

# Open Benchmarking for Click-Through Rate Prediction

Jieming Zhu\*  
Huawei Noah's Ark Lab  
Shenzhen, China  
jmzhu@ieee.org

Jinyang Liu\*  
The Chinese University of Hong Kong  
Hong Kong, China  
jyliu@cse.cuhk.edu.hk

Shuai Yang  
Peking University  
Beijing, China  
ethanyang@pku.edu.cn

Qi Zhang  
Huawei Noah's Ark Lab  
Beijing, China  
zhangqi193@huawei.com

Xiuqiang He  
Huawei Noah's Ark Lab  
Shenzhen, China  
hexiuqiang1@huawei.com

## ABSTRACT

Click-through rate (CTR) prediction is a critical task for many applications, as its accuracy has a direct impact on user experience and platform revenue. In recent years, CTR prediction has been widely studied in both academia and industry, resulting in a wide variety of CTR prediction models. Unfortunately, there is still a lack of standardized benchmarks and uniform evaluation protocols for CTR prediction research. This leads to non-reproducible or even inconsistent experimental results among existing studies, which largely limit the practical value and potential impact of their research. In this work, we aim to perform open benchmarking for CTR prediction and present a rigorous comparison of different models in a reproducible manner. To this end, we ran over 7,000 experiments for more than 12,000 GPU hours in total to re-evaluate 24 existing models on multiple dataset settings. Surprisingly, our experiments show that with sufficient hyper-parameter search and model tuning, many deep models have smaller differences than expected. The results also reveal that making real progress on the modeling of CTR prediction is indeed a very challenging research task. We believe that our benchmarking work could not only allow researchers to gauge the effectiveness of new models conveniently but also make them fairly compare with the state of the arts. We have publicly released the benchmarking tools, evaluation protocols, and experimental settings of our work to promote reproducible research in this field.

## CCS CONCEPTS

• **Information systems** → **Recommender systems**; **Computational advertising**.

## KEYWORDS

Recommender systems; sponsored search; computational advertising; CTR prediction; reproducible research; benchmarking

\*Both authors contributed equally to this work.

CIKM '21, November 1–5, 2021, Virtual Event, QLD, Australia

© 2021 Copyright held by the owner/author(s). Publication rights licensed to ACM. This is the author's version of the work. It is posted here for your personal use. Not for redistribution. The definitive Version of Record was published in *Proceedings of the 30th ACM International Conference on Information and Knowledge Management (CIKM '21)*, November 1–5, 2021, Virtual Event, QLD, Australia, <https://doi.org/10.1145/3459637.3482486>.

## ACM Reference Format:

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 and Knowledge Management (CIKM '21)*, November 1–5, 2021, Virtual Event, QLD, Australia. ACM, New York, NY, USA, 11 pages. <https://doi.org/10.1145/3459637.3482486>

## 1 INTRODUCTION

In many applications such as recommender systems, online advertising, and product search, click-through rate (CTR) is a key factor in business valuation. **CTR prediction, which aims to accurately predict the probability of a user clicking or interacting with a candidate item, is an important task in industry.** For applications with a large user base, even a small improvement in prediction accuracy can potentially lead to a large increase in the overall revenue. For example, existing studies from Google [8, 47] and Microsoft [28] reveal that an absolute improvement of 1% in logloss (or AUC) is considered as practically significant in real CTR prediction problems. **In contrast to other data types such as images and texts, CTR prediction problems usually involve data of large scale and high sparsity and comprise many categorical features of different fields** (e.g., billions of samples with millions of features in app recommendation of Google Play [8]). Therefore, it is a great challenge to make significant accuracy improvements in CTR prediction.

The importance and unique challenges of CTR prediction have attracted a lot of research attention from both academia and industry. CTR prediction models have evolved from simple logistic regression (LR) [33, 40], factorization machines (FM) [22, 39] and decision trees [6, 18], to deep neural networks (DNN) [8, 10]. Notably, numerous deep models have been proposed and shown remarkable performance gains in industrial CTR prediction problems, such as **Wide&Deep [8], DeepFM [16], DCN [47], xDeepFM [27], FiBiNET [21], DIN [56], and so on.**

Despite the success of these studies, there is still a lack of standardized benchmarks and uniform evaluation protocols for CTR prediction tasks. As a consequence, even though some common datasets (e.g., Criteo [1] and Avazu [2]) are used, existing studies often perform their own data partitions (e.g., with unknown train-test splitting or using unknown random seeds) and preprocessing steps (regarding how to cope with numeric features and how to filter rare categorical features). This leads to the non-reproducible and even inconsistent experimental results among these studies, since their nonstandard data preprocessing makes the results of any two different papers not comparable. Every work claims to

achieve the best results with a significant improvement on their own data partition, yet no one knows what it would be like if the same evaluation protocol had been applied for fair comparisons. Due to the lack of open benchmarking results for reference, readers may doubt whether the baseline models in a paper are implemented correctly and tuned with rigour, but none of these studies report the details or open the source code of their baseline implementations. In some cases, the official or third-party source code of some popular models is available (e.g., DeepCTR [4]), but we found that the training details about hyper-parameter settings, data loading, and early stopping are usually missing, which makes it difficult to reproduce the existing results. Such non-reproducibility and inconsistency issues largely limit the practical value and potential impact of the research in this field. Moreover, due to the lack of reusable and comparable benchmarking results in the literature, researchers need to re-implement all the baseline and re-evaluate them on their own data partition when publishing a new model. This is a tedious yet redundant effort, heavily increasing the burden of researchers to develop new models.

Inspired by the success of **ImageNet benchmark [12] in the CV domain and GLUE benchmark [46] in the NLP domain**, in this paper, we propose to perform open benchmarking for CTR prediction. Our work not only standardizes the open benchmarking pipeline for CTR prediction, but also makes a rigorous comparison of different models for reproducible research. To this end, we ran over 7,000 experiments for more than 12,000 GPU hours with a uniform setup to re-evaluate 24 existing models on multiple dataset settings of the two widely-used datasets, including Criteo [1] and Avazu [2]. Our experiments show somewhat surprising results. After sufficient hyper-parameter search and model tuning, many recent models have smaller differences than expected and sometimes are even inconsistent with what was reported in the literature. A similar study [11] also performed a re-evaluation of multiple representative papers on recommender systems, raising reproducibility issues of the results and concerns about lacking sufficient optimization of the baselines used for comparison. In contrast to this work, we take one step further to build an open CTR benchmark, which aims to set up a standardized benchmarking pipeline and to publicly release the most comprehensive benchmarking results as well as the reproducing steps to foster reproducible research.

We believe that our benchmarking work could be beneficial to multiple different groups of readers.

- **Researchers:** The benchmark could not only help researchers analyze strengths and bottlenecks of existing models, but also allow them to gauge the effectiveness of new models conveniently. Moreover, our benchmark demonstrates some good practices to fairly compare with the state of the arts for future research.
- **Practitioners:** The availability of benchmarking code and results can help industrial practitioners assess the applicability of new research models in their own problems and allow them to try new models with little efforts on their own datasets.
- **Competitors:** Leveraging our source code and hyper-parameters, competitors can easily implement high-performance baselines and ensembles in related competitions.

- **Beginners:** For beginners in this field, especially for students, our benchmarking code and the detailed reproducing steps can serve as a guidebook to learn the model implementations and model tuning tricks for CTR prediction. It is also valuable to apply our project for educational purposes.

In summary, our work makes the following main contributions:

- To the best of our knowledge, our work serves as the first step towards open benchmarking for CTR prediction.
- We release all the benchmarking code, settings, and results<sup>1</sup> to foster reproducible research on CTR prediction.
- Our results reveal the non-reproducibility and inconsistency issues in existing studies, and call for openness and rigour in future research on CTR prediction.

The remainder of this paper is organized as follows. Section 2 introduces the overview of CTR prediction. Section 3 describes our evaluation protocols and benchmarking results. Further, some discussions are made in Section 4 and the related work is reviewed in Section 5. Finally, we conclude the paper in Section 6.

## 2 CTR PREDICTION

In this section, we provide an overview of CTR prediction and then briefly review some of the representative models.

### 2.1 Overview

The objective of CTR prediction is to predict the probability that a user will click a given item. Yet, how to improve the accuracy of CTR prediction is a challenging research problem. **In contrast to other data types, such as images and texts, data in CTR prediction problems are typically in tabular format, comprising either numerical, categorical, or multi-valued (or sequence) features of multiple different fields.** The sample size is often large, yet the feature space is highly sparse. For example, app recommendation in Google Play [8] involves billions of samples with millions of features. In general, a CTR prediction model consists of the following key parts.

**2.1.1 Feature Embedding.** Input instances for CTR prediction generally contain three groups of features, *i.e.*, user profile, item profile, and context information. Each group has a number of fields as follows:

- *User profile:* age, gender, city, occupation, interests, etc.
- *Item profile:* item ID, category, tags, brand, seller, price, etc.
- *Context:* weekday, hour, position, slot id, etc.

**Features in each field may be categorical, numeric, or multi-valued** (e.g., multiple tags of a single item). Since most features are very sparse, leading to high-dimensional feature space after one-hot or multi-hot encoding, it is common to apply feature embedding to mapping these features into low-dimensional dense vectors. We summarize the embedding process of the three types of features in the following.

- *Categorical:* For a categorical feature field  $i$ , given a one-hot feature vector  $x_i$ , we have its embedding as  $e_i = V_i x_i$ , where the embedding matrix  $V_i \in \mathbb{R}^{d \times n}$  has vocabulary size  $n$  and embedding dimension  $d$ .

<sup>1</sup><https://openbenchmark.github.io/ctr-prediction>

- **Numeric:** For a numeric feature field  $j$ , there are multiple choices for feature embedding: 1) one can **bucketize** numeric values into discrete features, either by designing manually (e.g., grouping age 13~19 as teenager) or through training decision trees over numeric features (e.g., GBDT [19]), and then **embed them as categorical features**; 2) Given a *normalized* scalar value  $x_j$ , we set its embedding as  $e_j = v_j x_j$ , where  $v_j \in \mathbb{R}^d$  is the shared embedding vector of all features in field  $j$ ; 3) Instead of bucketizing each value into one bucket or assigning one single vector to each numeric field, one can apply AutoDis [15], **a numeric feature embedding** method, to dynamically bucketizing the numeric feature and computing the embedding from a meta embedding matrix.
- **Multi-valued:** For a multi-valued field  $h$ , each feature can be represented as a sequence. We obtain its embedding as  $e_h = V_h [x_{h1}, x_{h2}, \dots, x_{hk}] \in \mathbb{R}^{dxk}$ , given  $x_{hk}$  as a one-hot encoded vector of the sequence element and  $k$  denoting the maximal length of the sequence. Then the embedding  $e_h$  can be further aggregated to a  $d$ -dimensional vector, e.g., through mean/sum pooling. A further potential improvement is to apply sequential models, such as target attention in DIN [56] and GRU in DIEN [55], to aggregate the multi-valued behaviour sequence features.

**2.1.2 Feature Interaction.** It is straightforward to apply any classification model for CTR prediction after feature embedding. Nevertheless, for CTR prediction tasks, **interactions between features (a.k.a., feature conjunctions) are central to boost classification performance**. In factorization machines (FM) [39], inner products are shown as a simple yet effective way to capture pairwise feature interactions. Since the success of FM, a large body of research has been devoted to capturing interactions among features in different manners. Typical examples include inner product and outer product layers in PNN [37], Bi-interaction in NFM [17], cross network in DCN [47], compressed interaction in xDeepFM [27], convolution in FGCNN [29], circular convolution in HFM [45], bilinear interaction in FiBiNET [21], self-attention in AutoInt [42], graph neural network in FiGNN [26], hierarchical attention in InterHAt [25], just to name a few. Furthermore, **most current work investigates the way to combine both explicit and implicit features interactions with vanilla fully-connected networks** (i.e., MLPs).

**2.1.3 Loss Function.** The binary cross-entropy loss is widely used in CTR prediction tasks, which is defined as follows:

$$\mathcal{L} = -\frac{1}{N} \sum_{\mathbb{D}} \left( y \log \hat{y} + (1 - y) \log(1 - \hat{y}) \right), \quad (1)$$

where  $\mathbb{D}$  is the training set with  $N$  samples.  $y$  and  $\hat{y}$  denote the ground truth and the estimated click probability, respectively. We define  $\hat{y} = \sigma(\phi(x))$ , where  $\phi(x)$  represents the model function given input features  $x$  and  $\sigma(\cdot)$  is the sigmoid function to map  $\hat{y}$  to  $[0, 1]$ . The core of CTR prediction modeling lies in how to construct the model  $\phi(x)$  and learn its parameters from training data.

## 2.2 Representative Models

In this section, we summarize the representative models that we have evaluated and benchmarked in this work. Note that although

we enumerate only a part of the existing models above, they have covered a wide spectrum of studies on CTR prediction.

**2.2.1 Shallow Models.** It is common that industrial CTR prediction tasks have large-scale data. Therefore, shallow models have been in widespread use due to their simplicity and efficiency. Even today, LR [40] and FM [39] are still two strong baseline models deployed in industry. We describe the shallow models as follows:

- **LR.** Logistic regression (LR) is a simple baseline model for CTR prediction [40]. With the online learning algorithm, FTRL [33], proposed by Google, LR has been widely adopted in industry.
- **FM.** While LR fails to capture non-linear feature interactions, Rendle et al. propose factorization machine (FM) [39] that embeds features into dense vectors and models pairwise feature interactions as inner products of the corresponding embedding vectors. Notably, FM also has a linear time complexity in terms of the number of features.
- **FFM.** Field-aware factorization machine (FFM) [22] is an extension of FM that considers field information for feature interactions. It was a winner model in several Kaggle contests on CTR prediction.
- **HOFM.** Since FM only captures second-order feature interactions, HOFM [5] aims to extend FM to higher-order factorization machines. However, it results in exponential feature combinations that consume huge memory and take a long running time.
- **FwFM.** Recently, Pan et al. [35] extends FM by considering the field-wise weights of features interactions. Compared with FFM, it reports comparable performance but uses much fewer model parameters.
- **LorentzFM.** LorentzFM [50] has recently been proposed to embed features into a hyperbolic space and model feature interactions via triangle inequality of Lorentz distance.

**2.2.2 Deep Models.** Nowadays, deep neural networks have been widely studied and applied for CTR prediction. Compared with shallow models, deep models are more powerful in capturing sophisticated high-order feature interactions with non-linear activation functions, which usually yield better performance. Yet, efficiency has become the major bottleneck to scale deep models in practice.

- **DNN.** DNN is a straightforward deep model reported in [10], which applies a fully-connected network (termed DNN) after the concatenation of feature embeddings for CTR prediction.
- **CCPM.** CCPM [30] reports the first attempt to use convolution for CTR prediction, where feature embeddings are aggregated hierarchically through convolution networks.
- **Wide&Deep.** Wide&Deep [8] is a general learning framework proposed by Google that combines a wide (or shallow) network and deep network to achieve the advantages of both.
- **IPNN.** PNN [37] is a product-based network that feeds the inner (or outer) products of features embeddings as the input of DNN. Due to the huge memory requirement of pairwise outer products, we choose the inner product version, IPNN.
- **DeepCross.** Inspired by residual networks, [41] propose deep crossing to add residual connections between layers of DNNs.
- **NFM.** Similar to PNN, NFM [17] proposes a Bi-interaction layer that pools the pairwise feature interactions to a vector and then feed it to a DNN for CTR prediction.

**Table 1: Reproducibility requirements met by existing studies (✓ | - | × means totally | partially | not met, respectively)**

Reproducibility requirements	xDeepFM	FGCNN	AutoInt+	FiGNN	ONN	FiBiNET	LorentzFM	AFN+	InterHAt	Ours
Data preprocessing	×	×	✓	×	×	×	×	✓	✓	✓
Model source code	-	-	✓	✓	-	-	×	✓	✓	✓
Model hyper-parameters	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓
Baseline source code	×	×	×	×	×	×	×	×	×	✓
Baseline hyper-parameters	✓	✓	-	-	×	-	-	-	×	✓

- **AFM.** Instead of treating all feature interactions equally as in FM, AFM [49] learns the weights of feature interactions via attentional networks. Different from FwFM, AFM adjusts the weights dynamically according to the input data sample.
- **DeepFM.** DeepFM [16] is an extension of Wide&Deep that substitutes LR with FM to explicitly model second-order feature interactions.
- **DCN.** In DCN [47], a cross network is proposed to perform high-order feature interactions in an explicit way. In addition, it also integrates a DNN network following the Wide&Deep framework.
- **xDeepFM.** While high-order feature interactions modeled by DCN are bit-wise, xDeepFM [27] proposes to capture high-order feature interactions in a vector-wise way via a compressed interaction network (CIN).
- **HFM+.** HFM [45] proposes holographic representation and computes compressed outer products via circular convolution to model pairwise feature interactions. HFM+ further integrates a DNN network with HFM.
- **FGCNN.** FGCNN [29] applies convolution networks and recombination layers to generate additional combinatorial features to enrich existing feature representations.
- **AutoInt+.** AutoInt [42] leverages self-attention networks to learn high-order features interactions. AutoInt+ integrates AutoInt with a DNN network.
- **FiGNN.** FiGNN [26] leverages the message passing mechanism of graph neural networks to learn high-order features interactions.
- **ONN.** ONN (a.k.a., NFFM) [52] is a model built on FFM. It feeds the interaction outputs from FFM to a DNN network for CTR prediction.
- **FiBiNET.** FiBiNET [21] leverages squeeze-excitation network to capture important features, and proposes bilinear interactions to enhance feature interactions.
- **AFN+.** AFN [9] applies logarithmic transformation layers to learn adaptive-order feature interactions. AFN+ further integrates AFN with a DNN network.
- **InterHAt.** InterHAt [25] employs hierarchical attention networks to model high-order feature interactions in an efficient manner.

### 3 OPEN CTR BENCHMARK

In this section, we first summarize the key reproducibility requirements, and then present our evaluation protocols and open-source toolkit for benchmarking CTR prediction models. Finally, we report and analyze our benchmarking results.

#### 3.1 Reproducibility Requirements

CTR prediction has been widely studied in recent years. Unfortunately, there is a lack of rigour on reproducibility in many studies.

This may lead to non-reproducible and inconsistent results reported in existing studies, thereby hindering the development and evolution of new techniques in this field. In this work, we highlight the following **five key requirements to ensure reproducible research**. Yet, many current studies fail to meet all of these reproducibility requirements, as shown in Table 1. Note that we use ✓ | - | × to denote whether each requirement is totally, partially, or not met, respectively. It also indicates whether the artifacts are totally, partially, or not available for reproducing each step. For example, "-" indicates that FiBiNET has only unofficial model source code.

- **Data preprocessing:** Most work splits the training, validation, and testing data randomly, but others often cannot repeat the same data splitting due to the lack of scripts or random seeds used. Even more, some preprocessing details (e.g., how to handle numeric features and what threshold is used to filter rare categorical features) may be missing or incomplete. In such cases, researchers have to perform their own data splitting and preprocessing, leading to uncomparable results. It is worth noting that the authors of AutoInt [42], AFN [9], and InterHAt [25] have made a good starting point to share the data processing code or the preprocessed data.
- **Model source code:** In the spirit of open source, some studies have released the source code of their models on GitHub. Some popular models also have unofficial implementations available online from third-party libraries (e.g., DeepCTR [4]). But in many cases, we found that the source code is not ready for reproducibility studies, because it may either lack the training code (e.g., loading data, early stopping, etc.) or miss some key hyper-parameters on the given dataset.
- **Model hyper-parameters:** Most studies specify the detailed hyper-parameters of their own models in the papers. But without access to the original data preprocessed by the authors, it is inappropriate for others to use the same hyper-parameters on new data splits. It may only attain sub-optimal performance and needs to be retuned. Such practices lead to inconsistent results among existing papers.
- **Baseline source code:** Many studies report the details of their own models, but fail to clarify how they apply the baseline models. We noticed that existing studies seldom open source the baseline models, or tell which implementation is used for comparison. Model performance depends heavily on the quality of their code implementations. Bad implementations may introduce biases and make unfair model comparisons. However, this aspect is often overlooked by existing studies, making their performance improvements difficult to reproduce.
- **Baseline hyper-parameters:** It is desirable to exhaustively tune the hyper-parameters of baseline models to fairly compare model performance. Yet, this has not been guaranteed due to the lack of



**Table 2: Dataset statistics**

Dataset	#Instances	#Fields	#Features	Positive Ratio
Criteo	46M	39	5.55M	26%
Avazu	40M	24	8.37M	17%

open benchmarking. Most existing studies usually report inconsistent results of the same baseline models due to their unknown data preprocessing and baseline implementations.

To enable reproducible research and fairness of comparison, in this work, we aim to set up standardized evaluation protocols as well as provide the most comprehensive open benchmarking results for CTR prediction.

## 3.2 Evaluation Protocol

**3.2.1 Datasets.** In this work, we mainly use two real-world datasets for our evaluation: Criteo [1] and Avazu [2]. We provide more benchmarking results on other datasets in the Open-CTR-Benchmark website. Both of them are open datasets released by two leading ad companies, and have been widely used in the previous work (e.g., [16, 27, 42, 47]). We choose them because they are collected or sampled from real click logs in production, and both have tens of millions of samples, making the benchmarking results meaningful to industrial practitioners. Table 2 summarizes the data statistics information.

**3.2.2 Data splitting.** As with most of the existing studies, we follow them to randomly split both Criteo and Avazu into 8:1:1 as the training set, validation set, and test set, respectively. To make it exactly reproducible and easy to compare with existing work, we reuse the code provided by AutoInt [42] and control the random seed (i.e., seed=2018) for splitting. We mark these two data splits as *Criteo\_x4* and *Avazu\_x4*, respectively.

**3.2.3 Data preprocessing.** We mostly follow the steps in [42] to preprocess the features. However, we make some modifications that enable significant improvements and also fix a defect in [42].

- **Criteo.** The Criteo dataset consists of ad click data over a week. It comprises 26 categorical feature fields and 13 numerical feature fields. We create two different evaluation settings, denoted as *Criteo\_x4\_001* and *Criteo\_x4\_002*. Concretely, in *Criteo\_x4\_001*, instead of normalizing numeric values as in [42], we follow the winner solution of the Criteo contest to discretize each numeric value  $x$  to  $\lfloor \log^2(x) \rfloor$ , if  $x > 2$ ;  $x = 1$  otherwise, which yields much better performance. For categorical features, we replace infrequent features (min\_count=10) with a default "OOV" token. We also fix the feature embedding dimension to 16 in this setting. *Criteo\_x4\_002* differs in that we set min\_count=2 for categorical features, and the embedding dimension equals to 40 after tuning.
- **Avazu.** Avazu contains 10 days of click logs. It has a total of 23 fields with categorical features including app id, app category, device id, etc. We also create two evaluation settings, namely *Avazu\_x4\_001* and *Avazu\_x4\_002*. In *Avazu\_x4\_001*, we remove the id field that has a unique value in each data sample, which should be useless for CTR prediction. But it is retained in [42], resulting in a defect. In addition, we transform the timestamp field

into three new fields: hour, weekday, and is\_weekend. For all categorical features, we replace infrequent features (min\_count=2) with a default "OOV" token. We further fix the feature embedding dimension to 16 as in [42]. *Avazu\_x4\_002* is different in that we set min\_count=1 and fix the embedding dimension to 40. We emphasize that for *Avazu\_x4*, we set min\_count=1 or 2 after tuning because it shows much better performance than the original setting (min\_count=10) in [42], which is a counter-intuitive fact.

**3.2.4 Evaluation metrics.** We employ two most commonly-used metrics, AUC and logloss, for benchmarking.

- **AUC** (i.e., Area Under the ROC Curve) is a common metric to measure the probability that a randomly chosen positive sample is ranked higher than a randomly chosen negative sample. Higher AUC indicates better CTR prediction performance.
- **Logloss**, also known as the logistic loss or binary cross-entropy loss, is defined as Equation 1. Lower logloss indicates better CTR prediction performance.

It is worth noting that, when considering a large user base, an improvement of AUC at 0.001 level is generally considered as practically significant for an industrial CTR prediction task, as indicated by several existing studies from Google [8, 47], Microsoft [28], and Huawei [16]. In addition, although the metric gAUC [56] is also widely employed in practice, we cannot report the results because user\_id is missing in both Criteo and Avazu datasets.

**3.2.5 Toolkit Implementation of FuxiCTR.** While many open-source projects exist for CTR prediction, they mostly implement a few models in an ad-hoc manner and lack a complete workflow for benchmarking. In particular, DeepCTR [4] provides a good package with uniform implementations of many CTR prediction models. Nevertheless, our benchmarking needs a complete workflow, including data preprocessing, batch loading, model training, early stopping, learning rate decay, hyper-parameter search, and most importantly, seeding and logging for reproducibility.

In this work, we build the open-source FuxiCTR toolkit for benchmarking CTR prediction models, providing stunning features about configurability, tunability, and reproducibility. Since its inception in 2018, the project has evolved for over three years. The code of FuxiCTR consists of the following parts: 1) The data preprocessing part reads the raw data from CSV files, transforms all numeric, categorical, and sequential features, and outputs the transformed data into HDF5 files. 2) Tens of models have already been implemented in a uniform way in Pytorch. 3) The training part is implemented to read batches of data, compute forward and backward passes, and perform learning rate decay and early stopping if necessary. 4) The seeding and logging utilities are also specially designed for reproducibility, recording detailed running logs (including the used hyper-parameters) for each benchmarking experiment. 4) The hyper-parameter tuning part provides a configurable interface to allow grid search of hyper-parameters specified by users. We integrate all these parts as a complete benchmarking framework, allowing researchers to easily reuse our code, build new models, or adding new datasets. The goal of FuxiCTR is to provide an easy-to-use framework for future research on CTR prediction.

**3.2.6 Training details and hyper-parameters tuning.** During training, we apply the Reduce-LR-on-Plateau scheduler by default to

reduce learning rate by a factor of 10 when the given metric stops improving. To avoid overfitting, **early stopping is employed when the metric on the validation set stops improving in 2 or 3 consecutive epochs**. The default learning rate is  $1.e - 3$ . The batch size is initially set to 10000 and then decreases gradually using [5000, 2000, 1000] if an OOM error occurs in GPUs. We found that for CTR prediction models trained on millions of samples, using a large batch size usually makes the model run much faster. Given the large number of features, feature embeddings usually take up most of the model parameters. For fairness of comparison, we fix the embedding size to 16 or 40 in two different settings. We also found **that regularization weight makes a large effect on the model performance**. Thus, we carefully tuned it in the range of 0~1, usually at a multiplicative ratio of 10x. The model size (e.g., number of layers and units) is highly data-dependent, so we exhaustively tune these hyper-parameters. We also carefully tune some other hyper-parameters (e.g., whether use bath normalization), if any, to attain the best result for each model. To avoid exponential combination space, we usually tune important hyper-parameters first and then the other ones group by group. On average, we run 73 experiments for each model to obtain the best result. All the experiments were run on a shared GPU cluster with P100 GPUs, each with 16GB GPU memory.

**3.2.7 Reproducibility.** For reproducibility, we keep the md5sum values of each data split. We explicitly set the random seed for each experiment and record the data settings and model hyper-parameters into configure files. In particular, we choose to implement the models in Pytorch, since it offers better ability than Tensorflow to avoid non-determinism when running models on a GPU device. We have open the source code of FuxiCTR together with all the benchmarking settings and results to the community to foster more reproducible research in the future.

### 3.3 Results Analysis

In this section, we report our benchmarking results of 24 models, as shown in Table 3. To be specific, the "Best Reported" column shows the best results that we select from those reported by existing studies on Criteo and Avazu datasets. We report our benchmarking results w.r.t. logloss and AUC on four dataset settings, i.e., Criteo\_x4\_001, Criteo\_x4\_002, Avazu\_x4\_001, and Avazu\_x4\_002. For model efficiency, we also report training time, in terms of time per epoch and number of epochs, for Criteo\_x4\_002 and Avazu\_x4\_002. In addition, "#Params" denotes the number of parameters used in each model. "#Runs" records the number of experiments we run with grid search for model tuning. Note that the "#Runs" values typically depend on the number of hyper-parameters to tune in a model. The large number of runs (~73 on average) reveals that the models have been well tuned in our benchmarking. Moreover, we run the experiments on Criteo\_x4\_002 and Avazu\_x4\_002 first, so we perform a smaller number of experiments to tune models on Criteo\_x4\_001 and Avazu\_x4\_001. From the benchmarking results in Table 3, we have the following surprising observations:

- The best results reported by existing studies show certain inconsistency. For example, InterHAT has worse performance than LR on both datasets; DeepCross performs worse than LR on Avazu

as well. They are largely due to the fact that different data splitting or preprocessing steps are usually applied even on the same datasets. This reveals that standardized data splitting and preprocessing is necessary to make the results directly comparable among models.

- After model retuning on our dataset settings, we mostly obtain better performance than the best reported results. While their data settings may be different from what we described in Section 3.2, these results could serve as a good reference for our study. Considering that we follow the same data splitting as [42], the large improvements indicate that our modified data preprocessing is indeed more appropriate for model evaluation, resulting in stronger baselines results to build. This is also more practical, conforming to the real case where data preprocessing is usually heavily tuned.
- Our benchmarking follows the same evaluation protocol to make the results comparable. Yet, after exhaustive retuning, we found that the differences among the state-of-the-art models become small. For example, IPNN, DeepFM, DCN, xDeepFM, and ONN all achieve the same level of accuracy (~0.814 AUC) on Criteo, while DNN, DeepFM, DCN and xDeepFM attain comparable performance on Avazu. We run many experiments with different hyper-parameters, but do not obtain sufficiently distinct results. Especially, we run DCN for 544 and 496 experiments on Criteo\_x4\_002 and Avazu\_x4\_002 respectively, yet only make indistinguishable differences compared to DeepFM. The same happens for xDeepFM. Moreover, some recent models, such as InterHAT, AFN+, and LorentzFM obtain even worse results than some previous state-of-the-arts. Our results highlight the value of open benchmarking, which aims to provide reproducible and strong baseline results for future research.
- We also make a clear comparison between our benchmarking results and those reported by the other papers in Figure 1. For each dataset, we plot the AUC results from 8~9 existing papers. We can see that the results vary a lot among different papers due to unknown data splits and preprocessing. The X-axis is arranged sequentially according to the publication year, but there is no obvious pattern on model performance improvement. Some recent models only obtain diminished improvements and sometimes even lead to performance drops. Remarkably, our benchmarking presents the best results reported so far for all models.
- Memory consumption and model efficiency are another two important aspects to industrial CTR prediction tasks. As shown in the table, some models run very slowly (hours per epoch) due to the use of convolution networks (e.g., CCPM, FGCNN, HFM+), fieldwise interactions (e.g., FFM, ONN), graph neural networks (e.g., FiGNN), and so on. Some others have more parameters, such as FFM, ONN, FGCNN, etc. These drawbacks might hinder their practical use in industry.

**3.3.1 Model retuning.** To further demonstrate the necessity to retune baseline models, in Table 4, we present the results of four representative models in three settings: 1) The *Reported* setting denotes the results reported by the corresponding papers. 2) The *Rerun* setting indicates the cases that we rerun the experiments on our data splits according to the hyper-parameters given in the original papers. 3) The *Retuned* setting shows the results achieved

**Table 3: Benchmarking results on Criteo and Avazu. We highlight the top-5 best results in each dataset setting.**

Criteo												
Year	Model	Best Reported		Criteo_x4_001				Criteo_x4_002				
		Logloss	AUC	Logloss	AUC	#Params	#Runs	Logloss	AUC	#Params	#Runs	Time×Epochs
2007	LR	0.4474	0.7858	0.4568	0.7934	0.9M	6	0.4566	0.7936	5.5M	12	7m × 12
2010	FM	0.4464	0.7933	0.4431	0.8086	15.5M	9	0.4445	0.8078	227.5M	20	18m × 5
2015	CCPM	—	—	0.4415	0.8104	1.7M	9	0.4440	0.8077	222.0M	24	2h33m × 1
2016	HOFM	0.4508	0.8005	0.4411	0.8107	30.1M	8	0.4404	0.8115	255.3M	17	1h42m × 26
2016	FFM	0.4525	0.8001	0.4407	0.8113	139.3M	18	0.4409	0.8111	638.2M	20	3h59m × 5
2016	DNN	0.4491	0.7993	0.4380	0.8140	19.2M	51	0.4407	0.8112	226.5M	64	13m × 2
2016	WideDeep	0.4453	0.8062	<b>0.4377(5)</b>	<b>0.8142(5)</b>	20.1M	30	0.4389	0.8129	231.1M	80	11m × 4
2016	IPNN	0.4532	0.8038	0.4378	<b>0.8142(5)</b>	16.9M	24	0.4388	0.8132	258.5M	38	32m × 2
2016	DeepCross	0.4425	0.8009	0.4384	0.8135	17.1M	30	<b>0.4380(5)</b>	0.8139	284.4M	138	26m × 3
2017	AFM	0.4541	0.7965	0.4455	0.8060	15.5M	15	0.4443	0.8073	227.5M	10	31m × 22
2017	NFM	0.4537	0.7968	0.4424	0.8093	17.5M	24	0.4443	0.8072	229.6M	64	18m × 2
2017	DeepFM	0.4445	0.8085	<b>0.4376(2)</b>	<b>0.8143(3)</b>	20.1M	30	<b>0.4378(2)</b>	<b>0.8141(3)</b>	229.1M	128	18m × 9
2017	DCN	0.4419	0.8067	<b>0.4376(2)</b>	<b>0.8144(2)</b>	19.2M	81	<b>0.4378(2)</b>	<b>0.8141(3)</b>	245.1M	544	14m × 9
2018	FwFM	—	—	0.4408	0.8112	15.5M	18	0.4419	0.8098	222.0M	14	19m × 1
2018	xDeepFM	0.4418	0.8091	<b>0.4376(2)</b>	<b>0.8143(3)</b>	20.5M	106	<b>0.4375(1)</b>	<b>0.8144(1)</b>	232.3M	156	1h9m × 8
2019	HFM+	—	—	0.4392	0.8127	29.3M	34	0.4391	0.8127	260.2M	74	1h13m × 2
2019	FGCNN	—	—	0.4398	0.8121	61.1M	56	0.4381	<b>0.8142(2)</b>	317.4M	87	3h9m × 6
2019	AutoInt+	0.4434	0.8083	0.4390	0.8132	20.2M	75	0.4385	0.8134	285.4M	120	16m × 1
2019	FiGNN	0.4453	0.8062	0.4383	0.8138	14.7M	45	<b>0.4379(4)</b>	<b>0.8141(3)</b>	222.7M	40	1h40m × 14
2019	ONN	0.4358	0.8123	<b>0.4372(1)</b>	<b>0.8148(1)</b>	287.2M	30	0.4381	<b>0.8141(3)</b>	436.7M	108	2h13m × 8
2019	FiBiNET	0.4423	0.8103	0.4387	0.8131	71.1M	96	0.4386	0.8134	482.5M	78	2h21m × 2
2020	LorentzFM	—	—	0.4434	0.8083	14.6M	14	0.4413	0.8106	222.0M	8	2h12m × 20
2020	AFN+	0.4451	0.8074	0.4384	0.8138	56.0M	79	0.4387	0.8134	238.1M	114	38m × 6
2020	InterHAt	0.4577	0.7845	0.4414	0.81042	15.6M	106	0.4401	0.8117	222.1M	90	18m × 20

Avazu												
Year	Model	Best Reported		Avazu_x4_001				Avazu_x4_002				
		Logloss	AUC	Logloss	AUC	#Params	#Runs	Logloss	AUC	#Params	#Runs	Time×Epochs
2007	LR	0.3868	0.7676	0.3815	0.7775	3.8M	10	0.3799	0.7804	8.4M	90	5m × 26
2010	FM	0.3740	0.7793	0.3754	0.7887	63.8M	25	0.3736	0.7909	343.3M	187	18m × 1
2015	CCPM	0.3800	0.7812	0.3745	0.7892	60.6M	70	0.3721	0.7932	335.0M	24	57m × 1
2016	HOFM	0.3756	0.7701	0.3754	0.7891	123.8M	10	0.3733	0.7914	385.2M	31	1h47m × 1
2016	FFM	0.3781	0.7831	0.3719	0.7933	693.9M	19	0.3711	0.7948	778.7M	32	1h54m × 1
2016	DNN	—	—	0.3722	0.7928	63.4M	44	0.3705	0.7959	338.9M	252	12m × 1
2016	WideDeep	0.3744	0.7749	0.3720	0.7929	76.5M	27	0.3703	0.7958	345.9M	52	14m × 1
2016	IPNN	0.3737	0.7868	<b>0.3712(4)</b>	<b>0.7944(3)</b>	62.7M	18	<b>0.3684(4)</b>	<b>0.7989(4)</b>	336.1M	36	11m × 1
2016	DeepCross	0.3889	0.7643	0.3721	0.7930	63.1M	36	0.3700	0.7962	342.6M	45	12m × 1
2017	AFM	0.3766	0.7740	0.3792	0.7825	63.8M	29	0.3781	0.7840	343.3M	67	20m × 2
2017	NFM	0.3761	0.7708	0.3743	0.7894	66.8M	48	0.3721	0.7934	346.4M	54	16m × 1
2017	DeepFM	0.3742	0.7836	0.3719	0.7930	76.5M	24	0.3702	0.7962	373.2M	234	16m × 1
2017	DCN	0.3721	0.7681	0.3719	0.7931	68.8M	139	0.3699	0.7965	336.9M	496	7m × 1
2018	FwFM	0.3988	0.7406	0.3744	0.7907	60.0M	15	0.3724	0.7925	334.9M	15	11m × 1
2018	xDeepFM	0.3737	0.7855	0.3718	0.7933	64.6M	180	0.3697	0.7967	344.1M	288	31m × 1
2019	HFM+	—	—	<b>0.3714(5)</b>	<b>0.7944(3)</b>	66.2M	39	<b>0.3683(3)</b>	<b>0.7992(3)</b>	355.3M	92	43m × 1
2019	FGCNN	0.3746	0.7883	<b>0.3711(3)</b>	<b>0.7944(3)</b>	146.7M	52	<b>0.3696(5)</b>	<b>0.7971(5)</b>	374.5M	84	2h10m × 1
2019	AutoInt+	0.3811	0.7774	0.3746	0.7902	66.4M	114	0.3709	0.7953	337.9M	75	15m × 1
2019	FiGNN	0.3825	0.7762	0.3736	0.7915	60.1M	69	0.3711	0.7944	335.0M	64	2h25m × 1
2019	ONN	0.3945	0.7513	<b>0.3683(1)</b>	<b>0.7992(1)</b>	723.7M	18	<b>0.3677(2)</b>	<b>0.8001(2)</b>	406.3M	120	1h52 × 1
2019	FiBiNET	0.3786	0.7832	<b>0.3705(2)</b>	<b>0.7953(2)</b>	89.5M	72	<b>0.3675(1)</b>	<b>0.8003(1)</b>	395.9M	54	36m × 1
2020	LorentzFM	0.3828	0.7775	0.3756	0.7885	60.0M	5	0.3742	0.7912	334.9M	9	46m × 17
2020	AFN+	0.3718	0.7555	0.3726	0.7929	141.7M	20	0.3700	0.7965	363.5M	190	35m × 1
2020	InterHAt	0.3910	0.7582	0.3749	0.7882	60.1M	66	0.3722	0.7927	335.1M	56	17m × 1

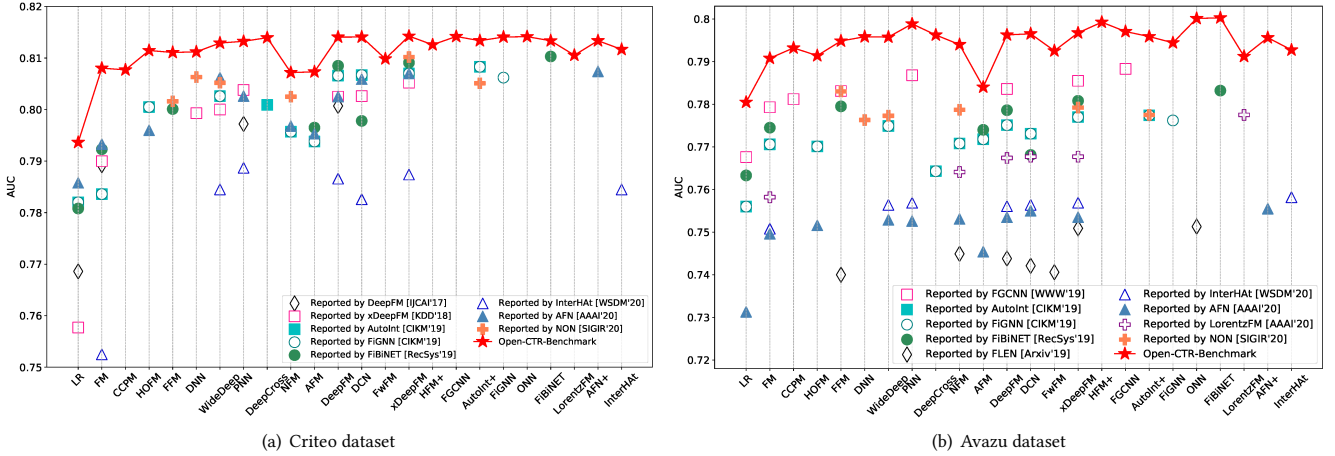


Figure 1: Comparison between our benchmarking results and those reported by the other papers

Table 4: Results before and after retuning hyper-parameters

Model	Setting	Criteo		Avazu	
		Logloss	AUC	Logloss	AUC
AutoInt+	Reported	0.4434	0.8083	0.3811	0.7774
	Rerun	0.4425	0.8092	0.3732	0.7913
	Retuned	<b>0.4385</b>	<b>0.8134</b>	<b>0.3709</b>	<b>0.7953</b>
FiGNN	Reported	0.4453	0.8062	0.3825	0.7762
	Rerun	0.4428	0.8089	0.3744	0.7893
	Retuned	<b>0.4379</b>	<b>0.8141</b>	<b>0.3711</b>	<b>0.7944</b>
FiBiNET	Reported	0.4423	0.8103	0.3786	0.7832
	Rerun	0.4431	0.8088	0.3724	0.7937
	Retuned	<b>0.4386</b>	<b>0.8134</b>	<b>0.3675</b>	<b>0.8003</b>
AFN+	Reported	0.4451	0.8074	0.3718	0.7555
	Rerun	0.4432	0.8084	0.3715	0.7939
	Retuned	<b>0.4387</b>	<b>0.8134</b>	<b>0.3700</b>	<b>0.7965</b>

after extensive hyper-parameter tuning. We can see that even on the same datasets, directly reusing the original hyper-parameters to rerun the experiments could bring large performance gaps in new data settings (i.e., Criteo\_x4\_002 and Avazu\_x4\_002 in our case). After model retuning, we achieve quite large improvements (up to 5%) over the original hyper-parameters, considering that 1% improvement is deemed to be significantly practical. This indicates that it is necessary to retune the hyper-parameters when testing a model on a new data split (even for the same dataset). However, it is not uncommon to find some studies (e.g., [21]) that choose to follow the baseline hyper-parameters used in the papers for fairness of comparison, but they experiment on a different data split. A common benchmark is thus desired to alleviate this issue.

### 3.4 Key Factors to Performance Tuning

During our benchmarking work, we also identify some key factors that are critical for performance tuning.

- **Data preprocessing.** Data often determine the upper bound of a model. However, existing work seldom tunes the min\_counts

thresholds for category features during data preprocessing. In our work, we set an appropriate threshold for infrequent feature filtering, which yields much better performance.

- **Batch size.** We observe that a large batch size usually results in faster training and better performance. For example, we set it to 10000 if the GPU does not raise an OOM error.
- **Embedding size.** While existing work usually set it to 10 or 16 in the experiments, we also experiment with other settings by using a larger embedding size (e.g., 40) within the GPU memory constraints.
- **Regularization weight and dropout rate.** Regularization and dropout are two key techniques to reduce model overfitting. They have a large impact on the performance of CTR prediction models. We exhaustively search the optimal value within a range.
- **Batch normalization.** In some cases, adding batch normalization between hidden layers of a DNN model can further boost prediction performance.

## 4 DISCUSSION

In this section, we discuss the limitations and potential directions for further exploration.

**More datasets:** In this work, we evaluate and benchmark existing CTR prediction models on two widely-used datasets, Criteo and Avazu. However, both datasets are anonymized and there is a lack of explicit user field and item field information. Therefore, they are able to be used to benchmark some models that require explicit user-item interaction (e.g., FLEN [7]) and user behavior sequence information (e.g., DIN [56]). We plan to extend more datasets from industrial-scale applications to make it a more comprehensive open benchmark for CTR prediction.

**Data splits:** To keep it consistent with most existing research, we split the datasets randomly to benchmark CTR prediction models. We do so with the following consideration. With randomly splitting, the data distributions among train, validation, and test sets are more consistent. This helps to better reveal the effectiveness of a CTR prediction model on capturing feature interactions, because in production it is necessary to perform CTR calibration after



prediction if the train-test distributions vary largely. As part of future work, we will evaluate the models by splitting data sequentially over time and also perform CTR calibration as necessary.

**Efficiency benchmarking:** Current models for CTR prediction has become more and more complex in structure, after using components such as attention [25, 42], convolution [29], and graph neural network [26]. In this version, we mainly evaluate the efficiency of these models through their training time. It is also desirable to test their inference time in future. Due to the strict latency constraints of CTR prediction in real-time applications, efficiency benchmarking would not only help practitioners choose an appropriate model, but also facilitate researchers to design effective yet efficient models.

**Auto-tuning of hyper-parameters:** As the experimental results shown in Section 3, hyper-parameter tuning is critical to the performance of CTR prediction models. How to quickly find the optimal hyper-parameter for a given model remains an open research problem. When data evolve with time, model hyper-parameters re-tuning is also required to adapt to the new data distribution. In our benchmark, we mainly apply grid search to find the best hyper-parameters of each model. It is highly expected to explore some advanced AutoML techniques (e.g., bayesian optimization in NNI [3]) to further boost hyper-parameter tuning process in future.

## 5 RELATED WORK

### 5.1 CTR Prediction

During the last decade, CTR prediction models have been widely studied and evolved through several generations from linear models [33], to factorization machines [39], and to deep learning-based models [10]. We have introduced some representative models in Section 2. Here, we present a review of more related studies on CTR prediction, which are summarized into the following categories.

**Feature interaction learning.** While simple linear models such as LR [40] and FTRL [33] have been widely used due to their simplicity and efficiency, they have difficulty capturing non-linear feature mappings and conjunctions. He et al. [18] propose the GBDT + LR approach that applies Gradient Boosting Decision Tree (GBDT) to extract meaningful feature conjunctions.

FM [39] is an effective model that captures pairwise feature interactions via inner products of feature vectors. Due to its success, many follow-up models have been proposed from different aspects, such as field awareness (e.g., FFM [22], FwFM [35]), importance of feature interactions (e.g., AFM [49], IFM [20]), outer-products based interaction (HFM [45]), robustness (RFM [36]), and interpretability (SEFM [23]). However, it is non-trivial for these models to capture high-order feature interactions in practice [5].

Recently, deep learning has become a popular technique in recommender systems [54], yielding an abundance of deep models for CTR prediction, including YoutubeDNN [10], Wide&Deep [8], PNN [37], DeepFM [16], DistillCTR [57], etc. Some of them aim to capture different orders of feature interactions explicitly (e.g., DCN [47], xDeepFM [27]). Some other models explore the use of convolutional networks (e.g., CCPM [30], FGCNN [29]), recurrent networks (e.g., [13, 55]), or attention networks (e.g., AutoInt [42], FiBiNET [21]) to learn implicit feature interactions.

**Behaviour sequence modeling.** The history behaviours of users have a large effect on predicting the click probability on the

next item. To better capture such history behaviours (e.g., item purchase sequences), some recent studies propose user interests modeling for CTR prediction via attention, LSTM, GRU, and memory networks. Typical examples include DIN [56], DIEN [55], DSIN [13], HPMN [38]) and DSTN [34].

**Multi-task learning.** In many recommender systems, users may have diverse behaviors beyond clicks, such as browsing, favorite, add-to-cart, and purchase. To improve the performance of CTR prediction, it is desirable to leverage other types of user feedbacks to enrich the supervision signals for CTR prediction. Towards this goal, some work proposes multi-task learning models to learn task relationships among different user behaviors, such as ESMM [32], MMoE [31] and PLE [44].

**Multi-modal learning.** Nowadays, multi-modal information contents are prevalent in recommender systems, such as e-commerce, news feeds, micro-blogs, and micro-videos. How to employ the rich multi-modal information (e.g., text, image, and video) of items to enhance CTR prediction models is an important research problem that needs more exploration. Some pioneer work (e.g., [14, 24, 48, 51, 53]) demonstrates the effectiveness to incorporate multi-modal content features in CTR prediction.

### 5.2 Benchmarking and Reproducibility

With the prevalence of deep learning, new models are emerging at an increasingly rapid pace. There is a high demand for an open benchmark to fairly compare against baseline models. Open benchmarking is valuable to promote research progress. For example, ImageNet [12] and GLUE [46] are two well-known benchmarks that contribute much to the progress in computer vision and natural language processing, respectively. In recommender systems, some datasets (e.g., Criteo and Avazu) are widely used. However, there is still a lack of standardized evaluation protocols, which results in the inconsistency and non-reproducibility issues of existing studies [11]. Notably, a concurrent work [43] also reports the benchmarking results on some classic recommendation models. Yet, their work fails to provide detailed configurations and hyper-parameter settings to allow reproducibility. In this work, we take an important step towards reproducible research by building the first open benchmark for CTR prediction as well as releasing the benchmarking results for over 20 models. More importantly, we provide all the intermediate artifacts (e.g., reproducing steps, running logs) to ensure reproducibility of our results.

## 6 CONCLUSION

In this paper, we present the first open benchmark for CTR prediction. We aim to alleviate the issues of non-reproducible and inconsistent results raised in current studies. We standardize the evaluation protocols and evaluate 24 existing models by running over 7,000 experiments for more than 12,000 GPU hours on two widely-used real-world datasets. We provide the most comprehensive benchmarking results that compare existing models in a rigorous manner. The results show that the difference between many models is smaller than expected and that inconsistent results exist in existing papers. We believe that our benchmark would not only drive more reproducible research but also help new beginners to learn the state-of-the-art CTR prediction models.

## REFERENCES

- [1] 2014. The Criteo Dataset. <https://www.kaggle.com/c/criteo-display-ad-challenge>
- [2] 2015. The Avazu Dataset. <https://www.kaggle.com/c/avazu-ctr-prediction>
- [3] 2020. Neural Network Intelligence — An open source AutoML toolkit. <https://github.com/Microsoft/nni>
- [4] 2021. The DeepCTR Package. <https://github.com/shenweichen/DeepCTR>
- [5] Mathieu Blondel, Akinori Fujino, Naonori Ueda, and Masakazu Ishihata. 2016. Higher-Order Factorization Machines. In *Annual Conference on Neural Information Processing Systems (NeurIPS)*. 3351–3359.
- [6] Tianqi Chen and Carlos Guestrin. 2016. XGBoost: A Scalable Tree Boosting System. In *Proceedings of the 22nd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining (KDD)*. 785–794.
- [7] Wenqiang Chen, Lizhang Zhan, Yuanlong Ci, and Chen Lin. 2019. FLEN: Leveraging Field for Scalable CTR Prediction. *CoRR* abs/1911.04690 (2019).
- [8] Heng-Tze Cheng, Levent Koc, Jeremiah Harmsen, Tal Shaked, et al. 2016. Wide & Deep Learning for Recommender Systems. In *Proceedings of the 1st Workshop on Deep Learning for Recommender Systems (DLRS@RecSys)*. 7–10.
- [9] Weiyu Cheng, Yanyan Shen, and Linpeng Huang. 2020. Adaptive Factorization Network: Learning Adaptive-Order Feature Interactions. In *The Thirty-Fourth AAAI Conference on Artificial Intelligence (AAAI)*. 3609–3616.
- [10] Paul Covington, Jay Adams, and Emre Sargin. 2016. Deep Neural Networks for YouTube Recommendations. In *Proceedings of the 10th ACM Conference on Recommender Systems (RecSys)*. 191–198.
- [11] Maurizio Ferrari Dacrema, Paolo Cremonesi, and Dietmar Jannach. 2019. Are we really making much progress? A worrying analysis of recent neural recommendation approaches. In *Proceedings of the 13th ACM Conference on Recommender Systems (RecSys)*. 101–109.
- [12] Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Fei-Fei Li. 2009. ImageNet: A large-scale hierarchical image database. In *IEEE Computer Society Conference on Computer Vision and Pattern Recognition (CVPR)*. 248–255.
- [13] 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. In *Proceedings of the Twenty-Eighth International Joint Conference on Artificial Intelligence (IJCAI)*. 2301–2307.
- [14] Tiezheng Ge, Liqin Zhao, Guorui Zhou, Keyu Chen, Shuying Liu, Huiming Yi, Zelin Hu, Bochao Liu, Peng Sun, Haoyu Liu, Pengtao Yi, Sui Huang, Zhiqiang Zhang, Xiaoqiang Zhu, Yu Zhang, and Kun Gai. 2018. Image Matters: Visually Modeling User Behaviors Using Advanced Model Server. In *Proceedings of the 27th ACM International Conference on Information and Knowledge Management (CIKM)*. 2087–2095.
- [15] Huifeng Guo, Bo Chen, Ruiming Tang, Weinan Zhang, Zhenguo Li, and Xiuqiang He. 2021. An Embedding Learning Framework for Numerical Features in CTR Prediction. In *The 27th ACM SIGKDD Conference on Knowledge Discovery and Data Mining (KDD)*. 2910–2918.
- [16] Huifeng Guo, Ruiming Tang, Yunming Ye, Zhenguo Li, and Xiuqiang He. 2017. DeepFM: A Factorization-Machine based Neural Network for CTR Prediction. In *International Joint Conference on Artificial Intelligence (IJCAI)*. 1725–1731.
- [17] Xiangnan He and Tat-Seng Chua. 2017. Neural Factorization Machines for Sparse Predictive Analytics. In *Proceedings of the 40th ACM SIGIR Conference on Research and Development in Information Retrieval (SIGIR)*. 355–364.
- [18] Xinran He, Junfeng Pan, Ou Jin, Tianbing Xu, et al. 2014. Practical Lessons from Predicting Clicks on Ads at Facebook. In *Proceedings of the Eighth International Workshop on Data Mining for Online Advertising (ADKDD)*. 5:1–5:9.
- [19] Xinran He, Junfeng Pan, Ou Jin, Tianbing Xu, Bo Liu, Tao Xu, Yanxin Shi, Antoine Atallah, Ralf Herbrich, Stuart Bowers, and Joaquin Quiñero Candela. 2014. Practical Lessons from Predicting Clicks on Ads at Facebook. In *Proceedings of the Eighth International Workshop on Data Mining for Online Advertising (ADKDD)*. 5:1–5:9.
- [20] Fuxing Hong, Dongbo Huang, and Ge Chen. 2019. Interaction-Aware Factorization Machines for Recommender Systems. In *The Thirty-Third AAAI Conference on Artificial Intelligence (AAAI)*. 3804–3811.
- [21] Tongwen Huang, Zhiqi Zhang, and Junlin Zhang. 2019. FiBiNET: combining feature importance and bilinear feature interaction for click-through rate prediction. In *Proceedings of ACM Conference on Recommender Systems (RecSys)*. 169–177.
- [22] Yu-Chin 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 (RecSys)*. 43–50.
- [23] Liang Lan and Yu Geng. 2019. Accurate and Interpretable Factorization Machines. In *The AAAI Conference on Artificial Intelligence (AAAI)*. 4139–4146.
- [24] Xiang Li, Chao Wang, Jiwei Tan, Xiaoyi Zeng, Dan Ou, and Bo Zheng. 2020. Adversarial Multimodal Representation Learning for Click-Through Rate Prediction. In *The Web Conference (WWW)*. 827–836.
- [25] Zeyu Li, Wei Cheng, Yang Chen, Haifeng Chen, and Wei Wang. 2020. Interpretable Click-Through Rate Prediction through Hierarchical Attention. In *The International Conference on Web Search and Data Mining (WSDM)*. 313–321.
- [26] Zekun Li, Zeyu Cui, Shu Wu, Xiaoyu Zhang, and Liang Wang. 2019. Fi-GNN: Modeling Feature Interactions via Graph Neural Networks for CTR Prediction. In *Proceedings of the 28th ACM International Conference on Information and Knowledge Management (CIKM)*. 539–548.
- [27] Jianxun Lian, Xiaohuan Zhou, Fuzheng Zhang, et al. 2018. xDeepFM: Combining Explicit and Implicit Feature Interactions for Recommender Systems. In *Proceedings of the 24th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining (KDD)*. 1754–1763.
- [28] Xiaoliang Ling, Weiwei Deng, Chen Gu, Hucheng Zhou, Cui Li, and Feng Sun. 2017. Model Ensemble for Click Prediction in Bing Search Ads. In *Proceedings of the 26th International Conference on World Wide Web Companion (WWW)*. 689–698.
- [29] Bin Liu, Ruiming Tang, Yingzhi Chen, Jinkai Yu, Huifeng Guo, and Yuzhou Zhang. 2019. Feature Generation by Convolutional Neural Network for Click-Through Rate Prediction. In *The World Wide Web Conference, (WWW)*. 1119–1129.
- [30] Qiang Liu, Feng Yu, Shu Wu, and Liang Wang. 2015. A Convolutional Click Prediction Model. In *Proceedings of the 24th ACM International Conference on Information and Knowledge Management (CIKM)*. 1743–1746.
- [31] Jiaqi Ma, Zhe Zhao, Xinyang Yi, Jilin Chen, Lichan Hong, and Ed H. Chi. 2018. 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 (KDD)*. 1930–1939.
- [32] Xiao Ma, Liqin Zhao, Guan Huang, Zhi Wang, Zelin Hu, Xiaoqiang Zhu, and Kun Gai. 2018. Entire Space Multi-Task Model: An Effective Approach for Estimating Post-Click Conversion Rate. In *Proceedings of the 41st International Conference on Research & Development in Information Retrieval (SIGIR)*. 1137–1140.
- [33] H. Brendan McMahan, Gary Holt, David Sculley, Michael Young, Dietmar Ebner, Julian Grady, Lan Nie, Todd Phillips, Eugene Davydov, Daniel Golovin, Sharat Chikkerur, Dan Liu, Martin Wattenberg, Arnar Mar Hrafnkelsson, Tom Boulos, and Jeremy Kubica. 2013. Ad click prediction: a view from the trenches. In *Proceedings of the 19th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining (KDD)*. 1222–1230.
- [34] Wentao Ouyang, Xiuwu Zhang, Li Li, Heng Zou, Xin Xing, Zhaojie Liu, and Yanlong Du. 2019. Deep Spatio-Temporal Neural Networks for Click-Through Rate Prediction. In *Proceedings of the 25th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining (KDD)*. 2078–2086.
- [35] 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 on World Wide Web, WWW 2018, Lyon, France, April 23-27, 2018*. 1349–1357.
- [36] Surabhi Punjabi and Priyanka Bhatt. 2018. Robust Factorization Machines for User Response Prediction. In *Proceedings of the 2018 World Wide Web Conference on World Wide Web (WWW)*. 669–678.
- [37] 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 *Proceedings of the IEEE 16th International Conference on Data Mining (ICDM)*. 1149–1154.
- [38] Kan Ren, Jiarui Qin, Yuchen Fang, Weinan Zhang, Lei Zheng, Weijie Bian, Guorui Zhou, Jian Xu, Yong Yu, Xiaoqiang Zhu, and Kun Gai. 2019. Lifelong Sequential Modeling with Personalized Memorization for User Response Prediction. In *Proceedings of the 42nd International ACM SIGIR Conference on Research and Development in Information Retrieval (SIGIR)*. 565–574.
- [39] Steffen Rendle. 2010. Factorization Machines. In *Proceedings of the 10th IEEE International Conference on Data Mining (ICDM)*. 995–1000.
- [40] 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 (WWW)*. 521–530.
- [41] Ying Shan, T. Ryan Hoens, Jian Jiao, Haijing Wang, Dong Yu, and J. C. Mao. 2016. Deep Crossing: Web-Scale Modeling without Manually Crafted Combinatorial Features. In *Proceedings of the 22nd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining (KDD)*. 255–262.
- [42] 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 (CIKM)*. 1161–1170.
- [43] Zhu Sun, Di Yu, Hui Fang, Jie Yang, Xinghua Qu, Jie Zhang, and Cong Geng. 2020. Are We Evaluating Rigorously? Benchmarking Recommendation for Reproducible Evaluation and Fair Comparison. In *Proceedings of the Fourteenth ACM Conference on Recommender Systems (RecSys)*. 23–32.
- [44] Hongyan Tang, Junling Liu, Ming Zhao, and Xudong Gong. 2020. Progressive Layered Extraction (PLE): A Novel Multi-Task Learning (MTL) Model for Personalized Recommendations. In *Proceedings of the Fourteenth ACM Conference on Recommender Systems (RecSys)*. 269–278.
- [45] Yi Tay, Shuai Zhang, Anh Tuan Luu, Siu Cheung Hui, Lina Yao, and Tran Dang Quang Vinh. 2019. Holographic Factorization Machines for Recommendation. In *The Thirty-Third AAAI Conference on Artificial Intelligence (AAAI)*. 5143–5150.
- [46] Alex Wang, Amanpreet Singh, Julian Michael, Felix Hill, Omer Levy, and Samuel R. Bowman. 2019. GLUE: A Multi-Task Benchmark and Analysis Platform for Natural Language Understanding. In *7th International Conference on Learning*

*Representations (ICLR).*

- [47] Ruoxi Wang, Bin Fu, Gang Fu, and Mingliang Wang. 2017. Deep & Cross Network for Ad Click Predictions. In *Proceedings of the 11th International Workshop on Data Mining for Online Advertising (ADKDD)*. 12:1–12:7.
- [48] Yu Wang, Jixing Xu, Aohan Wu, Mantian Li, Yang He, Jinghe Hu, and Weipeng P. Yan. 2018. Telepath: Understanding Users from a Human Vision Perspective in Large-Scale Recommender Systems. In *Proceedings of the Thirty-Second AAAI Conference on Artificial Intelligence (AAAI)*. 467–474.
- [49] 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 Twenty-Sixth International Joint Conference on Artificial Intelligence (IJCAI)*. 3119–3125.
- [50] Canran Xu and Ming Wu. 2020. Learning Feature Interactions with Lorentzian Factorization Machine. In *The Thirty-Fourth AAAI Conference on Artificial Intelligence (AAAI)*. 6470–6477.
- [51] Jiahao Xun, Shengyu Zhang, Zhou Zhao, Jieming Zhu, Qi Zhang, Jingjie Li, Xiuqiang He, Xiaofei He, Tat-Seng Chua, and Fei Wu. 2021. Why Do We Click: Visual Impression-aware News Recommendation. In *Proceedings of the 29th ACM International Conference on Multimedia (MM)*.
- [52] Yi Yang, Baile Xu, Shaofeng Shen, Furao Shen, and Jian Zhao. 2020. Operation-aware Neural Networks for user response prediction. *Neural Networks* 121 (2020), 161–168.
- [53] Qi Zhang, Jingjie Li, Qinglin Jia, Chuyuan Wang, Jieming Zhu, Zhaowei Wang, and Xiuqiang He. 2021. UNBERT: User-News Matching BERT for News Recommendation. In *Proceedings of the Thirtieth International Joint Conference on Artificial Intelligence (IJCAI)*. 3356–3362.
- [54] Shuai Zhang, Lina Yao, Aixin Sun, and Yi Tay. 2019. Deep Learning Based Recommender System: A Survey and New Perspectives. *ACM Comput. Surv.* 52, 1 (2019), 5:1–5:38.
- [55] Guorui Zhou, Na Mou, Ying Fan, Qi Pi, Weijie Bian, Chang Zhou, Xiaoqiang Zhu, and Kun Gai. 2018. Deep Interest Evolution Network for Click-Through Rate Prediction. *CoRR* abs/1809.03672 (2018).
- [56] Guorui Zhou, Xiaoqiang Zhu, Chengru Song, Ying Fan, Han Zhu, Xiao Ma, Yanghui Yan, Junqi Jin, Han Li, and Kun Gai. 2018. Deep Interest Network for Click-Through Rate Prediction. In *Proceedings of the SIGKDD International Conference on Knowledge Discovery & Data Mining (KDD)*. 1059–1068.
- [57] Jieming Zhu, Jinyang Liu, Weiqi Li, Jincai Lai, Xiuqiang He, Liang Chen, and Zibin Zheng. 2020. Ensembled CTR Prediction via Knowledge Distillation. In *Proceedings of the 29th ACM International Conference on Information and Knowledge Management (CIKM)*. 2941–2958.