and xCrossNet [64] employ a stack structure to enhance MLP-based CTR models. They aim to introduce explicit product-based feature interaction operations before using embeddings as inputs to the MLP, thereby breaking through the performance bottleneck of the MLP. From a semantic space segmentation perspective, this explicit product operation further enriches the feature interaction information within the current semantic space, resulting in improved performance. Wide & Deep [5], DeepFM [14], DeepLight [8], FinalMLP [34], xDeepFM [30], and DCN [53] employ a parallel structure. These models aim to introduce explicit feature interactions in a parallel manner to the simple MLP model. They achieve this by incorporating a fusion layer to capture both explicit and implicit feature interaction information simultaneously. While this parallel strategy of capturing feature interaction information in different semantic spaces has yielded promising results, it still fails to address the three issues we have identified, ultimately resulting in sub-optimal performance.

5.2 Contrastive Learning for CTR Prediction

To the best of our knowledge, until now, there has been very limited work combining contrastive learning with CTR prediction tasks based on feature interactions. This can occur because users

often have a propensity for multiple interests in items, and it becomes challenging to definitively distinguish between positive and negative samples based solely on the user's click behavior. Therefore, traditional contrastive learning based on alignment and uniformity principles cannot be directly applied to CTR. With the rapid advancement of self-supervised learning in the fields of Natural Language Processing (NLP) [27, 33] and Computer Vision (CV) [4, 13, 19], due to the similarities between click-through rate prediction models based on user behavior sequences and NLP, they initially incorporate contrastive learning [21, 50, 65]. MISS [15] analyzes user behavior sequences and employs contrastive loss to enhance user interest representations at the feature level, thereby improving model performance. AQCL [37] attempts to alleviate the problem of learning difficulties in representing the user's click history feature sequences in cold-start scenarios by introducing AQCL loss. CL4CTR [52] introduced contrastive learning for the first time into feature interaction-based CTR prediction tasks. It proposes feature alignment and field uniformity for the feature field concept of CTR to enhance the quality of feature representation. However, it does not incorporate InfoNCE [35] loss and fails to address the issues of diversity and homogeneity from an architectural perspective.

6 CONCLUSION AND FUTURE WORK

In this paper, we revisited the problem of effectively capturing feature interaction information from multiple semantic spaces, and composed the Simple Multi-Head Network (simMHN) with multiple Key-Value Blocks as parallel subcomponents. We then further enhanced simMHN around the two complementary principles of diversity and homogeneity, thereby proposing a new simple and effective CTR model, called the Contrast-enhanced Through Network (CETN). CETN builds upon simMHN by leveraging contrastive learning and through connections to further capture high-quality feature interaction information. Experimental results on four benchmark datasets validate the effectiveness of our proposed model. As we look toward future work, we are interested in refining the architecture of this model, seeking opportunities to make it even more simple and efficient.

REFERENCES

- [1] Alexey Borisov, Ilya Markov, Maarten De Rijke, and Pavel Serdyukov. 2016. A neural click model for web search. In *Proceedings of the 25th International Conference on World Wide Web*. 531–541.
- [2] Erika Bourguignon and Lenora Greenbaum. 1973. Diversity and homogeneity in world societies. (1973).
- [3] Bo Chen, Yichao Wang, Zhirong Liu, Ruiming Tang, Wei Guo, Hongkun Zheng, Weiwei Yao, Muyu Zhang, and Xiuqiang He. 2021. Enhancing explicit and implicit feature interactions via information sharing for parallel deep CTR models. In *Proceedings of the 30th ACM International Conference on Information & Knowledge Management*. 3757–3766.
- [4] Ting Chen, Simon Kornblith, Mohammad Norouzi, and Geoffrey Hinton. 2020. A simple framework for contrastive learning of visual representations. In *International Conference on Machine Learning*. PMLR, 1597–1607.
- [5] Heng-Tze Cheng, Levent Koc, Jeremiah Harmsen, Tal Shaked, Tushar Chandra, Hrishi 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.
- [6] Weiyu Cheng, Yanyan Shen, and Linpeng Huang. 2020. Adaptive factorization network: Learning adaptive-order feature interactions. In *Proceedings of the AAAI Conference on Artificial Intelligence*, Vol. 34. 3609–3616.
- [7] 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.
- [8] Wei Deng, Junwei Pan, Tian Zhou, Deguang Kong, Aaron Flores, and Guang Lin. 2021. Deeplight: Deep lightweight feature interactions for accelerating CTR predictions in ad serving. In *Proceedings of the 14th ACM International Conference on Web Search and Data Mining*. 922–930.
- [9] Stefan Elfwing, Eiji Uchibe, and Kenji Doya. 2018. Sigmoid-weighted linear units for neural network function approximation in reinforcement learning. *Neural Networks* 107 (2018), 3–11.
- [10] Yufei Feng, Fuyu Lv, Weichen Shen, Menghan Wang, Fei Sun, Yu Zhu, and Keping Yang. 2019. Deep session interest network for click-through rate prediction. *arXiv preprint arXiv:1905.06482* (2019).

- [11] Jianfeng Gao, Wei Yuan, Xiao Li, Kefeng Deng, and Jian-Yun Nie. 2009. Smoothing clickthrough data for web search ranking. In *Proceedings of the 32nd International ACM SIGIR Conference on Research and Development in Information Retrieval*. 355–362.
- [12] Zhabiz Gharibshah and Xingquan Zhu. 2021. User response prediction in online advertising. *ACM Computing Surveys (CSUR)* 54, 3 (2021), 1–43.
- [13] Spyros Gidaris, Praveer Singh, and Nikos Komodakis. 2018. Unsupervised representation learning by predicting image rotations. *arXiv preprint arXiv:1803.07728* (2018).
- [14] Hui Feng 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 (Melbourne, Australia) (IJCAI'17)*. AAAI Press, 1725–1731.
- [15] Wei Guo, Can Zhang, Zhicheng He, Jiarui Qin, Hui Feng Guo, Bo Chen, Ruiming Tang, Xiuqiang He, and Rui Zhang. 2022. Miss: Multi-interest self-supervised learning framework for click-through rate prediction. In *2022 IEEE 38th International Conference on Data Engineering (ICDE)*. IEEE, 727–740.
- [16] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. 2016. Deep residual learning for image recognition. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*. 770–778.
- [17] Xiangnan He and Tat-Seng Chua. 2017. Neural factorization machines for sparse predictive analytics. In *Proceedings of the 40th International ACM SIGIR Conference on Research and Development in Information Retrieval*. 355–364.
- [18] Xiangnan He, Kuan Deng, Xiang Wang, Yan Li, Yongdong Zhang, and Meng Wang. 2020. Lightgcn: Simplifying and powering graph convolution network for recommendation. In *Proceedings of the 43rd International ACM SIGIR conference on research and development in Information Retrieval*. 639–648.
- [19] Olivier J Hénaff, Skanda Koppula, Jean-Baptiste Alayrac, Aaron Van den Oord, Oriol Vinyals, and Joao Carreira. 2021. Efficient visual pretraining with contrastive detection. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*. 10086–10096.
- [20] HuaWei. 2021. An open-source CTR prediction library. <https://fuxictr.github.io>.
- [21] Mengyuan Jing, Yanmin Zhu, Tianzi Zang, and Ke Wang. 2023. Contrastive Self-Supervised Learning in Recommender Systems: A Survey. *ACM Transactions on Information Systems (TOIS)* 42, 2, Article 59 (nov 2023), 39 pages. <https://doi.org/10.1145/3627158>
- [22] Yuchin Juan, Yong Zhuang, Wei-Sheng Chin, and Chih-Jen Lin. 2016. Field-aware factorization machines for CTR prediction. In *Proceedings of the 10th ACM Conference on Recommender Systems*. 43–50.
- [23] Martin Kaloev and Georgi Krastev. 2021. Comparative analysis of activation functions used in the hidden layers of deep neural networks. In *2021 3rd International Congress on Human-Computer Interaction, Optimization and Robotic Applications (HORA)*. IEEE, 1–5.
- [24] Farhan Khawar, Xu Hang, Ruiming Tang, Bin Liu, Zhenguo Li, and Xiuqiang He. 2020. Autofeature: Searching for feature interactions and their architectures for click-through rate prediction. In *Proceedings of the 29th ACM International Conference on Information & Knowledge Management*. 625–634.
- [25] Diederik P Kingma and Jimmy Ba. 2014. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980* (2014).
- [26] Annelies Knoppers, Inge Claringbould, and Marianne Dortants. 2015. Discursive managerial practices of diversity and homogeneity. *Journal of Gender Studies* 24, 3 (2015), 259–274.
- [27] Zhenzhong Lan, Mingda Chen, Sebastian Goodman, Kevin Gimpel, Piyush Sharma, and Radu Soricut. 2019. Albert: A lite bert for self-supervised learning of language representations. *arXiv preprint arXiv:1909.11942* (2019).
- [28] Zekun Li, Zeyu Cui, Shu Wu, Xiaoyu Zhang, and Liang Wang. 2019. FiGNN: Modeling feature interactions via graph neural networks for ctr prediction. In *Proceedings of the 28th ACM International Conference on Information and Knowledge Management*. 539–548.
- [29] Zekun Li, Shu Wu, Zeyu Cui, and Xiaoyu Zhang. 2022. GraphFM: Graph factorization machines for feature interaction modeling. *arXiv preprint arXiv:2105.11866* (2022).
- [30] 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*. 1754–1763.
- [31] Patricia W Linville. 1998. The heterogeneity of homogeneity. (1998).
- [32] Bin Liu, Chenxu Zhu, Guilin Li, Weinan Zhang, Jincai Lai, Ruiming Tang, Xiuqiang He, Zhenguo Li, and Yong Yu. 2020. AutoFIS: Automatic feature interaction selection in factorization models for click-through rate prediction. In *Proceedings of the 26th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining*. 2636–2645.
- [33] Yinhan Liu, Myle Ott, Naman Goyal, Jingfei Du, Mandar Joshi, Danqi Chen, Omer Levy, Mike Lewis, Luke Zettlemoyer, and Veselin Stoyanov. 2019. Roberta: A robustly optimized bert pretraining approach. *arXiv preprint arXiv:1907.11692* (2019).

- [34] Kelong Mao, Jieming Zhu, Liangcai Su, Guohao Cai, Yuru Li, and Zhenhua Dong. 2023. FinalMLP: An Enhanced Two-Stream MLP Model for CTR Prediction. *Proceedings of the AAAI Conference on Artificial Intelligence*, 37(4), 4552–4560. (2023).
- [35] Aaron van den Oord, Yazhe Li, and Oriol Vinyals. 2018. Representation learning with contrastive predictive coding. *arXiv preprint arXiv:1807.03748* (2018).
- [36] Junwei Pan, Jian Xu, Alfonso Lobos Ruiz, Wenliang Zhao, Shengjun Pan, Yu Sun, and Quan Lu. 2018. Field-weighted factorization machines for click-through rate prediction in display advertising. In *Proceedings of the 2018 World Wide Web Conference*. 1349–1357.
- [37] Yujie Pan, Jiangchao Yao, Bo Han, Kunyang Jia, Ya Zhang, and Hongxia Yang. 2021. Click-through rate prediction with auto-quantized contrastive learning. *arXiv preprint arXiv:2109.13921* (2021).
- [38] Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan, Trevor Killeen, Zeming Lin, Natalia Gimelshein, Luca Antiga, et al. 2019. Pytorch: An imperative style, high-performance deep learning library. *Advances in Neural Information Processing Systems* 32 (2019).
- [39] Katherine W Phillips and Robert B Lount. 2007. The affective consequences of diversity and homogeneity in groups. In *Affect and Groups*. Vol. 10. Emerald Group Publishing Limited, 1–20.
- [40] Yanru Qu, Han Cai, Kan Ren, Weinan Zhang, Yong Yu, Ying Wen, and Jun Wang. 2016. Product-based neural networks for user response prediction. In *2016 IEEE 16th International Conference on Data Mining (ICDM)*. IEEE, 1149–1154.
- [41] Yanru Qu, Bohui Fang, Weinan Zhang, Ruiming Tang, Minzhe Niu, Huifeng Guo, Yong Yu, and Xiuqiang He. 2018. Product-based neural networks for user response prediction over multi-field categorical data. *ACM Transactions on Information Systems (TOIS)* 37, 1 (2018), 1–35.
- [42] Steffen Rendle. 2010. Factorization machines. In *2010 IEEE International Conference on Data Mining*. IEEE, 995–1000.
- [43] Steffen Rendle, Walid Krichene, Li Zhang, and John Anderson. 2020. Neural collaborative filtering vs. matrix factorization revisited. In *Proceedings of the 14th ACM Conference on Recommender Systems*. 240–248.
- [44] Matthew Richardson, Ewa Dominowska, and Robert Ragno. 2007. Predicting clicks: estimating the click-through rate for new ads. In *Proceedings of the 16th International Conference on World Wide Web*. 521–530.
- [45] Shai Shalev-Shwartz, Ohad Shamir, and Shaked Shammah. 2017. Failures of gradient-based deep learning. In *International Conference on Machine Learning*. PMLR, 3067–3075.
- [46] Yanyan Shen, Lifan Zhao, Weiyu Cheng, Zibin Zhang, Wenwen Zhou, and Lin Kangyi. 2023. RESUS: Warm-Up Cold Users via Meta-Learning Residual User Preferences in CTR Prediction. *ACM Transactions on Information Systems* 41, 3 (2023), 1–26.
- [47] Weiping Song, Chence Shi, Zhiping Xiao, Zhijian Duan, Yewen Xu, Ming Zhang, and Jian Tang. 2019. AutoInt: Automatic feature interaction learning via self-attentive neural networks. In *Proceedings of the 28th ACM International Conference on Information and Knowledge Management*. 1161–1170.
- [48] Christian Szegedy, Wei Liu, Yangqing Jia, Pierre Sermanet, Scott Reed, Dragomir Anguelov, Dumitru Erhan, Vincent Vanhoucke, and Andrew Rabinovich. 2015. Going deeper with convolutions. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*. 1–9.
- [49] Zhen Tian, Ting Bai, Wayne Xin Zhao, Ji-Rong Wen, and Zhao Cao. 2023. EulerNet: Adaptive Feature Interaction Learning via Euler’s Formula for CTR Prediction. In *Proceedings of the 46th International ACM SIGIR Conference on Research and Development in Information Retrieval*. 1376–1385.
- [50] Chenyang Wang, Weizhi Ma, Chong Chen, Min Zhang, Yiqun Liu, and Shaoping Ma. 2023. Sequential recommendation with multiple contrast signals. *ACM Transactions on Information Systems* 41, 1 (2023), 1–27.
- [51] Fangye Wang, Hansu Gu, Dongsheng Li, Tun Lu, Peng Zhang, and Ning Gu. 2023. Towards Deeper, Lighter and Interpretable Cross Network for CTR Prediction. In *Proceedings of the 32nd ACM International Conference on Information and Knowledge Management*. 2523–2533.
- [52] Fangye Wang, Yingxu Wang, Dongsheng Li, Hansu Gu, Tun Lu, Peng Zhang, and Ning Gu. 2023. CL4CTR: A Contrastive Learning Framework for CTR Prediction. In *Proceedings of the Sixteenth ACM International Conference on Web Search and Data Mining*. 805–813.
- [53] 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.
- [54] Ruoxi Wang, Rakesh Shivanna, Derek Cheng, Sagar Jain, Dong Lin, Lichan Hong, and Ed Chi. 2021. DCNv2: Improved deep & cross network and practical lessons for web-scale learning to rank systems. In *Proceedings of the Web Conference 2021*. 1785–1797.
- [55] 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).
- [56] Jiancan Wu, Xiang Wang, Fuli Feng, Xiangnan He, Liang Chen, Jianxun Lian, and Xing Xie. 2021. Self-supervised graph learning for recommendation. In *Proceedings of the 44th International ACM SIGIR Conference on Research and Development in Information Retrieval*. 726–735.

- [57] Jiancan Wu, Xiang Wang, Fuli Feng, Xiangnan He, Liang Chen, Jianxun Lian, and Xing Xie. 2021. Self-supervised graph learning for recommendation. In *Proceedings of the 44th International ACM SIGIR Conference on Research and Development in Information Retrieval*. 726–735.
- [58] Jun Xiao, Hao Ye, Xiangnan He, Hanwang Zhang, Fei Wu, and Tat-Seng Chua. 2017. Attentional factorization machines: learning the weight of feature interactions via attention networks. In *Proceedings of the 26th International Joint Conference on Artificial Intelligence*. 3119–3125.
- [59] Bing Xu, Naiyan Wang, Tianqi Chen, and Mu Li. 2015. Empirical evaluation of rectified activations in convolutional network. *arXiv preprint arXiv:1505.00853* (2015).
- [60] Yanwu Yang and Panyu Zhai. 2022. Click-through rate prediction in online advertising: A literature review. *Information Processing & Management* 59, 2 (2022), 102853.
- [61] Tiansheng Yao, Xinyang Yi, Derek Zhiyuan Cheng, Felix Yu, Ting Chen, Aditya Menon, Lichan Hong, Ed H Chi, Steve Tjoa, Jieqi Kang, et al. 2021. Self-supervised learning for large-scale item recommendations. In *Proceedings of the 30th ACM International Conference on Information & Knowledge Management*. 4321–4330.
- [62] Junliang Yu, Xin Xia, Tong Chen, Lizhen Cui, Nguyen Quoc Viet Hung, and Hongzhi Yin. 2023. XSimGCL: Towards Extremely Simple Graph Contrastive Learning for Recommendation. *IEEE Transactions on Knowledge and Data Engineering* (2023), 1–14. <https://doi.org/10.1109/TKDE.2023.3288135>
- [63] Junliang Yu, Hongzhi Yin, Xin Xia, Tong Chen, Lizhen Cui, and Quoc Viet Hung Nguyen. 2022. Are graph augmentations necessary? simple graph contrastive learning for recommendation. In *Proceedings of the 45th International ACM SIGIR Conference on Research and Development in Information Retrieval*. 1294–1303.
- [64] Runlong Yu, Yuyang Ye, Qi Liu, Zihan Wang, Chunfeng Yang, Yucheng Hu, and Enhong Chen. 2021. Xcrossnet: Feature structure-oriented learning for click-through rate prediction. In *Pacific-Asia Conference on Knowledge Discovery and Data Mining*. Springer, 436–447.
- [65] Tianzi Zang, Yanmin Zhu, Ruohan Zhang, Chunyang Wang, Ke Wang, and Jiadi Yu. 2023. Contrastive Multi-View Interest Learning for Cross-Domain Sequential Recommendation. *ACM Transactions on Information Systems* (nov 2023). <https://doi.org/10.1145/3632402> Just Accepted.
- [66] Yi Zhang, Yiwen Zhang, Dengcheng Yan, Shuiguang Deng, and Yun Yang. 2023. Revisiting graph-based recommender systems from the perspective of variational auto-encoder. *ACM Transactions on Information Systems* 41, 3 (2023), 1–28.
- [67] Zhao-Yu Zhang, Xiang-Rong Sheng, Yujing Zhang, Biye Jiang, Shuguang Han, Hongbo Deng, and Bo Zheng. 2022. Towards Understanding the Overfitting Phenomenon of Deep Click-Through Rate Models. In *Proceedings of the 31st ACM International Conference on Information & Knowledge Management*. 2671–2680.
- [68] 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.
- [69] Chenxu Zhu, Bo Chen, Weinan Zhang, Jincal Lai, Ruiming Tang, Xiuqiang He, Zhenguo Li, and Yong Yu. 2023. AIM: Automatic Interaction Machine for Click-Through Rate Prediction. *IEEE Transactions on Knowledge and Data Engineering* 35, 4 (2023), 3389–3403. <https://doi.org/10.1109/TKDE.2021.3134985>
- [70] Jieming Zhu, Jinyang Liu, Shuai Yang, Qi Zhang, and Xiuqiang He. 2021. Open benchmarking for click-through rate prediction. In *Proceedings of the 30th ACM International Conference on Information & Knowledge Management*. 2759–2769.

Received XXXXXXXXXX; revised XXXXXXXXXX; accepted XXXXXXXXXX