

Cross-domain Augmentation Networks for Click-Through Rate Prediction

Xu Chen, Zida Cheng, Shuai Xiao✉, Xiaoyi Zeng, Weilin Huang

Abstract—Data sparsity is an important issue for click-through rate (CTR) prediction, particularly when user-item interactions is too sparse to learn a reliable model. Recently, many works on cross-domain CTR (CDCTR) prediction have been developed in an effort to leverage meaningful data from a related domain. However, most existing CDCTR works have an impractical limitation that requires homogeneous inputs (*i.e.* shared feature fields) across domains, and CDCTR with heterogeneous inputs (*i.e.* varying feature fields) across domains has not been widely explored but is an urgent and important research problem. In this work, we propose a cross-domain augmentation network (CDAnet) being able to perform knowledge transfer between two domains with *heterogeneous inputs*. Specifically, CDAnet contains a designed translation network and an augmentation network which are trained sequentially. The translation network is able to compute features from two domains with heterogeneous inputs separately by designing two independent branches, and then learn meaningful cross-domain knowledge using a designed cross-supervised feature translator. Later the augmentation network encodes the learned cross-domain knowledge via feature translation performed in the latent space and fine-tune the model for final CTR prediction. Through extensive experiments on two public benchmarks and one industrial production dataset, we show CDAnet can learn meaningful translated features and largely improve the performance of CTR prediction. CDAnet has been conducted online A/B test in image2product retrieval at Taobao app over 20 days, bringing an absolute **0.11 point** CTR improvement and a relative **1.26%** GMV increase.

Index Terms—feature translation, cross-domain CTR prediction, heterogeneous input features

1 INTRODUCTION

CLICK-THROUGH rate (CTR) prediction which estimates the probability of a user clicking on a candidate item, has a vital role in online services like recommendation, retrieval, and advertising [1], [2], [3]. For example, Taobao, as one of the largest e-commercial platforms in the world, has a list of application domains such as text2product retrieval and image2product retrieval. Each domain has its own CTR prediction model, which is termed as single-domain CTR prediction. It has been pointed out that data sparsity is a key issue that significantly limits the improvement of single-domain CTR models [4], and recent effort has been devoted to improving the CTR models by leveraging data from the other related domains [4], [5], [6], [7], [8].

Various cross-domain CTR prediction (CDCTR) methods have recently been developed, which can be roughly categorized into two groups: joint training and pre-training & fine-tuning. The joint training approach is developed by combining multiple CTR objectives from different domains into a single optimization process. It usually has shared network parameters to build connections, and transfer the learned knowledge across different domains like MiNet [5], DDTCDR [9] and STAR [8]. However, since two domains usually have different objectives in optimization, these methods usually suffer from an optimization conflict problem, which might lead to a negative transfer result [10], [11]. To deal with this issue, a number of recent approaches



(a) text2product retrieval (b) image2product retrieval

Figure 1: An example of text2product retrieval and image2product retrieval at Taobao app.

- Xu Chen, Zida Cheng, Shuai Xiao, Xiaoyi Zeng and Weilin Huang are with Alibaba Group. E-mail: {huaisong.cx, chengzida.cz, shuai.xsh, weilin.hwl}@alibaba-inc.com, yuanhan@taobao.com. ✉ indicates the corresponding author

have been proposed for jointly training multiple CTR objectives [12], [13], [14]. On the other hand, the pre-training & fine-tuning methods often have two stages, by training a CTR model sequentially in a source domain and then in a

target domain, where the performance in the target domain can be generally improved by leveraging model parameters pre-trained from the source domain. Notice that only one objective is utilized for optimization in each training stage, which significantly reduces the impact of negative transfer. Therefore, this method has been widely applied at Taobao platform [6], [7].

Importantly, most of recent CDCTR methods have been developed to explore additional cross-domain data with *homogeneous input features* [5], [6], [7], [8], [12], [13], which means the feature fields across domains are shared. For example, both MiNet [5] and MMOE [12] have a shared embedding layer of inputs to transfer knowledge across domains. Particularly, recent pre-training & fine-tuning methods presented in [6], [7] also require homogeneous inputs that enables them to apply the data from a target domain directly to the pre-trained source model, and thus is able to generate additional features for fine-tuning the target model. However, the requirement with homogeneous inputs might largely limit their applications in practice. Many cross-domain scenarios often exist *heterogeneous input features* which means two domains have varying feature fields. For instance, for image2product retrieval in Taobao, image is an important query information, while text plays as the key query information in text2product retrieval. An image retrieval system relies on image features a lot, while a text retrieval system focuses on text features heavily. At Taobao app, text2product retrieval has developed for many years and generated an order of magnitude more data than that of image2product retrieval. It is desirable to perform knowledge transfer that learns user behaviors from text2product domain, and then apply them for improving the performance on image2product retrieval.

This inspired us to develop a new CDCTR approach allowing for *heterogeneous input features*, which has not been widely studied yet but is challenging and urgent to be solved in industrial applications. A straightforward solution is to build on recent works [9], [12], [13], by replacing the shared-embedding layers with multiple domain-specific layers which process heterogeneous inputs separately, but it would inevitably break the key module designed specifically for transferring knowledge, resulting in unsatisfied performance. The up-to-date STAR [8] can work with heterogeneous inputs, but it mainly relies on parameter-sharing to transfer knowledge and simultaneously involves multiple CTR objectives in learning, which may cause ineffective knowledge transfer.

In this paper, we propose novel cross-domain augmentation networks (CDAnet) consisting of a translation network and an augmentation network, which are performed sequentially by first learning cross-domain knowledge and then implicitly encoding the learned knowledge into the target model via cross-domain augmentation. Specifically, a translation network is designed to process the inputs from different domains separately (which allows for heterogeneous input features), and knowledge translation is learned in the latent feature space via a designed cross-supervised feature translator. Then cross-domain augmentation is performed in the augmentation network by augmenting target domain samples in latent space with their additional translated latent features. This implicitly encodes the knowl-

edge learned from the source domain, providing diverse yet meaningful additional information for improving the fine-tuning on the target model. A learning comparison between CDAnet and other methods is shown in Figure 2. Through experiments, we demonstrate that CDAnet can largely improve the performance of CTR prediction, and it has been conducted online A/B test in image2product retrieval at Taobao app, bringing an absolute **0.11 point** CTR improvement, with a relative **1.26%** GMV increase. It has been successfully deployed online, serving billions of consumers. In a nutshell, the contributions of this work are summarized as follows:

- We identify a new cross-domain CTR prediction with *heterogeneous inputs*, which allows the target model to learn meaningful additional knowledge from a different domain. This addresses the issue of data sparsity efficiently on CTR prediction, and also set it apart from most existing cross-domain methods only allowing for homogeneous inputs.
- We propose cross-domain augmentation networks (CDAnet) consisting of a designed translation network and an augmentation network. The translation network is able to learn cross-domain knowledge from heterogeneous inputs by computing the features of two domain separately. Then it works as an efficient domain translator that encodes the learned knowledge implicitly in the latent space via feature augmentation, during the fine-tuning of target model.
- Extensive experiments on various datasets demonstrate the effectiveness of the proposed method. Our CDAnet has been deployed in image2product retrieval at Taobao app, and achieved obvious improvements on CTR and GMV. In addition, through empirical studies, we show that CDAnet can learn meaningful translated features for boosted CTR improvement.

2 RELATED WORK

2.1 CTR Prediction

Single-domain CTR: Click-through rate (CTR) prediction plays a vital role in various online services, such as modern search engines, recommendation systems and online advertising. Previous works usually combine logistic regression [15] and feature engineering for CTR prediction. These methods often lack the ability to model sophisticated feature interactions, and heavily rely on human labor of designing features. With the excellent feature learning ability of deep neural networks (DNN), deep learning approaches have been extensively studied on CTR prediction, and recent works focus on applying DNN for learning feature interactions, such as Wide&Deep [1], DeepFM [16] and DCN [17]. For example, in [1], Cheng *et al.* combined shallow linear models and deep non-linear networks to capture both low and high-order features, while the power of factorization machine [18] and deep networks are combined for CTR prediction in [16]. In DCN [17], Wang *et al.* designed a deep & cross network to learn bounded-degree feature interactions. Furthermore, deep models also demonstrate a strong capability for modeling richer information for CTR tasks. For example, DIN [2] and DIEN [19] were proposed to capture user interests based on historical click behaviors,

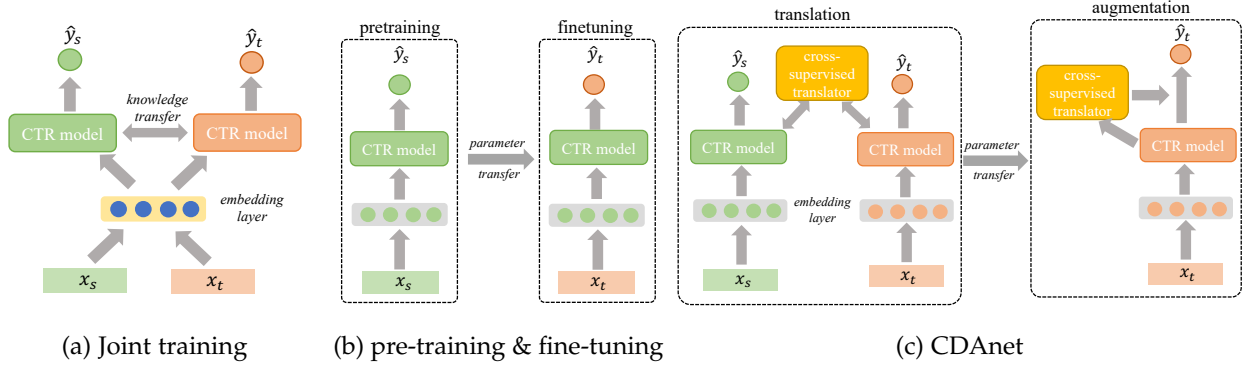


Figure 2: The comparison of joint training, pre-training & fine-tuning and our CDAnet. a) is the joint learning scheme where the embedding layer is usually shared and knowledge transfer techniques are employed between two CTR domains. b) shows the pre-training & fine-tuning learning style where a source domain CTR model will first be trained. Then the target domain model will load the pre-trained source domain model parameters and fine-tune itself with the target domain objective. c) shows the learning style of CDAnet. First, the translation network learns how to map the latent features from target domain to source domain by a cross-supervised translator. Then the augmentation network reuses the pre-trained parameters and employs the translated latent features as additional information to perform cross-domain augmentation for target domain CTR prediction.

and DSTN [3] takes the contextual ads when modeling user behaviors.

Cross-domain CTR: Single-domain CTR prediction suffers from a data sparsity issue because user behaviors in a real-world system are usually extremely sparse. Accordingly, cross-domain CTR prediction is developed to leverage user behaviors in a relevant but data-rich domain to facilitate learning in the current domain. A joint learning method was recently developed and has become a representative approach for cross-domain CTR. For example, STAR [8] is a star topology model that contains a centered network shared by different domains, with multiple domain-specific networks tailored for each domain. In MiNet [5], Ouyang *et al.* attempted to explore auxiliary data (e.g. historical user behaviors and ad title) from a source domain to improve the performance of a target domain. Meanwhile, a number of cross-domain recommendation (CDR) methods have been developed [9], [20], [21], [22], which can be naturally introduced to cross-domain CTR problems. For instance, a deep cross connection network is introduced in [20], and it is able to transfer user rating patterns across different domains. In DDTCDR [9], Li *et al.* proposed a deep dual transfer network that can bi-directionally transfer information across domains in an iterative style. However, such models have two different CTR objectives during joint learning, which might lead to the problem of gradient interference [10], [23], [24]. To handle the limitations, MMOE [12] was developed to learn an adaptive feature selection for different tasks, by using a shared mixture-of-expert model with task-specific gating networks. It explicitly models task relationships dynamically, allowing for automatically allocating model parameters which alleviates task conflicts in optimization. To decouple learning task-specific and task-shared information more explicitly, PLE [13] separates the network of task-shared components and task-specific components, and then adopts a progressive routing mechanism able to extract and separate deeper semantic knowledge gradually.

Apart from joint training, pre-training & fine-tuning is a widely-applied two-stage learning paradigm. In the pre-training stage, a model is first trained in a source domain. Then in fine-tuning stage, a target model would load the pre-trained model parameters, and then fine-tune itself for target domain CTR prediction. In each stage, only one objective is used for optimization, and thus the gradient interference issue can be alleviated to some extent. It has been shown that this method can be more efficient and effective for industrial systems [25]. Recently, Zhang *et al.* proposed KEEP [7] which is a two-stage framework that consists of a supervised pre-training knowledge extraction module performing on web-scale and long-time data, and a plug-in network that incorporates the extracted knowledge into the downstream fine-tuning model. In CTNet [6], Liu *et al.* focus on the CDCTR problem in a time-evolving scenario. In addition, a number of recent works such as [26], [27] investigate different layers of the network with various transfer manners (e.g. linear probing and fine-tuning), to perform knowledge transfer more efficiently.

Most of these CDCTR works focus on transferring knowledge between different user behavior distributions while ignoring the difference of input features, so they usually have a clear limitation on cross-domain homogeneous inputs [6], [7], [12], [13]. By contrast, the proposed CDAnet considers both the difference of user behavior distributions and input features in learning. It allows heterogeneous inputs and designs novel translation and augmentation networks to perform more effective knowledge transfer for this scenario.

2.2 Translation Learning

Translation is a task that translates the data from a source domain to a target domain while preserving the content information, which has two hot research topics: image2image translation [28] and neural machine translation [29]. As neural machine translation works on sequential data and

is not quite related to our work, we will mainly introduce image2image translation here. In earlier works of image2image translation, researchers usually use aligned image pairs to train the translation model. Pix2pix [30] is the first unified image2image translation framework based on conditional GANs [31]. However, the paired images are hard to be obtained and then the unpaired image2image translation networks are proposed [32], [33], [34], [35]. In unpaired image2image translation, some works propose the constraint of preserving the image properties of the source domain data such as semantic features [36] and class labels [37]. Another well-known and effective constraint is the cycle-consistency loss [33], [34], [35] in which if an input image is translated to the target domain and back, we can obtain the original input image. Despite the various works on image2image translation, an indispensable part among them is the cross translation part which translates the image of the source domain to the image in the target domain. Our CDAnet also has a similar cross-supervised translator. Whereas, the key difference is CDAnet does not aim to obtain the translated result in data space but to learn the mapped latent features to augment the CTR prediction in the target domain.

3 PROBLEM FORMULATION

In CDCTR prediction with heterogeneous input features, given source domain \mathcal{S} , we have its training samples (x_s, y_s) where $x_s \in \mathbb{R}^{F_s \times 1}$ denotes the input features and $y_s \in \{0, 1\}$ is the click label (i.e. 1 means click while 0 means non-click). Similarly, we have the target domain \mathcal{T} and its training samples (x_t, y_t) in which $y_t \in \{0, 1\}$ is the click label in target domain and $x_t \in \mathbb{R}^{F_t \times 1}$ is the input features. F_s and F_t denote the input feature dimension of source and target domain, respectively. We define x_s and x_t are heterogeneous when they have varying feature fields. For instance, the text query is important and widely used in text2product retrieval. While in image2product retrieval, we do not have the text query but have image query and the image feature is quite different from the text. When two domains have heterogeneous input features, many recent works cannot well handle this and it is challenging to bridge the heterogeneous gap for efficient knowledge transfer.

4 METHOD

4.1 Overview

The key of modeling in CDCTR with heterogeneous input features is to make sure that the model can embed the heterogeneous inputs and meanwhile have a good ability of transferring knowledge even without the widely used shared embedding layer technique appeared in recent works [5], [6], [12]. Accordingly, we propose our CDAnet which consists of two sequentially learned networks: translation network and augmentation network. First, the translation network embeds the heterogeneous input features with decoupled embedding layers and learns how the latent features are translated mutually by a designed cross-supervised translator. Then, the pre-trained parameters of translation network are transferred to the augmentation network. Next, the augmentation network will combine the

translated latent features and the original latent features of target domain samples together to fine-tune the target domain model. The model architecture is shown in Figure 3. Details about each module are demonstrated in the following parts.

4.2 Translation Network

Unlike image2image translation which aims to translate images in data space, our translation network learns to translate the latent features between domains. It has four parts including decoupled embedding layer, MMOE-based feature extractor, cross-supervised translator and prediction tower.

4.2.1 Decoupled Embedding Layer

Since the input features are heterogeneous, it is necessary to decouple the embedding layer of two domains so that the embedding layer can adapt to the input patterns in its own domain. Given the inputs $x_s \in \mathbb{R}^{F_s \times 1}$ of source domain and $x_t \in \mathbb{R}^{F_t \times 1}$ of the target domain, the embedding layer H_s and H_t , we have:

$$h_s = H_s(x_s), h_t = H_t(x_t) \quad (1)$$

where $h_s \in \mathbb{R}^{d \times 1}$ and $h_t \in \mathbb{R}^{d \times 1}$ are the embedding features of x_s and x_t , respectively. d is the latent feature dimension. H_s and H_t are domain-specific embedding layers.

4.2.2 MMOE-based Feature Extractor

In translation network, the two CTR objectives of source domain and target domain are jointly optimized. They have different optimization directions because of the different data distributions. Ignoring this may cause the optimization conflict problem and lead to negative transfer. Here, we place a shared MMOE-based network to extract useful features for each domain so that each domain can automatically have its own parameters. We denote F_i as the i -th expert with L -layer non-linear MLP, then we have:

$$f_s^i = F_i(h_s), f_t^i = F_i(h_t) \quad (2)$$

where $f_s^i \in \mathbb{R}^{d \times 1}$ and $f_t^i \in \mathbb{R}^{d \times 1}$ are the output of the i -th expert in source domain and target domain, respectively. Let G_s and G_t be the source domain gate and target domain gate that are non-linear MLP of L layers with $\mathbb{R}^d \rightarrow \mathbb{R}^d$, then we have:

$$g_s = \text{softmax}(W_s^{\text{gate}} G_s(h_s)), g_t = \text{softmax}(W_t^{\text{gate}} G_t(h_t)) \quad (3)$$

where $W_s^{\text{gate}} \in \mathbb{R}^{K \times d}$, $W_t^{\text{gate}} \in \mathbb{R}^{K \times d}$ are the affine transformation and K means the number of experts. $g_s \in \mathbb{R}^{K \times 1}$ and $g_t \in \mathbb{R}^{K \times 1}$ are the outputs of source domain gate and target domain gate. With Eq. 3, the source and target CTR objective can automatically select its own parameters for optimization. To this end, we can obtain the latent features of two domains as $z_s \in \mathbb{R}^{d \times 1}$ and $z_t \in \mathbb{R}^{d \times 1}$:

$$z_s = \sum_{i=1}^K g_s^i f_s^i, z_t = \sum_{i=1}^K g_t^i f_t^i \quad (4)$$

With this MMOE-based feature extractor, on one hand, the parameter-sharing can help knowledge transfer. On the

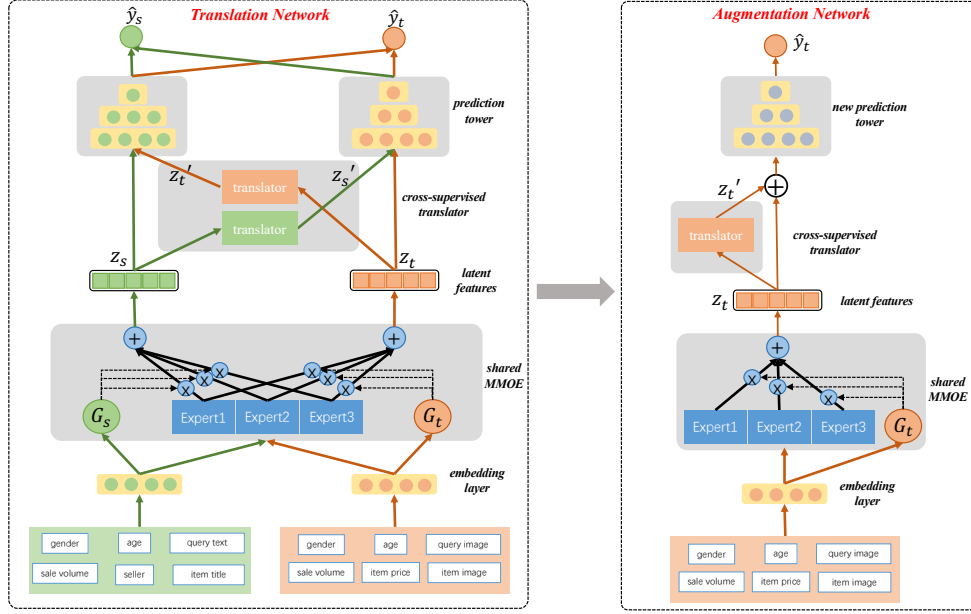


Figure 3: The architecture of our cross-domain augmentation networks (CDAnet). First, the translation network mainly focuses on learning the translator by cross supervision. Then the translation network parameters except the tower layer are transferred to the augmentation network. After this, the augmentation network takes the target domain samples and their translated latent features together to augment the fine-tuning of the target domain model.

other hand, its architecture can effectively alleviate the optimization conflict issue [12]. This MMOE-based network is placed just behind the embedding layer because we believe the features are richer and generalized in low layers while more compact and specialized in high levels. The former one has more transferable patterns that can benefit downstream tasks.

4.2.3 Cross-supervised Translator

After obtaining the latent features of each domain, we propose a cross-supervised translator to learn the latent feature translation, which is inspired by image2image translation. In image2image translation, images from target domain are translated as the images in the source domain, supervised by the true images in source domain. Whereas in CDCTR, the sample (e.g. combined user feature and item feature) in the target domain may not share the user behavior label with a sample in the source domain, because these two samples do not have paired relationship. Instead, a widely held belief in CDCTR is that if a user favors an item in the target domain, the favor behavior is preserved if the user and item are mapped into the source domain as corresponding features. For example, if a user likes science movies, he or she will also tend to love science novels.

In other words, given the latent feature z_t of target domain, we aim to translate it into source domain as z'_s which can preserve the semantics of z_t . Then, z'_s can be taken into the prediction tower of source domain while supervised by the target domain label y_t . Let $W_t^{tran} \in \mathbb{R}^{d \times d}$ be the translator of target domain and BCE be the binary cross entropy loss, then the cross-supervised translator of target domain is optimized by:

$$\min \mathcal{L}_t^{cross} = BCE(\sigma(R_s(z'_s)), y_t), \quad z'_s = W_t^{tran} z_t \quad (5)$$

where σ is the sigmoid function and R_s is the prediction tower of source domain. To stabilize the training, we have the symmetric formulation for source domain like Eq. 5 as:

$$\min \mathcal{L}_s^{cross} = BCE(\sigma(R_t(z'_t)), y_s), \quad z'_t = W_s^{tran} z_s \quad (6)$$

where R_t is the prediction tower of target domain and $W_s^{tran} \in \mathbb{R}^{d \times d}$ is the translator of source domain. Note that the network architecture of R_s and R_t could be different architectures according to the domain's own characteristics.

In order to better conduct the latent feature translation, apart from the above cross supervision, we add an orthogonal mapping constraint on the translators W_s^{tran} and W_t^{tran} as:

$$\min \mathcal{L}^{orth} = \|\mathbb{T}(W_s^{tran})(W_s^{tran} z_s) - z_s\|_F^2 + \|\mathbb{T}(W_t^{tran})(W_t^{tran} z_t) - z_t\|_F^2 \quad (7)$$

where \mathbb{T} is the transpose operation. The orthogonal transformation in mathematics has a characteristic that can preserve lengths and angles between vectors. Therefore, the orthogonal mapping constraint in Eq. 7 can help the latent features z_t of different samples preserve the similarity and avoid a case where multiple z_t collapse to a single point after translation.

4.2.4 Objective Function in Translation

Apart from the above objective of learning translator, we still need the vanilla objective for optimizing the CTR task in each domain. Namely, the vanilla CTR objectives of source and target domain are formulated as:

$$\min \mathcal{L}_s^{vani} = BCE(\sigma(R_s(z_s)), y_s), \mathcal{L}_t^{vani} = BCE(\sigma(R_t(z_t)), y_t) \quad (8)$$

These two objectives help the model learn the network parameters (e.g. the tower network) targeted for CTR prediction. Meanwhile, when the tower networks are optimized

for each domain’s CTR prediction, the translators can adapt to the tower networks and learn feature translation in a meaningful and right direction. To sum up, the objective function of translation network is:

$$\min \mathcal{L}_{trans} = \underbrace{\mathcal{L}_s^{vani} + \mathcal{L}_t^{vani}}_{\mathcal{L}^{vani}} + \alpha \underbrace{(\mathcal{L}_s^{cross} + \mathcal{L}_t^{cross})}_{\mathcal{L}^{cross}} + \beta \mathcal{L}^{orth} \quad (9)$$

where α and β are hyper-parameters on loss weights.

4.3 Augmentation Network

After the translation learning, we would first transfer the network parameters including the embedding layer, the shared MMOE module and the translator to the augmentation network. In order to give enough learning flexibility for augmentation network, as shown in Figure 3, the prediction tower is newly initialized rather than coming from the translation network.

4.3.1 Cross-domain Augmentation

Meanwhile, the translator learned in translation network can map the latent feature z_t into another space as z_t' and provide more useful information. That is to say, we can use the translated features of target domain samples as additional information to augment the target domain model training. The augmented latent feature of target domain is formulated as:

$$z_t^{aug} = z_t \oplus z_t' \quad (10)$$

where \oplus denotes the concatenation operation. When caring about the augmentation in source domain, the augmented feature can be obtained in a similar formula.

4.3.2 Objective Function in Augmentation

With the augmented latent feature z_t^{aug} , we would feed it into the new prediction tower R_t^{aug} and use the vanilla CTR objective for fine-tuning. The objective is defined as:

$$\mathcal{L}_{aug} = BCE(\sigma(R_t^{aug}(z_t^{aug})), y_t) \quad (11)$$

In this augmentation stage, the model has only one CTR objective and can avoid the optimization conflict problem in multi-objective models. When focusing on the performance of source domain, its augmentation network can be optimized in a similar way. The constraint in Eq. 7 also works here with the same β as in Eq. 9.

Considering the technique in knowledge transfer, the joint training works [8], [12], [13], [20] and fine-tuning methods [5], [6], [7] mainly rely on parameter-sharing, while our CDAnet proposes a novel translation then augmentation idea. In addition, CDAnet supports heterogeneous input features while most existing CDCTR models require homogeneous inputs.

5 EXPERIMENTS AND ANALYSIS

5.1 Experiment Setup

5.1.1 Datasets

We first conduct our experiments on two public benchmarks whose two domains have heterogeneous input fea-

Table 1: The statistics of datasets.

dataset	Amazon (movie and book)		Taobao (ad and rec)	
domain	movie	book	ad	rec
#users	29,680	52,690	141,917	186,731
#items	16,494	47,302	165,689	379,817
#input feature dim	8,044	24,466	44,897	5,004
#train samples	1,685,836	7,133,107	3,576,414	12,168,878
#val samples	210,730	891,638	447,052	1,521,110
#test samples	210,730	891,639	447,052	1,511,110
#positive samples	351,216	1,486,064	234,736	855,362

tures: Amazon¹ and Taobao²³. We choose the two largest domains—movie and book of Amazon to conduct the experiments. In the movie domain, we have the user ID, movie ID, movie genre, movie director and movie name information. While in the book domain, we have the user ID, book ID, book category, book writer and book name information, which are varying from those in the movie domain. Also, each domain has its own user behaviors. The original user behaviors are 0-5 ratings and we process the scores larger than 3 as positive feedback and other scores as negative feedback for CTR prediction. The user and item ID are both embedded as 64-dim features. Both movie and book name are processed as vectors by a word2vec Glove-6B model⁴. The movie genre, movie director, book category and book writer are mapped as one-hot features, respectively for each domain. Taobao dataset contains the user-item interactions of the advertisement (ad) and recommendation (rec) domain. We consider the buy behavior type as positive feedback and other types as negative feedback. In the ad domain, we get the user ID, age, gender, occupation and some other user profile information for the user side. As for the ad side, we have the ad ID, ad category and ad brand information. In the rec domain, we get user ID, item ID and item category information, which shows quite different features from the ad domain. For Taobao dataset, the user and item ID are both embedded as 64-dim features. Other discrete features are processed as one-hot or multi-hot features. For both datasets, we also apply a k -core filtering to guarantee each user or item has at least k interactions. k is 5 and 10 for movie and book domain of Amazon, respectively. For Taobao, k also equals 5 and 10 for ad and rec domain, respectively. The dataset statistics are summarized in Table 1.

Furthermore, since we address the real-world CDCTR with heterogeneous input features at Taobao, we also evaluate our model on Alibaba production data and online experiments on Taobao mobile app. In particular, we collect six-month user behaviors in text2product retrieval as the source domain data and one-year user behaviors in image2product retrieval as the target domain data. Both domains contain hundreds of billions of samples and rich user-side and item-side input features that are used in the online production system. The input features in these two domains are quite different due to the characteristic of these scenarios.

1. <https://jmcauley.ucsd.edu/data/amazon/>

2. <https://tianchi.aliyun.com/dataset/56>

3. <https://tianchi.aliyun.com/dataset/649>

4. <https://nlp.stanford.edu/projects/glove/>

5.1.2 Baselines

Existing pre-training & fine-tuning methods (e.g. Keep [7], CTNet [6]) require shared feature fields for knowledge transfer, so we do not include them as the baselines. Instead, we adopt various methods for comparison, including single domain and joint training methods:

- MLP. A deep multi-layer perception model is a common and efficient ranking model in online search and recommendation systems. Here, MLP serves as a single-domain CTR model for each domain.
- ShareMiddle. Considering the input features across domains are heterogeneous, we separate the embedding layers while sharing the middle layer just after the embedding layer of each domain inspired by the Share-Bottom [12] technique.
- STAR [8]. STAR is a star topology model that trains a single model to serve all domains by leveraging the data from all domains simultaneously.
- DDTCDR [9]. It is a deep dual transfer learning model that transfers knowledge between related domains in an iterative manner.
- MMOE [12]. MMOE implicitly models task relationships for multi-task learning by a shared mixture-of-experts module and task-specific gates. We utilize MMOE here for CDCTR with domain-specific embedding layers.
- PLE [13]. PLE is a multi-task learning model that separates shared components and task-specific components explicitly and adopts a progressive routing mechanism to extract and transfer knowledge.

5.1.3 Parameter Settings

In our experiments on public benchmarks, we split the data into train, validation, and test sets with the common 8:1:1 setting according to chronological order. The experiments are conducted multiple times and the mean value is taken as the model performance. We set the latent feature size as 64 for all models. The number of training epochs is set as 200 which can ensure the model’s convergence. To make a fair comparison, we conduct a grid search of hyper-parameters and the number of layers for all models. Considering CDAnet, for both domains of Amazon dataset, α is 0.01, β equals 0.1, the number of experts is 2 and each expert has 2 layers, the prediction tower has 2 layers including the logit mapping layer. For both domains of Taobao dataset, α is 0.03, β equals 0.1, the number of experts is 2, each expert has 2 layers and the prediction tower has 3 layers including the logit mapping layer. On our production dataset, the prediction tower of text2product domain is a DCN-v2 [38] model⁵ while that of our image2product domain is an online serving model. Based on the experience on public benchmarks, we also set the number of experts is 2, each expert has 2 layers, α and β are both 0.01 for our production dataset⁶.

5. We have no access to the online production model of text2product domain, so we empirically choose DCN-v2 as the prediction tower network.

6. We do not tune the hyper-parameters on this production dataset since it would cost too much computation resource in a limited time. Different hyper-parameters would be tuned in future experiments.

Table 2: The AUC comparison of different models on Amazon and Taobao. The results of both domains are listed.

Dataset	Amazon		Taobao	
Domain	Movie	Book	Ad	Rec
MLP	0.6595	0.7604	0.6161	0.6865
ShareMiddle	0.6604	0.7603	0.6160	0.5020
STAR	0.6503	0.7626	0.6149	0.6795
DDTCDR	0.6636	0.7807	0.6162	0.6756
MMOE	0.7025	0.7899	0.6177	0.6930
PLE	0.6993	0.7846	0.6164	0.6933
CDAnet	0.7225	0.7811	0.6200	0.7034

Table 3: The AUC comparison results of image2product domain on our production dataset. Here we consider text2product as the source domain and our image2product domain as the target domain.

Task	CTR
Base	0.7845
CDAnet	0.7866(+0.21%)

5.2 Overall Comparison

5.2.1 Offline Comparison

In this part, we show the comparison results on both domains of the datasets. By following popular works [6], [8], we take AUC as the evaluation metric. The results are shown in Table 2 and Table 3.

From Table 2, we observe that CDAnet generally can improve the CTR performance compared to various models. The joint training models (e.g. DDTCDR, MMOE, PLE) originally are not supportive for CDCTR with heterogeneous input features, while we revise them with domain-specific embedding layers so that they can work. Their inferior performance to CDAnet indicates that simply revising existing models with domain-specific embedding layers cannot well transfer the knowledge across domains. Although STAR can work with heterogeneous inputs, it mainly relies on a centered and domain-shared network to transfer knowledge, which is not so effective according to our results in Table 2. Instead, our translation then augmentation idea is superior and can provide better performance. In addition, considering the results in Table 3, we see CDAnet benefits CTR task even with so extremely large-scale data. It brings an absolute 0.21% AUC increase for CTR. A possible reason for this is that CDAnet can well transfer the knowledge of large-scale and high-quality data in text2product domain to image2product domain.

5.2.2 Online A/B test

We further conduct online experiments in an A/B test framework of image2product retrieval at Taobao app over 20 days. The baseline model is the online serving model trained on only image2product data. The online evaluation metrics are real CTR and GMV. The online A/B test shows that our CDAnet leads to an absolute 0.11 point CTR increase and a relative 1.26% GMV increase. In addition, we compare some online ranking examples between the baseline model and our CDAnet. The result is shown in Figure 4. In this figure, we give two examples to illustrate the better ranking ability of our CDAnet. In (a), if the user uploads a picture of Wangzai milk, we see that the online base model ranks a bottle of Wangzai candy that is quite

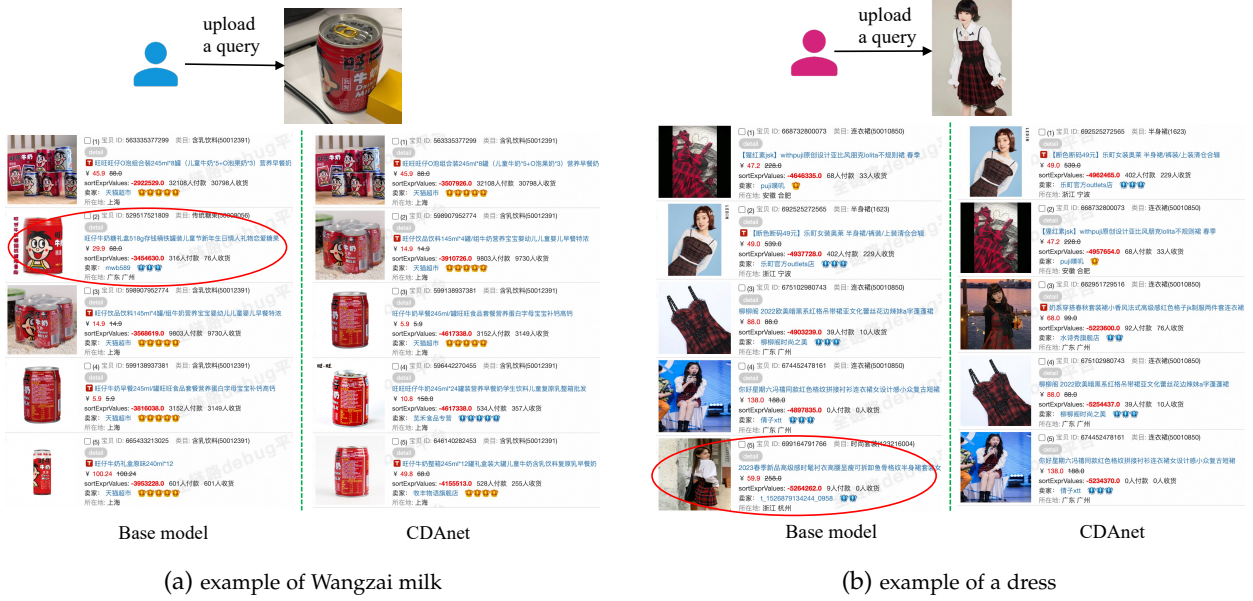


Figure 4: The ranking examples of the online base model and our proposed CDAnet. Given a query with the same user, we list the top5 ranking items of different models. (a) is a example of Wangzai milk and (b) is an example of a dress.

visually similar in position 2. Instead, our CDAnet can exclude this case and the top5 items are all about Wangzai milk. A possible reason may be that the query of milk and candy in text domain are quite different. The CTR model in text2product retrieval can well capture this while the CTR model in image2product domain may easily be confused by the visual patterns. Our CDAnet can transfer the knowledge of text2product domain into our image2product domain and improve the ranking results. In (b), we show a ranking example of a dress. Although our database does not have the the same item as the query, CDAnet has better ability of ranking similar items as the dress, while the online base model has a skirt in position 5. These examples can help us better understand the superiority of CDAnet.

5.2.3 Different Sparse Data

In order to investigate how CDAnet performs under more sparse conditions, we conduct an experiment to see the model performance at different sparsity levels of the training data. In particular, fixing the validation and test set, we vary the ratio of the original train data to the new train set.

The results are shown in Figure 5, where we can see that CDAnet can achieve better results at almost all sparsity levels compared to other cross-domain CTR methods. CDAnet has a more advanced knowledge transfer style that considers both the heterogeneous inputs and optimization conflict problem, so that is can show more generalized performance.

5.3 The Semantics of Translated Features

CDAnet learns how to translate the latent features between domains and exploits these translated latent features as additional information to boost the target domain CTR prediction, so it is curious to see whether the translated features are meaningful for target domain CTR prediction. In this part, we design an experiment to see whether the latent

features have related semantics before and after translation. To be specific on Amazon dataset, given a user that has training instances sharing positive labels in both domains⁷, we can obtain the latent feature matrix $Z_b \in \mathbb{R}^{N_b \times d}$ of book domain and $Z_m \in \mathbb{R}^{N_m \times d}$ of movie domain. Then, we can get the translated features of Z_m by the learned translator and it is denoted as $Z'_m \in \mathbb{R}^{N_m \times d}$. Next, for a row in the the translated feature Z'_m , we find its k -nearest neighbours in Z_b . In this case, we denote the book names of these neighbors as the corresponding information of the translated feature. Finally, for a row in Z'_m , we can check whether the book name of the neighbors in book space have related semantics with the original movie name in movie space. The results are summarized in Table 4.

From this Table, we find that the items of 5-nearest neighbours from book domain have close semantics to the corresponding item in the movie domain. For example, given UserId “21778”, the interacted movie is “8329” and the movie name is “The Addams Family” which tells an emotional and love story. Meanwhile, for the corresponding translated features, the items of 5-nearest neighbors from book domain are novels also about emotion and love. It shows that the translated features can capture related semantics from movie domain to book domain. This result matches a common belief in CDCTR that if a user likes an item in one domain, the user may also love the items sharing related semantics in other domains [6], [9], [22]. Notice that our model does not use any alignment information of items across domains, but it can automatically capture this and provide more useful information to boost the CTR

7. Notice our method does not require a user to have training instances in both domains. We make this requirement here just to better analyze the translated features. We let the instances in two domains share the positive label here for analysis because the “negative” instances in CTR are usually sampled and have less confidence than positives.

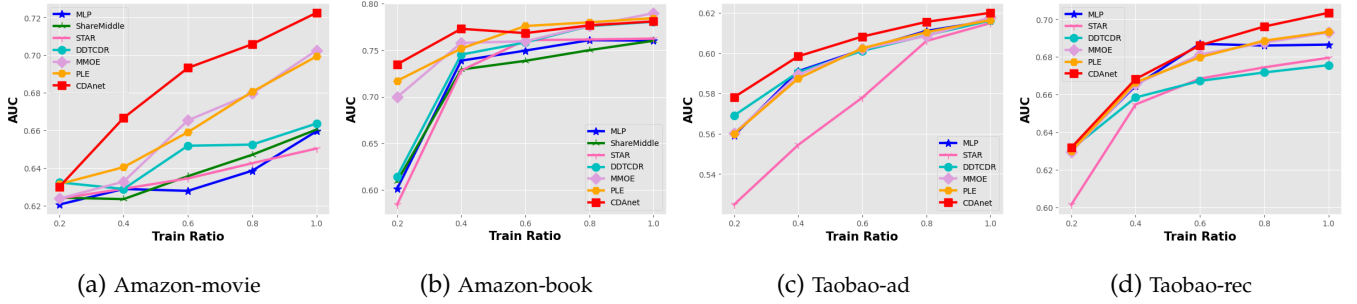


Figure 5: The effects of different sparsity levels on Amazon and Taobao. Train ratio means the ratio of the original train data.

Table 4: Results to show that the translated features in book space have related item semantics as its original item semantics in movie space. The title is the movie name or book name. Bolded book titles mean close semantics between the movie and book.

items in movie domain			items of 5-nearest neighbours from Book domain		
UserID	ItemID	Title	ItemId	Title	item genre
21778	8329	The Addams Family	39040	Where My Heart Breaks	emotion, love
			33024	The Program	
			29368	Lick (A Stage Dive Novel)	
			37926	Snowed Over	
			44092	Where the Stars Still Shine	
4842	13887	Mystery Science Theater 3000	10505	A Wild Sheep Chase: A Novel	mystery, thrillers
			7713	Wyrd Sisters	
			8205	Polar Star	
			135765	Harry Potter and the Chamber of Secrets	
			6374	The Ambler Warning	
2078	4977	Stargate SG-1 Season 2	43606	Countess So Shameless	fantasy, romance
			16339	Crimson City	
			35940	Not Your Ordinary Wolf Girl	
			10612	The Resisters	
			12671	Dagger-Star (Epic of Palins, Book 1)	

Table 5: The ablation study on different model parts of CDAnet. The figures are AUC results. “w/o” means without the corresponding model part. “w/o translation network” means we directly train the augmentation network which is random initialized. “w/o augmentation network” indicates we just use the trained translation network for evaluation.

Dataset	Amazon		Taobao	
Domain	movie	book	ad	rec
w/o MMOE	0.7186	0.7648	0.6199	0.7012
w/o \mathcal{L}^{orth}	0.7153	0.7734	0.6190	0.7022
w/o \mathcal{L}^{cross}	0.7188	0.7757	0.6191	0.7003
w/o translation network	0.6464	0.7649	0.6075	0.7001
w/o augmentation network	0.6818	0.7686	0.6168	0.6811
CDAnet	0.7225	0.7811	0.6200	0.7034

prediction performance. This result further demonstrates the value of our translation network.

5.4 Ablation Study

5.4.1 The Impacts of Different Model Parts

In order to assess the consequences of varying model components, we undertook an experiment to analyze the effect on performance upon the removal of the associated module. The results are shown in Table 5.

Comparing the result between CDAnet and w/o MMOE, we see the MMOE-based feature extractor has positive effects on boosting model performance. This MMOE module may help translation network alleviate optimization

conflict issue and better transfer knowledge. Considering the result of w/o \mathcal{L}^{orth} and w/o \mathcal{L}^{cross} , we see either removing the orthogonal constraint loss \mathcal{L}^{orth} or removing the cross-supervision loss \mathcal{L}^{cross} would cause deteriorated performance. The orthogonal constraint loss \mathcal{L}^{orth} can help CDAnet keep the similarity among latent features after translation and avoid mode collapse problem. The cross-supervision loss \mathcal{L}^{cross} plays a vital role in learning the two translators W_s^{trans} and W_t^{trans} . Without \mathcal{L}^{cross} , the translators cannot learn how to translate latent features into another space.

Further, either without translation network or without augmentation network would cause rather poor performance. The translation network learns how to transfer knowledge between domains. When removing it, we cannot have useful knowledge for later augmentation network. The augmentation network reuses the pre-trained parameters of translation network and employs the additional translated latent features for final CTR prediction goal. When removing the augmentation network, we only have the translation network whose goal is translation and cannot guarantee good CTR prediction performance.

5.4.2 Parameter Sensitivity

In CDAnet, α and β control the loss weight on cross-supervision loss and orthogonal constraint loss, respectively. K controls the number of experts in the shared MMOE-based feature extractor. Here, we study the parameter sen-



Figure 6: The effects of hyper-parameters on different datasets. α and β are loss weights on different loss parts. K denotes the number of experts.

sitivity of these hyper-parameters to investigate their effects on model performance. The results are shown in Figure 6.

In this figure, α controls the strength of cross-supervision on learning the translators. Too large α may dominate the learning of \mathcal{L}^{vani} and cause bad effects on learning cross-supervised translators. β controls the weight of \mathcal{L}^{orth} and too large β may limit the learning flexibility of the translators. K is the number of experts in the shared MMOE-based feature extractor and $K = 2$ shows the best performance. A possible reason may be too large K involves too many parameters and the data does not support the model’s complexity.

6 CONCLUSION AND FUTURE WORK

Cross-domain click-through rate (CDCTR) prediction with heterogeneous input features is an important and practical problem in real-world systems. In this paper, we propose a novel model named CDAnet that contains a translation network and augmentation network for effective knowledge transfer in CDCTR with heterogeneous input features. Through extensive experiments, we show CDAnet is able to learn meaningful translated latent features and boost the CTR prediction performance. Results on the large-scale production dataset and online system at Taobao app show its superiority in real-world applications.

Although CDAnet has achieved better performance than existing models, it still has inadequacies in transferring knowledge. The latent features of two domains are usually not fully overlapped, which means some ingredients of the latent features may not be translated and may have negative

effects in translation. In the future, we will study a more effective technique for latent feature translation and boost the target domain CTR prediction.

REFERENCES

- [1] H.-T. Cheng, L. Koc, J. Harmsen, T. Shaked, T. Chandra, H. Aradhye, G. Anderson, G. Corrado, W. Chai, M. Ispir *et al.*, “Wide & deep learning for recommender systems,” in *Proceedings of the 1st workshop on deep learning for recommender systems*, 2016, pp. 7–10.
- [2] G. Zhou, X. Zhu, C. Song, Y. Fan, H. Zhu, X. Ma, Y. Yan, J. Jin, H. Li, and K. Gai, “Deep interest network for click-through rate prediction,” in *Proceedings of the 24th ACM SIGKDD international conference on knowledge discovery & data mining*, 2018, pp. 1059–1068.
- [3] W. Ouyang, X. Zhang, L. Li, H. Zou, X. Xing, Z. Liu, and Y. Du, “Deep spatio-temporal neural networks for click-through rate prediction,” in *Proceedings of the 25th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining*, 2019, pp. 2078–2086.
- [4] C. Li, Y. Lu, Q. Mei, D. Wang, and S. Pandey, “Click-through prediction for advertising in twitter timeline,” in *Proceedings of the 21th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, 2015, pp. 1959–1968.
- [5] W. Ouyang, X. Zhang, L. Zhao, J. Luo, Y. Zhang, H. Zou, Z. Liu, and Y. Du, “Minet: Mixed interest network for cross-domain click-through rate prediction,” in *Proceedings of the 29th ACM international conference on information & knowledge management*, 2020, pp. 2669–2676.
- [6] L. Liu, Y. Wang, T. Wang, D. Guan, J. Wu, J. Chen, R. Xiao, W. Zhu, and F. Fang, “Continual transfer learning for cross-domain click-through rate prediction at taobao,” *arXiv preprint arXiv:2208.05728*, 2022.
- [7] Y. Zhang, Z. Chan, S. Xu, W. Bian, S. Han, H. Deng, and B. Zheng, “Keep: An industrial pre-training framework for online recommendation via knowledge extraction and plugging,” in *Proceedings of the 31st ACM International Conference on Information & Knowledge Management*, 2022, pp. 3684–3693.
- [8] X.-R. Sheng, L. Zhao, G. Zhou, X. Ding, B. Dai, Q. Luo, S. Yang, J. Lv, C. Zhang, H. Deng *et al.*, “One model to serve all: Star topology adaptive recommender for multi-domain ctr prediction,” in *Proceedings of the 30th ACM International Conference on Information & Knowledge Management*, 2021, pp. 4104–4113.
- [9] P. Li and A. Tuzhilin, “Ddtcd: Deep dual transfer cross domain recommendation,” *International conference on web search and data mining*, 2020.
- [10] O. Sener and V. Koltun, “Multi-task learning as multi-objective optimization,” *Advances in neural information processing systems*, vol. 31, 2018.
- [11] S. Vandenhende, S. Georgoulis, W. Van Gansbeke, M. Proesmans, D. Dai, and L. Van Gool, “Multi-task learning for dense prediction tasks: A survey,” *IEEE transactions on pattern analysis and machine intelligence*, 2021.
- [12] J. Ma, Z. Zhao, X. Yi, J. Chen, L. Hong, and E. H. Chi, “Modeling task relationships in multi-task learning with multi-gate mixture-of-experts,” in *Proceedings of the 24th ACM SIGKDD international conference on knowledge discovery & data mining*, 2018, pp. 1930–1939.
- [13] H. Tang, J. Liu, M. Zhao, and X. Gong, “Progressive layered extraction (ple): A novel multi-task learning (mtl) model for personalized recommendations,” in *Fourteenth ACM Conference on Recommender Systems*, 2020, pp. 269–278.
- [14] X. Lin, H. Chen, C. Pei, F. Sun, X. Xiao, H. Sun, Y. Zhang, W. Ou, and P. Jiang, “A pareto-efficient algorithm for multiple objective optimization in e-commerce recommendation,” in *Proceedings of the 13th ACM Conference on recommender systems*, 2019, pp. 20–28.
- [15] M. Richardson, E. Dominowska, and R. Ragno, “Predicting clicks: estimating the click-through rate for new ads,” in *Proceedings of the 16th international conference on World Wide Web*, 2007, pp. 521–530.
- [16] H. Guo, R. Tang, Y. Ye, Z. Li, and X. He, “Deepfm: A factorization-machine based neural network for ctr prediction,” in *Proceedings of the 26th International Joint Conference on Artificial Intelligence*, ser. IJCAI’17. AAAI Press, 2017, p. 1725–1731.
- [17] R. Wang, B. Fu, G. Fu, and M. Wang, “Deep & cross network for ad click predictions,” in *Proceedings of the ADKDD’17*, 2017, pp. 1–7.
- [18] S. Rendle, “Factorization machines,” in *2010 IEEE International conference on data mining*. IEEE, 2010, pp. 995–1000.

- [19] G. Zhou, N. Mou, Y. Fan, Q. Pi, W. Bian, C. Zhou, X. Zhu, and K. Gai, "Deep interest evolution network for click-through rate prediction," in *Proceedings of the AAAI conference on artificial intelligence*, vol. 33, no. 01, 2019, pp. 5941–5948.
- [20] G. Hu, Y. Zhang, and Q. Yang, "Conet: Collaborative cross networks for cross-domain recommendation," in *International conference on information and knowledge management*. ACM, 2018, pp. 667–676.
- [21] F. Yuan, L. Yao, and B. Benatallah, "Darec: Deep domain adaptation for cross-domain recommendation via transferring rating patterns," *International joint conference on artificial intelligence*, 2019.
- [22] X. Chen, Y. Zhang, I. W. Tsang, Y. Pan, and J. Su, "Towards equivalent transformation of user preferences in cross domain recommendation," *ACM Transactions on Information Systems (TOIS)*, 2020.
- [23] Z. Wang, Y. Tsvetkov, O. Firat, and Y. Cao, "Gradient vaccine: Investigating and improving multi-task optimization in massively multilingual models," *arXiv preprint arXiv:2010.05874*, 2020.
- [24] T. Yu, S. Kumar, A. Gupta, S. Levine, K. Hausman, and C. Finn, "Gradient surgery for multi-task learning," *Advances in Neural Information Processing Systems*, vol. 33, pp. 5824–5836, 2020.
- [25] L. Chen, F. Yuan, J. Yang, X. He, C. Li, and M. Yang, "User-specific adaptive fine-tuning for cross-domain recommendations," *IEEE Transactions on Knowledge and Data Engineering*, 2021.
- [26] U. Evci, V. Dumoulin, H. Larochelle, and M. C. Mozer, "Head2toe: Utilizing intermediate representations for better transfer learning," in *International Conference on Machine Learning*. PMLR, 2022, pp. 6009–6033.
- [27] Y.-L. Sung, J. Cho, and M. Bansal, "Lst: Ladder side-tuning for parameter and memory efficient transfer learning," *arXiv preprint arXiv:2206.06522*, 2022.
- [28] Y. Pang, J. Lin, T. Qin, and Z. Chen, "Image-to-image translation: Methods and applications," *IEEE Transactions on Multimedia*, 2021.
- [29] D. W. Otter, J. R. Medina, and J. K. Kalita, "A survey of the usages of deep learning for natural language processing," *IEEE transactions on neural networks and learning systems*, vol. 32, no. 2, pp. 604–624, 2020.
- [30] P. Isola, J.-Y. Zhu, T. Zhou, and A. A. Efros, "Image-to-image translation with conditional adversarial networks," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2017, pp. 1125–1134.
- [31] M. Mirza and S. Osindero, "Conditional generative adversarial nets," *arXiv preprint arXiv:1411.1784*, 2014.
- [32] M.-Y. Liu, T. Breuel, and J. Kautz, "Unsupervised image-to-image translation networks," *Advances in neural information processing systems*, vol. 30, 2017.
- [33] Z. Yi, H. Zhang, P. Tan, and M. Gong, "Dualgan: Unsupervised dual learning for image-to-image translation," in *Proceedings of the IEEE international conference on computer vision*, 2017, pp. 2849–2857.
- [34] J.-Y. Zhu, T. Park, P. Isola, and A. A. Efros, "Unpaired image-to-image translation using cycle-consistent adversarial networks," in *Proceedings of the IEEE international conference on computer vision*, 2017, pp. 2223–2232.
- [35] T. Kim, M. Cha, H. Kim, J. K. Lee, and J. Kim, "Learning to discover cross-domain relations with generative adversarial networks," in *International conference on machine learning*. PMLR, 2017, pp. 1857–1865.
- [36] Y. Taigman, A. Polyak, and L. Wolf, "Unsupervised cross-domain image generation," *arXiv preprint arXiv:1611.02200*, 2016.
- [37] K. Bousmalis, N. Silberman, D. Dohan, D. Erhan, and D. Krishnan, "Unsupervised pixel-level domain adaptation with generative adversarial networks," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2017, pp. 3722–3731.
- [38] R. Wang, R. Shivanna, D. Cheng, S. Jain, D. Lin, L. Hong, and E. Chi, "Dcn v2: Improved deep & cross network and practical lessons for web-scale learning to rank systems," in *Proceedings of the Web Conference 2021*, 2021, pp. 1785–1797.