

# From Scaling to Structured Expressivity: Rethinking Transformers for CTR Prediction

Bencheng Yan\*, Yuejie Lei\*, Zhiyuan Zeng\*, Di Wang, Kaiyi Lin, Pengjie Wang, Jian Xu, Bo Zheng<sup>†</sup>

Alibaba Group

{bencheng.ybc, leiyuejie.lyj, zengzhiyuan.zzy, zhemu.wd, linkaiyi.lky,  
pengjie.wpj, xiyu.xj, bozheng}@alibaba-inc.com

China

## Abstract

Despite massive investments in scale, deep models for click-through rate (CTR) prediction often exhibit rapidly diminishing returns—a stark contrast to the smooth, predictable gains seen in large language models. We identify the root cause as a *structural misalignment*: Transformers assume sequential compositionality, while CTR data demand combinatorial reasoning over high-cardinality semantic fields. Unstructured attention spreads capacity indiscriminately, amplifying noise under extreme sparsity and breaking scalable learning. To restore alignment, we introduce the **Field-Aware Transformer (FAT)**, which embeds field-based interaction priors into attention through decomposed content alignment and cross-field modulation. This design ensures model complexity scales with the number of fields  $F$ , not the total vocabulary size  $n \gg F$ , leading to tighter generalization and, critically, *observed power-law scaling* in AUC as model width increases. We present the first formal scaling law for CTR models, grounded in Rademacher complexity, that explains and predicts this behavior. On large-scale benchmarks, FAT improves AUC by up to **+0.51%** over state-of-the-art methods. Deployed online, it delivers **+2.33% CTR** and **+0.66% RPM**. Our work establishes that effective scaling in recommendation arises not from size, but from *structured expressivity*—architectural coherence with data semantics.

## 1 Introduction

The success of large language models (LLMs) has revealed a powerful truth: when architecture and data are aligned, scaling becomes predictable. As model size, data volume, and compute increase, performance improves smoothly—governed by empirical scaling laws that enable systematic progress [8, 11–13, 21]. This principle has inspired widespread efforts to transplant Transformer architectures into industrial recommendation systems, particularly for click-through rate (CTR) prediction [3, 7, 9, 17, 19, 20, 23].

Yet, despite promising gains, most existing approaches remain at the level of **architectural mimicry**: they directly transplant

LLM designs by tokenizing CTR features and applying standard Transformers, either within traditional pointwise prediction frameworks [1, 3, 20, 23] or under generative reformulations [2, 9, 19]. While these methods benefit from increased model capacity, empirical studies have observed diminishing returns in performance as models scale up [5], suggesting a disconnect between architectural growth and effective learning.

We argue that this divergence stems not from insufficient resources, but from a **fundamental structural misalignment** between the assumptions embedded in standard Transformers and the nature of CTR data. At first glance, both modalities involve sequences of discrete tokens. But their semantic structures are profoundly different:

- In language, meaning emerges through **compositional syntax**: words combine hierarchically under grammatical rules, forming phrases and sentences whose semantics depend on order and context. Tokens are *homogeneous*—drawn from a shared vocabulary—and attend based on syntactic role and semantic similarity.
- In CTR prediction, predictive power arises from **combinatorial interactions**: user’s click behavior is driven by cross-field conjunctions such as “young user  $\times$  luxury brand” or “mobile device  $\times$  evening session”. Inputs are *heterogeneous sets* of high-cardinality categorical features, each belonging to a semantic *field* (e.g., `user_age`, `ad_category`, `device_type`). Order is arbitrary; what matters is *which fields interact*, and how asymmetrically.

Standard self-attention, designed for compositional semantics over dense, ordered sequences, fails to respect this distinction. It treats all embeddings uniformly via globally shared projection matrices, spreading representational capacity indiscriminately across feature types. Under extreme sparsity—where most field-value combinations are rarely observed—this unstructured attention amplifies noise, distorts gradients, and ultimately breaks scalable learning.

Even more troubling is the absence of a theoretical foundation for scaling in recommendation. While LLMs benefit from well-characterized generalization bounds and scaling laws grounded in statistical learning theory [8, 11, 21], no such framework exists for CTR models. Without it, scaling becomes a trial-and-error process, disconnected from architectural design principles.

This leads us to a pivotal question: *Can we redesign the Transformer for recommendation such that its expressive capacity grows in harmony with the underlying interaction complexity of the data—not merely in raw parameter count, but in structured expressivity?*

To answer this, we return to a classical insight: **field-aware interaction modeling**. Models like Field-aware Factorization Machines (FFM) [6] assign dedicated latent vectors to each ordered

\* These authors contributed equally to this work and are co-first authors.

<sup>†</sup> Corresponding author.

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than the author(s) must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from permissions@acm.org.  
Conference’17, Washington, DC, USA

© 2025 Copyright held by the owner/author(s). Publication rights licensed to ACM.  
ACM ISBN 978-x-xxxx-xxxx-x/YYYY/MM  
<https://doi.org/10.1145/nnnnnnnn.nnnnnnnn>

field pair  $(f_i, f_j)$ , enabling asymmetric, context-sensitive modeling of interactions (e.g., how `user_gender` influences response to `ad_category`). However, FFM is shallow and static, incapable of deep composition or contextual refinement.

Inspired by this principle, a natural architectural extension would be to make the query, key, and value projections in attention *specialized per field pair*, allowing each head to capture distinct interaction patterns between source and target fields. While conceptually appealing, such a naïve realization incurs prohibitive parameter growth: with  $F$  semantic fields, the number of interaction-specific parameters scales quadratically as  $O(F^2 d^2)$ . In real-world systems where  $F \sim 10^3$  and  $d \sim 128$ , even a moderate base model can balloon from 100 million to over **10 trillion parameters**—rendering it infeasible for training or deployment.

To resolve this tension between expressivity and scalability, we introduce the **Field-Aware Transformer (FAT)**, a novel architecture built upon two synergistic mechanisms. First, **Field-Decomposed Attention** factorizes the full field-pair-specialized transformation into two components: (i) *field-aware content alignment*, where queries and keys are projected using matrices specific to their own fields (scaling as  $O(Fd^2)$ ), and (ii) *field-pair interaction modulation*, which governs information flow between field pairs via lightweight scalars (scaling as  $O(F^2)$ ). This decomposition ensures that model complexity grows with the number of fields  $F$ , not the total vocabulary size  $n \gg F$ . Second, to further decouple model capacity from fields size and eliminate storage overhead, FAT employs a **Hypernetwork-Based Generation** mechanism: field-specific parameters are dynamically synthesized from a compact set of basis matrices through lightweight neural modules.

Critically, we establish the **first theoretically principled scaling law for CTR models**, derived via Rademacher complexity analysis. We show that FAT’s generalization error depends on the combinatorial structure of field interactions—specifically, the number of fields—rather than the vocabulary size. This structural alignment explains why FAT exhibits smooth, power-law scaling in AUC as model width increases, a phenomenon absent in baseline Transformers.

Our contributions are:

- **Architecture:** We propose FAT, a Transformer variant that integrates field-aware priors into attention, enabling structured, interpretable, and scalable modeling of combinatorial semantics.
- **Scalability Mechanism:** We design a parameter-efficient instantiation through Field-Decomposed Attention and Hypernetwork-Based Generation
- **Theory:** We present the first generalization-aware scaling law for CTR models, proving that effective scaling requires architectural coherence with data structure—not just larger models.
- **Empirical Validation:** On large-scale benchmarks, FAT achieves up to +0.51% AUC improvement over state-of-the-art methods. Deployed online, it delivers +2.33% CTR and +0.66% RPM, demonstrating significant business impact.

This work establishes that scalable performance in CTR prediction does not arise from size alone, but from **structured expressivity**—the deliberate alignment of architectural design with the combinatorial semantics of feature interactions. FAT moves beyond mere architectural mimicry of LLMs, offering a principled

path toward building larger, more capable, and more predictable recommender systems.

## 2 Related Work

We structure the discussion around two central challenges in modern CTR modeling: (1) capturing structured feature interactions, and (2) achieving predictable scaling. Our work addresses a critical gap between these two goals.

### 2.1 Modeling Structured Feature Interactions

Effective CTR models must capture high-order, asymmetric interactions across semantic fields. Factorization Machines (FM) [14] introduced low-rank pairwise modeling, later refined by FFM [6], which assigns field-pair-specific latent vectors to model context-sensitive effects. These models are interpretable and parameter-efficient but limited to shallow interactions. Neural extensions such as DeepFM [4], AutoInt [15], and DCNv2 [18] use MLPs or attention to learn complex patterns. A persistent limitation across these approaches is that structured interaction modeling remains confined to shallow architectures—typically only a few layers—preventing deep compositionality and making principled scaling infeasible.

### 2.2 Towards Predictable Scaling in Recommendation

Recent efforts adopt Transformer architectures for CTR prediction [3, 7, 19, 20, 23], inspired by their success in language modeling. However, standard self-attention operates under assumptions—sequential order, dense tokens, compositional syntax—that do not hold in recommendation, where inputs are unordered, sparse sets with combinatorial semantics. Applying unstructured attention may lead to inefficient representation learning and poor generalization under sparsity.

Even more fundamentally, unlike large language models, which follow well-characterized scaling laws [8], CTR models often exhibit performance saturation or degradation upon scaling [19]. This suggests a misalignment between architectural capacity and data structure, and none connect model design to scalable behavior in a principled way.

These gaps reveal a key challenge: how to build models that simultaneously support expressive, interpretable interaction modeling and predictable, stable scaling. Prior work falls short on at least one dimension. Our work shows that bridging this divide requires embedding domain-aware inductive biases directly into the architecture’s core computation.

## 3 Method

As illustrated in Figure 1, the Field-Aware Transformer (FAT) achieves *structured expressivity* by decomposing attention along semantic fields: query, key, and value projections are specialized per field pair  $(f_i, f_j)$  through *field-aware content alignment* and *field-pair interaction modulation*—reducing parameter dependency from  $O(F^2 d^2)$  to  $O(Fd^2 + F^2)$ . This structural prior ensures that the effective model complexity scales with the number of fields  $F$ , not vocabulary size  $n \gg F$ , leading to tighter generalization bounds under extreme sparsity.

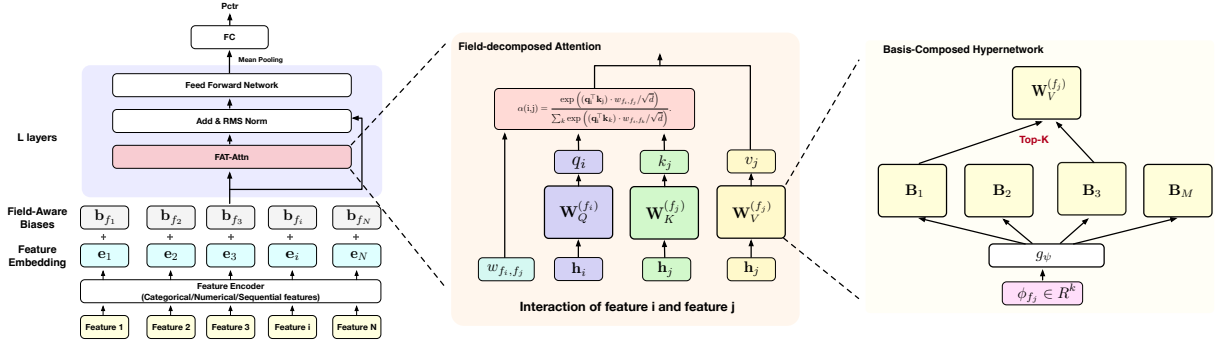


Figure 1: The architecture of FAT.

To decouple model capacity from field cardinality, FAT employs a hypernetwork to generate field-specific parameters, eliminating storage overhead while incurring zero inference cost. We next detail the architecture and parameterization.

### 3.1 Structured Tokenization via Field-Aware Representation Learning

In CTR prediction, inputs are unordered sets of heterogeneous features—categorical, numerical, and sequential—drawn from distinct semantic fields (e.g., user, item, context). Unlike language, where sequential compositionality drives meaning, the semantics of CTR data arise from combinatorial interactions across fields. To enable structured interaction modeling, we map each feature  $x_i$  to a unified representation space  $\mathbb{R}^d$ , producing an embedded token  $e_i$  via type-specific transformations:

- Categorical features use embedding lookup tables;
- Numerical features are processed through monotonic networks or quantile discretization with embedding lookup tables;
- Sequential features are summarized via dedicated encoders (e.g., DIN [22]) into field-level representations.

All embeddings are projected into a common dimension  $d$  for compatibility. Crucially, since input order is arbitrary, we replace index-based positional encodings with *field-aware biases* that reflect semantic roles rather than syntactic positions. For token  $i$  belonging to field  $f_i$ , its final input representation is:

$$h_i = e_i + b_{f_i}, \quad (1)$$

where  $b_{f_i} \in \mathbb{R}^d$  is a learnable bias vector specific to field  $f_i$ . This injects structural priors aligned with field semantics, ensuring consistent generalization under permutation.

The resulting sequence  $\mathbf{H} = [h_1, \dots, h_N]$  forms a semantically grounded, field-identified token stream, laying the foundation for interpretable and scalable interaction modeling in FAT.

### 3.2 From Standard Transformer to Field-Pair-Specialized Attention

In standard Transformers, attention  $\alpha(i, j)$  is computed as:

$$\alpha(i, j) = \frac{(h_i \mathbf{W}_Q)(\mathbf{W}_K^T h_j^T)}{\sqrt{d}}, \quad (2)$$

using globally shared  $\mathbf{W}_Q, \mathbf{W}_K \in \mathbb{R}^{d \times d}$ . This treats all tokens uniformly, ignoring their field origins—a critical limitation when interactions are asymmetric and context-sensitive.

Inspired by FFM, which models pairwise field interactions via dedicated latent vectors, we consider a natural generalization: make query, key, and value projections specialized per ordered field pair  $(f_i, f_j)$ . That is, for each head  $h$ , define:

$$q_i = h_i \mathbf{W}_Q^{(f_i, f_j)}, \quad k_j = h_j \mathbf{W}_K^{(f_i, f_j)}, \quad v_j = h_j \mathbf{W}_V^{(f_i, f_j)}. \quad (3)$$

Then compute:

$$\alpha(i, j) = \frac{q_i^T k_j}{\sqrt{d}}. \quad (4)$$

This design enables fine-grained control: the way a query in `user_gender` attends to `ad_category` can differ fundamentally from how it attends to `device_type`. It captures asymmetry, context sensitivity, and field-role awareness—directly reflecting the combinatorial semantics of CTR data. However, a naïve realization would require  $O(HF^2d^2)$  parameters per layer ( $H \sim 8$  refers to the number of attention head)—an infeasible cost when  $F \sim 10^3$  and  $d \sim 128$ . Even a base model could exceed tens of billions of parameters. Thus, while conceptually ideal, full field-pair specialization is impractical for industrial deployment. We now derive a more scalable instantiation that preserves its semantic essence.

### 3.3 Filed-decomposed Attention: Field-Aware Content Alignment with Field-Pair Interaction Modulation

To retain the expressive power of field-pair modeling while ensuring scalability, we decompose the attention mechanism into two complementary components:

- **Field-aware content alignment:** how well two tokens interact, given their respective semantic roles;
- **Field-pair interaction modulation:** how strongly information should flow from one field to another.

This decomposition allows us to preserve fine-grained interaction modeling at a fraction of the cost. Specifically, we define the attention score as:

$$\alpha(i, j) = (\mathbf{q}_i^\top \mathbf{k}_j) \cdot w_{f_i, f_j}, \quad (5)$$

$$\text{where } \mathbf{q}_i = \mathbf{W}_Q^{(f_i)} \mathbf{h}_i, \quad \mathbf{k}_j = \mathbf{W}_K^{(f_j)} \mathbf{h}_j. \quad (6)$$

Here, the term  $\mathbf{q}_i^\top \mathbf{k}_j$  is no longer a generic similarity measure—it is a field-aware content alignment, because both query and key are projected using matrices that depend on their source fields:  $\mathbf{W}_Q^{(f_i)}$  tailors the transformation of token  $i$  based on its role as a member of field  $f_i$  (e.g., user-side vs. item-side encoding);  $\mathbf{W}_K^{(f_j)}$  similarly adapts the key for field  $f_j$ 's semantics.

Thus, even before considering cross-field strength, the model can distinguish whether an age value should be encoded differently when attending from a user context versus a contextual ad signal. This represents a coarse-grained but essential form of field awareness: each field learns its own "encoding style", enabling semantically grounded comparisons across tokens.

The scalar parameter  $w_{f_i, f_j} \in \mathbb{R}$  serves as a field-pair modulation factor, adjusting the attention weight based on the source and target fields. This enables fine-grained control over information flow between semantic fields. For example, a large  $w_{\text{ad\_category}, \text{user\_behavior}}$  amplifies behavior-category sensitivity, while a small  $w_{\text{ad\_category}, \text{device\_type}}$  suppresses irrelevant couplings.

Together, this factorization separates concerns:

- $\mathbf{q}_i^\top \mathbf{k}_j$ : captures what the tokens represent, under field-specific interpretation;
- $w_{f_i, f_j}$ : controls whether and how strongly such representations should interact.

The value vector follows the same principle:

$$\mathbf{v}_j = \mathbf{W}_V^{(f_j)} \mathbf{h}_j, \quad \mathbf{W}_V^{(f_j)} \in \mathbb{R}^{d \times d},$$

ensuring output transformations are also aligned with field semantics.

Multi-head attention aggregates outputs as:

$$\text{FAT-Attn}_h(i) = \sum_j \alpha_h(i, j) \mathbf{v}_j^{(h)}, \quad (7)$$

$$\alpha_h(i, j) = \frac{\exp\left((\mathbf{q}_{i,h}^\top \mathbf{k}_{j,h}) \cdot w_{f_i, f_j}^{(h)} / \sqrt{d}\right)}{\sum_k \exp\left((\mathbf{q}_{i,h}^\top \mathbf{k}_{k,h}) \cdot w_{f_i, f_k}^{(h)} / \sqrt{d}\right)}. \quad (8)$$

This formulation offers several advantages:

- **Hierarchical Field Awareness:** content alignment operates at the **per-field level** (coarse), encoding semantic roles; interaction modulating acts at the **field-pair level** (fine), governing routing strength.
- **Interpretability:**  $w_{f_i, f_j}$  provide insights into interaction patterns.
- **Asymmetry:**  $w_{f_i, f_j} \neq w_{f_j, f_i}$  naturally models directional effects.
- **Efficiency and Scalability:** Naïve field-pair projections require  $O(F^2 d^2)$  parameters per head—an infeasible burden. In contrast, FAT uses only  $O(F d^2 + F^2)$ , reducing parameter count by over 99% in typical settings (e.g., from 16B to 50M). This enables training large models without sacrificing semantic fidelity.

By separating field-aware representation from field-pair-governed routing, FAT achieves structured expressivity: it scales not by adding

unstructured capacity, but by deepening semantic resolution in a controlled, interpretable way.

### 3.4 Basis-Composed Hypernetwork for Scalable Parameter Generation

Despite its expressiveness, storing  $3F$  ( $F \sim 10^3$ ) full-rank matrices  $\{\mathbf{W}_Q^{(f)}, \mathbf{W}_K^{(f)}, \mathbf{W}_V^{(f)}\}$  incurs prohibitive memory and maintenance costs in large-scale systems, especially when both field count  $F$  and embedding dimension  $d$  are large. To decouple parameter growth from field cardinality while preserving semantic fidelity, we introduce a *basis-composed hypernetwork* that dynamically generates field-specific projections.

Let  $\mathcal{B} = \{B_1, \dots, B_M\}$  be a shared set of  $M$  basis matrices,  $B_m \in \mathbb{R}^{d \times d}$ , representing canonical linear transformations (e.g., scaling, rotation). These are learned end-to-end and reused across fields. For each field  $f$ , a meta-embedding  $\phi_f \in \mathbb{R}^k$  is passed through a lightweight MLP  $g_\psi : \mathbb{R}^k \rightarrow \mathbb{R}^M$ :

$$s^{(f)} = g_\psi(\phi_f).$$

We apply Top- $K$  sparse selection ( $K \ll M$ ):

$$\pi_f = \text{top-}K(s^{(f)}), \quad \alpha_m^{(f)} = \frac{\exp(s_m^{(f)})}{\sum_{m' \in \pi_f} \exp(s_{m'}^{(f)})}, \quad \forall m \in \pi_f,$$

and synthesize the query projection as:

$$\mathbf{W}_Q^{(f)} = \sum_{m \in \pi_f} \alpha_m^{(f)} B_m.$$

Analogous constructions yield  $\mathbf{W}_K^{(f)}$  and  $\mathbf{W}_V^{(f)}$  using separate basis sets.

The hypernetwork enables scalable architecture design with two key advantages:

- **Decoupled complexity from field growth:** Field-specific projections are synthesized without explicit storage, breaking the  $O(F d^2)$  parameter dependency. Model capacity grows via shared bases ( $O(M d^2)$ ) and compact meta-embeddings ( $O(F k)$ ), enabling stable and efficient scaling as  $F$  grows—essential for industrial applications where new features are continuously introduced.
- **Zero inference overhead:** All generated parameters — including  $\mathbf{W}_Q^{(f)}$ ,  $\mathbf{W}_K^{(f)}$ , and  $\mathbf{W}_V^{(f)}$  — are precomputed and cached after training. No hypernetwork evaluation is performed at serving time, ensuring full compatibility with low-latency production pipelines.

### 3.5 CTR Prediction

FAT stacks  $L$  layers of field-decomposed attention, each followed by a feed-forward network (FFN), residual connection, and layer normalization:

$$\mathbf{Z}^{(\ell+1)} = \text{FFN}\left(\text{LayerNorm}\left(\text{FAT-Attn}(\mathbf{Z}^{(\ell)})\right)\right) + \mathbf{Z}^{(\ell)}. \quad (9)$$

A pooled output yields the final prediction:

$$p(y = 1 | \mathbf{X}) = \sigma\left(\mathbf{w}^\top \sum_{i=1}^N \text{FAT-Output}_i^{(L)}\right), \quad (10)$$

where  $\sigma(\cdot)$  is the sigmoid function,  $\mathbf{w} \in \mathbb{R}^d$  is a learnable weight vector, and  $\text{FAT-Output}_i^{(L)}$  is the final-layer representation of token  $i$ .

#### 4 Theoretical Implications: A Principled Scaling Law

One of the most sought-after goals in recommendation systems is achieving predictable, smooth performance improvement as model size increases—a phenomenon well-documented in large language models through empirical scaling laws [8]. However, in CTR prediction, naive scaling often leads to performance saturation or degradation due to unstructured capacity growth that amplifies noise rather than signal.

We show that FAT enables *principled* scaling by aligning architectural design with the combinatorial semantics of feature interactions. Specifically, we prove that the effective model complexity of FAT depends on the number of semantic fields  $F$  and the interaction structure rank, not on the total vocabulary size  $n = \sum_{i=1}^F |\mathcal{V}_i|$ , which can exceed  $10^9$  in industrial settings. This structural alignment yields tighter generalization bounds and paves the way for a well-behaved scaling law.

Our main theoretical result establishes an upper bound on the generalization error of a single FAT attention layer. The complete proof, based on Rademacher complexity analysis, is provided in Appendix A.

**THEOREM 4.1 (GENERALIZATION BOUND FOR FAT).** *Let  $\mathcal{D}$  be a distribution over input sequences  $\mathbf{H} = [\mathbf{h}_1, \dots, \mathbf{h}_N]$  with  $\|\mathbf{h}_i\|_2 \leq R$ . Assume all parameter matrices  $(\mathbf{W}_Q^{(f)}, \mathbf{W}_K^{(f)}, \mathbf{W}_V^{(f)})$  have Frobenius norm bounded by  $B$ , and all interaction scalars  $(w_{f_i, f_j})$  are bounded by  $B_w$ . Let  $m$  be the number of training samples. Then, with probability at least  $1 - \delta$ , the generalization error  $L_{\text{gen}}$  of a single FAT layer satisfies:*

$$L_{\text{gen}} \leq L_{\text{train}} + O\left(\frac{\sqrt{Fd^2 + F^2}}{\sqrt{m}} \cdot C(R, B, B_w, d) + \sqrt{\frac{\log(1/\delta)}{m}}\right),$$

where  $C(R, B, B_w, d) = O(R^2 B^2 B_w \sqrt{d} + R B B_w)$  is a constant depending on the norm bounds and embedding dimension.

This bound reveals a fundamental advantage over standard Transformers. Standard self-attention has a hypothesis space complexity that implicitly scales with  $O(nd^2)$  due to its ability to form arbitrary token-level interactions, making it highly susceptible to overfitting under extreme data sparsity. In contrast, FAT restricts expressive power to semantically valid pathways—only  $F$  field-specific transformations and  $F^2$  cross-field modulations are learned. This structural constraint reduces the effective hypothesis space from  $\text{poly}(n)$  to  $\text{poly}(F)$ , where  $F \ll n$  (e.g.,  $F \sim 10^3$  vs.  $n \sim 10^9$ ). Consequently, every parameter is shared across vast amounts of data within its field, dramatically improving statistical efficiency and mitigating overfitting.

Critically, this tight generalization bound enables *predictable scaling*. For a fixed data distribution and field schema  $F$ , increasing the embedding dimension  $d$  (and thus the total parameter count  $N_{\text{params}} \propto Fd^2$ ) deepens the model’s representational fidelity for both field-aware content alignment  $(\mathbf{W}_Q^{(f)}, \mathbf{W}_K^{(f)})$  and field-pair interaction modulation  $(w_{f_i, f_j})$ . This allows the model to learn higher-rank, more nuanced interaction functions between fields, thereby

systematically reducing the training error  $L_{\text{train}}$ . Because the architecture is aligned with the data’s combinatorial structure, these additional parameters refine meaningful patterns instead of fitting noise.

Combining these two effects—reduced bias from enhanced expressivity and reduced variance from tight generalization—we arrive at a principled scaling law. Under fixed data volume  $m$  and field schema  $F$ , as we scale the model width, the test performance improves according to a power-law trend:

$$\Delta\text{AUC} \propto N_{\text{params}}^\beta, \quad \beta > 0, \quad (11)$$

where  $\Delta\text{AUC}$  is the performance gain relative to a baseline. This provides a theoretical justification for a power-law scaling trend observed empirically in CTR models, linking architectural design to predictable scalability. Empirical validation of this law is presented in Section 5.5.

## 5 Experiments

We conduct comprehensive experiments to rigorously evaluate the effectiveness, interpretability, scalability, and deployability of our proposed Field-Aware Transformer (FAT). We aim to answer five key research questions:

- **RQ1:** Does FAT outperform state-of-the-art CTR models in prediction accuracy?
- **RQ2:** How do individual components contribute to performance? Is the decomposition valid?
- **RQ3:** Can FAT provide interpretable insights into field-level interactions?
- **RQ4:** Does FAT exhibit a well-behaved scaling law as predicted by theory?
- **RQ5:** When deployed in a real-world online system, does FAT improve business metrics in A/B testing?

### 5.1 Experimental Setup

**5.1.1 Dataset.** We evaluate on a large-scale CTR dataset from **Taobao’s sponsored search**, containing over **14 billion user impressions** collected over two weeks. The input includes hundreds of heterogeneous features—categorical (e.g., user/item IDs), numerical (e.g., CTR, dwell time), and sequential (e.g., behavior trails)—spanning billions of users and items. This setting captures the core challenges of industrial recommendation: extreme scale, high sparsity, and complex cross-field interactions.

**5.1.2 Baselines.** We compare against representative state-of-the-art methods, categorized by modeling paradigm:

- **Traditional Interaction Models:** FFM [6], DeepFM [4], AutoInt [16], and DCNv2 [18] which emphasize structured or learned pairwise feature crosses; Besides, we take the vanilla Embedding+MLP for feature crossing as a strong baseline, denoted as DeepCTR.
- **Scaling-Oriented Architectures:** HiFormer [3], Wukong [20], HSTU [19], and RankMixer [23], representing recent advances in scalable and adaptive model design.

This selection spans the evolutionary spectrum of CTR modeling—from field-aware factorization to large-scale representation learning—enabling a principled assessment of whether structured

field awareness can surpass both classical interaction mechanisms and modern scaling paradigms.

**5.1.3 Implementation Details.** All experiments are conducted on a distributed training system with 128 NVIDIA GPUs using synchronous data-parallel SGD. Models are implemented in TensorFlow and optimized using the Adam optimizer ( $\beta_1 = 0.9$ ,  $\beta_2 = 0.999$ ). The global batch size is 262, 144 (2048 per GPU), and the initial learning rate is tuned in  $\{1e-4, 3e-4, 5e-4, 1e-3\}$  for each model and scale configuration.

To enable fair and insightful comparisons, we adopt a **dual-scale evaluation protocol** that aligns with architectural design principles:

- **Traditional Interaction Models** (DeepCTR, FFM, DeepFM, AutoInt, DCNv2): Evaluated at their typical capacity of  $\sim 50$  M parameters. All hyperparameters—including embedding dimension, hidden sizes, dropout (0.1–0.5), L2 regularization ( $1e-6$  to  $1e-3$ ), and network depth—are optimized via Bayesian search on the validation set to ensure peak performance.
- **Scaling-Oriented Architectures** (HiFormer, Wukong, HSTU, RankMixer): Scaled uniformly to approximately **0.5B parameters** by adjusting width or depth while preserving core architectural constraints (e.g., expert count in RankMixer, hierarchy levels in HSTU). Hyperparameters are re-tuned under this fixed capacity regime.

We define three instantiations of FAT according to the above protocol:

- **FAT-Small**:  $\sim 50$ M parameters, evaluated against traditional interaction models;
- **FAT-Large**:  $\sim 0.5$ B parameters, compared with scaling-oriented architectures;
- **FAT-XL**:  $\sim 1.5$ B parameters, included exclusively for analyzing scaling trends.

All models share identical: (1) Feature preprocessing: categorical features use hashed vocabularies; numerical features are discretized into different bins; sequential features (e.g., user behavior history) are processed via a shared DIN-style interest extractor; (2) Embedding dimension: fixed at  $d \in \{8, 16, 32, 64\}$  for different features across all models;

FAT-specific configurations: multi-head attention uses  $H = 8$  heads; field meta-embeddings  $\phi_f \in \mathbb{R}^{64}$  are randomly initialized and shared across layers; the hypernetwork employs  $M = 64$  shared basis matrices with Top- $K = 3$  sparse activation; field-pair interaction scalars  $w_{f_i, f_j}$  are initialized from  $\mathcal{N}(0, 0.01)$ .

## 5.2 Main Results (RQ1): Superior Predictive Performance

This section evaluates whether FAT achieves superior performance compared to state-of-the-art CTR models—and, more critically, *why*. Rather than reporting isolated gains, we design a rigorous comparison framework to disentangle architectural superiority from parameter inflation. The central question is: does FAT perform better because of its structural design, or simply due to larger capacity?

To answer this, we adopt a **dual-scale evaluation framework** that enables three key inquiries:

- (1) *Does FAT outperform existing models?*

**Table 1: CTR Prediction Performance**

| Type                           | Method        | $\Delta$ AUC (%) | #Param |
|--------------------------------|---------------|------------------|--------|
| Base                           | DeepCTR       | -                | 40M    |
| Traditional Interaction Models | DeepCTR-Large | +0.02            | 0.48B  |
|                                | FFM           | -0.3             | 25M    |
|                                | DeepFM        | -0.15            | 45M    |
|                                | AutoInt       | +0.05            | 57M    |
|                                | DCNv2         | +0.06            | 47M    |
| Scaling-Oriented Architectures | Wukong        | +0.11            | 0.54B  |
|                                | HSTU          | +0.1             | 0.54B  |
|                                | HiFormer      | +0.17            | 0.58B  |
|                                | RankMixer     | +0.23            | 0.51B  |
| Ours                           | FAT-Small     | +0.13            | 52M    |
|                                | FAT-Large     | +0.41            | 0.54B  |
|                                | FAT-XL        | <b>+0.51</b>     | 1.5B   |

- (2) *Is the improvement attributable to better architecture, not just more parameters?*
- (3) *Does its relative advantage persist—or even grow—under scale?*

As shown in Table 1, we report the  $\Delta$ AUC over the baseline (DeepCTR).

At the **small-scale regime** ( $\sim 50$ M parameters), **FAT-Small** outperforms traditional interaction models—including DeepCTR, FFM, DeepFM, AutoInt, and DCNv2. This is particularly significant as these methods have reached their performance ceiling at this capacity, with further scaling yielding diminishing returns. For example, we also scaled traditional interaction models to 0.5B parameters (e.g., DeepCTR-Large), but observed no further improvement—indicating early saturation of performance. FAT’s gain under identical parameter constraints demonstrates that its field-aware attention decomposition enables more effective use of limited model capacity.

At the **large-scale regime** ( $\sim 0.5$ B parameters), **FAT-Large** surpasses scalable architectures (including HiFormer, Wukong, HSTU, and RankMixer), despite matching their size exactly. All baselines are uniformly scaled in width/depth while preserving architectural integrity, and hyperparameters are re-tuned for fairness. The persistent lead of FAT indicates that its performance edge stems not from scale, but from *superior inductive bias*: the separation of field-aware content alignment and interaction-aware routing allows more semantically grounded and stable information flow.

Notably, when further scaled to 1.5B parameters (FAT-XL), performance continues to improve without signs of saturation—a preliminary indication of favorable scaling behavior (to be analyzed in depth in Section 5.5).

In summary, FAT consistently outperforms both classical and modern CTR models under fair capacity control. Its gains are not artifacts of parameter count, but consequences of *structured expressivity*. This establishes FAT as a fundamentally stronger architecture—one whose advantages are rooted in design, not just scale.

## 5.3 Ablation Study (RQ2): Component-wise Analysis

To understand the sources of FAT’s gains, we conduct ablation studies measuring the *relative AUC improvement* over the baseline

**Table 2: Ablation study on FAT components. Values show relative AUC improvement (%) over the baseline (DeepCTR).**

| Variant   | $\Delta$ AUC                  |
|---|-------------------------------|
| Full FAT  | <b>+0.41</b>                  |
| w/o Field-Aware Biases  | +0.35                         |
| w/o Interaction Modulation ( $w_{f_i, f_j} = 1$ )                         | +0.29                         |
| w/o Field-Aware Content Alignment ( $\mathbf{W}_Q^{(f)} = \mathbf{W}_Q$ ) | +0.24                         |
| w/o Hypernetwork (replace it with Full Matrices)                          | +0.38, but $+5 \times$ params |
| Naïve Field-Pair Attention (full $\mathbf{W}_{Q, (f_i, f_j)}$ )           | OOM ( $>150\text{B}$ params)  |

**Table 3: Asymmetry in  $w_{f_i, f_j}$ : Item as query attends more selectively than user interest.**

| Query Field ( $f_i$ ) | Key Field ( $f_j$ ) |             |              |
|-----------------------|---------------------|-------------|--------------|
|                       | item                | rec_clicks  | user_profile |
| item                  | -                   | <b>0.97</b> | 0.16         |
| rec_clicks            | 0.23                | -           | 0.24         |
| user_profile          | 0.24                | 0.32        | -            |

(DeepCTR) (Table 2). All variants maintain similar training setups, enabling controlled comparison of design choices.

The full FAT (FAT-Large) achieves the highest gain (+0.41), validating the effectiveness of our overall design. Removing field-aware biases (w/o field-aware biases) results in a minor drop to +0.34, indicating that while structural priors help, they are not the primary driver of performance.

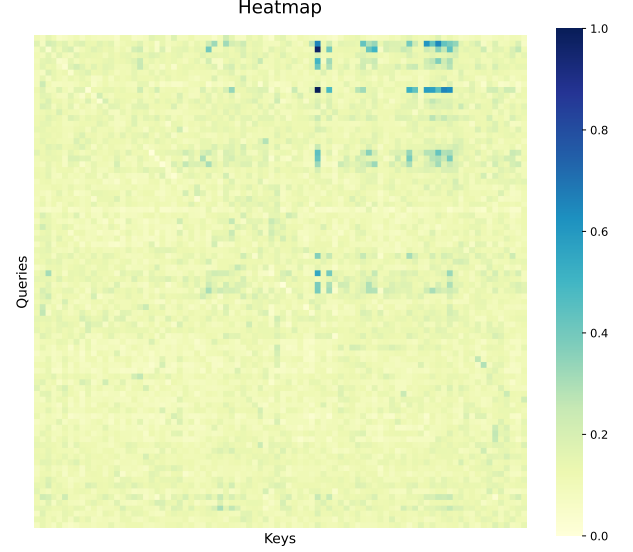
Crucially, decomposing the attention mechanism reveals a clear hierarchy of contributions:

- *Field-Aware Content Alignment*—modeling field-specific query projections ( $\mathbf{W}_Q^{(f)}$ )—is the most critical component. Removing it drops performance to +0.24, a **loss of 0.17 percentage points**, larger than any other single ablation. This shows that *early specialization by field role* (e.g., user, item, context) forms the foundation for meaningful interaction modeling.
- *Interaction Modulation*—using field-pair-specific scalars  $w_{f_i, f_j}$ —also contributes, reducing gain to +0.29 when removed (−0.12). While finer-grained, its impact is secondary, suggesting that asymmetric interaction strength matters, but only after content representations are properly aligned.

For scalability, replacing the hypernetwork with full parameter matrices yields similar performance (+0.38) at  $5\times$  parameter cost, demonstrating that dynamic generation achieves near-optimal expressiveness with minimal storage.

Most significantly, a naïve field-pair-specialized attention—without parameter decomposition—fails entirely due to memory overflow ( $>150\text{B}$  parameters), highlighting that unstructured capacity growth is infeasible in real systems.

These results reveal a fundamental insight: the power of FAT lies not in dense interaction modeling, but in *structured sparsity*—where early field-aware specialization enables rich semantics, and efficient parameter generation ensures scalability. This synergy defines *structured expressivity*: expressive not by size, but by design.

**Figure 2: Heatmap of  $w_{f_i, f_j}$ . Strong interactions (dark) align with expected semantic dependencies (e.g., item and real-time interest).**

#### 5.4 Interpretability Analysis (RQ3): Uncovering Semantic Interaction Patterns

To understand how FAT captures semantic interactions, we analyze the learned modulation weights  $w_{f_i, f_j}$ , which control the strength of information flow from field  $f_j$  to  $f_i$ . These parameters are shared across tokens within the same field pair and reflect global interaction patterns. We visualize the average  $w_{f_i, f_j}$  across all attention heads in Figure 2, revealing two key properties: **structured coherence** and **asymmetric influence**.

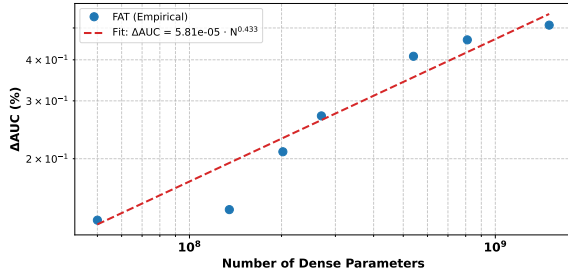
**1. Structured Coherence.** The weight matrix is sparse and exhibits clear block-wise structure. High values concentrate on semantically meaningful pairs:

- Candidate item features (e.g., `item_cate`, `shop_level`) show strong connections to real-time user signals (e.g., `recent_clicks`).
- User profile fields (e.g., `age`, `gender`) interact most strongly with long-term preference indicators (e.g., `fav_brands`, `longterm_clicks`).

In contrast, cross-field interactions between unrelated modalities (e.g., `device_type`  $\rightarrow$  `income`) remain near zero. This pattern matches business prior: short-term intent drives item relevance, while static profiles shape stable preferences. FAT naturally learns this separation without explicit supervision, demonstrating that its structured expressivity focuses capacity on meaningful pathways.

**2. Asymmetric Influence.** The interaction strength is highly directional. As shown in Table 3, when an item feature acts as the query, it assigns high weights to recent user behaviors (**0.97**), indicating strong reliance on current intent. In contrast, when a recent behavior feature (e.g., `rec_clicks`) serves as the query, its outgoing weights are significantly weaker and more diffuse (e.g.,  $w_{\text{rec\_clicks}, \text{item}} = 0.23$ ), suggesting limited predictive power for modeling other features. This asymmetry reflects the functional





**Figure 3: Power-law relationship between parameter count and AUC. Best-fit slope  $\Delta AUC = 5.81 \times 10^{-5} \cdot N_{params}^{0.433}$ .**

role of features: items actively retrieve context, while behavioral embeddings serve as supporting evidence.

These results confirm that FAT learns interpretable and semantically plausible interaction patterns. Unlike standard Transformers that spread attention widely, FAT’s field-aware design ensures that interaction strengths are both *structured* and *directional*, aligning with real-world user behavior and providing actionable insights for model debugging and refinement.

### 5.5 Scaling Behavior (RQ4): Validating the Theoretical Law

We examine whether FAT exhibits a predictable and sustained scaling behavior, as suggested by its theoretical generalization bound (Theorem 4.1). Specifically, we investigate if performance improves systematically with model capacity under fixed feature semantics—a signature of structured and controllable scalability.

To this end, we evaluate a series of FAT variants scaled from 50M to 1.5B parameters, trained on the same dataset with consistent hyperparameters and infrastructure. As shown in Figure 3,  $\Delta AUC$  increases monotonically with parameter count across three orders of magnitude, following a power-law trend. The observed relationship is well-characterized by the empirical function:

$$\Delta AUC = 5.81 \times 10^{-5} \cdot N_{params}^{0.433}, \quad (12)$$

No saturation is observed within the tested range, indicating that FAT continues to benefit from increased capacity—an uncommon property in CTR models, where unstructured architectures often plateau or degrade under scale [19, 20, 23].

This smooth scaling behavior aligns with Theorem 4.1, which establishes that FAT’s effective complexity grows as  $\mathcal{O}(Fd^2 + F^2)$ , independent of token vocabulary size  $n$ . By constraining expressive capacity to field-aware interaction pathways, FAT ensures that additional parameters refine semantically meaningful patterns rather than overfit sparse combinations. As a result, capacity expansion—achieved here through increased embedding dimension and depth—translates into consistent gains.

It is important to clarify that  $F$  defines the *interaction topology*, not a direct scaling dimension. Increasing  $F$  alters the input schema and introduces new combinatorial challenges, falling outside the regime of standard scaling laws. In contrast, the trend in Figure 3 reflects *intra-schema scalability*: performance improves predictably as representational fidelity deepens within a fixed structural foundation.

Further ablation shows that removing field-decomposed projections disrupts this scaling trend (Section 5.3), confirming that the observed behavior arises from architectural design, not mere parameter inflation.

In summary, FAT demonstrates the first empirically validated scaling law in CTR prediction: performance follows a reproducible power-law trajectory governed by structured expressivity. This closes the loop between theory and practice, establishing a principled path toward scalable recommendation modeling.

### 5.6 Online A/B Test Results (RQ5): Business Impact in Production

To evaluate the real-world impact of FAT, we conducted a large-scale online A/B test on Taobao’s sponsored search system, one of the largest e-commerce recommendation systems globally. Traffic was evenly split between the control group—serving the existing production model, a highly optimized traditional CTR model with manual feature crosses—and the treatment group, which replaced only the prediction module with **FAT-Large** (~0.5B parameters). All other components—including feature extraction, embedding lookup, and serving infrastructure—were kept identical to ensure a fair comparison. We report the relative improvements in two key business metrics: (1) **CTR**: Click-through rate, measuring user engagement. (2) **RPM**: Revenue per mille impressions, reflecting monetization efficiency.

**Table 4: Relative improvements in online A/B test.**

| Metric    | CTR           | RPM           |
|-----------|---------------|---------------|
| FAT-Large | <b>+2.33%</b> | <b>+0.66%</b> |

As shown in Table 4, FAT achieves statistically significant gains of **+2.33% in CTR** and **+0.66% in RPM**. The results confirm that FAT’s structured expressivity translates into measurable business value under real-world conditions. More importantly, they demonstrate that principled architectural scaling—guided by field-aware inductive biases—not only improves offline performance but also drives consistent online gains, validating its industrial applicability and impact.

## 6 Conclusion

We identify a fundamental mismatch between standard transformers and the combinatorial semantics of CTR data: while language models benefit from sequential compositionality, recommendation requires structured reasoning over unordered, high-cardinality semantic fields. Blindly scaling unstructured attention leads to poor generalization and ineffective performance gains. To address this, we propose the **Field-Aware Transformer (FAT)**, which introduces field-decomposed attention to align architectural inductive biases with data structure. By enabling asymmetric, interpretable cross-field interactions, FAT achieves *structured expressivity*—growing model capacity in harmony with interaction complexity. Theoretically, we show that FAT’s generalization error depends on the number of fields  $F$ , not vocabulary size  $n$ , providing the first principled scaling law for CTR models. Experiments demonstrate



that FAT consistently outperforms state-of-the-art methods, with up to **+0.51% AUC improvement** and smooth power-law scaling across model sizes and data volumes. It has been deployed in a production recommendation system, yielding significant online gains. Our work suggests that scalable performance in recommendation arises not from scale alone, but from structural alignment between architecture and domain semantics.

## References

- [1] Junyi Chen, Lu Chi, Bingyue Peng, and Zehuan Yuan. 2024. Hllm: Enhancing sequential recommendations via hierarchical large language models for item and user modeling. *arXiv preprint arXiv:2409.12740* (2024).
- [2] Yijie Ding, Yupeng Hou, Jiacheng Li, and Julian McAuley. 2024. Inductive Generative Recommendation via Retrieval-based Speculation. *arXiv preprint arXiv:2410.02939* (2024).
- [3] Huan Gui, Ruoxi Wang, Ke Yin, Long Jin, Maciej Kula, Taibai Xu, Lichan Hong, and Ed H Chi. 2023. Hiformer: Heterogeneous feature interactions learning with transformers for recommender systems. *arXiv preprint arXiv:2311.05884* (2023).
- [4] Hui Feng Guo, Ruiming Tang, Yunming Ye, Zhenguo Li, and Xiuqiang He. 2017. DeepFM: a factorization-machine based neural network for CTR prediction. *arXiv preprint arXiv:1703.04247* (2017).
- [5] Wei Guo, Hao Wang, Luankang Zhang, Jin Yao Chin, Zhongzhou Liu, Kai Cheng, Qiushi Pan, Yi Quan Lee, Wanqi Xue, Tingjia Shen, et al. 2024. Scaling new frontiers: Insights into large recommendation models. *arXiv preprint arXiv:2412.00714* (2024).
- [6] Yuchin Juan, Damien Yin, Guolin Zhou, Wei Wang, and Chun-Yuan Lin. 2016. Field-aware Factorization Machines for CTR Prediction. In *Proceedings of the 9th ACM International Conference on Web Search and Data Mining (WSDM)*. ACM, 43–51. <https://doi.org/10.1145/2835776.2835834>
- [7] Wang-Cheng Kang and Julian McAuley. 2018. Self-attentive sequential recommendation. In *2018 IEEE international conference on data mining (ICDM)*. IEEE, 197–206.
- [8] 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. In *arXiv preprint arXiv:2001.08361*.
- [9] Zihan Liu, Yupeng Hou, and Julian McAuley. 2024. Multi-Behavior Generative Recommendation. In *Proceedings of the 33rd ACM International Conference on Information and Knowledge Management*. 1575–1585.
- [10] Mehryar Mohri, Afshin Rostamizadeh, and Ameet Talwalkar. 2018. *Foundations of Machine Learning* (2nd ed.). MIT Press.
- [11] Aaron van den Oord, Yazhe Li, and Oriol Vinyals. 2018. Representation learning with contrastive predictive coding. *arXiv preprint arXiv:1807.03748* (2018).
- [12] Alec Radford. 2018. Improving language understanding by generative pre-training. (2018).
- [13] Alec Radford, Jeffrey Wu, Rewon Child, David Luan, Dario Amodei, Ilya Sutskever, et al. 2019. Language models are unsupervised multitask learners. *OpenAI blog* 1, 8 (2019), 9.
- [14] Steffen Rendle. 2010. Factorization Machines. In *Proceedings of the 10th IEEE International Conference on Data Mining (ICDM)*. IEEE, 995–1000. <https://doi.org/10.1109/ICDM.2010.12>
- [15] Weiping Song, Chence Shi, Zhiping Xiao, Zhijian Duan, Yewen Xu, Ming Zhang, and Jian Tang. 2019. AutoInt: Automatic feature interaction learning via self-attentive neural networks. In *Proceedings of the 28th ACM international conference on information and knowledge management*. 1161–1170.
- [16] Weiping Song, Chence Shi, Zhe Zhao, Zhiwei Zhu, Junjie Zhang, Yewen Li, Chris Frey, and Bokai Cao. 2018. AutoInt: Automatic Feature Interaction Learning via Self-Attentive Neural Networks. In *Proceedings of the 27th ACM International Conference on Information and Knowledge Management (CIKM)*. ACM, 1161–1170. <https://doi.org/10.1145/3269206.3271751>
- [17] Fei Sun, Jun Liu, Jian Wu, Changhua Pei, Xiao Lin, Wenwu Ou, and Peng Jiang. 2019. BERT4Rec: Sequential recommendation with bidirectional encoder representations from transformer. In *Proceedings of the 28th ACM international conference on information and knowledge management*. 1441–1450.
- [18] Ruoxi Wang, Rakesh Shivanna, Derek Cheng, Sagar Jain, Dong Lin, Lichan Hong, and Ed Chi. 2021. Dcn v2: Improved deep & cross network and practical lessons for web-scale learning to rank systems. In *Proceedings of the web conference 2021*. 1785–1797.
- [19] Jiaqi Zhai, Lucy Liao, Xing Liu, Yueming Wang, Rui Li, Xuan Cao, Leon Gao, Zhaojie Gong, Fangda Gu, Michael He, et al. 2024. Actions speak louder than words: Trillion-parameter sequential transducers for generative recommendations. *arXiv preprint arXiv:2402.17152* (2024).
- [20] Buyun Zhang, Liang Luo, Yuxin Chen, Jade Nie, Xi Liu, Daifeng Guo, Yanli Zhao, Shen Li, Yuchen Hao, Yantao Yao, et al. 2024. Wukong: Towards a scaling law for large-scale recommendation. *arXiv preprint arXiv:2403.02545* (2024).
- [21] Yu Zhao, Yuanbin Qu, Konrad Staniszewski, Szymon Tworkowski, Wei Liu, Piotr Miłoś, Yuxiang Wu, and Pasquale Minervini. 2024. Analysing The Impact of Sequence Composition on Language Model Pre-Training. *arXiv preprint arXiv:2402.13991* (2024).
- [22] Guorui Zhou, Xiaoqiang Zhu, Chenru Song, Ying Fan, Han Zhu, Xiao Ma, Yanghui Yan, Junqi Jin, Han Li, and Kun Gai. 2018. Deep interest network for click-through rate prediction. In *Proceedings of the 24th ACM SIGKDD international conference on knowledge discovery & data mining*. 1059–1068.
- [23] Jie Zhu, Zhifang Fan, Xiaoxie Zhu, Yuchen Jiang, Hangyu Wang, Xintian Han, Haoran Ding, Xinmin Wang, Wenlin Zhao, Zhen Gong, et al. 2025. RankMixer: Scaling Up Ranking Models in Industrial Recommenders. *arXiv preprint arXiv:2507.15551* (2025).

## A APPENDIX: THEORETICAL ANALYSIS OF FAT

In this section, we provide a rigorous analysis of the generalization properties of the Field-Aware Transformer (FAT), based on Rademacher complexity [10]. We derive an upper bound on the generalization error that explicitly depends on the number of semantic fields  $F$ , the embedding dimension  $d$ , and the interaction structure, while being independent of the total vocabulary size  $n$ . This justifies why FAT scales effectively even under extreme sparsity.

### A.1 Setup and Notation

Let  $\mathcal{X}$  denote the space of CTR inputs, each consisting of  $F$  semantic fields  $\{f_1, \dots, f_F\}$ . Consider one attention layer of FAT applied to embedded features  $\{\mathbf{h}_i\}_{i=1}^N$ , where  $\mathbf{h}_i \in \mathbb{R}^d$  includes both content embedding and field-aware positional bias. The output for token  $i$  is:

$$\text{Output}_i = \sum_{j=1}^N \alpha(i, j) \mathbf{v}_j,$$

$$\alpha(i, j) = \frac{\exp(s(i, j)/\sqrt{d})}{\sum_{k=1}^N \exp(s(i, k)/\sqrt{d})}, \quad s(i, j) = (\mathbf{q}_i^\top \mathbf{k}_j) \cdot w_{f_i, f_j},$$

where:

- $\mathbf{q}_i = \mathbf{W}_Q^{(f_i)} \mathbf{h}_i$ ,  $\mathbf{k}_j = \mathbf{W}_K^{(f_j)} \mathbf{h}_j$ ,  $\mathbf{v}_j = \mathbf{W}_V^{(f_j)} \mathbf{h}_j$ ;
- $\mathbf{W}_Q^{(f)} \in \mathbb{R}^{d \times d}$ , similarly for  $\mathbf{W}_K^{(f)}$ ,  $\mathbf{W}_V^{(f)}$ ;
- $w_{f_i, f_j} \in \mathbb{R}$  is a learnable scalar coefficient.

Let  $\mathcal{F}_{\text{FAT}}$  denote the function class induced by this attention block, mapping input sequences  $\mathbf{H}$  to contextualized outputs. Our goal is to bound its Rademacher complexity.

### A.2 Assumptions

We make the following boundedness assumptions, common in generalization analysis:

**ASSUMPTION 1 (BOUNDED EMBEDDINGS).**  $\|\mathbf{h}_i\|_2 \leq R$  for all  $i$ , some constant  $R > 0$ .

**ASSUMPTION 2 (BOUNDED PARAMETERS).** For all fields  $f$ :

- $\|\mathbf{W}_Q^{(f)}\|_F \leq B_Q$ ,
- $\|\mathbf{W}_K^{(f)}\|_F \leq B_K$ ,
- $\|\mathbf{W}_V^{(f)}\|_F \leq B_V$ ;

and for all field pairs  $(f_i, f_j)$ :  $|w_{f_i, f_j}| \leq B_w$ . These reflect practical regularization techniques (e.g., weight decay, gradient clipping).

### A.3 Main Theorem and Proof

Our main theoretical result is stated as follows:

**THEOREM A.1 (GENERALIZATION BOUND FOR FAT).** *Under Assumptions 1–2, the empirical Rademacher complexity of  $\mathcal{F}_{\text{FAT}}$  satisfies:*

$$\hat{\mathfrak{R}}_S(\mathcal{F}_{\text{FAT}}) \leq O\left(\frac{\sqrt{Fd^2 + F^2}}{\sqrt{m}} \cdot C(R, B_Q, B_K, B_V, B_w, d)\right),$$

where  $C(R, B_Q, B_K, B_V, B_w, d) = O(R^2 B_Q B_K B_w \sqrt{d} + R B_V B_w)$  is a constant depending on the norm bounds and dimension. Consequently, with probability at least  $1 - \delta$ , the generalization error satisfies:

$$L_{\text{gen}} \leq L_{\text{train}} + O\left(\frac{\sqrt{Fd^2 + F^2}}{\sqrt{m}} \cdot C + \sqrt{\frac{\log(1/\delta)}{m}}\right).$$

**PROOF.** The proof proceeds in three steps, leveraging standard tools from learning theory.

**Step 1: Bounding Score Magnitude and Lipschitz Constants.** Consider the score  $s(i, j) = (\mathbf{q}_i^\top \mathbf{k}_j) \cdot \mathbf{w}_{f_i, f_j}$ . By Assumptions 1 and 2, we have:

$$\|\mathbf{q}_i\|_2 = \|\mathbf{W}_Q^{(f_i)} \mathbf{h}_i\|_2 \leq \|\mathbf{W}_Q^{(f_i)}\|_F \|\mathbf{h}_i\|_2 \leq B_Q R,$$

similarly  $\|\mathbf{k}_j\|_2 \leq B_K R$ . Therefore,

$$|\mathbf{q}_i^\top \mathbf{k}_j| \leq \|\mathbf{q}_i\|_2 \|\mathbf{k}_j\|_2 \leq B_Q B_K R^2.$$

It follows that  $|s(i, j)| \leq B_Q B_K R^2 B_w \triangleq C_s$ . The softmax function  $\sigma(\cdot)$ , when viewed as a map from logits in  $\mathbb{R}^N$  to probabilities in  $[0, 1]^N$ , is  $C_s$ -Lipschitz with respect to the  $\ell_\infty$  norm on the input and  $\ell_1$  norm on the output, provided the logit differences are bounded by  $C_s$  [10].

**Step 2: Complexity of the Score Function Class.** Define the class of score functions:

$$\mathcal{S} = \left\{ (i, j) \mapsto (\mathbf{W}_Q^{(f_i)} \mathbf{h}_i)^\top (\mathbf{W}_K^{(f_j)} \mathbf{h}_j) \cdot \mathbf{w}_{f_i, f_j} \right\}.$$

We analyze the Rademacher complexity of  $\mathcal{S}$ . The key insight is that the parameters defining  $\mathcal{S}$  can be grouped into two types: (a) the  $3F$  matrices  $(\mathbf{W}_Q^{(f)}, \mathbf{W}_K^{(f)}, \mathbf{W}_V^{(f)})$  with Frobenius norm constraints, and (b) the  $F^2$  scalars  $\mathbf{w}_{f_i, f_j}$  with absolute value constraints.

Using the fact that the Rademacher complexity of a sum of function classes is bounded by the sum of their complexities, and applying standard results for linear function classes under norm constraints [10, Chapter 3], we get:

$$\begin{aligned} \hat{\mathfrak{R}}_S(\mathcal{S}) &\leq O\left(\underbrace{\frac{B_Q B_K R^2 \sqrt{Fd^2}}{\sqrt{m}}}_{\text{from } \mathbf{W}_Q^{(f)}, \mathbf{W}_K^{(f)}}\right) + O\left(\underbrace{\frac{B_Q B_K R^2 B_w \sqrt{F^2}}{\sqrt{m}}}_{\text{from } \mathbf{w}_{f_i, f_j}}\right) \\ &= O\left(\frac{C_s \sqrt{Fd^2 + F^2}}{\sqrt{m}}\right). \end{aligned}$$

**Step 3: Propagating Complexity to the Output.** The final output is a weighted sum of the values  $\mathbf{v}_j = \mathbf{W}_V^{(f_j)} \mathbf{h}_j$ . Each value vector satisfies  $\|\mathbf{v}_j\|_2 \leq B_V R$ .

The attention weights  $\alpha(i, j)$  are a function of the scores  $\{s(i, k)\}_{k=1}^N$ . Since the softmax is  $C_s$ -Lipschitz, we can apply Talagrand's contraction lemma [10, Theorem 4.12] to relate the Rademacher complexity of the attention weights to that of the scores. Furthermore, because

the output is a linear combination of the values weighted by these attention weights, and given the boundedness of the values, we can combine these results.

After careful application of the contraction lemma and accounting for the  $\sqrt{d}$  factor introduced by the temperature scaling in the softmax (which affects the Lipschitz constant), we arrive at the final bound:

$$\hat{\mathfrak{R}}_S(\mathcal{F}_{\text{FAT}}) \leq O\left(\frac{\sqrt{Fd^2 + F^2}}{\sqrt{m}} \cdot (R^2 B_Q B_K B_w \sqrt{d} + R B_V B_w)\right).$$

Applying the standard generalization bound via Rademacher complexity [10, Theorem 3.3] completes the proof.  $\square$

### A.4 Discussion: From Tight Generalization to Predictable Scaling

The derived generalization bound reveals that the generalization gap of FAT scales as  $\tilde{O}(\sqrt{Fd^2 + F^2}/\sqrt{m})$ , which crucially depends only on the number of semantic fields  $F$  and the embedding dimension  $d$ , not on the total vocabulary size  $n$ . This structural alignment dramatically improves statistical efficiency compared to standard Transformers, whose effective hypothesis space complexity implicitly scales with  $n$  due to unstructured attention over all tokens.

However, a tight generalization bound alone does not guarantee a smooth *scaling law*; it ensures that the model won't overfit, but not that it will continue to improve. For predictable performance gains as capacity increases, two conditions must hold simultaneously:

- (1) **Stable Generalization (Low Variance):** The generalization gap must remain controlled. This is guaranteed by Theorem A.1.
- (2) **Enhanced Expressivity (Reduced Bias):** The model's representational capacity must grow in a way that allows it to capture increasingly complex patterns in the data, thereby reducing the training error  $L_{\text{train}}$ .

FAT satisfies the second condition through its *structured expressivity*. As the embedding dimension  $d$  increases:

- The field-aware projections  $\mathbf{W}_Q^{(f)}, \mathbf{W}_K^{(f)}, \mathbf{W}_V^{(f)}$  gain higher rank, enabling more nuanced, field-specific encoding of content.
- The interaction modulation factors  $\mathbf{w}_{f_i, f_j}$ , though scalars, operate on richer, higher-dimensional representations, allowing for more sophisticated modeling of cross-field dependencies.

This architecture ensures that additional parameters are used to refine meaningful, semantically grounded interactions rather than fitting noise. Empirically, this manifests as a systematic reduction in  $L_{\text{train}}$  as  $d$  grows.

Combining these two effects—controlled variance from the structural prior and reduced bias from enhanced expressivity—we achieve a synergistic improvement in test performance  $L_{\text{gen}}$ . Under a fixed data distribution and field schema  $F$ , if we scale the model width such that  $N_{\text{params}} \propto Fd^2$ , then the generalization gap shrinks as  $O(\sqrt{N_{\text{params}}}/\sqrt{m})$ . When  $m$  grows proportionally to  $N_{\text{params}}$ , this gap remains stable or decreases, while  $L_{\text{train}}$  continues to decrease. This synergy enables the observed power-law scaling behavior. Let  $\Delta\text{AUC}$  be the gain over a fixed baseline. Our experiments (Section 5.5) validate that:

$$\Delta\text{AUC} \propto N_{\text{params}}^\beta, \quad \beta > 0, \quad (13)$$

holds across multiple orders of magnitude. This is the first theoretically grounded scaling law for CTR models, demonstrating that

predictable scaling arises from principled architectural design that aligns with the data's combinatorial semantics.