

GRAB: An LLM-Inspired Sequence-First Click-Through Rate Prediction Modeling Paradigm

Shaopeng Chen¹ Chuyue Xie¹ Huimin Ren¹ Shaozong Zhang¹ Han Zhang¹ Ruobing Cheng¹
Zhiqiang Cao¹ Zehao Ju¹ Yu Gao¹ Jie Ding¹ Xiaodong Chen¹ Xuewu Jiao¹ Shuanglong Li¹ Lin Liu¹

Abstract

Traditional Deep Learning Recommendation Models (DLRMs) face increasing bottlenecks in performance and efficiency, often struggling with generalization and long-sequence modeling. Inspired by the scaling success of Large Language Models (LLMs), we propose Generative Ranking for Ads at Baidu (GRAB), an end-to-end generative framework for Click-Through Rate (CTR) prediction. GRAB integrates a novel Causal Action-aware Multi-channel Attention (CamA) mechanism to effectively capture temporal dynamics and specific action signals within user behavior sequences. Full-scale online deployment demonstrates that GRAB significantly outperforms established DLRMs, delivering a 3.05% increase in revenue and a 3.49% rise in CTR. Furthermore, the model demonstrates desirable scaling behavior: its expressive power shows a monotonic and approximately linear improvement as longer interaction sequences are utilized.

1. Introduction

For a long time, Deep Learning Recommendation Models (DLRMs) (Naumov et al., 2019) have remained the mainstream choice in industrial recommender systems, especially for advertising Click-Through Rate (CTR) prediction (Mudigere et al., 2022; Zhou et al., 2018; Cheng et al., 2016; Guo et al., 2017; Bai et al., 2025; Wang et al., 2017; Ma et al., 2018a), due to their strong capability to process high-cardinality sparse features and to model feature interactions with expressive neural networks. However, as user behavior data grows exponentially, traditional DLRMs face increasing bottlenecks in both *performance* and *efficiency* (detailed discussions in Appendix A.1). Fundamentally, DLRMs rely on rule-based feature engineering and suffer from the inherent flaw of “strong memory, weak reasoning” (Cheng

et al., 2016; Wu et al., 2024). They often fail to generalize to new ads or scenarios that require logical inference, and their gains exhibit diminishing returns: further improvements typically demand exponentially increasing computational costs, rendering long-term deployment and iteration economically unsustainable (Zhang et al., 2024b; Mudigere et al., 2022).

Departing from the structural constraints of DLRMs, the rise of Large Language Models (LLMs) has been driven by scaling laws (Kaplan et al., 2020; Zhang et al., 2024a;b), where performance predictably improves with increased parameters, data, and compute. This success has inspired the extension of scaling laws to recommendation systems, fostering the LLMs for Recommendation (LLM4Rec) paradigm (Wu et al., 2024; Li et al., 2024a)(see Appendix A.2 for a taxonomy). A key innovation within this framework is Generative Recommendation (GR)(Li et al., 2024a; Rajput et al., 2023). Representative works like HSTU (Zhai et al., 2024) formulate recommendation as autoregressive sequence prediction, effectively modeling long user sequences to capture temporal dynamics (Zhou et al., 2018; 2019; Kang & McAuley, 2018; Sun et al., 2019). Crucially, GR exhibits scaling properties similar to LLMs, offering a practical path to transcend the performance bottlenecks of traditional DLRMs (Naumov et al., 2019; Zhai et al., 2024; Zhang et al., 2024a).

Despite these theoretical advancements, deploying GR models in high-throughput industrial systems remains challenging due to strict online serving and optimization constraints. The primary obstacle is *computational efficiency*. Standard Transformer training requires extensive padding for variable-length sequences, resulting in significant computational waste (Vaswani et al., 2017; Krell et al., 2021). While sequence packing—a common Natural Language Processing (NLP) technique for concatenating multiple short sequences—effectively mitigates this issue (Krell et al., 2021), its straightforward application to recommendation systems triggers a more subtle yet damaging failure mode: *Distribution Skew* (Baylor et al., 2017; Polyzotis et al., 2019; Sculley et al., 2015; Han et al., 2025).

In recommendations, packing a user’s full history creates mini-batches with excessive intra-user correlation, which violates the i.i.d. assumption typically relied on by SGD-

¹Baidu Inc., Beijing, China. Correspondence to: Shuanglong Li <lishuanglong@baidu.com>.

style optimization (Doan et al., 2020). This skew (details in Appendix D.1) causes sparse parameters (i.e., embedding tables) to overfit specific users, hindering the generalization of dense parameters (e.g., Transformer weights responsible for inference) (Naumov et al., 2019; Li et al., 2024b). This reveals a fundamental tension: sparse parameters require diverse, uncorrelated samples for robust “memorization”, whereas dense parameters benefit from long, coherent contexts for sequential “reasoning” (Cheng et al., 2016; Kang & McAuley, 2018; Sun et al., 2019). This misalignment implies that standard synchronous training on packed sequences may lead to suboptimal convergence due to the conflicting gradient requirements of the sparse and dense components (Yu et al., 2020).

Meanwhile, existing GR models typically ignore data heterogeneity, resulting in performance limitations (see Appendix A.3 for detailed discussion). To overcome these challenges, we propose Generative Ranking for Ads at Baidu (GRAB), an end-to-end sequential training and inference framework for industrial-grade CTR prediction. GRAB introduces three core innovations to reconcile the demands for performance, efficiency, and training stability:

- **End-to-End Framework:** We introduce GRAB, an end-to-end framework that combines the strengths of DLRMs and GR. Specifically, it fuses the large-scale sparse feature engineering inherent to DLRMs with the sequential inference capabilities of GR, thereby achieving a balance between explicit memorization and implicit reasoning.
- **Causal Action-aware Multi-channel Attention (CamA) mechanism:** We propose CamA, a multi-channel, action-aware mechanism designed to boost model performance by modeling both user exposure and interaction signals, improving generalization and robustness across tasks and scenarios.
- **Sequence-Then-Sparse (STS) Training:** To address distribution skew from sequence packing, we propose STS, a training strategy that decouples the optimization of dense parameters and sparse embeddings. This resolves their gradient conflict, stabilizes training, and improves convergence without extra compute, enabling high-throughput industrial deployment of GR.

We have completed a comprehensive evaluation of GRAB in Baidu’s commercial advertising CTR ranking business. In offline comparisons, GRAB outperformed mainstream industrial DLRMs as well as emerging GR models (achieving 0.19% relative improvement over the best baseline). Compared to the production DLRM baseline, GRAB achieved an AUC uplift of approximately 2 basis points in online A/B testing, resulting in a 3.05% increase in CPM and a 3.49% increase in CTR. Furthermore, scaling analysis demonstrates

that the model’s AUC improves monotonically with both model capacity and the length of behavior sequences, indicating that the architecture can stably benefit from modeling longer behavior chains without saturation.

2. Related Works

DLRM-based industrial CTR prediction. Industrial CTR prediction has long been dominated by DLRMs, which embed high-cardinality categorical fields and model feature interactions via MLP/Cross-style modules (Naumov et al., 2019; Cheng et al., 2016; Guo et al., 2017; Wang et al., 2017). To incorporate user histories, production systems often attach explicit behavior encoders to DLRMs, e.g., target-attention/memory based models such as DIN, DIEN, MIMN, and SIM (Zhou et al., 2018; 2019; Pi et al., 2019; 2020), as well as stronger industrial variants like TWIN (Si et al., 2024). Despite their effectiveness, these approaches still heavily rely on hand-crafted statistics and engineered cross features (He et al., 2014; Cheng et al., 2016), and typically compress long histories into fixed-size vectors, making it difficult to scale to long sequences and heterogeneous action signals (Pi et al., 2019).

GR. Recent GR work models recommendation as causal Transformer-based sequential prediction, enabling long-context modeling and exhibiting favorable scaling behavior (Zhai et al., 2024; Chai et al., 2025; Petrov & Macdonald, 2023). However, deploying GR in the industrial advertising CTR stack still presents challenges in the following aspects: (i) bridging large-scale sparse feature engineering with tokenized sequential modeling (Han et al., 2025; Chai et al., 2025), (ii) modeling heterogeneous action semantics often discarded by naive homogeneous serialization (Zhai et al., 2024), and (iii) training instability introduced by sequence packing (distribution skew) under strict optimization constraints (Krell et al., 2021). Detailed discussion of GR and comparisons between GRAB and related work are given in Appendix B.

3. Methodology

3.1. DLRMs

The traditional DLRM architecture, as shown in Fig. 1, follows a modular processing pipeline for CTR prediction, handling raw features from users, candidate ads, and contextual signals. The pipeline involves: (a) expanding categorical features into fixed fields via feature engineering, (b) mapping these fields through hashing to obtain discrete ID vectors for embedding lookup in a Sparse Parameter Server Table(PSTable), and (c) concatenating and normalizing the retrieved embeddings to form a fixed-length flattened vector. This unified representation is then fed into an MLP, typically enhanced with a gating network, to model high-order

feature interactions and generate the final CTR prediction.

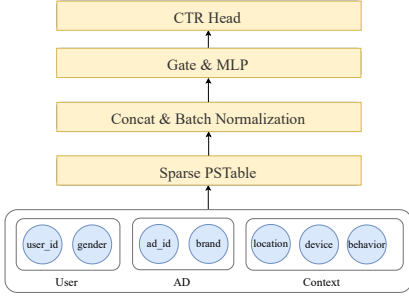


Figure 1. The traditional DLRM architecture: sparse features are hashed to IDs and embedded via PTable, and then concatenated into a fixed-length flattened vector for CTR prediction.

3.2. Overall Architecture of GRAB

GRAB, with the overall architecture shown in Fig. 2, is designed to model user behavior history sequences in an end-to-end manner, as applied in scenarios like CTR prediction. GRAB follows a three-stage pipeline: (i) *sparse feature layer*; (ii) *dense tokenizer*; and (iii) *sequence modeling layer*. Given raw behavior logs, GRAB first converts heterogeneous categorical signals into sparse IDs at the event level, then tokenizes each event into a dense representation, and finally applies a sequence model to estimate the click probability of candidate ads. GRAB uses its dense representation calculated from the dense tokenizer to bridge DLRM-style sparse feature engineering and GR-style sequential modeling, enabling end-to-end training and inference along a single, unified computation path from input to output, thereby improving CTR prediction performance through end-to-end sequential modeling of event-level user behaviors.

Sparse Feature Layer. The sparse feature layer (details in Appendix C.1) processes raw logs into time-ordered event sequences. Each event’s categorical fields are converted into sparse IDs using standard DLRM feature engineering (Section 3.1), yielding a structured sequence of events annotated with field-wise IDs.

Dense Tokenizer. Unlike DLRM, which collapses field embeddings into a fixed-length, order-agnostic vector for pointwise processing, GRAB preserves the temporal event structure. It aggregates per-event field embeddings and projects them into $\mathbb{R}^{d_{\text{model}}}$ to form sequential event tokens (Appendix C.2), resulting in a time-ordered token sequence. This sequence serves as the input to a subsequent Transformer, thereby enabling the modeling of long-range dependencies and interest drift.

Autoregressive-like Sequence Modeling Layer. Built on sequence packing (Section 3.3.1), heterogeneous tokens

(Section 3.3.2), and action-aware relative attention bias (Section 3.3.3), our core contribution is the **CamA** mechanism (Section 3.3.4). CamA integrates a multi-channel design for parallel processing of diverse behaviors and inherits action-aware contextualization from RAB, providing a unified and efficient framework for modeling complex user interest patterns across scenarios.

3.3. Autoregressive-like Sequence Modeling Layer

Following the dense tokenizer, this layer is designed to capture the temporal dependencies and dynamic evolution of user interests, which takes the sequence of dense event tokens generated by the preceding layer as input (as described in Appendix C.3). Formally, for a user u , the input sequence consists of the behavior history $\mathbf{E}^{\text{beh}} = \{\mathbf{e}_t^{\text{beh}}\}_{t=1}^{T_u}$ and the candidate advertisements $\mathbf{E}^{\text{ad}} = \{\mathbf{e}_i^{\text{ad}}\}_{i=1}^{N_u}$, where $\mathbf{e}_t^{\text{beh}}, \mathbf{e}_i^{\text{ad}} \in \mathbb{R}^{d_{\text{model}}}$ are the dense embeddings of the t -th behavior event and the i -th candidate ad, respectively, T_u is the behavior history length, and N_u is the number of candidate ads.

3.3.1. SEQUENCE PACKING AND USER-ISOLATED CAUSAL MASK

In industrial training logs, as shown in the left image of Fig. 3a., a mini-batch is typically formed by sampling B_{ins} impression instances. Each instance contains a variable-length token sequence composed of (i) a subsequence of the user’s historical behavior tokens and (ii) target advertisement tokens to be scored. A straightforward batching strategy pads every instance to a fixed length L_{max} , yielding a dense tensor with dimensions $B_{\text{ins}} \times L_{\text{max}} \times d_{\text{model}}$,

which introduces substantial computational waste when most instances are much shorter than L_{max} .

To eliminate such padding overhead while preserving the temporal semantics, GRAB performs sequence packing by grouping tokens by user. Specifically, tokens from multiple impression instances belonging to the same user u are merged into a single contiguous token segment, while segments of different users are strictly separated. Within each user segment, all tokens are stably sorted by timestamp so that the packed segment forms a single timeline for sequential modeling. After packing, the batch is represented as one long packed tensor $H = \text{Pack}(\mathbf{E}^{\text{beh}}, \mathbf{E}^{\text{ad}}) \in \mathbb{R}^{1 \times L \times d_{\text{model}}}$, where L denotes the total packed length across all users in the mini-batch.

For convenience, we associate each packed position $p \in \{1, \dots, L\}$ with (i) a *segment id* $\sigma(p) \in U_B$ indicating which user it belongs to, and (ii) a *local time index* $\ell(p) \in \{1, \dots, L_{\sigma(p)}\}$ within that user segment.

User-isolated causal mask. On the packed tensor H , we construct an additive attention mask $M^{\text{pack}} \in \mathbb{R}^{L \times L}$ that

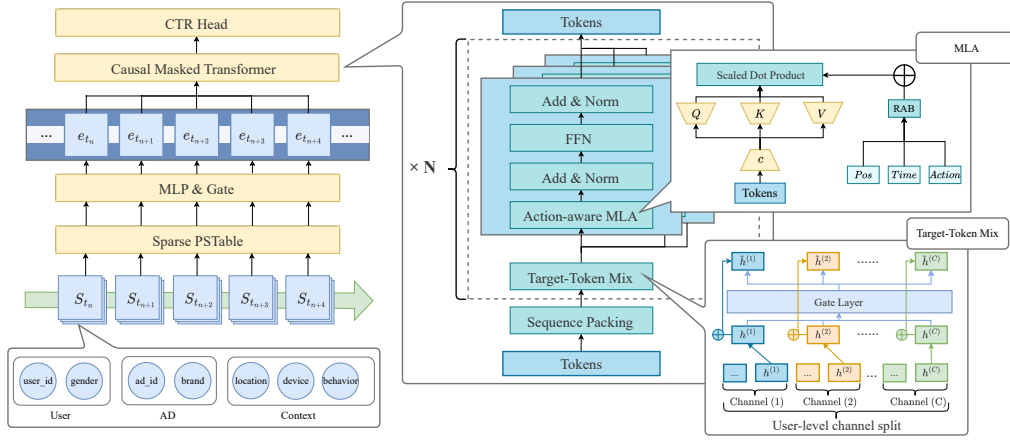


Figure 2. Overview of GRAB’s end-to-end CTR prediction pipeline: (1) Tokenizing raw fields via a sparse PSTable and fusing them into event tokens. (2) Packing tokens per user with causal and heterogeneous masks. (3) Processing through N Transformer layers equipped with the Causal Action-aware Multi-channel Attention (CamA) mechanism. (4) Final CTR prediction from the output representations.

enforces two constraints: (1) *user isolation* (no cross-user attention), and (2) *causality* within each user’s timeline (no future leakage). Formally, for query position p and key position q ,

$$M_{p,q}^{\text{pack}} = \begin{cases} 1, & \text{if } \sigma(p) = \sigma(q) \text{ and } \ell(q) \leq \ell(p), \\ 0, & \text{otherwise.} \end{cases} \quad (1)$$

This yields a block-diagonal lower-triangular structure(as shown in Fig. 3b), where each block corresponds to one user segment.

3.3.2. HETEROGENEOUS BEHAVIOR TOKENS AND HETEROGENEOUS VISIBILITY MASK

After sequence packing, for each user u , we obtain a user-isolated, time-ordered packed stream with its causal mask M^{pack} . To further reduce redundancy in the packed history while preserving the information needed for scoring the current candidate, we instantiate two token views at each packed timestamp t : **the partial token (history)** $h_t \in \mathbb{R}^{d_{\text{model}}}$, which retains only time-varying information that is useful for representing history and discards static or highly repetitive fields (e.g., user_id) that would otherwise be duplicated across historical steps and could lead to overfitting; and **the full token (candidate)** $h'_t \in \mathbb{R}^{d_{\text{model}}}$, which retains the complete information required to score the candidate at time t , including the static fields omitted from the partial history view. We then interleave them to form the heterogeneous packed sequence: $H_u = [\mathbf{h}_1, \mathbf{h}'_1, \mathbf{h}_2, \mathbf{h}'_2, \dots, \mathbf{h}_{T_u}, \mathbf{h}'_{T_u}]$.

Heterogeneous Visibility Mask. On top of the user-isolated causal constraint encoded by M^{pack} , we apply a mask-rewriting operator $\mathcal{R}(\cdot)$ to obtain the heterogeneous visibility mask M^{het} . Concretely, $\mathcal{R}(\cdot)$ rewrites the visibility pattern according to the token types in the following way:

(i) partial (\mathcal{P}) tokens only attend to partial history tokens; and (ii) full (\mathcal{F}) tokens attend to partial history tokens and themselves, but never attend to other full tokens. Formally, index positions in H_u by $n \in \{1, \dots, 2T_u\}$, we define the time index $\tau(n) = \lceil n/2 \rceil$ and token type $\kappa(n) = \mathcal{P}$ if n is odd, otherwise $\kappa(n) = \mathcal{F}$. Then the heterogeneous mask (as shown in Fig. 4) is

$$M_{p,q}^{\text{het}} = \begin{cases} 1, & \kappa(p) = \mathcal{P}, \kappa(q) = \mathcal{P}, \tau(q) \leq \tau(p), \\ 1, & \kappa(p) = \mathcal{F}, \kappa(q) = \mathcal{P}, \tau(q) \leq \tau(p), \\ 1, & \kappa(p) = \mathcal{F}, p = q, \\ 0, & \text{otherwise.} \end{cases} \quad (2)$$

3.3.3. ACTION-AWARE ATTENTION: RELATIVE ENCODING AND EFFICIENT COMPUTATION

On top of the heterogeneous behavior tokens and the heterogeneous visibility mask M^{het} , we further adopt a action-aware RAB(i.e., relative attention bias) causal attention mechanism. It augments standard multi-head self-attention with three designs: a causal mask to prevent future leakage, a dual sliding-window visibility constraint to support streaming-style training, and a query-aware relative bias that enables the query to directly interact with relative position/time/action signals.

Action-aware relative attention logits. Given a query q_i and a key k_j , the attention logit is computed as

$$w_{i,j} = q_i^\top \cdot (k_j + \text{Pos}_{i,j} + \text{Action}_{i,j} + \text{Time}_{i,j}), \quad (3)$$

where $\text{Pos}_{i,j}$, $\text{Action}_{i,j}$, and $\text{Time}_{i,j}$ are learnable embeddings derived from relative position, relative action, and relative time, respectively. For continuous or large-range

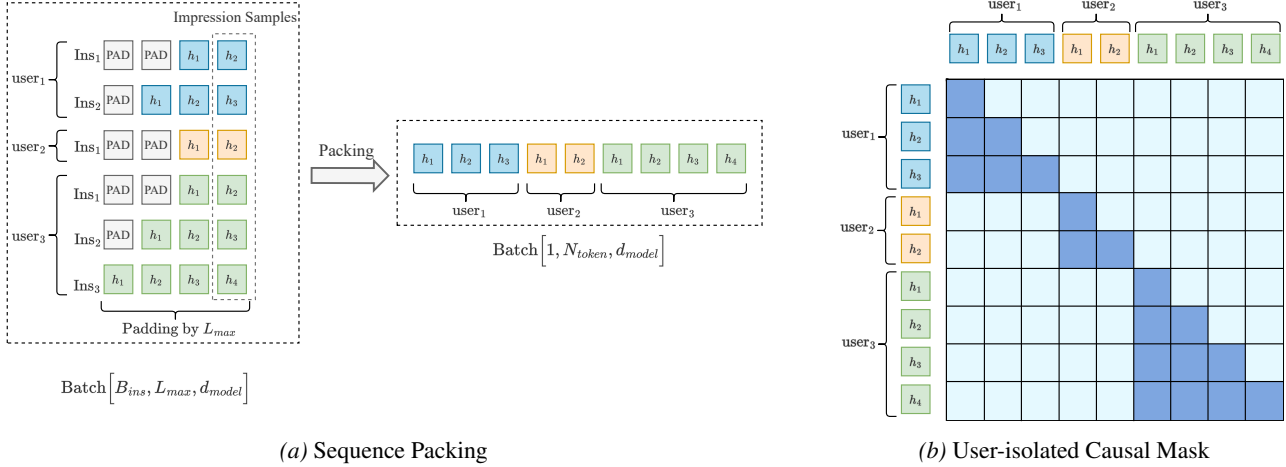


Figure 3. Sequence packing and user-isolated causal masking in GRAB. (a) Instead of padding each impression instance to a fixed length L_{\max} , tokens from multiple impressions are concatenated within each user and different users are kept in disjoint segments, yielding a single packed sequence of length N_{token} for compute-efficient batching. (b) The user-isolated causal mask exhibits a block-diagonal lower-triangular pattern, so each token can only attend to past tokens within the same user segment, enforcing both user isolation and temporal causality.

signals (e.g., action statistics or play durations), we first discretize them into buckets and then perform embedding lookup.

Compared with a query-agnostic relative bias (e.g., $w_{i,j} = q_i^\top k_j + \text{Pos}_{i,j} + \dots$), Eq. 3 makes the relative signals action-aware via the interaction $q_i^\top \text{Pos}_{i,j}$, $q_i^\top \text{Action}_{i,j}$, and $q_i^\top \text{Time}_{i,j}$, allowing the model to adaptively emphasize different contextual relations under different queries (i.e., target ads).



Figure 4. Heterogeneous behavior tokens and heterogeneous visibility mask M^{het} (blue entries). Partial tokens attend only to partial-history tokens up to the current time, while full tokens attend to partial-history tokens up to their time index and to themselves, but never to other full tokens, preventing duplicated static information from propagating along time.

Causal mask with dual sliding windows. We enforce causality and further restrict attention using combined time and length windows. The mask is defined as $M_{p,q}^{\text{rab}} = 1$ if $q \leq p$ and the distance $p - q$ does not exceed the length sliding-window limit L_w ; otherwise $M_{p,q}^{\text{rab}} = 0$.

This serves two key industrial purposes: (1) it bounds per-token computation, guaranteeing stable throughput/latency over growing behavior histories; (2) it matches the online training paradigm—events arrive incrementally, and the model updates attention context on the fly without reprocessing the full sequence, boosting training efficiency and serving practicality.

Efficient computation. The naive implementation of Eq. 3 would yield an $O(L^2 d_{\text{model}})$ intermediate tensor, which is prohibitively memory-intensive in practice. We adopt the optimization in (Golovneva et al., 2024) to re-order the computation. We define codebooks $B^{\text{pos/act/time}} \in \mathbb{R}^{N_* \times d_{\text{model}}}$ and bucketized indices $p_{i,j}, a_{i,j}, t_{i,j}$. Then Eq. 3 can be equivalently written as:

$$w_{i,j} = q_i^\top k_j + (s_i^{\text{pos}})[p_{i,j}] + (s_i^{\text{act}})[a_{i,j}] + (s_i^{\text{time}})[t_{i,j}]. \quad (4)$$

where $s_i^{\text{pos}} = q_i^\top B^{\text{pos}}$, $s_i^{\text{act}} = q_i^\top B^{\text{act}}$, and $s_i^{\text{time}} = q_i^\top B^{\text{time}}$. In practice, we first compute the projection vectors s_i^* , then obtain relative terms via fast gather operations. This completely avoids the large $L \times L \times d_{\text{model}}$ tensor, dramatically reducing peak memory and improving computational efficiency.

3.3.4. MULTI-CHANNEL ATTENTION

While the action-aware RAB attention (Section 3.3.3) enhances each individual attention logit with relative position/action/time signals, it still treats the packed stream as a single mixed sequence. However, in industrial logs, user behaviors are highly heterogeneous (e.g., spanning different time windows or encompassing different behavior types),

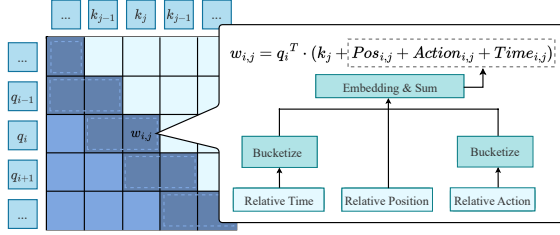


Figure 5. Action-aware relative attention bias (RAB) with efficient computation. Left: a causal mask with dual sliding windows, which limits each query to attend only to recent past tokens visible within the sliding-window. Right: the action-aware relative encoding pipeline: relative time, position, and action signals are bucketized (as needed), embedded, summed, and injected to the attention logits.

and different behavioral subsets often exhibit distinct temporal dynamics and predictive value. A straightforward design is to flatten all tokens into a single sequence and apply causal self-attention, yet this couples heterogeneous sources into one interaction graph and incurs a quadratic cost (e.g., $O((n+m)^2)$ for two sources with lengths n and m). To improve both modeling effectiveness and efficiency, we further introduce the **Causal Action-aware Multi-channel Attention (CamA)** mechanism, which integrates a multi-channel design, conceptually analogous to multi-head attention but with channel-specific visibility constraints. We therefore model each channel with an independent causal self-attention stack, and fuse the channel-wise representations via a lightweight gated mixing module. Let $\mathcal{C} = \{1, \dots, C\}$ denote the channel set. For each user, channel c provides a token sequence $\mathbf{X}^{(c)} \in \mathbb{R}^{T_c \times d}$, and we append the shared target token $\mathbf{x}^{ad} \in \mathbb{R}^d$:

$$\mathbf{S}^{(c)} = [\mathbf{X}^{(c)}; \mathbf{x}^{tar}] \in \mathbb{R}^{(T_c+1) \times d}, \quad t^* = T_c + 1. \quad (5)$$

Each channel is equipped with its own causal visibility mask $\mathbf{M}^{(c)}$, and is encoded independently:

$$\begin{aligned} \mathbf{H}^{(c, \ell+1)} &= \text{Layer}_{\ell}^{(c)}(\tilde{\mathbf{H}}^{(c, \ell)}; \mathbf{M}^{(c)}), \\ \mathbf{H}^{(c, 0)} &= \mathbf{S}^{(c)}, \quad c \in \mathcal{C}. \end{aligned} \quad (6)$$

Target-token gated mixing. To enable cross-channel information sharing while keeping computation lightweight, we perform mixing only on the target position t^* at each layer. The mixed representation $\tilde{\mathbf{h}}^{(c, \ell)}$ is obtained by first computing channel-wise gating weights $\beta^{(c, \ell)}$ and then aggregating information from all other channels:

$$\tilde{\mathbf{h}}^{(c, \ell)} = \mathbf{h}^{(c, \ell)} + \sum_{i \in \mathcal{C} \setminus \{c\}} \beta^{(i, \ell)} \odot \mathbf{h}^{(i, \ell)}. \quad (7)$$

This updated representation replaces $\mathbf{h}^{(c, \ell)}$ at position t^* , forming the updated channel representation $\tilde{\mathbf{H}}^{(c, \ell)}$ used in

(6). Finally, the concatenated last-layer target representations from all channels are used for CTR prediction.

3.4. Sequence Then Sparse Training

While sequence packing (Section 3.3.1) significantly enhances computational efficiency, it introduces a critical challenge: distribution skew. Since samples within a packed mini-batch belong to the same user, the high intra-user correlation leads to redundant updates for specific sparse IDs, causing the model to overfit to specific user-ad interactions, rather than learning generalizable patterns. To mitigate this, we propose the Sequence Then Sparse (STS) training paradigm (detailed discussions in Appendix D), a two-stage decoupled optimization strategy that balances long-range sequential modeling with robust sparse feature learning.

3.4.1. STAGE I: SEQUENCE MODELING (SEQUENCE PHASE)

The first stage focuses on capturing the evolution of user interests and temporal dependencies. We perform end-to-end autoregressive-like learning on the packed user sequences \mathbf{Z} , which include candidate tokens and their historical trajectories. In this phase, we optimize the dense tokenizer and the causal Transformer, while keeping the Sparse Embedding Table Φ frozen. By freezing Φ , we stabilize the token space, forcing the Transformer to focus exclusively on the relational dynamics between events rather than over-memorizing specific ID features.

3.4.2. STAGE II: SPARSE FEATURE LEARNING (SPARSE PHASE)

The second stage is designed to refine the discrete representations, particularly for long-tail IDs. In this phase, we revert to a non-sequential format, treating each sample as an independent user-ad exposure to break the distribution skewness. This stage optimizes the sparse embeddings Φ , which act as a robust corrector for the gradient accumulation amplified by sequence packing. It ensures that the basic feature representations remain accurate and unbiased across the entire traffic distribution.

4. System Deployment

GRAB has been successfully deployed in a large-scale feed ad ranking system, handling billions of daily requests. Unlike conventional memory-bound DLRMs, GR is markedly compute-bound due to the quadratic complexity of Transformer self-attention. To satisfy stringent latency requirements, we implemented a co-designed hardware-software architecture. Due to space constraints, we provide the comprehensive system overview (Fig. 8) and detailed deployment optimizations in Appendix E.

5. Experiment

5.1. Overall Performance Comparison

We first compared the performance of GRAB against state-of-the-art recommendation models on the Baidu real-world industrial dataset. The training data, derived from the Baidu real recommendation advertising scene, contains billions of users, exposure logs, and click logs. The test set includes millions of users, billions of exposure logs, and millions of click logs. The baselines encompass both DL-RMs and GR models, including: DIN (Zhou et al., 2018), which models short-term user behavior with target attention; SIM(Soft) (Pi et al., 2020), a sequential model that uses soft-search to encode user interests; TWIN (Si et al., 2024), which extends multi-head target attention from ESU to GSU; HSTU (Zhai et al., 2024), an efficient model for long-sequence behavior modeling; and LONGER (Chai et al., 2025), a Transformer-based architecture designed for ultra-long behavior sequences. Experimental results are presented in Table 1: GRAB outperforms all other baselines, achieving a 0.19% relative improvement over the most competitive model. Meanwhile, Fig. 6a illustrates the performance of different models across varying lengths of user behavior sequences. GRAB surpasses other recommendation models at all sequence lengths, with its performance gains becoming more pronounced as the sequence length increases.

Table 1. Overall performance in industrial settings

Model	AUC
DIN	0.83309
SIM Soft	0.83520
TWIN	0.83556
HSTU	0.83590
LONGER	0.83615
GRAB-small	0.83661
GRAB-standard	0.83772

5.2. Scaling Analysis

We evaluate model performance across different capacity scales by independently scaling the number of Transformer blocks(n_{layer}), the number of attention heads(n_{head}), and the feature dimension of the model(d_{model}) in Table 2, Fig. 6b presents the test-set performance of the GRAB model under varying configurations (i.e., n_{layer} , n_{head} and d_{model}). These results demonstrate that increasing model capacity effectively improves model performance. We also found that as the model capacity increases, the performance improvement on longer user behavior sequences becomes more significant. Moreover, no significant saturation trend is observed within the current range of configurations, which also confirms the strong scalability of the GRAB model.

Table 2. Comparison of models with different settings

Model	Params	Setting
GRAB _{2l-2h-64d}	6.51M	$n_{layer}=2, n_{head}=2, d_{model}=64$
GRAB _{4l-2h-64d}	6.67M	$n_{layer}=4, n_{head}=2, d_{model}=64$
GRAB _{6l-2h-64d}	6.83M	$n_{layer}=6, n_{head}=2, d_{model}=64$
GRAB _{2l-4h-64d}	6.48M	$n_{layer}=2, n_{head}=4, d_{model}=64$
GRAB _{4l-4h-64d}	6.63M	$n_{layer}=4, n_{head}=4, d_{model}=64$
GRAB _{4l-4h-128d}	7.05M	$n_{layer}=4, n_{head}=4, d_{model}=128$
GRAB _{4l-4h-256d}	8.13M	$n_{layer}=4, n_{head}=4, d_{model}=256$
GRAB _{4l-4h-512d}	11.27M	$n_{layer}=4, n_{head}=4, d_{model}=512$

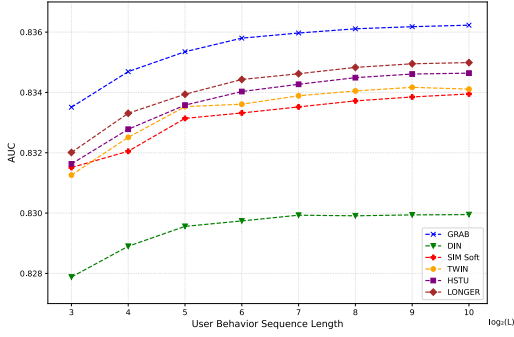
5.3. Ablation Study

Heterogeneous Tokens. We conduct ablation studies on heterogeneous representations with three configurations: GRAB with heterogeneous, only partial, or only full tokens (Table 3). Results show that heterogeneous representations achieve the best performance. Using only partial tokens leads to significant degradation, confirming that full feature representations are more beneficial for target scoring. Notably, using only full tokens also degrades performance, suggesting that artificially designed statistical features can introduce confusion and impair sequence modeling.

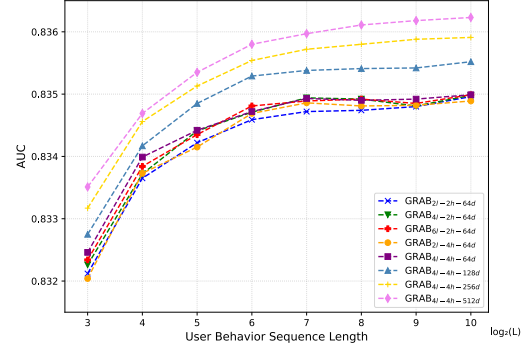
Table 3. Ablation studies of GRAB

Model	AUC
GRAB	0.83772
GRAB w/ Partial Token	0.83492
GRAB w/ Full Token	0.83749
GRAB w/o relative pos	0.83768
GRAB w/o relative time	0.83743
GRAB w/o relative action	0.83724
GRAB w/o Multi-channel	0.83743
GRAB w/o Target-token mix	0.83768
GRAB_{sparse}	0.83614
GRAB _{sparse} w/o STS	0.83549

Action-aware Attention. We ablate three components of GRAB’s Action-aware Attention: relative position, time, and action. The results (Table 3) show that removing any of these components degrades performance. The decline is more pronounced for time and action than for position, indicating that historical sequences are more sensitive to behavioral and temporal signals. We also analyze the attention weight distribution across buckets defined by relative position/time differences (smaller values denote more recent tokens). As shown in Figure 7, weights decrease as bucket values increase, confirming that more recent behaviors better reflect user interest and receive higher weights. For relative action, we compare positive (click) and negative (non-click) labels. The weight distribution is highly skewed: positive labels account for 88% of the total weight, versus only 12%



(a) Overall Performance



(b) Scaling Performance

Figure 6. DLRMs vs. GRs across different user behavior sequence lengths (a), with a +0.1% improvement in AUC, indicating a significant enhancement. GRABs comparison in different parameter scale(b)

for negative labels. This suggests that incorporating more positive feedback could further improve sequence modeling.

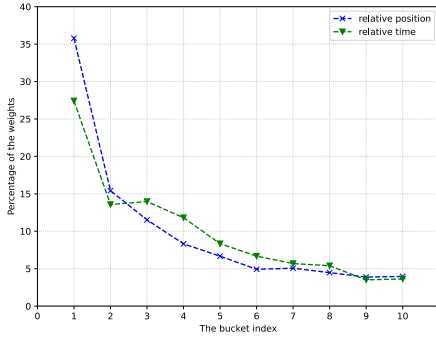


Figure 7. The weight distribution of action-aware attention in relative position and relative time.

Multi-channel Attention. To verify the effectiveness of multi-channel attention in sequence modeling, we conduct the following settings: 1) the GRAB model without multi-channel attention, that is, using a single channel for sequence modeling, 2) remain the multi-channel attention and only remove the target token mix component. As shown in Table 3, both configurations have varying degrees of performance degradation, indicating that each component is indispensable. In terms of performance, multi-channel attention is crucial, and adding the target token mix component can further improve performance.

STS Training. We evaluate the STS paradigm by comparing GRAB’s second-stage training with and without sequence modeling for sparse feature learning. With STS, sparse embeddings are updated through sequence modeling on packed user behavior sequences; without STS, the same batch data is treated as independent exposures. Results (as shown in Table 3) show that STS brings significant accuracy

gains in sparse feature learning, confirming the efficacy of the two-stage training. This demonstrates that STS alleviates the distribution skew and overfitting caused by direct sequence-packed training.

5.4. Online A/B Test

To assess the online performance of GRAB, we deployed it in Baidu home feed scenario of Baidu and compared its performance with the current online DLRM model. The experiment used 10% of the main traffic and remained online for about a month. Online evaluation shows that GRAB delivered 3.49% improvement in CTR and 3.05% improvement in CPM, which indicates that GRAB achieves more accurate advertising estimation and brings considerable revenue increments. Notably, GRAB has already been fully deployed on Baidu, and the online inference costs on par with the previous online DLRM model.

6. Conclusion

We propose GRAB, an end-to-end generative ranking framework that integrates a novel CamA mechanism to effectively capture temporal dynamics and specific action signals within user behavior sequences. On Baidu billion-scale industrial dataset, GRAB establishes a new state-of-the-art, outperforming DLRM and other GR baselines. Ablation studies validate the necessity of its key components, and our proposed STS training paradigm effectively mitigates distribution shift. Scaling analysis indicates continued gains from larger models and longer sequences. Finally, full online A/B testing in Baidu home feed ads shows that GRAB boosts CTR by 3.49% and CPM by 3.05%, leading to full production deployment. Further discussion of this work can be found in the Appendix F.

References

- Agarwal, S., Yan, C., Zhang, Z., and Venkataraman, S. Bagpipe: Accelerating deep recommendation model training. In *Proceedings of the 29th Symposium on Operating Systems Principles*, pp. 348–363, 2023.
- Bai, J., Geng, X., Deng, J., Xia, Z., Jiang, H., Yan, G., and Liang, J. A comprehensive survey on advertising click-through rate prediction algorithm. *The Knowledge Engineering Review*, 40:e3, 2025.
- Bao, K., Zhang, J., Zhang, Y., Wang, W., Feng, F., and He, X. Tallrec: An effective and efficient tuning framework to align large language model with recommendation. In *Proceedings of the 17th ACM conference on recommender systems*, pp. 1007–1014, 2023.
- Baylor, D., Breck, E., Cheng, H.-T., Fiedel, N., Foo, C. Y., Haque, Z., Haykal, S., Ispir, M., Jain, V., Koc, L., et al. Tfx: A tensorflow-based production-scale machine learning platform. In *Proceedings of the 23rd ACM SIGKDD international conference on knowledge discovery and data mining*, pp. 1387–1395, 2017.
- Cao, Y., Mehta, N., Yi, X., Keshavan, R. H., Heldt, L., Hong, L., Chi, E., and Sathiamoorthy, M. Aligning large language models with recommendation knowledge. In *Findings of the Association for Computational Linguistics: NAACL 2024*, pp. 1051–1066, 2024.
- Chai, Z., Ren, Q., Xiao, X., Yang, H., Han, B., Zhang, S., Chen, D., Lu, H., Zhao, W., Yu, L., et al. Longer: Scaling up long sequence modeling in industrial recommenders. In *Proceedings of the Nineteenth ACM Conference on Recommender Systems*, pp. 247–256, 2025.
- Chen, J., Chi, L., Peng, B., and Yuan, Z. Hllm: Enhancing sequential recommendations via hierarchical large language models for item and user modeling. *arXiv preprint arXiv:2409.12740*, 2024.
- Cheng, H.-T., Koc, L., Harmsen, J., Shaked, T., Chandra, T., Aradhye, H., Anderson, G., Corrado, G., Chai, W., Ispir, M., et al. Wide & deep learning for recommender systems. In *Proceedings of the 1st workshop on deep learning for recommender systems*, pp. 7–10, 2016.
- Covington, P., Adams, J., and Sargin, E. Deep neural networks for youtube recommendations. In *Proceedings of the 10th ACM conference on recommender systems*, pp. 191–198, 2016.
- Di Palma, D., Biancofiore, G. M., Anelli, V. W., Narducci, F., Di Noia, T., and Di Sciascio, E. Evaluating chatgpt as a recommender system: A rigorous approach. *arXiv preprint arXiv:2309.03613*, 2023.
- Doan, T. T., Nguyen, L. M., Pham, N. H., and Romberg, J. Finite-time analysis of stochastic gradient descent under markov randomness. *arXiv preprint arXiv:2003.10973*, 2020.
- Geng, B., Huan, Z., Zhang, X., He, Y., Zhang, L., Yuan, F., Zhou, J., and Mo, L. Breaking the length barrier: Llm-enhanced ctr prediction in long textual user behaviors. In *Proceedings of the 47th International ACM SIGIR Conference on Research and Development in Information Retrieval*, pp. 2311–2315, 2024.
- Golovneva, O., Wang, T., Weston, J., and Sukhbaatar, S. Contextual position encoding: Learning to count what’s important. *arXiv preprint arXiv:2405.18719*, 2024.
- Guo, H., Tang, R., Ye, Y., Li, Z., and He, X. Deepfm: a factorization-machine based neural network for ctr prediction. *arXiv preprint arXiv:1703.04247*, 2017.
- Han, R., Yin, B., Chen, S., Jiang, H., Jiang, F., Li, X., Ma, C., Huang, M., Li, X., Jing, C., et al. Mtgr: Industrial-scale generative recommendation framework in meituan. In *Proceedings of the 34th ACM International Conference on Information and Knowledge Management*, pp. 5731–5738, 2025.
- He, X., Pan, J., Jin, O., Xu, T., Liu, B., Xu, T., Shi, Y., Atallah, A., Herbrich, R., Bowers, S., et al. Practical lessons from predicting clicks on ads at facebook. In *Proceedings of the eighth international workshop on data mining for online advertising*, pp. 1–9, 2014.
- He, Z., Xie, Z., Jha, R., Steck, H., Liang, D., Feng, Y., Majumder, B. P., Kallus, N., and McAuley, J. Large language models as zero-shot conversational recommenders. In *Proceedings of the 32nd ACM international conference on information and knowledge management*, pp. 720–730, 2023.
- Hou, Y., Mu, S., Zhao, W. X., Li, Y., Ding, B., and Wen, J.-R. Towards universal sequence representation learning for recommender systems. In *Proceedings of the 28th ACM SIGKDD conference on knowledge discovery and data mining*, pp. 585–593, 2022.
- Hou, Y., He, Z., McAuley, J., and Zhao, W. X. Learning vector-quantized item representation for transferable sequential recommenders. In *Proceedings of the ACM Web Conference 2023*, pp. 1162–1171, 2023.
- Jia, J., Wang, Y., Li, Y., Chen, H., Bai, X., Liu, Z., Liang, J., Chen, Q., Li, H., Jiang, P., et al. Learn: Knowledge adaptation from large language model to recommendation for practical industrial application. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 39, pp. 11861–11869, 2025.

- Kang, W.-C. and McAuley, J. Self-attentive sequential recommendation. In *2018 IEEE international conference on data mining (ICDM)*, pp. 197–206. IEEE, 2018.
- Kaplan, J., McCandlish, S., Henighan, T., Brown, T. B., Chess, B., Child, R., Gray, S., Radford, A., Wu, J., and Amodei, D. Scaling laws for neural language models. *arXiv preprint arXiv:2001.08361*, 2020.
- Krell, M. M., Kosec, M., Perez, S. P., and Fitzgibbon, A. Efficient sequence packing without cross-contamination: Accelerating large language models without impacting performance. *arXiv preprint arXiv:2107.02027*, 2021.
- Li, L., Zhang, Y., Liu, D., and Chen, L. Large language models for generative recommendation: A survey and visionary discussions. In *Proceedings of the 2024 Joint International Conference on Computational Linguistics, Language Resources and Evaluation (LREC-COLING 2024)*, pp. 10146–10159, 2024a.
- Li, R., Deng, W., Cheng, Y., Yuan, Z., Zhang, J., and Yuan, F. Exploring the upper limits of text-based collaborative filtering using large language models: Discoveries and insights. In *Proceedings of the 34th ACM International Conference on Information and Knowledge Management*, pp. 1643–1653, 2025.
- Li, S., Guo, H., Tang, X., Tang, R., Hou, L., Li, R., and Zhang, R. Embedding compression in recommender systems: A survey. *ACM Computing Surveys*, 56(5):1–21, 2024b.
- Lin, J., Dai, X., Xi, Y., Liu, W., Chen, B., Zhang, H., Liu, Y., Wu, C., Li, X., Zhu, C., et al. How can recommender systems benefit from large language models: A survey. *ACM Transactions on Information Systems*, 43(2):1–47, 2025.
- Lin, Z., Ding, H., Hoang, N. T., Kveton, B., Deoras, A., and Wang, H. Pre-trained recommender systems: A causal debiasing perspective. In *Proceedings of the 17th ACM International Conference on Web Search and Data Mining*, pp. 424–433, 2024.
- Liu, J., Liu, C., Zhou, P., Lv, R., Zhou, K., and Zhang, Y. Is chatgpt a good recommender? a preliminary study. *arXiv preprint arXiv:2304.10149*, 2023.
- Luo, S., He, B., Zhao, H., Shao, W., Qi, Y., Huang, Y., Zhou, A., Yao, Y., Li, Z., Xiao, Y., et al. Recranker: Instruction tuning large language model as ranker for top-k recommendation. *ACM Transactions on Information Systems*, 43(5):1–31, 2025.
- Ma, J., Zhao, Z., Yi, X., Chen, J., Hong, L., and Chi, E. H. 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*, pp. 1930–1939, 2018a.
- Ma, X., Zhao, L., Huang, G., Wang, Z., Hu, Z., Zhu, X., and Gai, K. Entire space multi-task model: An effective approach for estimating post-click conversion rate. In *The 41st International ACM SIGIR Conference on Research & Development in Information Retrieval*, pp. 1137–1140, 2018b.
- Mudigere, D., Hao, Y., Huang, J., Jia, Z., Tulloch, A., Sridharan, S., Liu, X., Ozdal, M., Nie, J., Park, J., et al. Software-hardware co-design for fast and scalable training of deep learning recommendation models. In *Proceedings of the 49th Annual International Symposium on Computer Architecture*, pp. 993–1011, 2022.
- Naumov, M., Mudigere, D., Shi, H.-J. M., Huang, J., Sundaraman, N., Park, J., Wang, X., Gupta, U., Wu, C.-J., Azzolini, A. G., et al. Deep learning recommendation model for personalization and recommendation systems. *arXiv preprint arXiv:1906.00091*, 2019.
- Ning, L., Liu, L., Wu, J., Wu, N., Berlowitz, D., Prakash, S., Green, B., O’Banion, S., and Xie, J. User-llm: Efficient llm contextualization with user embeddings. In *Companion Proceedings of the ACM on Web Conference 2025*, pp. 1219–1223, 2025.
- Petrov, A. V. and Macdonald, C. Generative sequential recommendation with gptrec. *arXiv preprint arXiv:2306.11114*, 2023.
- Pi, Q., Bian, W., Zhou, G., Zhu, X., and Gai, K. Practice on long sequential user behavior modeling for click-through rate prediction. In *Proceedings of the 25th ACM SIGKDD international conference on knowledge discovery & data mining*, pp. 2671–2679, 2019.
- Pi, Q., Zhou, G., Zhang, Y., Wang, Z., Ren, L., Fan, Y., Zhu, X., and Gai, K. Search-based user interest modeling with lifelong sequential behavior data for click-through rate prediction. In *Proceedings of the 29th ACM International Conference on Information & Knowledge Management*, pp. 2685–2692, 2020.
- Polyzotis, N., Zinkevich, M., Roy, S., Breck, E., and Whang, S. Data validation for machine learning. *Proceedings of machine learning and systems*, 1:334–347, 2019.
- Rajput, S., Mehta, N., Singh, A., Hulikal Keshavan, R., Vu, T., Heldt, L., Hong, L., Tay, Y., Tran, V., Samost, J., et al. Recommender systems with generative retrieval. *Advances in Neural Information Processing Systems*, 36: 10299–10315, 2023.

- Sculley, D., Holt, G., Golovin, D., Davydov, E., Phillips, T., Ebner, D., Chaudhary, V., Young, M., Crespo, J.-F., and Dennison, D. Hidden technical debt in machine learning systems. *Advances in neural information processing systems*, 28, 2015.
- Sheng, X.-R., Gao, J., Cheng, Y., Yang, S., Han, S., Deng, H., Jiang, Y., Xu, J., and Zheng, B. Joint optimization of ranking and calibration with contextualized hybrid model. In *Proceedings of the 29th ACM SIGKDD Conference on Knowledge Discovery and Data Mining*, pp. 4813–4822, 2023.
- Shin, K., Kwak, H., Kim, S. Y., Ramström, M. N., Jeong, J., Ha, J.-W., and Kim, K.-M. Scaling law for recommendation models: Towards general-purpose user representations. In *Proceedings of the AAAI conference on artificial intelligence*, volume 37, pp. 4596–4604, 2023.
- Si, Z., Guan, L., Sun, Z., Zang, X., Lu, J., Hui, Y., Cao, X., Yang, Z., Zheng, Y., Leng, D., et al. Twin v2: Scaling ultra-long user behavior sequence modeling for enhanced ctr prediction at kuaishou. In *Proceedings of the 33rd ACM International Conference on Information and Knowledge Management*, pp. 4890–4897, 2024.
- Sun, F., Liu, J., Wu, J., Pei, C., Lin, X., Ou, W., and Jiang, P. Bert4rec: Sequential recommendation with bidirectional encoder representations from transformer. In *Proceedings of the 28th ACM international conference on information and knowledge management*, pp. 1441–1450, 2019.
- Sun, Z., Si, Z., Zang, X., Zheng, K., Song, Y., Zhang, X., and Xu, J. Large language models enhanced collaborative filtering. In *Proceedings of the 33rd ACM International Conference on Information and Knowledge Management*, pp. 2178–2188, 2024.
- Tan, J., Xu, S., Hua, W., Ge, Y., Li, Z., and Zhang, Y. Idgenrec: Llm-recsys alignment with textual id learning. In *Proceedings of the 47th international ACM SIGIR conference on research and development in information retrieval*, pp. 355–364, 2024.
- Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J., Jones, L., Gomez, A. N., Kaiser, Ł., and Polosukhin, I. Attention is all you need. *Advances in neural information processing systems*, 30, 2017.
- Wang, R., Fu, B., Fu, G., and Wang, M. Deep & cross network for ad click predictions. In *Proceedings of the ADKDD’17*, pp. 1–7, 2017.
- Wu, L., Zheng, Z., Qiu, Z., Wang, H., Gu, H., Shen, T., Qin, C., Zhu, C., Zhu, H., Liu, Q., et al. A survey on large language models for recommendation. *World Wide Web*, 27(5):60, 2024.
- Xu, L., Zhang, J., Li, B., Wang, J., Chen, S., Zhao, W. X., and Wen, J.-R. Tapping the potential of large language models as recommender systems: A comprehensive framework and empirical analysis. *ACM Transactions on Knowledge Discovery from Data*, 19(5):1–51, 2025.
- Yu, T., Kumar, S., Gupta, A., Levine, S., Hausman, K., and Finn, C. Gradient surgery for multi-task learning. *Advances in neural information processing systems*, 33: 5824–5836, 2020.
- Yuan, Z., Yuan, F., Song, Y., Li, Y., Fu, J., Yang, F., Pan, Y., and Ni, Y. Where to go next for recommender systems? id-vs. modality-based recommender models revisited. In *Proceedings of the 46th International ACM SIGIR Conference on Research and Development in Information Retrieval*, pp. 2639–2649, 2023.
- Yue, Z., Rabhi, S., Moreira, G. d. S. P., Wang, D., and Oldridge, E. Llamarec: Two-stage recommendation using large language models for ranking. *arXiv preprint arXiv:2311.02089*, 2023.
- Zhai, J., Liao, L., Liu, X., Wang, Y., Li, R., Cao, X., Gao, L., Gong, Z., Gu, F., He, M., et al. Actions speak louder than words: Trillion-parameter sequential transducers for generative recommendations. *arXiv preprint arXiv:2402.17152*, 2024.
- Zhang, B. et al. Wukong: Towards a scaling law for large-scale recommendation, 2024a.
- Zhang, J., Xie, R., Hou, Y., Zhao, X., Lin, L., and Wen, J.-R. Recommendation as instruction following: A large language model empowered recommendation approach. *ACM Transactions on Information Systems*, 43(5):1–37, 2025.
- Zhang, Y. et al. Scaling law of large sequential recommendation models, 2024b.
- Zhao, Z., Fan, W., Li, J., Liu, Y., Mei, X., Wang, Y., Wen, Z., Wang, F., Zhao, X., Tang, J., et al. Recommender systems in the era of large language models (llms). *IEEE Transactions on Knowledge and Data Engineering*, 36(11):6889–6907, 2024.
- Zhou, G., Zhu, X., Song, C., Fan, Y., Zhu, H., Ma, X., Yan, Y., Jin, J., Li, H., and Gai, K. Deep interest network for click-through rate prediction. In *Proceedings of the 24th ACM SIGKDD international conference on knowledge discovery & data mining*, pp. 1059–1068, 2018.
- Zhou, G., Mou, N., Fan, Y., Pi, Q., Bian, W., Zhou, C., Zhu, X., and Gai, K. Deep interest evolution network for click-through rate prediction. In *Proceedings of the AAAI conference on artificial intelligence*, volume 33, pp. 5941–5948, 2019.

Zhu, Y., Wu, L., Guo, Q., Hong, L., and Li, J. Collaborative large language model for recommender systems. In *Proceedings of the ACM Web Conference 2024*, pp. 3162–3172, 2024.

A. Extended Background

A.1. The Performance–Efficiency Trade-off in Industrial CTR Prediction

The overall design of industrial-grade recommendation systems and their recommendation models almost always revolves around two goals: **performance** and **efficiency** (Covington et al., 2016; Naumov et al., 2019; Mudigere et al., 2022; Agarwal et al., 2023; Bai et al., 2025). Performance is not only reflected in the model’s fitting capabilities as measured by metrics such as AUC and PCOC, but also in its ability to sensitively capture user interest drift and changes in content distribution under varying traffic patterns (Sheng et al., 2023; Pi et al., 2020). Efficiency, on the other hand, is comprehensively reflected in the computational power consumption, memory/bandwidth usage, and online inference speed during the training and inference phases. Among these, training and inference costs are prerequisites for the long-term deployment and continuous iteration of the model in a real production environment (Naumov et al., 2019; Mudigere et al., 2022; Agarwal et al., 2023).

Although DLRMs have achieved considerable success, it faces bottlenecks in both performance and efficiency (Bai et al., 2025; Zhang et al., 2024a; Han et al., 2025). On the one hand, DLRMs rely on an experience- and rule-based feature system, which suffers from the inherent flaw of “strong memory, weak reasoning.” This makes the DLRMs insufficiently generalizable when dealing with new advertisements, new users, or scenarios requiring logical inference (Cheng et al., 2016; He et al., 2014; Ma et al., 2018b;a; Wu et al., 2024). At the same time, with the exponential growth of user behavior, traditional DLRMs suffer from significant information loss in ultra-long sequence modeling and has poor adaptability to different scenarios (Zhou et al., 2018; Pi et al., 2020; 2019; Zhang et al., 2024b). On the other hand, as the network design of DLRMs become increasingly complex, the performance improvement of the model shows diminishing returns. To achieve further performance improvements, it often requires exponentially increased computational costs, making the long-term deployment and continuous iteration of the model in a real production environment problematic (Zhang et al., 2024a;b; Mudigere et al., 2022; Han et al., 2025).

A.2. A Taxonomy of LLM-based Recommendation Research

LLMs have recently emerged as a promising direction for recommendation systems, giving rise to a growing line of research commonly referred to as LLM4Rec (Wu et al., 2024; Lin et al., 2025; Zhao et al., 2024; Li et al., 2024a). The motivation behind this paradigm shift lies in the inherent limitations of traditional ID-based recommendation models, which often struggle with semantic understanding, cold-start problems, and cross-domain generalization (Yuan et al., 2023; Li et al., 2025). LLMs offer the potential to introduce extensive world knowledge, robust reasoning capabilities, and high-quality textual generation into the recommendation pipeline (Zhang et al., 2025; He et al., 2023). However, integrating LLMs into industrial-scale systems presents unique challenges, primarily the “ID-Text dilemma”—where high-cardinality sparse IDs do not map naturally to the continuous token space of LLMs (Tan et al., 2024; Rajput et al., 2023)—and the prohibitive inference latency of decoder-only architectures in real-time scoring (Yue et al., 2023). Based on recent literature and industrial practices, LLM4Rec approaches can be systematically categorized into three distinct paradigms: *LLM as Recommender*, *LLM for Representation*, and *Generative Sequential Modeling*.

LLM as Recommender. This category explores the direct application of LLM capabilities—such as memory, reasoning, and zero-shot generalization—to core recommendation tasks including retrieval and ranking (Wu et al., 2024; Lin et al., 2025; Xu et al., 2025). Methods in this domain typically adapt recommendation data into natural language prompts, leveraging techniques like Instruction Tuning to align the LLM with recommendation objectives (Zhu et al., 2024; Zhang et al., 2025; Bao et al., 2023; Luo et al., 2025). While these methods demonstrate promise in explainability and conversational recommendation, their performance on traditional metrics (e.g., CTR) often falls short of specialized ID-based models (Liu et al., 2023; Di Palma et al., 2023; Cao et al., 2024). In recommendation scenarios, user behavior is heavily influenced by implicit feedback and specific context rather than the explicit semantic logic found in natural language; consequently, general world-knowledge reasoning does not necessarily translate effectively to modeling complex user–item interaction patterns (Bao et al., 2023; Cao et al., 2024; Zhu et al., 2024). Furthermore, the inference latency remains a significant bottleneck for real-time industrial deployment (Xu et al., 2025).

LLM for Representation. In this paradigm, LLMs function as sophisticated feature encoders (Lin et al., 2025; Wu et al., 2024). Instead of performing the ranking directly, the intermediate layers or final output embeddings of the LLM are extracted and utilized as semantic features to augment the input of traditional recommendation models (Sun et al., 2024; Jia et al., 2025; Geng et al., 2024; Chen et al., 2024; Ning et al., 2025). This approach aims to enhance the model’s semantic understanding without bearing the full cost of LLM inference during the serving phase. LLM-derived representations significantly mitigate the limitations of discrete feature models, particularly regarding the generalization capability for

long-tail items and cold-start users/ads (Hou et al., 2022; 2023). However, this methodology faces notable limitations. There is typically a limited gain on warm items, as the strong collaborative filtering signals derived from abundant historical interactions often outweigh the semantic benefits provided by the LLM (Hou et al., 2023; Lin et al., 2024). Furthermore, employing large-scale models for representation learning introduces a high inference cost, which creates substantial latency and resource bottlenecks during both the offline feature extraction and online serving phases (Lin et al., 2025).

Generative Sequential Modeling. This category represents a structural adaptation rather than a direct semantic application. It borrows the architectural innovations underlying LLMs—specifically the Transformer architecture, Causal Masking, and Long-context modeling capabilities—to reconstruct recommendation systems (Vaswani et al., 2017; Kang & McAuley, 2018; Sun et al., 2019). These models (such as GR models) treat user history as a sequence and the next item prediction as a generative task, similar to next token prediction (Kang & McAuley, 2018; Petrov & Macdonald, 2023; Han et al., 2025). By employing generative sequential modeling techniques and combining them with discrete features that precisely characterize user historical behavior, these models have shown significant potential (Han et al., 2025). A key observation in this domain is the emergence of “scaling laws” within recommendation systems, where model performance metrics improve predictably as the sequence length and model capacity increase (Shin et al., 2023; Zhang et al., 2024b), mirroring the trajectory seen in NLP.

A comparative overview of the three LLM4Rec paradigms is presented in Table 4, highlighting their core mechanisms, key strengths, and primary limitations.

Table 4. Comparison of LLM4Rec Paradigms.

Category	Core Mechanism	Key Strength	Primary Limitation
LLM as Recommender	Instruction tuning / prompting for ranking	Explainability, Zero-shot generalization	High latency, ineffective interaction modeling
LLM for Representation	Feature Encoder (Extracting embeddings)	Cold-start handling, Semantic understanding	Limited gain on warm items, high inference cost
Generative Seq. Modeling	Transformer decoder, next token prediction	Scaling laws, high capacity	High inference latency, high resource consumption

A.3. Limitations of GR in Performance

Existing GR models often inherit NLP-style causal Transformers with minimal adaptation to recommender logs, implicitly assuming that user history can be represented as a homogeneous token stream. In practice, recommendation data are inherently *heterogeneous*: events comprise multiple fields (e.g., item, context, query, creator), and user trajectories interleave distinct behavior types (e.g., exposure, click, like, skip, dwell). A naive serialization pipeline typically collapses this structured record into a single sequence of item IDs (or flattened tokens), which *discards action semantics*—the critical distinction between what the user was shown and how the user responded.

This structural mismatch leads to a performance bottleneck: the model conflates semantically different interactions, dilutes supervision signals, and learns spurious correlations (e.g., treating exposures as implicit positives or mixing weak/strong feedback). As a result, even with larger models and longer contexts, GR may underperform in industrial CTR settings where fine-grained behavior semantics and cross-field interactions are decisive, highlighting the need for action-aware, heterogeneity-preserving sequence modeling rather than direct NLP-style tokenization.

B. Extended Related Work

B.1. Limitations of Emerging Generative Ranking Models

While GR models have successfully introduced the scaling laws of LLMs into recommendation systems, their direct application to industrial CTR prediction faces distinct structural and optimization challenges.

Mitigation of Distribution Skew in Sequence Packing. To improve training efficiency with variable-length user sequences, standard GR models often employ sequence packing techniques borrowed from NLP (e.g., concatenating multiple short sequences). However, unlike NLP where samples are generally Independent and Identically Distributed (I.I.D.), packing in

recommendation systems groups multiple interactions from the *same user* into a single training instance to maintain context. This creates a severe distribution skew, where a mini-batch is dominated by highly correlated samples from a few users. This correlation causes the model—especially the sparse embedding parameters—to overfit specific user identities rather than learning generalizable interaction patterns.

Action heterogeneity. Existing GR models often treat user history as a homogeneous token stream, neglecting the inherent heterogeneity of recommendation data. This reliance on naive serialization discards critical action semantics—distinguishing what was shown from how the user responded—thereby diluting supervision signals and limiting performance in complex industrial scenarios (as discussed in Appendix A.3).

Explicit Modeling of Relative Action Signals. Standard GR models rely on the vanilla Scaled Dot-Product Attention, often supplemented only by absolute or relative positional encodings. While effective for capturing sequential order (as emphasized in LONGER), this approach treats the nature of the interaction as implicit. It fails to explicitly differentiate between varying feedback signals (e.g., “clicked” vs. “viewed” vs. “purchased”) and their relative timing in a query-dependent manner.

B.2. Comparative Discussion with Existing Ranking Models

To position GRAB within the evolving landscape of recommendation systems, we provide a qualitative comparison against two primary categories of existing models: traditional DLRMs and emerging GR approaches in Table 5.

Table 5. Comparison of DLRM, HSTU, LONGER, and GRAB based on Key Dimensions from the Paper

Dimension	DLRMs	HSTU	LONGER	GRAB
Core Architecture	Embedding + MLP + Target Attention (DIN/SIM)	Pure Transformer (Generative Ranking)	Transformer-based (Optimized for Long Seq)	End-to-End Generative Framework
Feature Interaction	Manual Cross Features & MLP / Gating	Standard Self-Attention (Global)	Standard Self-Attention (Long Context)	Causal Action-aware Multi-channel Attention (CamA)
Heterogeneous Features	Concatenation & Flattening (Fixed-length vector)	Homogeneous Tokens (Often redundant static info)	Homogeneous Tokens (Similar to HSTU)	Heterogeneous Tokens (Partial History / Full Candidate)
Action Semantic Modeling	Implicit or via Engineered Features	Implicit (Learned via sequence order)	Implicit (Focus on sequence length)	Action-aware RAB (Explicit Relative Bias)
Training Strategy	Standard Supervised (Pointwise BCE)	Autoregressive / Next-Token (Suffers from Dist. Skew)	Autoregressive (Standard GR Training)	Sequence-Then-Sparse (STS) (Decoupled Optimization)
Scaling Potential	Diminishing Returns (Saturated Performance)	High (Follows Scaling Laws)	High (Scaling in Sequence Length)	High & Linear (No saturation observed)

C. Detailed Architecture and Data Flow of GRAB

C.1. Sparse Feature Layer

In the sparse feature layer, GRAB expands user u ’s raw behavior sequence and the candidate ad sequence into event-level representations, preserving their original temporal order to form a structured, time-ordered event sequence:

$$\{S_t^{beh}\}_{t=1}^{T_u}, \quad \{S_i^{ad}\}_{i=1}^{N_u} \quad (8)$$

where T_u denotes the length of user u ’s behavior history, N_u denotes the number of candidate advertisements for user u . GRAB takes user behavior history and candidate ads as input sequences. Specifically, the t -th behavior event S_t^{beh} consists of user attributes U , context C , and behavior attributes B , while the i -th candidate ad S_i^{ad} consists of context C and item attributes A , as follows:

$$S_t^{beh} = (U_u, C_t, B_t), \quad S_i^{ad} = (A_i, C_i) \quad (9)$$

Following the discrete feature engineering standard of DLRMs, we apply a structured expansion function Φ to transform each event into a fixed multi-field representation. Subsequently, each field value is mapped to a discrete ID via a sparse PSTable II. The event-level representations of the raw behavior sequence and the candidate ad sequence can be obtained as:

$$\begin{aligned} \mathbf{x}_t^{beh} &= \Pi(\Phi(S_t^{beh})), & t &= 1, \dots, T_u \\ \mathbf{x}_i^{ad} &= \Pi(\Phi(S_i^{ad})), & i &= 1, \dots, N_u \end{aligned} \quad (10)$$

C.2. Dense Tokenizer

Unlike the DLRM approach, which concatenates all field embeddings into a fixed-length flattened vector, GRAB preserves the event structure by aggregating field embeddings within each event into a single event token. This yields a time-ordered token sequence that is fed into a Transformer to capture long-range behavioral dependencies and interest drift. Given the structured discrete ID sequences from Section C.1, GRAB converts each event into a dense token for Transformer-based sequential modeling. Specifically, each event is first transformed into a vector through a field-wise embedding lookup followed by a multi-field fusion process, as follows:

$$v_t = \text{Emb}(x_t), \quad v_i = \text{Emb}(x_i) \quad (11)$$

Subsequently, the token representation for each event is generated by the GateMLP module, which consists of an MLP and a Gate layer, formulated as:

$$\mathbf{e} = \text{GateMLP}(v) \in \mathbb{R}^{d_{\text{model}}} \quad (12)$$

$$\mathbf{E}^{\text{beh}} = \{\mathbf{e}_t^{\text{beh}}\}_{t=1}^{T_u}, \quad \mathbf{E}^{\text{ad}} = \{\mathbf{e}_i^{\text{ad}}\}_{i=1}^{N_u}.$$

C.3. Autoregressive-like Sequence Modeling Layer

Taking the output (\mathbf{E}^{beh} and \mathbf{E}^{ad}) from the previous dense layer as input, we feed it into the sequence modeling layer (for more details about this layer, please see Section 3.3) to capture sequential dependencies. Let Z denote the final output representation of the sequence layer:

$$Z = \text{SeqLayer}(\mathbf{E}^{\text{beh}}, \mathbf{E}^{\text{ad}}). \quad (13)$$

Finally, the output of the sequence modeling layer is fed into a logistic head to yield the CTR prediction. The model is trained to minimize the binary cross-entropy loss:

$$\hat{y} = \sigma(\mathbf{w}^\top Z + b) \in (0, 1), \quad (14)$$

$$\mathcal{L}_{\text{BCE}} = -[y \log \hat{y} + (1 - y) \log(1 - \hat{y})].$$

D. In-Depth Analysis of STS Training

While sequence packing dramatically improves computational resource utilization by eliminating padding, it introduces a non-trivial optimization challenge known as Distribution Skew. In this section, we provide the theoretical justification for the proposed STS training paradigm, detail its mathematical formulation, and discuss how it reconciles the learning space inconsistency between stages.

D.1. Distribution Skew

Standard Stochastic Gradient Descent (SGD) relies on the assumption that samples within a mini-batch are I.I.D.. Formally, for a loss function \mathcal{L} , the gradient g computed on a batch \mathcal{B} is an unbiased estimator of the true gradient:

$$\mathbb{E}[g_{\mathcal{B}}] = \nabla \mathcal{L}, \quad \text{Var}(g_{\mathcal{B}}) = \frac{\sigma^2}{|\mathcal{B}|} \quad (15)$$

where $|\mathcal{B}|$ is the batch size.

In sequence packing, we form a packed mini-batch $\mathcal{B}_{\text{pack}}$ by concatenating multiple actions from the same user u into a single training instance (or a user-dominated batch) to avoid padding waste. While efficient, this construction makes samples within $\mathcal{B}_{\text{pack}}$ highly correlated: for $i, j \in \mathcal{B}_{\text{pack}}$, $\text{Cov}(x_i, x_j) \gg 0$, because they share the same user_id, static context, and long-term interests. As a result, the effective batch size is substantially reduced and the variance of the stochastic gradient estimator increases, yielding noisier and less stable updates.

This issue is most damaging for the sparse embedding table Φ . Since a packed batch repeatedly contains the same user features (e.g., user_id=123 appears L times along the packed sequence), the update for that single embedding vector is amplified by repeated contributions:

$$\Delta \Phi_u \propto \sum_{t=1}^L \nabla \mathcal{L}_t. \quad (16)$$

Such oversized, user-specific updates encourage Φ to memorize individual trajectories rather than learn generalizable interaction patterns. Meanwhile, the dense sequence model (e.g., Transformer) suffers from batch-to-batch distribution skew: consecutive packed batches may be dominated by different users (User A \rightarrow User B), causing abrupt shifts in inputs and gradients, which hinders stable convergence of sequential reasoning parameters.

D.2. Formalization of STS Stages

To mitigate the distribution skew, STS decouples the optimization into two orthogonal objectives: relational reasoning (Dense) and feature representation (Sparse). The algorithm flow of STS is shown in Algorithm 1.

Algorithm 1 Two-Stage Alternating Training Strategy

Require: Training dataset \mathcal{D} , Learning rate η

- 1: **Initialize:** Dense tokenizer θ_{cont} , Causal Transformer θ_{tr} , Sparse Embedding Table Φ
 - 2: **while** not converged **do**
 - 3: **// Stage I: Sequence Modeling (Sequence Phase)**
 - 4: Sample packed user sequences batch $Z = (x_{hist}, x_{cand}) \sim \mathcal{D}$
 - 5: Freeze Sparse Embedding Table Φ
 - 6: Compute sequence output:
 - 7: $\hat{y}_{seq} \leftarrow f_{seq}(Z; \theta_{cont}, \theta_{tr}, \Phi)$
 - 8: Compute sequence loss:
 - 9: $\mathcal{L}_{seq} \leftarrow \text{BCE}(\hat{y}_{seq}, y)$
 - 10: Update dense modules:
 - 11: $\{\theta_{cont}, \theta_{tr}\} \leftarrow \text{SGD}(\nabla_{\{\theta_{cont}, \theta_{tr}\}} \mathcal{L}_{seq})$
 - 12: **// Stage II: Sparse Feature Learning (Sparse Phase)**
 - 13: Sample independent user-ad batch $(x^{beh}, x^{ad}) \sim \mathcal{D}$
 - 14: Freeze Sequential modules $\{\theta_{cont}, \theta_{tr}\}$
 - 15: Compute aggregated features (breaking distribution skewness):
 - 16: $s \leftarrow \text{Agg}(\{\Phi(x_t^{beh})\}) \parallel \Phi(x^{ad})$
 - 17: Compute sparse phase prediction:
 - 18: $\hat{y}_{sp} \leftarrow f_{sp}(s)$
 - 19: Update sparse embeddings:
 - 20: $\Phi \leftarrow \text{SGD}(\nabla_{\Phi} \mathcal{L}_{sp})$
 - 21: **end while**
-

D.3. Discussion: A Pre-training & Transfer Perspective

The STS paradigm can be viewed through the lens of pre-training and transfer learning. Stage I serves as a sequential pre-training step that encodes interest evolution into the dense space, while Stage II transfers these insights back to the target sparse feature space for fine-tuning. Although this non-end-to-end mode might introduce a subtle objective inconsistency between stages, our results demonstrate that the benefit of resolving distribution skew far outweighs the cost of misalignment. STS ensures that the end-to-end sequence predictor (Stage I) consistently validates against the optimal embeddings refined in Stage II, providing a pragmatic path to balance efficiency and generalization in large-scale recommendation systems.

E. System Deployment Details

E.1. Platform Architecture

As illustrated in Fig. 8, the proposed system is implemented within a comprehensive Online Advertising System that operates as a closed-loop platform integration of Online Services and Offline Services. The architecture is designed to handle high-concurrency requests while maintaining model freshness through continuous updates.

E.1.1. ONLINE SERVING

The online serving component processes user interactions in real-time. The workflow initiates with a Page View (PV) Request, which sequentially passes through *matching* and *ranking* phases to select appropriate advertisements.

The core of the ranking mechanism involves a CTR Prediction module that relies on two primary inputs:

- **User Representation:** A user model processes historical tokens and maintains a KV-cache to efficiently manage state.
- **Ad Representation:** The system generates ad tokens corresponding to candidate advertisements.

This process results in the display of ads, where user interactions are captured in the online behavior log, consisting of impression logs and action logs.

E.1.2. OFFLINE TRAINING AND FEEDBACK LOOP

The offline component ensures the model evolves with user behavior. The process involves:

- **Data Collection and Storage:** Logs are collected and stored in behavior storage, which organizes data into User (user_id, gender), AD (ad_id, brand), and Context (location, device, behavior) categories.
- **Feature Engineering:** The system performs sparse feature engineering on the collected logs.
- **Training:** Training data is grouped by user ID (uid) to facilitate model offline training.
- **Deployment:** Updated models are pushed back to the online environment via an hourly release mechanism, updating the CTR prediction and user model components.

E.2. Deployment Constraints and Optimization

The system, referred to as GRAB, has been deployed in a feed ad ranking system handling billions of daily requests. The deployment addresses several critical engineering challenges distinguishable from conventional DLRMs:

- **Computational Complexity:** Unlike memory-bound DLRMs, GR in this context is compute-bound. This is primarily attributed to the quadratic complexity ($O(L^2)$) of Transformer self-attention mechanisms required for processing long sequences.
- **Latency Requirements:** The system’s performance is bound by critical latency thresholds defined in the Service Level Agreements (SLAs).

To meet these demands, the deployment utilizes a co-designed architecture incorporating data compression, hierarchical storage, and disaggregated serving.

E.3. Data Infrastructure and Training Optimization

User-Centric Data Layout and Compression. Traditional industrial pipelines use time-partitioned logs, necessitating expensive global shuffling. We transitioned to a User Slice storage architecture, where user behavior logs are pre-aggregated by User ID into contiguous physical file blocks. To further reduce I/O overhead, we upgraded from standard Gzip text storage to a Binary + LZ4 compression scheme. The binary format combined with LZ4 (which is highly efficient for repetitive user behavior patterns) reduced storage costs by 12% compared to Gzip and decreased decoding latency by 70%, enabling the system to stream complete user histories with near-linear scalability.

Hierarchical Parameter Server (PaddleBox). To handle terabyte-scale embedding tables, we utilized the training framework with a three-tier storage hierarchy:

- **L1 (GPU HBM):** Stores hot embeddings for the current micro-batch.
- **L2 (CPU DRAM):** Buffers warm parameters.
- **L3 (SSD):** Utilizes NVMe SSDs for massive long-tail feature embeddings. An intelligent prefetching engine asynchronously moves parameters between tiers, masking SSD I/O latency.

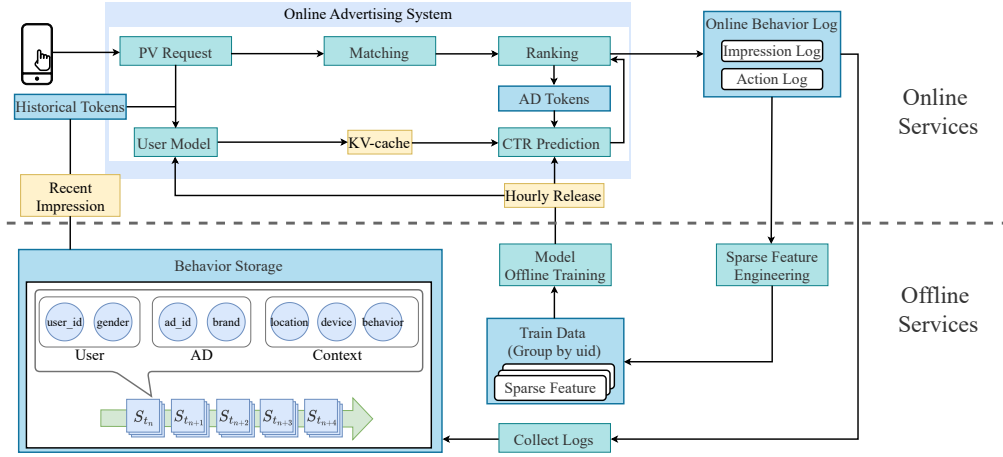


Figure 8. Overview of an online advertising CTR system with an online-offline closed loop. Online services handle PV requests via matching and ranking, and feed the CTR predictor with user-side historical tokens (maintained by a user model with KV-cache) and candidate ad tokens; user interactions are continuously logged as impression/action logs. Offline services collect these logs, apply sparse feature engineering, group training samples by user ID, and perform offline training; updated models are released (e.g., hourly) back to online.

Handling Long-Tail Sequences. Real-world user history lengths exhibit a heavy-tailed distribution, where the top 5% of sequences can cause “Out Of Memory” (OOM). We implemented an Inverse Sliding Window strategy during training. Instead of random slicing, sequences are sliced from the most recent action backwards. This prioritizes recent user interests and ensures that extreme long-tail data does not destabilize GPU memory usage.

E.4. High-Performance Inference Architecture

Disaggregated Serving and Parallelism. We adopted a disaggregated serving architecture using a User Interest Center (UIC). The UIC asynchronously computes and updates the Transformer’s Key-Value (KV) cache triggered by user actions. Crucially, we implemented Parallel Material Recall, where the user’s historical sequence encoding overlaps with the candidate generation (ad retrieval) phase. By the time candidate ads are retrieved, the user’s dense state is already computed, significantly hiding latency.

KV-Cache Reuse and M-FALCON. To avoid re-computing the user history for every candidate ad, we integrated the M-FALCON algorithm (Zhai et al., 2024). It utilizes a broadcast-attention mechanism where the fetched user KV cache is shared across a micro-batch of candidate items, reducing the marginal inference complexity per item from quadratic to linear.

Operator Fusion and Mixed Precision. To maximize throughput on GPUs, we employed aggressive operator fusion (e.g., fusing Gemm + Bias + LayerNorm), which reduced kernel launch overheads and improved latency by roughly 43%. Furthermore, rather than simple INT8 quantization which may degrade accuracy, we adopted a Mixed Precision strategy: utilizing TF32 for Transformer matrix operations to accelerate computation and FP16 for fully connected layers, achieving a balance between 28% performance gain and negligible precision loss.

E.5. Data Consistency

A major challenge in GRs is the Freshness Gap (Train-Serve Skew). We addressed this by implementing a streaming data pipeline based on Flink & TableStore. We utilized a Global Strictly Incremental ID mechanism to ensure strict ordering of user actions across distributed nodes. This allows the inference engine to fetch the exact state of the user corresponding to the training checkpoint, reducing data synchronization delay from minutes to seconds and ensuring the model always predicts based on the most consistent context.

F. Discussion

F.1. Limitations and Challenges

Operational Complexity of Two-Stage Training. A primary limitation of GRAB lies in the operational overhead introduced by the Sequence Then Sparse (STS) training paradigm. While STS effectively resolves the distribution skew caused by sequence packing and stabilizes the optimization of dense versus sparse parameters, it inherently complicates the model iteration pipeline. Unlike standard DLRMs that support continuous, single-stage online learning, STS requires a decoupled scheduling of sequence modeling (freezing sparse features) and sparse feature learning (freezing dense parameters). This increases the engineering maintenance cost and introduces latency in incorporating fresh feature embeddings into the dense sequential context, potentially affecting the model’s responsiveness to emerging trends in real-time environments.

Compute-Bound Hardware Constraints. The shift from DLRMs to Generative Ranking marks a fundamental transition from memory-bound to compute-bound workloads in recommendation infrastructure. Although we mitigated inference latency via optimizations like Action-aware RAB and KV-cache reuse (Appendix E.4), the quadratic complexity of attention mechanisms—even when bounded by sliding windows—remains computationally heavier than the MLP layers of traditional models. Scaling GRAB to significantly longer sequences (e.g., user lifetimes spanning months) or deploying it on edge devices with limited compute capacity poses a significant challenge, necessitating further research into linear attention mechanisms or more aggressive token pruning strategies specifically tailored for recommendation data.

Interpretability of Generative Signals. While the Multi-channel Attention mechanism allows us to inspect which behavior subsets contribute to a prediction, the end-to-end generative nature of GRAB can obscure the precise “why” behind a ranking decision compared to feature-engineered linear models. Understanding whether a prediction is driven by short-term intent (sequential reasoning) or long-term habit (sparse memorization) remains a non-trivial task, which is critical for debugging bad cases in commercial systems.

F.2. Future Directions

Towards Multimodal Generative Ranking. Currently, GRAB operates on discretized ID tokens derived from categorical features. However, the architecture’s Heterogeneous Token design (Section 3.3.2) is naturally extensible to other modalities. A promising future direction is to integrate raw modal tokens—such as image patches (Visual Tokens) or ad textual descriptions (Language Tokens)—directly into the interaction sequence. By leveraging GRAB’s strong sequence modeling capabilities, the model could learn semantic alignment between visual cues and user behaviors end-to-end, overcoming the information loss inherent in pre-extracted ID features and enabling true multimodal CTR prediction.

Unified Generative Representation Across Domains. Finally, the “pre-training & fine-tuning” paradigm common in NLP has yet to be fully realized in industrial recommendation. We envision extending GRAB to learn a Universal User Representation by pre-training on diverse behavior logs across multiple business scenarios (e.g., Home Feed, Search, and Short Video). A unified GRAB model could transfer learned sequential patterns from data-rich domains to cold-start scenarios, effectively solving the “data silo” problem prevalent in large-scale platforms.

Foundation for Agent-based Recommender Systems. GRAB’s ability to model the transition probabilities of user states ($s_t \rightarrow s_{t+1}$) positions it as a powerful “World Model” or User Simulator for future agent-based recommendation systems. By accurately predicting not just the next click, but the evolution of user interests over time, GRAB can serve as the environment model for Reinforcement Learning (RL) agents. This would allow the system to move beyond myopic CTR optimization toward maximizing Long-Term Value (LTV) or user satisfaction by simulating how current recommendations influence future user trajectories.