# STORE: Semantic Tokenization, Orthogonal Rotation and Efficient Attention for Scaling Up Ranking Models

Yi Xu
Alibaba Group
Beijing, China
xy397404@alibaba-inc.com

Chaofan Fan
Alibaba Group
Beijing, China
fanchaofan.fcf@alibaba-inc.com

Jinxin Hu*
Alibaba Group
Beijing, China
jinxin.hjx@alibaba-inc.com

Yu Zhang
Alibaba Group
Beijing, China
daoji@alibaba-inc.com

Xiaoyi Zeng
Alibaba Group
Beijing, China
yuanhan@taobao.com

Jing Zhang
Wuhan University
Wuhan, China
jingzhang.cv@gmail.com

## Abstract

Ranking models have become an important part of modern personalized recommendation systems. However, significant challenges persist in handling high-cardinality, heterogeneous, and sparse feature spaces, particularly regarding model scalability and efficiency. We identify two key bottlenecks: (i) Representation Bottleneck: Driven by the high cardinality and dynamic nature of features, model capacity is forced into sparse-activated embedding layers, leading to low-rank representations. This, in turn, triggers phenomena like "One-Epoch" and "Interaction-Collapse," ultimately hindering model scalability. (ii) Computational Bottleneck: Integrating all heterogeneous features into a unified model triggers an explosion in the number of feature tokens, rendering traditional attention mechanisms computationally demanding and susceptible to attention dispersion. To dismantle these barriers, we introduce STORE, a unified and scalable token-based ranking framework built upon three core innovations: (1) Semantic Tokenization fundamentally tackles feature heterogeneity and sparsity by decomposing high-cardinality sparse features into a compact set of stable semantic tokens; and (2) Orthogonal Rotation Transformation is employed to rotate the subspace spanned by low-cardinality static features, which facilitates more efficient and effective feature interactions; and (3) Efficient attention that filters low-contributing tokens to improve computional efficiency while preserving model accuracy. Across extensive offline experiments and online A/B tests, our framework consistently improves prediction accuracy(online CTR by 2.71%, AUC by 1.195%) and training effeciency (1.84× throughput).

## CCS Concepts

• **Information systems → Recommender systems**.

*Corresponding author

## Keywords

Recommendation System; Click-Through Rate Prediction; Semantic ID;

## 1 Introduction

Ranking models form the backbone of modern online services. Their core task is to model complex user behavior by processing a vast and heterogeneous collection of features. To handle this feature diversity, existing ranking models have evolved into a collection of specialized modules for feature interaction. While effective for accuracy, this intricate and fragmented design presents a major obstacle to scalability. Unlike Large Language Models (LLMs), where "Scaling Laws" provide a clear path to predictable performance gains, ranking models fail to exhibit similar scaling behavior[6, 14]. There are two fundamental challenges that prevent ranking models from benefiting from Scaling Laws:

(1) Representation Bottleneck: The high-cardinality features force model capacity into sparse-activated embedding layers over deep networks. This yields low-rank embeddings, triggers the "One-Epoch"[16] and "Interaction-Collapse"[5] problems. Ultimately, this hinders model scalability where adding depth or training epochs offers diminishing returns. This phenomenon leads to the loss of effective high-order feature interactions[5]. As model size grows, gains diminish rapidly, undermining capacity utilization and predictable scaling. (2) Computational Bottleneck: As the model scale to incorporate vast feature sets triggers an explosion in the number of feature tokens. This renders vanilla attention, with its $O(L^2)$ complexity, computationally prohibitive and simultaneously exacerbates attention dispersion, where vital signals are lost amidst a sea of irrelevant tokens.

To address these challenges, we introduce STORE (**S**emantic **T**okenization, **O**rthogonal **R**otation, and **E**fficient attention), a unified and scalable ranking framework built upon three synergistic components. (1) Semantic Tokenization efficiently decomposes high-cardinality features into a set of compact, orthogonal semantic IDs. This approach fundamentally mitigates feature heterogeneity and sparsity, thereby unlocking more efficient model scaling.
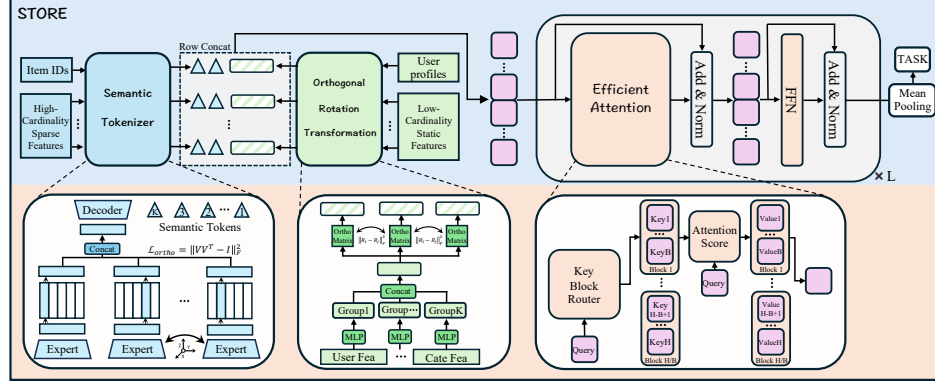
**Figure 1: Overview of the proposed STORE.**

(2) Orthogonal Rotation Transformation is employed to rotate the subspace spanned by low-cardinality static features, which facilitates more efficient and effective feature interactions in different high-dimensional spaces. (3) Efficient Attention adaptively utilizes sparsity mechanism to prune low-contributing tokens conditioned on the target item and context, reducing computational complexity and alleviating attention dispersion while preserving accuracy. Our key contributions are summarized as follows:

- We present STORE, a unified token-based ranking model framework that effectively tackles heterogeneity and sparsity in large and dynamic feature spaces. This paradigm mitigates long-standing scaling-law bottlenecks in large-scale recommender systems, supporting more predictable scaling.
- We design a synergistic trio of architectural innovations: semantic tokenization and orthogonal rotation are proposed to resolve the representation bottleneck, while efficient attention is proposed to mitigate the computational bottleneck by sparsifying attention and reducing quadratic costs.
- Extensive offline experiments and online A/B tests demonstrate the superiority of STORE in both effectiveness and efficiency, showing an improvement of 1.195% in AUC, 2.71% in CTR, and 1.84× higher training throughput.

## 2 Methodology

This section details the proposed STORE framework. To address the heterogeneity and sparsity of feature space, features are categorized into high-cardinality sparse features(e.g., item identifiers) and low-cardinality static features(e.g., category id, age, gender), employed wtih distinct processing strategies: Semantic Tokenizer for high-cardinality sparse features, detailed in section 2.1, and Orthogonal Rotation Transformation for low-cardinality static features, detailed in section 2.2. To migrate the computational efficiency bottleneck for ranking model, the Efficient Attention for Unified Feature Interaction is proposed, detailed in section 2.3.

### 2.1 Semantic Tokenizer

Using item IDs as a representative example of high-cardinality sparse features, we mitigate this issue by mapping raw item IDs into a more stable and structured semantic space via Semantic IDs (SIDs). We achieve this by quantizing powerful pre-trained embeddings (e.g., from SASRec) into a sequence of SIDs.

$$(SID_1, SID_2, \ldots, SID_K) = \mathcal{T}_{\text{item}}(\mathbf{e}_p \in \mathbb{R}^d) \qquad (1)$$

where $\mathbf{e}_p$ denotes the pre-trained item embedding and $\mathcal{T}_{\text{item}}$ denotes the item semantic tokenization function producing $K$ SIDs for each item. In this paper, the setting is $K = H$.

To efficiently encode high-cardinality IDs into a set of compact and parallel SIDs, we propose an **O**rthogonal, **P**arallel, **M**ulti-expert **Q**uantization network(**OPMQ**). For each item, the network utilizes $K$ experts to encode its pre-trained embedding into $K$ latent representations. The formulation is as follows.

$$\mathbf{z}_i = E_i(\mathbf{e}_p), \quad i \in \{1, \ldots, K\} \qquad (2)$$

$$c_i = \arg \min_{j \in \{1, \ldots, V\}} \|\mathbf{z}_i - \mathbf{s}_j\|_2^2 \qquad (3)$$

where $E_i$ is the $i$-th expert and the latent representation is $\mathbf{z}_i$. For each latent vector $\mathbf{z}_i$, we assigning it to the index of its nearest neighbor codeword $c_i$, the codeword vector $\mathbf{s}_i$. The entire OPMQ network is trained end-to-end by minimizing the reconstruction error between the original embedding and the output of a decoder that aggregates the quantized vectors.

$$\mathcal{L}_{recon} = ||\mathbf{e}_p - decoder[\sum_i^K (\mathbf{z_i} + sg(\mathbf{s}_i - \mathbf{z_i}))]||^2 \qquad (4)$$

To capture diverse and non-redundant aspects of the original item, the orthogonal regularization of SIDs is performed on the parameters of multi-experts. Formally, for the $i$-expert, we define the parameter vector $\mathbf{V_i} \in \mathbb{R}^{d_1 d_2}$ as the L2-normalized version of the flattened parameter matrix $\mathbf{W}_i \in \mathbb{R}^{d_1 \times d_2}$. The orthogonal regularization is performed on the set of $K$ parameter vectors, which is formulated as follows.

$$\mathcal{L}_{\text{orth}} = \left\|\mathbf{V}\mathbf{V}^\top - \mathbf{I}\right\|_F^2, \qquad (5)$$

where $\mathbf{I}$ is the identity matrix and $\|\cdot\|_F^2$ denotes the Frobenius norm.

### 2.2 Orthogonal Rotation Transformation

Unlike high-cardinality sparse features with heterogeneity and sparsity, for low-cardinality static features with controllable features sizes, we employ the original embeddings. For simplicity and efficiency, we perform manual grouping based on their semantic meanings of domain knowledge. These features are partitioned into $K$ semantic groups, with each group containing several features. For each feature group, a shallow MLP is employed for intra-group feature fusion. By concatenating all the semantically fused feature groups, we obtain an instance-wise feature block, denoted as $\mathbf{C}$. The formulation is as follows.

$$C = [MLP_1(g_1), \ldots, MLP_K(g_K)] \tag{6}$$

To facilitate efficient and effective feature interactions in high-dimensional spaces, the orthogonal rotation transformation is employed to rotate the instance-wise feature block. To obtain $K$ diverse instance-wise feature block, we rotate the $\mathbf{C}$ with $K$ group of orthogonal matrices. For the $i$-th rotation, the formulation is as follows.

$$\mathbf{O_i} = \mathbf{CR_i} \tag{7}$$

where $\mathbf{R_i}$ is an orthogonal matrix. To prevent the rotation matrices from collapsing during training (e.g., all becoming identical), we introduce a diversity regularization term. This works in concert with the orthogonality constraint to encourage a diverse set of learned transformations, formulated as follows.

$$\min_{\mathbf{R_1}, \ldots, \mathbf{R_k}} \quad -\lambda \sum_{i=1}^{K} \sum_{j=i+1}^{K} \|\mathbf{R_i} - \mathbf{R_j}\|_F^2 \tag{8}$$

$$\text{s.t.} \quad \mathbf{R_i}^T \mathbf{R_i} = \mathbf{I}, \quad \forall i \in \{1, \ldots, K\} \tag{9}$$

where $||\cdot||_F$ is the Frobenius norm. $\lambda$ is the hyperparameter, which is set to 0.1 in this paper. The rotation matrices and the parameters of the main network are alternative optimized.

## 2.3 Efficient Attention for Unified Feature Interaction

To efficiently capture feature interactions in a unified framework, we propose the efficient attention with instance-wise tokens. Following the distinct processing of high-cardinality sparse and low-cardinality static features, we concatenate the embedding of SIDs with the rotated feature block in the first layer, $\mathbf{X_0^i} = [\mathbf{s_i}, \mathbf{O_i}], \mathbf{X_0^i} \in \mathbb{R}^{H \times d}$. The input of efficient attention for feature interactions is the instance-wise token sequence $\mathbf{X_0} = [\mathbf{X_0^1}, \mathbf{X_0^2}, \ldots, \mathbf{X_0^H}]$, which construct the $\mathbf{Q}, \mathbf{K}, \mathbf{V}$.

The iterative unified efficient attention for feature interaction formulated as follows:

$$\mathbf{X_l} = \text{LN}(\text{EfficientAttention}(\mathbf{X_{l-1}}) + \mathbf{X_{l-1}}) \tag{10}$$

where $\mathbf{X_{l-1}}$ is the input to the l-th layer, LN denotes Layer Normalization. Vanilla self-attention's $O(H^2)$ computational complexity becomes prohibitive as the number of instance-wise tokens $H$ grows.

$$\text{MoBA}(\mathbf{Q}, \mathbf{K}, \mathbf{V}) = \text{Softmax}\left(\mathbf{Q}\mathbf{K}[Ind]^T\right)\mathbf{V}[Ind], \tag{11}$$

$$Ind_i = [(i-1) \times B + 1, i \times B] \tag{12}$$

To overcome this bottleneck, our framework incorporates an efficient attention mechanism. Specifically, we adopt MOBA [7], which employs a routing strategy for each query to attend to only a small subset of key-value pairs. As formulated in Eq 11,12, $Ind_i \subseteq \{1, \ldots, H\}$ is the dynamically selected set of indices of key-value pairs, the size of selective block is $B$. This approach reduces the complexity from quadratic significantly, a choice made viable by our framework's effective mitigation of feature heterogeneity and sparsity.

## 3 Experiments

In this paper, we conduct extensive experiments on both industrial and public datasets to evaluate the effectiveness of STORE with the following questions:

**RQ1**: How does STORE compare to SOTA ranking models?

**RQ2**: What is the contribution of each component in STORE?

**RQ3**: What is the scalability of STORE?

## 3.1 Experimental Setup

**Table 1: Overall performance comparison in CTR prediction on public and industrial datasets. "Improv." denotes the relative improvement of STORE over the best baseline. The best baseline performance score is denoted in <u>underline</u>.**

| Dataset | Avazu | | | Industrial | | |
|---|---|---|---|---|---|---|
| Method | AUC | GAUC | Logloss | AUC | GAUC | Logloss |
| FM | 0.7291 | 0.7248 | 0.4052 | 0.6711 | 0.6011 | 0.1144 |
| DNN | 0.7231 | 0.7211 | 0.4052 | 0.6721 | 0.6005 | 0.1148 |
| Wide&Deep | 0.7356 | 0.7329 | 0.3988 | 0.6720 | 0.6018 | 0.1144 |
| DeepFM | 0.7404 | 0.7375 | 0.3965 | 0.6707 | 0.5907 | 0.1152 |
| DCN | 0.7344 | 0.7310 | 0.4042 | 0.6734 | 0.6029 | 0.1141 |
| AutoInt | 0.7439 | 0.7408 | 0.3948 | 0.6728 | 0.6021 | 0.1142 |
| GDCN | 0.7370 | 0.7344 | 0.3989 | 0.6726 | 0.6022 | 0.1142 |
| MaskNet | 0.7426 | 0.7383 | 0.3942 | 0.6753 | 0.6054 | <u>0.1140</u> |
| PEPNet | 0.7411 | 0.7380 | 0.3961 | 0.6741 | 0.6039 | 0.1148 |
| RankMixer | 0.7450 | 0.7412 | 0.3951 | <u>0.6774</u> | 0.6053 | <u>0.1140</u> |
| OneTrans | <u>0.7461</u> | <u>0.7432</u> | <u>0.3943</u> | 0.6771 | <u>0.6058</u> | 0.1141 |
| **STORE** | **0.7479** | **0.7451** | **0.3912** | **0.6804** | **0.6064** | **0.1139** |
| **STORE-4 Epoch** | **0.7488** | **0.7463** | **0.3900** | **0.6855** | **0.6086** | **0.1134** |
| **Improv.** | *+0.362%* | *+0.417%* | *+0.913%* | *+1.195%* | *+0.462%* | *+0.526%* |

*3.1.1 Dataset.* To validate the effectiveness of our proposed framework, we conduct experiments on real-world large-scale datasets and public datasets.

**Avazu**: Avazu is a widely-used public benchmark for CTR prediction, consisting of 9 million of chronologically ordered ad click logs, 23 feature fields and 3437 site ids.

**Industrial Dataset**: This dataset contains 7 billion user interaction records from an international e-commerce advertising system, featuring diverse item features and user behavior sequences.

*3.1.2 Evaluation Metrics.* For the evaluation, we use the widely used AUC, GAUC and LogLoss for prediction accuracy. We use training TFlops/Batch(batch size=1024) to evaluate the training effeciency.

*3.1.3 Baselines.* To demonstrate the effectiveness of our framework, we evaluate the proposed framework with state-of-the-art CTR prediction models, including FM[8], DNN, Wide&Deep[2], DeepFM[4], DCN[11], AutoInt[9], GDCN[10], MaskNet[12], PEPNet[1], RankMixer[17] and OneTrans[15].

*3.1.4 Implementation Details.* We utilize the pre-trained item embeddings from a pre-trained SASRec model. The number of SIDs (K) and the codebook size are set to (3, 16) for the public dataset and (32, 300) for the industrial dataset, respectively.

## 3.2 Overall Performance(RQ1)

Table 1 presents the overall prediction performance of all methods on both industrial and public datasets, alongside with the relative improvement against the best baseline. The performance comparision on two datasets can demonstrate the effectiveness of STORE. It is noteworthy that while methods like Rankmixer and OneTrans project aggregated feature groups to mitigate feature heterogeneity, achieving gains over traditional networks. In contrast, our proposed method fundamentally resolves the issues by leveraging SIDs and rotation techniques. This leads to substantial improvements in prediction accuracy.

## 3.3 Ablation Study(RQ2)

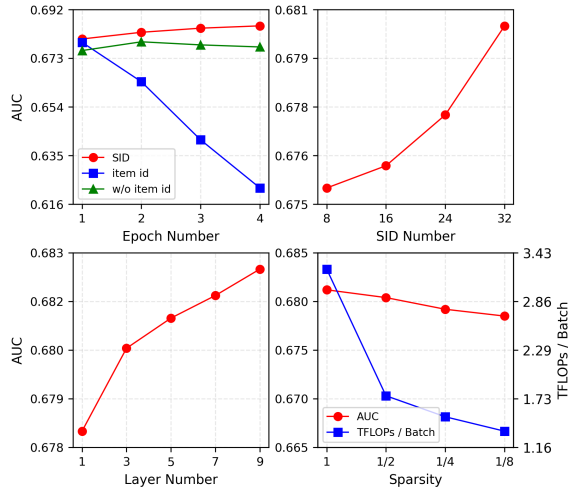As shown in Table 2, we conduct ablation experiments to evaluate the impact of each component in STORE.

- **Different Semantic Tokenizer**: We compare OPMQ against two widely-used tokenizers: RQ-VAE[13], and Optimized Product Quantization (OPQ)[3].
- **Orthogonal Rotation Transformation**: We conduct experiments on w/o orthogonal rotation transformation.
- **Efficient Attention**: Compared with vanilla attention, the efficient attention achieve comparable prediction accuracy with higher training effeciency.

**Table 2: Ablation study of STORE variants and components.**

| Variants | AUC | GAUC | Logloss | TFlops/Batch |
|---|---|---|---|---|
| **STORE-4 Epoch** | **0.6855** | **0.6086** | **0.1134** | **1.764** |
| **STORE** | **0.6804** | **0.6064** | **0.1139** | **1.763** |
| w OPQ | 0.6787 | 0.6045 | 0.1140 | 1.763 |
| w RQ-VAE | 0.6768 | 0.6047 | 0.1141 | 1.762 |
| w/o Orthogonal Rotation | 0.6780 | 0.6050 | 0.1140 | 1.760 |
| w Vanilla-Attention | 0.6812 | 0.6068 | 0.1137 | 3.240 |

## 3.4 Scaling Laws Study with Different Hyperparameters(RQ3)

In this section, we conducted experiments to demonstrate the scalability and efficiency of our proposed method. We will compare them from the following dimensions. As shown in Fig 2: (a) Epoch number: that demonstrate SIDs combats the "One-Epoch" phenomenon which models with ItemIDs show decreased performance over mult-epoch. (b) SID Number: The more SIDs, the better the effect. (c) Layer Number: The more Layers, the better the effect. (d) Sparsity: We explore the relationship between attention sparsity, the training effeciency and model accuracy has demonstrated STORE's ability to reduce computational cost with minimal impact on performance.



**Figure 2: Scaling Laws Study of (a) Epoch Number (b) SID Number (c) Layer Number (d) Sparsity.**

## 4 Online Experiments

We conducted a 15-day online A/B test on a large-scale e-commerce platform, comparing STORE with the production baseline. STORE achieved a relative CTR increase of 2.71%. In deployment, the OPMQ is set to K=32, codebook size=300 for SIDs. The sparsity of attention is set to 1/2, improving inference efficiency and response speed while maintaining performance.

## 5 Conclusion

We introduced STORE, a framework that resolves both representation and computational bottlenecks in recommenders. Through semantic tokenization, orthogonal rotation, and an efficient attention mechanism, STORE unlocks superior model scalability and computational efficiency. Extensive experiments validate STORE as a practical and effective path toward building more powerful large-scale ranking models.

## References

[1] Jianxin Chang, Chenbin Zhang, Yiqun Hui, Dewei Leng, Yanan Niu, Yang Song, and Kun Gai. 2023. PEPNet: Parameter and Embedding Personalized Network for Infusing with Personalized Prior Information. arXiv:2302.01115 [cs.IR] https://arxiv.org/abs/2302.01115

[2] Heng-Tze Cheng, Levent Koc, Jeremiah Harmsen, Tal Shaked, Tushar Chandra, Hrishi Aradhye, Glen Anderson, Greg Corrado, Wei Chai, Mustafa Ispir, Rohan Anil, Zakaria Haque, Lichan Hong, Vihan Jain, Xiaobing Liu, and Hemal Shah. 2016. Wide & Deep Learning for Recommender Systems. arXiv:1606.07792 [cs.LG] https://arxiv.org/abs/1606.07792

[3] Tiezheng Ge, Kaiming He, Qifa Ke, and Jian Sun. [n. d.]. Optimized Product Quantization. ([n. d.]).

[4] Huifeng Guo, Ruiming Tang, Yunming Ye, Zhenguo Li, and Xiuqiang He. 2017. DeepFM: A Factorization-Machine based Neural Network for CTR Prediction. *ArXiv* abs/1703.04247 (2017). https://api.semanticscholar.org/CorpusID:970388

[5] Xingzhuo Guo, Junwei Pan, Ximei Wang, Baixu Chen, Jie Jiang, and Mingsheng Long. 2023. On the Embedding Collapse when Scaling up Recommendation Models. *ArXiv* abs/2310.04400 (2023). https://api.semanticscholar.org/CorpusID:263831500

[6] Jared Kaplan, Sam McCandlish, Tom Henighan, Tom B. Brown, Benjamin Chess, Rewon Child, Scott Gray, Alec Radford, Jeffrey Wu, and Dario Amodei. 2020. Scaling Laws for Neural Language Models. *CoRR* abs/2001.08361 (2020). arXiv:2001.08361 https://arxiv.org/abs/2001.08361

[7] Enzhe Lu, Zhejun Jiang, Jingyuan Liu, Yulun Du, Tao Jiang, Chao Hong, Shaowei Liu, Weiran He, Enming Yuan, Yuzhi Wang, Zhiqi Huang, Huan Yuan, Suting Xu, Xinran Xu, Guokun Lai, Yanru Chen, Huabin Zheng, Junjie Yan, Jianlin Su, Yuxin Wu, Neo Y. Zhang, Zhilin Yang, Xinyu Zhou, Mingxing Zhang, and Jiezhong Qiu. 2025. MoBA: Mixture of Block Attention for Long-Context LLMs. arXiv:2502.13189 [cs.LG] https://arxiv.org/abs/2502.13189

[8] Steffen Rendle. 2010. Factorization Machines. In *Proceedings of the 2010 IEEE International Conference on Data Mining (ICDM '10)*. IEEE Computer Society, USA, 995–1000. doi:10.1109/ICDM.2010.127

[9] 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 '19)*. ACM, 1161–1170. doi:10.1145/3357384.3357925

[10] Fangye Wang, Hansu Gu, Dongsheng Li, Tun Lu, Peng Zhang, and Ning Gu. 2023. Towards Deeper, Lighter and Interpretable Cross Network for CTR Prediction. arXiv:2311.04635 [cs.IR] https://arxiv.org/abs/2311.04635

[11] Ruoxi Wang, Bin Fu, Gang Fu, and Mingliang Wang. 2017. Deep & Cross Network for Ad Click Predictions. arXiv:1708.05123 [cs.LG] https://arxiv.org/abs/1708.05123

[12] Zhiqiang Wang, Qingyun She, and Junlin Zhang. 2021. MaskNet: Introducing Feature-Wise Multiplication to CTR Ranking Models by Instance-Guided Mask. *ArXiv* abs/2102.07619 (2021). https://api.semanticscholar.org/CorpusID:231924665

[13] Neil Zeghidour, Alejandro Luebs, Ahmed Omran, Jan Skoglund, and Marco Tagliasacchi. 2021. SoundStream: An End-to-End Neural Audio Codec. arXiv:2107.03312 [cs.SD] https://arxiv.org/abs/2107.03312

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

[15] Zhaoqi Zhang, Haolei Pei, Jun Guo, Tianyu Wang, Yufei Feng, Hui Sun, Shaowei Liu, and Aixin Sun. 2025. OneTrans: Unified Feature Interaction and Sequence Modeling with One Transformer in Industrial Recommender. arXiv:2510.26104 [cs.IR]

[16] Zhao-Yu Zhang, Xiang-Rong Sheng, Yujing Zhang, Biye Jiang, Shuguang Han, Hongbo Deng, and Bo Zheng. 2022. Towards Understanding the Overfitting Phenomenon of Deep Click-Through Rate Prediction Models. arXiv:2209.06053 [cs.IR] https://arxiv.org/abs/2209.06053

[17] Jie Zhu, Zhifang Fan, Xiaoxie Zhu, Yuchen Jiang, Hangyu Wang, Xintian Han, Haoran Ding, Xinmin Wang, Wenlin Zhao, Zhen Gong, Huizhi Yang, Zheng Chai, Zhe Chen, Yuchao Zheng, Qiwei Chen, Feng Zhang, Xun Zhou, Peng Xu, Xiao Yang, Di Wu, and Zuotao Liu. 2025. RankMixer: Scaling Up Ranking Models in Industrial Recommenders. arXiv:2507.15551 [cs.IR] https://arxiv.org/abs/2507.15551