

An Efficient Large Recommendation Model: Towards a Resource-Optimal Scaling Law

Songpei Xu, Shijia Wang*, Da Guo, Xianwen Guo, Qiang Xiao[†], Fangjian Li, Chuanjiang Luo
 {xusongpei, wangshijia1, guoda, guoxianwen, hzxiaoqiang, hzlifangjian, luochuanjiang03}@corp.netease.com
 NetEase Cloud Music
 Hangzhou, China

Abstract

The pursuit of scaling up recommendation models confronts intrinsic tensions between expanding model capacity and preserving computational tractability. While prior studies have explored scaling laws for recommendation systems, their resource-intensive paradigms—often requiring tens of thousands of A100 GPU hours—remain impractical for most industrial applications. This work addresses a critical gap: achieving sustainable model scaling under strict computational budgets. We propose Climber, a resource-efficient recommendation framework comprising two synergistic components: the ASTRO model architecture for algorithmic innovation and the TURBO acceleration framework for engineering optimization. ASTRO (Adaptive Scalable Transformer for Recommendation) adopts two core innovations: (1) multi-scale sequence partitioning that reduces attention complexity from $O(n^2d)$ to $O(n^2d/N_b)$ via hierarchical blocks, enabling more efficient scaling with sequence length; (2) dynamic temperature modulation that adaptively adjusts attention scores for multimodal distributions arising from inherent multi-scenario and multi-behavior interactions. Complemented by TURBO (Two-stage Unified Ranking with Batched Output), a co-designed acceleration framework integrating gradient-aware feature compression and memory-efficient Key-Value caching, Climber achieves 5.15× throughput gains without performance degradation.

Comprehensive offline experiments on multiple datasets validate that Climber exhibits a more ideal scaling curve. To our knowledge, this is the first publicly documented framework where controlled model scaling drives continuous online metric growth (12.19% overall lift) without prohibitive resource costs. Climber has been successfully deployed on Netease Cloud Music, one of China’s largest music streaming platforms, serving tens of millions of users daily. These advancements establish a new paradigm for industrial recommendation systems, demonstrating that coordinated algorithmic-engineering innovation—not just brute-force scaling—unlocks sustainable performance growth under real-world computational constraints.

CCS Concepts

• Information systems → Retrieval models and ranking.

Keywords

Recommendation system, Transformer, Scaling law

ACM Reference Format:

Songpei Xu, Shijia Wang, Da Guo, Xianwen Guo, Qiang Xiao, Fangjian Li, Chuanjiang Luo. 2018. An Efficient Large Recommendation Model: Towards a Resource-Optimal Scaling Law. In *Proceedings of Make sure to enter the correct conference title from your rights confirmation email (Conference acronym 'XX)*. ACM, New York, NY, USA, 12 pages. <https://doi.org/XXXXXXX.XXXXXXX>

1 Introduction

Scaling laws, initially explored in language models [14, 19], establish predictable relationships between model performance and key factors such as model size, training data volume, and computational budget. For instance, Kaplan et al. [19] demonstrated that transformer-based language models follow power-law improvements in perplexity as model parameters and token counts increase. Similar trends have been observed in vision [9] and multimodal models [2], where scaling model dimensions and data diversity directly correlate with the performance of downstream tasks.

Recent research has shown that scaling laws also apply to recommendation systems, providing valuable insights into model design and resource allocation [1, 11]. The HSTU model [38] employs hierarchical self-attention mechanisms to model long-term user behavior sequences, achieving better performance than traditional Transformers. Similarly, the MARM model [26] introduces memory augmentation techniques to reduce computational complexity, enabling multi-layer sequence modeling with minimal inference costs. However, these approaches demand prohibitive computational resources (e.g., thousands of GPU hours or terabytes of memory), rendering them impractical for real-world deployment. Furthermore, the interplay between key scaling factors—sequence length, model depth, and heterogeneous user behaviors—remains underexplored, leading to suboptimal resource allocation and diminished returns on scaling efforts.

Inspired by the DeepSeek series [3, 22, 23], which has significantly enhanced the efficiency of large language model (LLM) development and reduced computational resource costs, we aim to address the following question: **how can we scale up the recommendation model at a substantially lower cost?** To gain insights, we conducted industrial-scale analysis on two mainstream models, the Deep Learning Recommendation Model (DLRM) [28] and the

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 acronym 'XX, June 03–05, 218, Woodstock, NY

© 2018 Copyright held by the owner/author(s). Publication rights licensed to ACM.

ACM ISBN 978-1-4503-XXXX-X/18/06

<https://doi.org/XXXXXXX.XXXXXXX>

*Corresponding Author

†Corresponding Author

Transformer model. In Figure 1(a), we present the scaling curves for both DLRM* and Transformer models, and an oracle curve is included to represent an ideal scaling curve, characterized by a higher starting point and a larger slope. In Figure 1(b), we present an AUC curves derived from simulations of various combinations of sequence lengths and layer number, and introduce the concept of "performance interval", which represents the range of AUC variation for the model under the equivalent FLOPs. However, our findings reveal that some issues still remain in recommendation systems:

- **Transformer's degraded performance under FLOPS constraints:** As shown in Figure 1(a), the FLOPs corresponding to the crossover point are $10^{8.2}$. Using this FLOPs value as the boundary, the performance comparison between DLRM and Transformer shows different trends. Transformer model outperforms traditional architectures such as DLRM when FLOPs exceed $10^{8.2}$ magnitude. However, when FLOPs less than $10^{8.2}$, Transformer model shows worse performance compared to DLRM. This highlights the pursuit for models with the oracle curve shown in Figure 1(a), which enables the model to achieve better performance even when the FLOPs are limited.
- **The impact of factor combinations on model performance under equivalent FLOPs:** In recommender systems, factors such as sequence length and layer number significantly influence FLOPs, and different combinations of these factors can lead to varying model performance. For instance, there is a obvious performance interval (nearly 1%) under different combinations at 10^9 FLOPs, as shown in Figure 1(b). Current researches lack comprehensive analysis of how factor combinations impact recommendation model performance, which hinders efficient scaling for models.
- **Huge resource consumption:** As depicted in Figure 1(a), to outperform the DLRM, Transformers require FLOPs on the order of 10^9 to 10^{10} . However, the FLOPs of the currently deployed DLRM models are only at the order of 10^7 , indicating a $100\times$ increase of FLOPs to achieve comparable performance, which requires tens of thousands of GPU hours in practice. These resource demands are already beyond the affordability of most industrial applications. Consequently, achieving resource-optimal scaling remains a significant challenge for industrial recommender systems.

Driven by the above insights, we introduce Climber, a novel framework that rethinks the scaling paradigm for recommendation systems. At its core, Climber integrates two complementary innovations: Adaptive Scalable Transformer for Recommendation (ASTRO) model architecture and Two-stage Unified Ranking with Batched Output (TURBO) acceleration framework. ASTRO redefines how recommendation systems handle heterogeneous user behaviors by introducing multi-scale sequence partitioning, which decomposes user behavior sequences into smaller, fine-grained subsequences. This approach not only reduces computational complexity, but also enables more precise modeling of user interests across diverse scenarios. In addition, ASTRO incorporates dynamic temperature modulation, a mechanism that adaptively adjusts attention scores to account for the varying importance of different

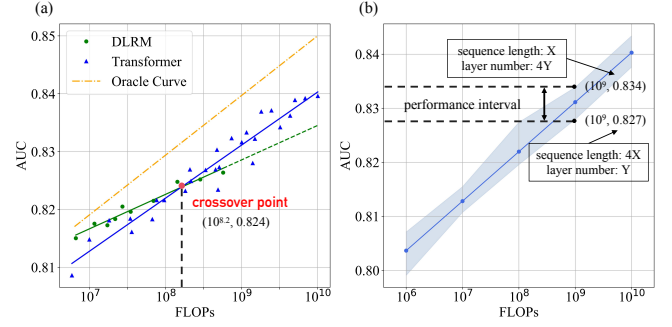


Figure 1: (a) Scalability: DLRM vs. Transformer. The left part of crossover point indicates that Transformer underperforms DLRM when FLOPs are limited. The right part of crossover point highlights the increasing computational demands as Transformer performance improves. **(b) Toy Example: Simulated AUC curve with performance interval.** This figure depicts an AUC curve derived from simulations of various combinations of sequence length and layer number.

behaviors and scenarios. On the engineering side, TURBO introduces a unified ranking framework that optimizes both training and inference efficiency. As a unified ranking framework, TURBO transforms traditional "Single User, Single Item" samples into a format that aligns with actual online requests, namely "Single User, Multiple Items". TURBO accelerates forward propagation during training and inference with encoder-level KV Cache, achieving significant efficiency improvements. Finally, we investigate the scalability of our proposed model and the impact of factor combinations on AUC under equivalent FLOPs. These insights enable more rational resource allocation and guide the identification of key factors for rapid model scaling.

Our contributions are mainly categorized as follows:

- We present the first industrial-scale study of scaling laws in recommendation systems from the computational resources-constrained aspects, explicitly quantifying the impact of factor combinations under equivalent FLOPs. This analysis reveals that balanced scaling—alternating sequence and depth expansions—yields both offline and online metric gains.
- We propose a novel Transformer variant, ASTRO, which resolves the scaling dilemmas in recommendation systems through multi-scale partitioning and adaptive temperature modulation. To our knowledge, the proposed method pioneers sustainable scaling—delivering +12.19% online metric growth, which is the highest annual improvement in our production system.
- A unified framework bridging training-inference gaps via dynamic feature compression and block-parallel KV caching. Deployed on Netease Cloud Music, TURBO sustains $5.15\times$ throughput gains, enabling $100\times$ model scaling with moderate increase in computational resources.

2 RELATED WORK

Wukong [39] explored parameter scaling in retrieval models but relied on strong assumptions about feature engineering. HSTU [38] reformulated recommendations as sequential transduction tasks,

*The baseline model that has been implemented for online deployment.

achieving trillion-parameter models with linear computational scaling via hierarchical attention and stochastic sequence sampling. However, HSTU’s focus on generative modeling left gaps in bridging traditional feature-based DLRMs. Concurrently, MARM [26] proposed caching intermediate attention results to reduce inference complexity from $O(n^2d)$ to $O(nd)$, empirically validating cache size as a new scaling dimension. While effective, MARM’s caching strategy assumes static user patterns, overlooking real-time behavior shifts.

Techniques to mitigate computational costs have been widely adopted. In NLP, KV caching [10, 25] avoids redundant attention computations during autoregressive inference. MARM adapted this idea to recommendations by storing historical attention outputs, enabling multi-layer target-attention with minimal FLOPs overhead. Similarly, HSTU introduced Stochastic Length to sparsify long sequences algorithmically, reducing training costs by 80% without quality degradation. For advertisement retrieval, Wang et al. [36] designed R/R^* , an eCPM-aware offline metric, to estimate online revenue scaling laws at low experimental costs. These works collectively highlight the importance of tailoring efficiency strategies to recommendation-specific constraints, such as high-cardinality features and millisecond-level latency requirements.

3 Method

To achieve efficient scaling, we first propose the ASTRO model, which reduces computational complexity through multi-scale sequence partitioning and adapts Transformer architecture to recommender systems by integrating multi-scenario and multi-behavior characteristics. Besides, we introduce a co-design acceleration framework, TURBO, which incorporates dynamic compression and encoder-level KV cache to enhance the efficiency of training and inference without performance degradation.

3.1 Scaling Dilemma in Recommender Systems

In Recommendation Systems(RS), model scaling mainly involves two approaches: feature scaling and model capacity scaling. Feature scaling, particularly the extension of user behavior sequences, has been a focal point, with extensive research demonstrating its effectiveness in enhancing model performance [6, 7, 29, 32]. Meanwhile, the integration of Transformer architectures has markedly enhanced the modeling capabilities of recommendation systems[8, 24, 33, 42]. However, combining these two approaches for synchronous scaling results in a quadratic increase in computational cost. This increase not only raises the resource requirements for training and inference but also puts more stress on real-time efficiency. Therefore, achieving a balance between performance improvement and computational efficiency has become an urgent challenge in recommendation systems.

In NLP, words form continuous sequences with explicit syntactic and semantic relationships. In contrast, user behavior sequences in RS are inherently fragmented, spanning multiple scenarios (e.g., browsing, purchasing) and heterogeneous behaviors (e.g., clicks, likes)[15, 16, 21, 40]. When fed into Transformer models, these fragmented sequences often lead to disordered attention distributions, as the model struggles to prioritize relevant behaviors within sparse and irregular patterns. This inefficiency limits the Transformer’s

scalability, particularly under computational constraints, where specialized models like DLRMs often outperform due to their tailored handling of sparse data.

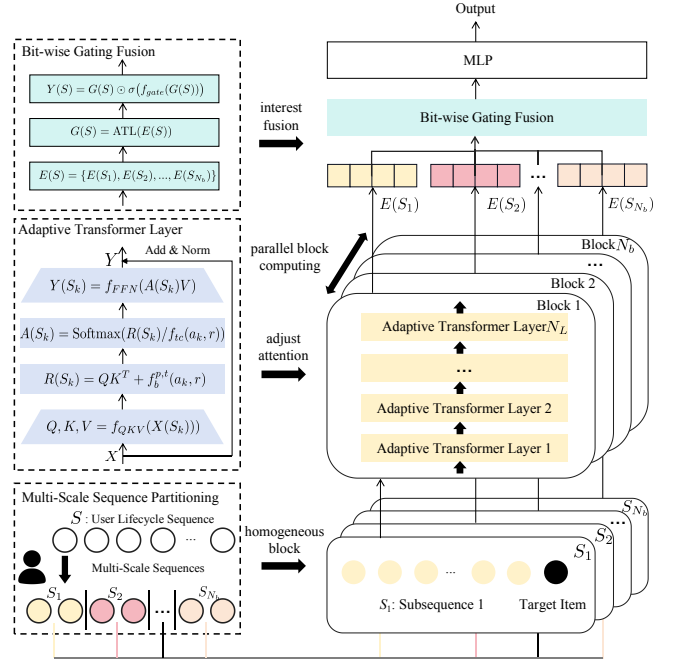


Figure 2: ASTRO Model Architecture.

3.2 ASTRO Model Architecture

3.2.1 Overall. To address the dual challenges of computational complexity and fragmented sequence structures in recommendation systems, we propose ASTRO model. This model integrates the multi-behavior and multi-scenario characteristics of recommendation systems into the Transformer architecture, while considering the resource demands associated with scaling. Our model comprises three modules: Multi-scale Sequence Partitioning (MSP), Adaptive Transformer Layer (ATL), and Bit-wise Gating Fusion (BGF). Specifically, MSP generates multi-scale sequences from user lifecycle sequence. These multi-scale sequences represent different types of subsequences (multi-scenario/multi-behavior/short and long-term behaviors). For each subsequence, we apply a corresponding transformer block composed of stacked ATLs to extract user interest. Additionally, we extend the time span of important subsequences to cover the user’s entire lifecycle. ATL uses an adaptive temperature coefficient in the activation function to adjust the attention distribution across different subsequences in various scenarios. Finally, BGF aggregates the representations from adaptive Transformer blocks corresponding to each subsequence to produce a unified output, which enables interest fusion between multi-scale sequences through ATL and a bit-wise gating mechanism. Figure 2 illustrates the detailed workflow of ASTRO model.

3.2.2 Multi-Scale Sequence Partitioning. The approach of MSP involves reorganizing user sequences based on different strategies,

which can be represented by the following formula:

$$S = \{x_1, x_2, \dots, x_{n_s}\} \quad (1)$$

$$S_k = \text{MSP}(S, a_k) = \{x_1^{a_k}, x_2^{a_k}, \dots, x_{n_k}^{a_k}\} \quad (2)$$

where S represents the user lifecycle sequence, x_i denotes i -th item ID from the entire item set \mathbf{X} . S_k represents the k -th subsequence extracted from the user lifecycle sequence S based on the extraction strategy a_k , and $x_j^{a_k}$ indicates the j -th item in the subsequence S_k . The n_s and n_k denote the length of S and S_k , respectively. We assume that there are a total of N_b extraction strategies and $\sum_{k=1}^{N_b} n_k = n$. In our practical application, $n \ll n_s$. This is because the extraction strategy typically retains only the user's positive behaviors. Consequently, the computational complexity under a single Transformer can be reduced from $O(n_s^2 d)$ to $O(n^2 d)$. We further improve the training process by employing the corresponding Transformer block for each subsequence S_k , thereby achieving a time complexity of $O(n_k^2 d)$. We design each extraction strategy to extract subsequences of equal length, such that $n_k = n/N_b$. In the case of fully serial operations, this results in a time complexity of $O(n^2 d/N_b)$. Therefore, even when $N_b = 2$, we can still achieve substantial training acceleration. In summary, MSP reduces the computational complexity from $O(n^2 d)$ to $O(n^2 d/N_b)$, and transform user lifecycle sequence into multi-scale sequences. These enhancements improve the efficiency and scalability of the recommendation system.

3.2.3 Adaptive Transformer Layer. The Softmax activation function, which normalizes the attention scores, plays a pivotal role in the Transformer architecture. Specifically, the attention mechanism is normalized by $\frac{QK^T}{\sqrt{d_k}}$ and subsequently multiplied by V . The division by $\sqrt{d_k}$ ensures that the distribution of the attention matrix aligns with that of Q and K [13]. However, we generate multi-scale sequences based on different extraction strategies from user lifecycle sequence, and apply corresponding Transformer block to each subsequence. A single scaling factor $\sqrt{d_k}$ is insufficient to accommodate the diverse requirements of all Transformer blocks corresponding to multi-scale sequences [12]. To further refine the attention distribution within each Transformer block, we introduce an adaptive temperature coefficient for each layer of each block. We refer to this change as Adaptive Transformer Layer (ATL), which can be mathematically expressed as follows:

$$\begin{aligned} Q, K, V &= f_{QKV}(X(S_k)), \\ R(S_k) &= QK^T + f_b^{p,t}(a_k, r), \\ A(S_k) &= \text{Softmax}(R(S_k)/f_{tc}(a_k, r)), \\ Y(S_k) &= f_{FFN}(A(S_k)V). \end{aligned} \quad (3)$$

Compared to conventional Transformer layer, we introduce an adaptive temperature coefficient and adjust relative attention bias from a recommendation perspective. Here, $X(S_k) \in \mathbb{R}^{s \times d}$, $R(S_k) \in \mathbb{R}^{h \times s \times s}$, $A(S_k) \in \mathbb{R}^{h \times s \times s}$ and $Y(S_k) \in \mathbb{R}^{s \times d}$ represent layer input, raw attention matrix, normalized attention matrix, and layer output, respectively. s, h, d represent sequence length, head number and feature dimension, respectively. $f_{QKV}(X(S_k))$ is used to derive query, key, and value matrices from input $X(S_k)$. $f_b^{p,t}(a_k, r)$ denotes relative attention bias [31, 38] that incorporates positional

(p) and temporal (t) information. $f_{tc}(a_k, r)$ denotes a function that derives the temperature coefficient. It is important to note that both extraction strategy a_k and recommendation scenario r influence the relative attention bias and the temperature coefficient. f_{FFN} processes the attention-weighted value matrix through a feedforward neural network (FFN) to produce the final output of the layer. This approach is inspired by the multi-scenario and multi-behavior characteristics of recommendation systems. Specifically, the adaptive temperature coefficient allows for more flexible attention weighting, addressing the limitations of the fixed scaling factor $\sqrt{d_k}$ in capturing the inheritance of diverse behaviors and scenarios.

3.2.4 Bit-wise Gating Fusion. However, when user lifecycle sequence is separated into multi-scale sequences, there is a lack of interaction between them [20, 37]. Therefore, we propose a bit-wise gating fusion module that integrates information across the N_b subsequences corresponding to our N_b blocks. Specifically, each block generates an output vector $E(S_k)$ through each adaptive Transformer block, and these N_b output vectors are concatenated as $E(S) \in \mathbb{R}^{N_b \times d}$. The concatenated vectors are then processed through a new ATL, followed by a sigmoid activation function to achieve bit-level gating. Finally, the vectors pass through a subsequent network to produce the final output score. The bit-wise gating fusion module can be formally expressed as:

$$\begin{aligned} E(S) &= \{E(S_1), E(S_2), \dots, E(S_{N_b})\}, \\ G(S) &= \text{ATL}(E(S)), \\ Y(S) &= G(S) \odot \sigma(f_{gate}(G(S))). \end{aligned} \quad (4)$$

where σ is sigmoid activation and f_{gate} represents a squeeze-and-excitation module [17], which ensure the identity between the input and output dimensions to dynamically adjusts the contribution of each element in $G(S) \in \mathbb{R}^{N_b \times d}$. $f_{gate}(G(S)) \in \mathbb{R}^{N_b \times d}$ and $Y(S) \in \mathbb{R}^{N_b \times d}$ represent the bit-wise attention matrix and the output of fusion module, respectively. ATL (Adaptive Transformer Layer in Equation 3) is used to calculate the similarity between different blocks to interact the different interest from subsequences. In contrast, the ATL in fusion function doesn't incorporate relative attention bias, and the temperature coefficient is solely determined by the recommendation scenario. The attention operation of ATL on N_b blocks can be regarded as field-wise interactions [17, 18], which facilitate information exchange at the feature level. Our method enhances interest fusion by adding bit-wise interactions. This allows the model to capture precise relationships among sequences. As a result, the model's ability to understand multi-scale sequences is significantly improved. It is worth mentioning that although this attention mechanism has a complexity of $O(N_b^2 d)$, the number of extraction strategy N_b is much smaller than the dimension d of the feature vector in the fusion stage. Thus the computational complexity of bit-wise gating fusion module is relatively low.

3.3 TURBO Acceleration Framework

A typical ranking model for recommendation systems involves elaborate features from one user and one item as model inputs. Probabilities or scores for specific tasks such as is-click or not are calculated as model outputs through a well-designed neural network. Such an arrangement for model inputs and outputs can

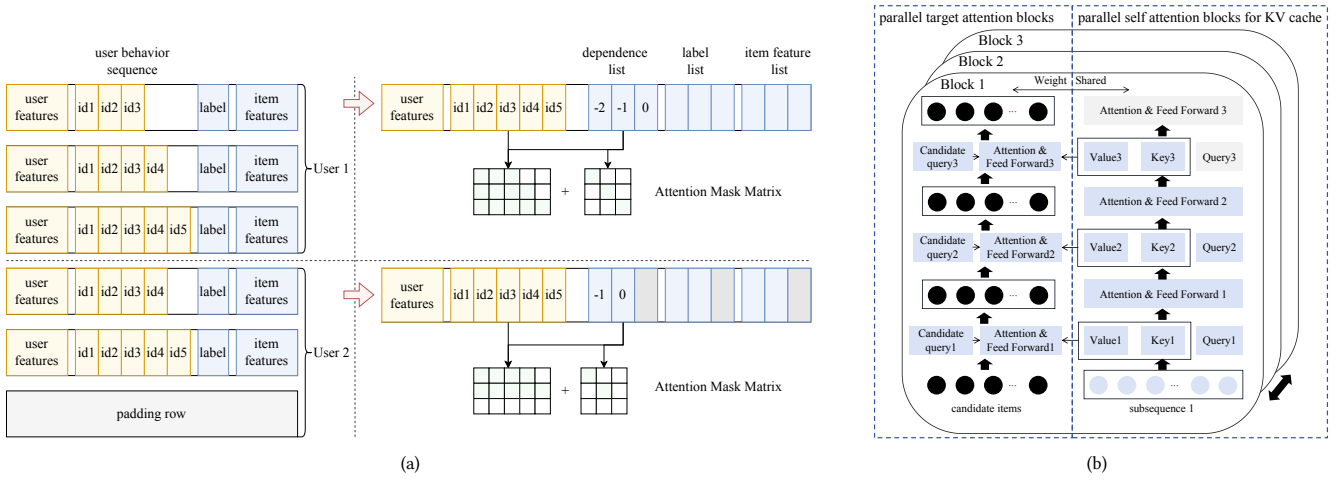


Figure 3: Two key features of TURBO: (a) Sketch of compression algorithm on training dataset; (b) Sketch of encoder-level KV cache with block-wise parallelism

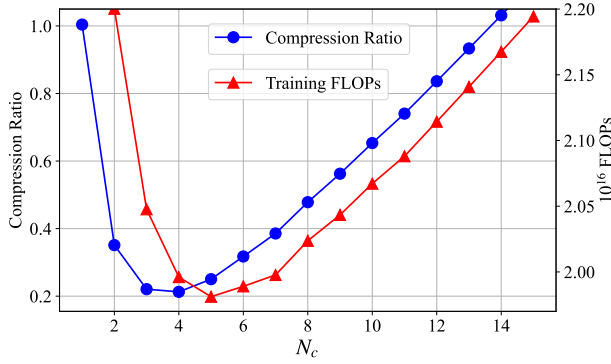


Figure 4: Compression Ratio and Training FLOPs based on our dataset

be summarized as "Single User, Single Item (SUSI)". It's instinctive to expand the first dimension of input tensor to batch size during training stage to achieve "Multiple Users, Multiple Items (MUMI)".

When it comes to the online serving system, "Single User, Multiple Items (SUMI)" is required to get scores for multiple candidates in a single request from one single user. Simply stacking user features to the number of candidates as "MUMI" seems like a common approach. However, "MUMI" might not be affordable for online serving due to the extra time costs caused by repeated computation of user feature. This redundant computation is almost unacceptable especially for transformer-dominated architectures with user features composed by long sequences. Zhai, et al.[38] proposed the M-FALCON algorithm to achieve "SUMI"-style batched inference. Online throughput performance can be boosted with the aid of kv-cache [30] and micro-batch parallelism.

Further, we propose TURBO, formulating these prefill-encode two stages into one forward pass. With re-arrangement of training dataset via dynamic compression algorithm explained below, TURBO serves as a unified framework compatible for both offline training and online inference.

3.3.1 Dynamic Compression on Training Data. To accommodate the training phase to the "SUMI" setting in TURBO, a fixed number of labels and features from the item side related to the same user need to be aggregated to the form of label list and feature list. The list can be padded to a constant size N_c across all users, and we call N_c as the degree of aggregation. It should be noted that user behavior might slowly accumulate over a period of time. If user behavior sequence grows up within an aggregated batch, the longest as well as the latest sequence should be kept to avoid loss of information. What's more, an extra dependence list need to be stored for each row of training sample as shown in Figure 3. With the help of this auxiliary list, we can build customized attention mask matrix for each row to prevent violation of causality.

The value of N_c should be determined carefully for a given training dataset. Either too small or too large value of N_c would be harmful to data compression and computation efficiency. Sample distribution $g(n)$ of a given dataset can be defined as the number of unique user who have the total number of record n . For instance, $g(2) = 1000$ means that there are 1000 unique users each has 2 row of records in a given dataset. Let N denote total number of record of training dataset, N_p denote total number of rows to be padded, D_{user} and D_{item} denote the average number of bytes in each row to store user features (including user behavior sequence) and item features respectively. The compression ratio α can be estimated by the following formulation 7. The term "sizeof(int)" shows the part of storage for dependence list composed by integer element.

$$p(n) = \left\lfloor \frac{n}{N_c} \right\rfloor \cdot N_c - n \quad (5)$$

$$N_p = \sum_n g(n) \cdot p(n) \quad (6)$$

$$\alpha = \frac{(N + N_p)D_{user}/N_c + (N + N_p)(D_{item} + \text{sizeof}(int))}{N(D_{user} + D_{item})} \quad (7)$$

3.3.2 Encoder-level KV-cache with Block-wise Parallelism. Since each block under ASTRO architecture is independent, calculations between blocks can be parallelized. We propose to fuse all transformer blocks with the same number of layers into a whole module. The first dimension in input tensors, weight and bias parameter is extended, which represents the block number in multi-head attention module and FFN module. Block-wise parallelism is actually achieved through scheduling of larger matrix operations at CUDA level.

To achieve lower computation costs, we followed the idea in M-FALCON to conduct a two-stage forward pass. The first stage is to get key/value cache vector for items in user behavior sequence. For the last layer in transformer module, results of kv-cache can be returned immediately without subsequent calculating like attention scores and FFN module. We named this trick "last layer short-circuit" to save a certain amount of computation. The second stage is to get scores for target items in a batched way with kv-cache obtained in the first stage. Computation cost in the second stage is proportional to the size of micro-batch, which is N_c during model training phase. The total training FLOPs F_{total} for one epoch can be estimated as summation of the two stages:

$$F_{total} = \frac{N + N_p}{N_c} \cdot (F_{stage1} + F_{stage2}(N_c)) \quad (8)$$

We plot the total training FLOPs for one epoch of our dataset against N_c in Figure 4, along with the compression ratio estimated through Equation 7. It can be shown that the optimal N_c to achieve maximum storage savings and the one to achieve minimum training cost are very close. Thus, we determined to use $N_c = 5$ as the base setting for the rest of the discussion.

4 Experiments

In this section, we detail the offline and online experiments conducted on real industrial data to evaluate our proposed method, addressing the following four research questions:

- RQ1: How does ASTRO model perform in offline evaluation compared to state-of-the-art (SOTA) models?
- RQ2: How does ASTRO model demonstrate superior scalability compared to DLRM and Transformer?
- RQ3: How can we allocate resources to scale up our model by considering the impact of different factor combinations on AUC under equivalent FLOPs?
- RQ4: How does Climber framework perform in industrial recommendation systems?

4.1 Experimental Setting

4.1.1 Dataset. To validate the effectiveness of our method in recommendation systems, we construct a dataset using real user behavior sequence as the primary feature. The dataset is used to predicts

different user actions on the candidate items based on historical interactions. Important user behaviors include full-play, like, share and comment. In order to protect data privacy, we only present statistical data for our recommendation scenarios. Additionally, we evaluated our model on three recommendation datasets, **Spotify**[4], **30Music**[34] and **Amazon-Book**[27]. A more detailed analysis of these datasets after preprocessing is presented in Appendix B.1.

4.1.2 Compared methods. In our comparative study, we select DLRM as benchmark model. DLRM is a model that leverages life-long sequences and complex feature interactions, which has been deployed in our online systems. In addition to DLRM, we incorporate several SOTA models to provide a comprehensive evaluation. These models include: DIN[41], TWIN[6], Transformer[35], and HSTU[38]. To systematically evaluate the model architectures, we focus on three key components: multi-scale sequence partitioning, adaptive Transformer layer and bit-wise gating fusion. We utilize the ASTRO-large model ($6\times$ layer number, and $4\times$ sequence length in Industrial setting) to validate the model's scalability. Throughout the experiments, we ensure fairness and consistency by adopting a unified configuration for both training hyper-parameters and feature representations across all models. More details about compared methods and experimental setting are presented in Appendix B.2.

4.2 Overall Performance (RQ1)

4.2.1 Performance Comparison. As shown in Table 1, the ASTRO model achieve the best performance in four recommendation datasets. Notably, the application scenarios of **Spotify** and **30Music** belong to the same domain as our industrial dataset, which is music recommendation systems. In contrast, **Amazon-Book** differs significantly from our scenario. However, our model still achieves favorable results on this dataset, indicating its potential adaptability to diverse applications. Next, we focus on comparing the AUC improvements of ASTRO relative to other methods. 1) as our primary online model, DLRM outperforms DIN and TWIN because DLRM includes a wide range of feature interaction structures beyond lifelong sequence and attention mechanisms. 2) The Transformer achieves +0.134% AUC improvement over TWIN. This is because Transformer computes the similarity between all historical items and target item in a single stage. 3) HSTU implements several enhancements to the Transformer architecture, achieving an AUC improvement of +0.036% on our dataset compared to Transformer. However, these enhancements also result in increased computational complexity. 4) ASTRO reduces computational complexity by sequence partitioning and adjusts attention distribution across multi-scenario and multi-behavior through adaptive temperature coefficients. This results in a +0.170% improvement over DLRM. Furthermore, ASTRO-large achieves a +2.21% AUC improvement by scaling up the model, achieving the largest offline gain in the past year.

4.2.2 Ablation Study. To systematically evaluate the contributions of each component in our ASTRO model, we conduct a comprehensive set of experiments across multiple datasets. For illustrative purposes, we focus on the detailed ablation study conducted on the Industrial dataset and select Transformer for comparison with ASTRO series. By incorporating MSP, ASTRO (-ATL, -BGF) transforms user lifecycle sequence into multi-scale subsequence blocks.

Table 1: Evaluation of methods on public/industrial datasets

	Spotify		Amazon-Book		30Music		Industrial	
	AUC	LogLoss	AUC	LogLoss	AUC	LogLoss	AUC	LogLoss
DLRM(Baseline)	0.7606	0.5761	0.7842	0.5541	0.8927	0.2012	0.8216	0.7067
DIN	0.7557	0.5803	0.7796	0.5580	0.8861	0.2044	0.8158	0.7109
TWIN	0.7589	0.5772	0.7831	0.5563	0.8903	0.2019	0.8203	0.7078
Transformer	0.7621	0.5735	0.7836	0.5557	0.8930	0.2007	0.8214	0.7074
HSTU	0.7626	0.5722	0.7869	0.5520	0.8938	0.1982	0.8217	0.7053
ASTRO (-ATL,-BGF)	0.7635	0.5710	0.7873	0.5518	0.8944	0.1978	0.8221	0.7045
ASTRO (-BGF)	0.7655	0.5694	0.7881	0.5510	0.8950	0.1972	0.8225	0.7034
ASTRO	0.7663	0.5687	0.7887	0.5501	0.8986	0.1957	0.8230	0.7029
ASTRO-large	0.7702	0.5666	0.7914	0.5472	0.9035	0.1916	0.8398	0.6911
%Improve	+1.26%	-1.64%	+0.91%	-1.24%	+1.20%	-4.77%	+2.21%	-2.20%

This enhancement results in a positive AUC gain of +0.085% compared to the Transformer model. ASTRO (-BGF) further improves the model by introducing an adaptive temperature coefficient. This component adjusts the attention distribution dynamically, leading to an AUC improvement of +0.134%. Finally, the incorporation of BGF brings about an AUC improvement of +0.195%. This module integrates user interests represented by different subsequences, emphasizing the importance of interest fusion in recommendation systems. In summary, our model demonstrates strong performance and adaptability based on the offline evaluations across various datasets.

4.3 Scalability (RQ2)

Before discussing model scalability, we formally define FLOPs:

$$C \propto s * l \quad (9)$$

where s represents the sequence length and l represents the layer number in the model. In large-scale scenarios, the quadratic computational complexity of attention mechanisms accounts for only a minor proportion of the model's overall FLOPs even with longer sequence [5, 14, 19]. Thus in our computational analysis, we focus solely on the linear part of sequence length s within $C \propto s * l$ relationship. The scaling curves of DLRM, Transformer, and ASTRO are shown in Figure 5(a). Although Transformer can achieve better performance than DLRM when FLOPs exceed 10^9 , its efficiency is notably lower than DLRM between 10^7 and 10^8 . Compared to the Transformer, ASTRO exhibits a more ideal scaling curve due to its higher starting point and larger slope. When FLOPs are below $10^{7.5}$, ASTRO's performance remains weaker than DLRM, but the crossover point shifts to the left, enabling the ASTRO model to achieve performance transformation more efficiently than Transformer. In this experiment, the two main factors affecting FLOPs are layer number and sequence length. For the ASTRO, we illustrate the relationship between model performance and the layer number in Figure 5(b) and sequence length in Figure 5(c), respectively. When the sequence length is fixed, model performance improves in a manner similar to a power-law with the increase in the layer number; when the layer number is fixed, there is a similar improvement in model performance as the sequence length increases. Thus, our proposed ASTRO model exhibits scaling curves in terms of FLOPs,

sequence length, and layer number, and has a more efficient scaling curve compared to Transformer.

4.4 Efficient Allocation (RQ3)

From Figure 5(b, c), it is evident that increasing the sequence length and the layer number can improve the model's AUC. However, the priority of these two parameters has not been extensively discussed. Table 2 presents the model's AUC for the same FLOPs with different layer number and sequence length. It is clear that under the same FLOPs, the combination of layer number and sequence length can lead to significant changes in the offline testing AUC. According to $C \propto s * l$, the product of layer number and sequence length remains constant for the same FLOPs. When the FLOPs is 4.11×10^8 , the model achieves the best performance of 0.8301 AUC on a $(400s \times 4l)$ model; When the FLOPs is 1.01×10^9 , the model achieves the best performance of 0.8335 AUC on a $(400s \times 8l)$ model. We observe that expanding a single factor may limit the model's development. Therefore, it is best to consider both layer number and sequence length equally when scaling up the model. For example, with a $(400s \times 4l)$ model, if we need to increase FLOPs by 4 times, we can choose between $(1600s \times 4l)$, $(800s \times 8l)$, and $(400s \times 16l)$. From Table 2, it can be found that the best choice is $(800s \times 8l)$, which jointly expands both factors to increase the model's AUC from 0.8301 to 0.8382. This conclusion also guides how to allocate resources online. In our practical recommendation systems, usually only one factor is chosen for each iteration. Therefore, when scaling up online, we alternate between increasing sequence length and layer number.

4.5 Online A/B Test (RQ4)

As shown in Table 3, the results of deploying Climber framework are presented. Consistent with the conclusion that sequence length and layer number are equally important, our model demonstrates online scaling curves for metric and FLOPs by only adjusting sequence length and layer number. Firstly, ASTRO model with 5.82×10^6 FLOPs exhibits a negative metric improvement. When the FLOPs value of the ASTRO (6x) model match that of DLRM (6x), there is only a slight negative metric, indicating that ASTRO model's efficiency is lower with fewer FLOPs compared with DLRM. Furthermore, when the FLOPs value of ASTRO (479x) model is 2.79×10^9 , the online metric improvement reaches +7.78%. Finally, when the

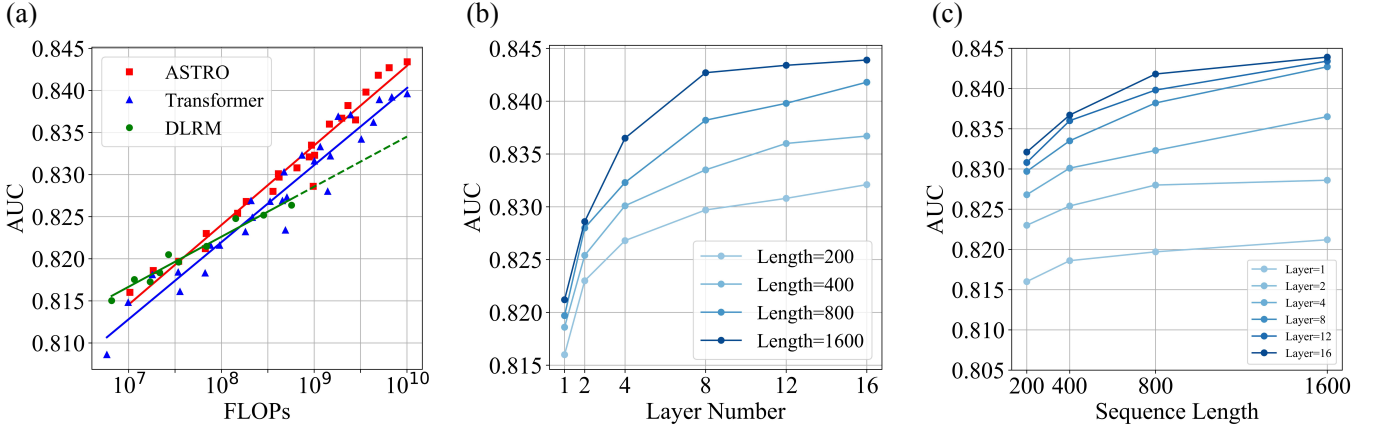


Figure 5: Scaling Curve. (a) Scalability: DLRM vs Transformer vs ASTRO in Industrial Dataset. (b) Model performance of scaling layer number with constant sequence length. (c) Model performance of scaling sequence length with constant layer number.

Table 2: Performance Comparison under Equivalent FLOPs

FLOPs	Sequence Length	Layer Number	AUC
4.11×10^8	1600	1	0.8212
	800	2	0.8280
	400	4	0.8301
	200	8	0.8297
1.01×10^9	1600	2	0.8286
	800	4	0.8323
	400	8	0.8335
	200	16	0.8321
2.55×10^9	1600	4	0.8365
	800	8	0.8382
	400	16	0.8367

Table 3: Online Metric Improvement Compared with DLRM

Method	FLOPs	Sequence Length	Layer Number	Online Metric
DLRM	3.45×10^7 (6×)	-	-	+0%
Climber	5.82×10^6 (1×)	100	1	-4.95%
	1.84×10^7 (3×)	400	1	-1.31%
	3.46×10^7 (6×)	800	1	-1.22%
	4.11×10^8 (71×)	400	4	+3.65%
	1.01×10^9 (174×)	800	4	+4.29%
	2.79×10^9 (479×)	1600	4	+7.78%
	2.31×10^9 (397×)	800	8	+10.68%
	3.61×10^9 (620×)	800	12	+12.19%

FLOPs value of ASTRO (620×) model is 3.61×10^9 , an online metric improvement of +12.19% is achieved. Despite a 100× increase in FLOPs compared to DLRM, the change in online inference efficiency is regarded as acceptable in practical online applications due to TURBO acceleration framework. Our TURBO framework

achieves a 5.15× acceleration during training and a 2.92× acceleration in online inference with the same FLOPs. This significantly boosts the efficiency of both training and inference stage. Currently, this model has been applied to serve all user traffic in our actual music recommendation system. It is noteworthy that there has been only a moderate increase in the CPU/GPU/Memory resources throughout this process. To the best of our knowledge, ASTRO is the first recommendation model to display both offline and online scaling curves while maintaining resource balance. Moreover, it achieves +12.19% metric improvement, representing the largest improvement in the past year.

5 Conclusion

The Climber method is designed to scale up recommendation models under resource constraints, which consists of two main components: the ASTRO model architecture and the TURBO acceleration framework. The ASTRO model integrates the multi-scenario and multi-behavior characteristics of recommendation systems into the Transformer architecture through multi-scale sequence partitioning, adaptive Transformer layer, and bit-wise gating fusion, while reducing the model's time complexity. This integration enables the model to exhibit superior scalability in offline evaluation compared to both DLRM and Transformer models. Furthermore, we introduce the TURBO acceleration framework, which employs dynamic compression and an encoder-level KV cache. This framework allows the deployment of models that are 100× more complex without increasing prohibitive computing resources. The Climber method demonstrates a scaling curve online and achieves a 12.19% improvement in online metric. Additionally, experimental findings have elucidated the impact of sequence length and layer number on model performance under equivalent FLOPs, emphasizing the significance of both factors in scaling up model. This insight contributes to the rational allocation of resources. Our findings demonstrate the superiority of the Climber method in efficiently implementing scaling laws while maintaining resource balance. This capability empowers recommendation systems to overcome resource constraints and explore larger models.

References

- [1] Newsha Ardalani, Carole-Jean Wu, Zeliang Chen, Bhargav Bhushanam, and Adnan Aziz. 2022. Understanding scaling laws for recommendation models. *arXiv preprint arXiv:2208.08489* (2022).
- [2] Jinze Bai, Shuai Bai, Shusheng Yang, Shijie Wang, Sinan Tan, Peng Wang, Junyang Lin, Chang Zhou, and Jingren Zhou. 2023. Qwen-vl: A frontier large vision-language model with versatile abilities. *arXiv preprint arXiv:2308.12966* (2023).
- [3] Xiao Bi, Deli Chen, Guanting Chen, Shanhuang Chen, Damai Dai, Chengqi Deng, Honghui Ding, Kai Dong, Qiusi Du, Zhe Fu, et al. 2024. Deepseek llm: Scaling open-source language models with longtermism. *arXiv preprint arXiv:2401.02954* (2024).
- [4] Brian Brost, Rishabh Mehrotra, and Tristan Jehan. 2019. The music streaming sessions dataset. In *The World Wide Web Conference*. 2594–2600.
- [5] Adam Casson. 2023. Transformer FLOPs. (2023). <https://adamcasson.com/posts/transformer-flops>
- [6] Jianxin Chang, Chenbin Zhang, Zhiyi Fu, Xiaoxue Zang, Lin Guan, Jing Lu, Yiqun Hui, Dewei Leng, Yanan Niu, Yang Song, et al. 2023. TWIN: Two-stage interest network for lifelong user behavior modeling in CTR prediction at kuaishou. In *Proceedings of the 29th ACM SIGKDD Conference on Knowledge Discovery and Data Mining*. 3785–3794.
- [7] Qiwei Chen, Changhua Pei, Shanshan Lv, Chao Li, Junfeng Ge, and Wenwu Ou. 2021. End-to-end user behavior retrieval in click-through rate prediction model. *arXiv preprint arXiv:2108.04468* (2021).
- [8] Qiwei Chen, Huan Zhao, Wei Li, Pipei Huang, and Wenwu Ou. 2019. Behavior sequence transformer for e-commerce recommendation in alibaba. In *Proceedings of the 1st international workshop on deep learning practice for high-dimensional sparse data*. 1–4.
- [9] Mehdi Cherti, Romain Beaumont, Ross Wightman, Mitchell Wortsman, Gabriel Ilharco, Cade Gordon, Christoph Schuhmann, Ludwig Schmidt, and Jenia Jitsev. 2023. Reproducible scaling laws for contrastive language-image learning. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 2818–2829.
- [10] Harry Dong, Xinyu Yang, Zhenyu Zhang, Zhangyang Wang, Yuejie Chi, and Beidi Chen. 2024. Get More with LESS: Synthesizing Recurrence with KV Cache Compression for Efficient LLM Inference. *arXiv preprint arXiv:2402.09398* (2024).
- [11] Wei Guo, Hao Wang, Luankang Zhang, Jin Yao Chin, Zhongzhou Liu, Kai Cheng, Qiusi 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).
- [12] Yu-Lin He, Xiao-Liang Zhang, Wei Ao, and Joshua Zhexue Huang. 2018. Determining the optimal temperature parameter for Softmax function in reinforcement learning. *Applied Soft Computing* 70 (2018), 80–85.
- [13] Geoffrey Hinton. 2015. Distilling the Knowledge in a Neural Network. *arXiv preprint arXiv:1503.02531* (2015).
- [14] Jordan Hoffmann, Sebastian Borgeaud, Arthur Mensch, Elena Buchatskaya, Trevor Cai, Eliza Rutherford, Diego de Las Casas, Lisa Anne Hendricks, Johannes Welbl, Aidan Clark, et al. 2022. Training compute-optimal large language models. *arXiv preprint arXiv:2203.15556* (2022).
- [15] Yupeng Hou, Zhankui He, Julian McAuley, and Wayne Xin Zhao. 2023. Learning vector-quantized item representation for transferable sequential recommenders. In *Proceedings of the ACM Web Conference 2023*. 1162–1171.
- [16] Yupeng Hou, Shanlei Mu, Wayne Xin Zhao, Yaliang Li, Bolin Ding, and Ji-Rong Wen. 2022. Towards universal sequence representation learning for recommender systems. In *Proceedings of the 28th ACM SIGKDD Conference on Knowledge Discovery and Data Mining*. 585–593.
- [17] Jie Hu, Li Shen, and Gang Sun. 2018. Squeeze-and-excitation networks. In *Proceedings of the IEEE conference on computer vision and pattern recognition*. 7132–7141.
- [18] Tongwen Huang, Zhiqi Zhang, and Junlin Zhang. 2019. FiBiNET: combining feature importance and bilinear feature interaction for click-through rate prediction. In *Proceedings of the 13th ACM conference on recommender systems*. 169–177.
- [19] 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. *arXiv preprint arXiv:2001.08361* (2020).
- [20] Chao Li, Zhiyuan Liu, Mengmeng Wu, Yuchi Xu, Huan Zhao, Pipei Huang, Guoliang Kang, Qiwei Chen, Wei Li, and Dik Lun Lee. 2019. Multi-interest network with dynamic routing for recommendation at Tmall. In *Proceedings of the 28th ACM international conference on information and knowledge management*. 2615–2623.
- [21] Jiacheng Li, Ming Wang, Jin Li, Jinmiao Fu, Xin Shen, Jingbo Shang, and Julian McAuley. 2023. Text is all you need: Learning language representations for sequential recommendation. In *Proceedings of the 29th ACM SIGKDD Conference on Knowledge Discovery and Data Mining*. 1258–1267.
- [22] Aixian Liu, Bei Feng, Bin Wang, Bingxuan Wang, Bo Liu, Chenggang Zhao, Chengqi Deng, Chong Ruan, Damai Dai, Daya Guo, et al. 2024. Deepseek-v2: A strong, economical, and efficient mixture-of-experts language model. *arXiv preprint arXiv:2405.04434* (2024).
- [23] Aixian Liu, Bei Feng, Bing Xue, Bingxuan Wang, Bochao Wu, Chengda Lu, Chenggang Zhao, Chengqi Deng, Chenyu Zhang, Chong Ruan, et al. 2024. Deepseek-v3 technical report. *arXiv preprint arXiv:2412.19437* (2024).
- [24] Chi Liu, Jiangxia Cao, Rui Huang, Kai Zheng, Qiang Luo, Kun Gai, and Guorui Zhou. 2024. Kuaiformer: Transformer-Based Retrieval at Kuaishou. *arXiv preprint arXiv:2411.10057* (2024).
- [25] Zichang Liu, Aditya Desai, Fangshuo Liao, Weitao Wang, Victor Xie, Zhaozhuo Xu, Anastasios Kyrillidis, and Anshumali Shrivastava. 2024. Scissorhands: Exploiting the persistence of importance hypothesis for llm kv cache compression at test time. *Advances in Neural Information Processing Systems* 36 (2024).
- [26] Xiao Lv, Jiangxia Cao, Shijie Guan, Xiaoyou Zhou, Zhiguang Qi, Yaqiang Zang, Ming Li, Ben Wang, Kun Gai, and Guorui Zhou. 2024. MARM: Unlocking the Future of Recommendation Systems through Memory Augmentation and Scalable Complexity. *arXiv preprint arXiv:2411.09425* (2024).
- [27] Julian McAuley, Christopher Targett, Qinfeng Shi, and Anton Van Den Hengel. 2015. Image-based recommendations on styles and substitutes. In *Proceedings of the 38th international ACM SIGIR conference on research and development in information retrieval*. 43–52.
- [28] Dheevatsa Mudigere, Yuchen Hao, Jianyu Huang, Zhihao Jia, Andrew Tulloch, Srinivas Sridharan, Xing Liu, Mustafa Ozdal, Jade Nie, Jongsoo Park, et al. 2022. 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*. 993–1011.
- [29] Qi Pi, Guorui Zhou, Yujing Zhang, Zhe Wang, Lejian Ren, Ying Fan, Xiaoqiang Zhu, and Kun Gai. 2020. 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*. 2685–2692.
- [30] Reiner Pope, Sholto Douglas, Aakanksha Chowdhery, Jacob Devlin, James Bradbury, Jonathan Heek, Kefan Xiao, Shivani Agrawal, and Jeff Dean. 2023. Efficiently scaling transformer inference. *Proceedings of Machine Learning and Systems* 5 (2023), 606–624.
- [31] Colin Raffel, Noam Shazeer, Adam Roberts, Katherine Lee, Sharan Narang, Michael Matena, Yanqi Zhou, Wei Li, and Peter J Liu. 2020. Exploring the limits of transfer learning with a unified text-to-text transformer. *Journal of machine learning research* 21, 140 (2020), 1–67.
- [32] Zihua Si, Lin Guan, ZhongXiang Sun, Xiaoxue Zang, Jing Lu, Yiqun Hui, Xingchao Cao, Zeyu Yang, Yichen Zheng, Dewei Leng, et al. 2024. 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*. 4890–4897.
- [33] 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.
- [34] Roberto Turrin, Massimo Quadrona, Andrea Condorelli, Roberto Pagano, Paolo Cremonesi, et al. 2015. 30Music Listening and Playlists Dataset. *RecSys Posters* 75 (2015).
- [35] A Vaswani. 2017. Attention is all you need. *Advances in Neural Information Processing Systems* (2017).
- [36] Yunli Wang, Zixuan Yang, Zhen Zhang, Zhiqiang Wang, Jian Yang, Shiyang Wen, Peng Jiang, and Kun Gai. 2024. Scaling Laws for Online Advertisement Retrieval. *arXiv preprint arXiv:2411.13322* (2024).
- [37] Zhibo Xiao, Luwei Yang, Wen Jiang, Yi Wei, Yi Hu, and Hao Wang. 2020. Deep multi-interest network for click-through rate prediction. In *Proceedings of the 29th ACM International Conference on Information & Knowledge Management*. 2265–2268.
- [38] 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).
- [39] 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).
- [40] Gaowei Zhang, Yupeng Hou, Hongyu Lu, Yu Chen, Wayne Xin Zhao, and Ji-Rong Wen. 2024. Scaling law of large sequential recommendation models. In *Proceedings of the 18th ACM Conference on Recommender Systems*. 444–453.
- [41] 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.
- [42] Pablo Zivic, Hernan Vazquez, and Jorge Sánchez. 2024. Scaling Sequential Recommendation Models with Transformers. In *Proceedings of the 47th International ACM SIGIR Conference on Research and Development in Information Retrieval*. 1567–1577.

A More Details about Method and Deployment

A.1 Rank Task Paradigm Statement

In recommendation systems, Click-Through Rate (CTR) prediction is a key task that estimates the probability of a user clicking on a recommended item. This task is crucial for enhancing user engagement and optimizing system performance. CTR prediction can be formulated as a supervised binary classification problem. Given a user u and a candidate item c , the goal is to predict the probability that the user will click on the item. This can be mathematically represented as:

$$\hat{y}_{u,c} = \sigma(f(X_u, X_c)) \quad (10)$$

where $y_{u,c}$ is a binary label indicating whether the user clicked on the item (1) or not (0), $\hat{y}_{u,c}$ is the prediction value of $y_{u,c}$. X_u and X_c are the feature vectors representing the user and item, respectively. $f(\cdot)$ is a model that estimates the probability of a click and $\sigma(\cdot)$ is the sigmoid function. The model is trained by minimizing the negative log-likelihood:

$$\mathcal{L} = -\frac{1}{N_s} \sum_{u,c} [y_{u,c} \log(\hat{y}_{u,c}) + (1 - y_{u,c}) \log(1 - \hat{y}_{u,c})] \quad (11)$$

where N_s is the total number of samples.

Traditional DLRM predicts users' next actions on candidate items based on their historical characteristics and elaborately engineered item features. This approach is typically represented as:

$$f(X_u, X_c) = f(u_{fea}, c_{fea}) \quad (12)$$

where u_{fea} and c_{fea} respectively denote the user features and item features generated through feature engineering. However, due to limitations in scalability and the complexity of feature engineering, we propose replacing handcrafted features with user behavior sequences. Specifically, we adopt a Transformer-based sequence model to scale up our recommendation system.

With sequence modeling, our task is to predict the user's behaviors on candidate items based on their historical behavior sequences. Specifically, we utilize the user's historical n items and their corresponding behaviors to predict the potential behavior on the candidate item c , which can be formulated as:

$$f(X_u, X_c) = f(\{x_1, x_2, \dots, x_n\}, x_c) \quad (13)$$

where x_i denotes a item ID from the entire music set X . Here, x_1 to x_n represent items on which the user has already taken action, while x_c is a candidate item obtained through recall and prerank stage.

ASTRO model employs multi-scale sequence partitioning to extract multi-scale subsequences with extraction strategy from user lifecycle behaviors. By integrating Equation 1 and 2, the task paradigm associated with this sequence structuring under the ASTRO model is characterized as follows:

$$f(X_u, X_c) = f(\{a_k, x_1^{a_k}, x_2^{a_k}, \dots, x_{n_k}^{a_k}\}_{k=1}^{N_b}, x_c) \quad (14)$$

where N_b represents the number of extraction strategy a_k , and n_k represents the length of subsequence extracted by a_k .

To ensure unified training and inference stage and align training data more closely with real online requests, TURBO organizes samples according to "SUMI" style. Therefore, we predict multiple

candidate items for one user simultaneously in the training stage, which can be expressed as:

$$f(X_u, \{X_{c,j}\}_{j=1}^{N_c}) = f(\{a_k, x_1^{a_k}, x_2^{a_k}, \dots, x_{n_k}^{a_k}\}_{k=1}^{N_b}, \{x_{c,j}\}_{j=1}^{N_c}) \quad (15)$$

where $x_{c,j}$ represents the j -th candidate item, and N_c represents the number of total candidate items.

The proposed recommendation model offers several notable advantages that significantly enhance its performance and applicability: (1) The model achieves a substantial reduction in computational complexity from $O(n^2d)$ to $O(n^2d/N_b)$, where N_b denotes the number of subsequence blocks. This improvement enables a marked acceleration in both training and inference processes, even when $N_b = 2$. (2) Our model overcomes temporal limitations for certain crucial behaviors. For highly sparse yet significant behaviors, we leverage users' entire lifecycle behaviors. Compared to the previous method that relies solely on data within a single time window, our method enables the incorporation of a broader range of items that more effectively represent user interests. (3) By extending important sparse behaviors to the entire lifecycle, the model gains a significant enrichment of information. Additionally, the block-wise approach alleviates the unfair allocation of attention scores between different subsequences, which will be further discussed below. Overall, these advantages contribute to the robustness and superiority of our recommendation model, thereby reinforcing its potential as a valuable approach in recommendation systems.

A.2 Deployment

As depicted in Figure 6, the implementation of the comprehensive recommendation system integrates online services, offline training, and model inference procedures. The system's performance is augmented by the TURBO acceleration framework and the ASTRO model. At the online service stage, the system initially conducts recall and prerank operations to identify the candidate item set. Subsequently, it utilizes the ASTRO model in conjunction with the TURBO framework to process these candidates, eventually generating a prediction list. In the offline training phase, user behavior logs are recorded and stored in a database. These logs are then dynamically compressed and transferred to an offline storage system. The log data are used to create a list of project labels, offering crucial data support for model training. During the training process, the system computes the loss and optimizes the model parameters via the backpropagation mechanism. The actual sorting process is divided into two stages. In the first stage, the model analyzes the user behavior sequence to generate a KV-cache vector, thus swiftly capturing user characteristics. In the second stage, the model capitalizes on these extracted features to score the candidate items and generate the final recommendation list. The collaborative design of the TURBO framework and the ASTRO model boosts the efficiency of the recommendation system and guarantees the accuracy of the recommendation results, thereby furnishing users with more personalized and satisfactory recommendation services.

B More Details about Experimental Setting

B.1 Dataset

Table 4 presents the number of users, items, and interactions for the four processed datasets. To protect data privacy, we have applied

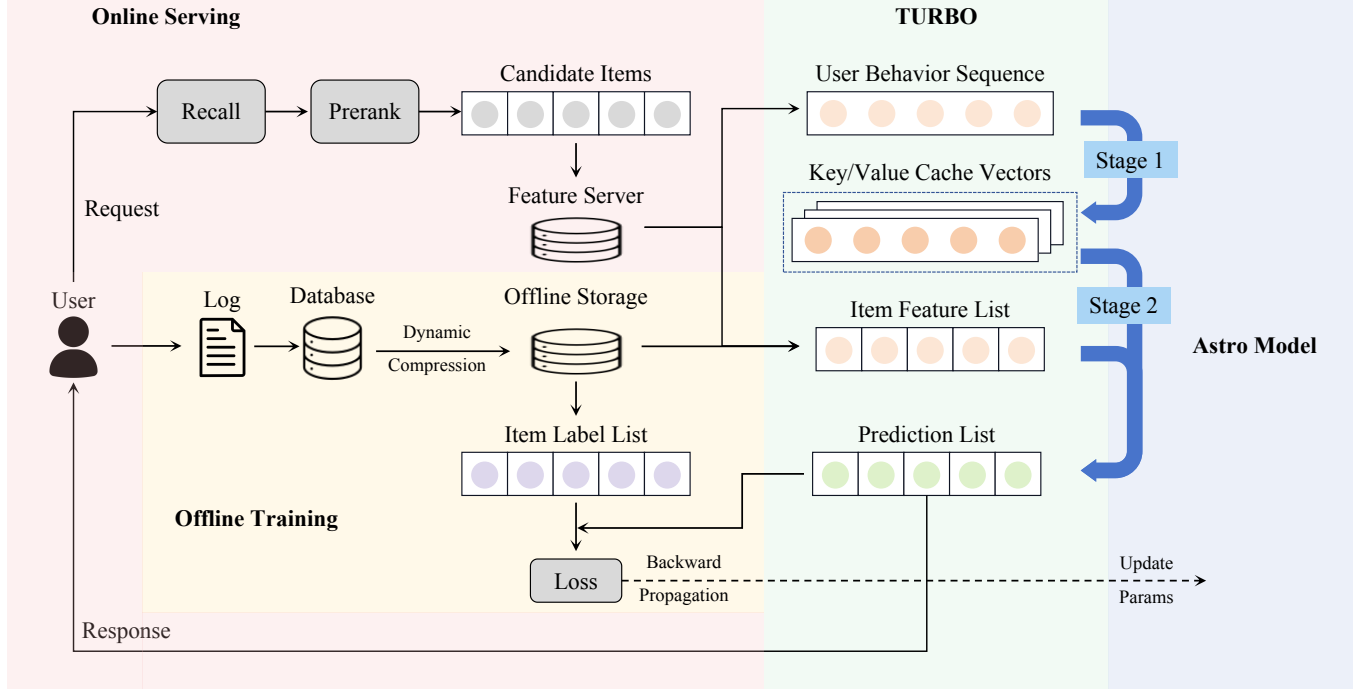


Figure 6: An Overview of our efficient large recommendation model deployment

Table 4: Dataset Statistics

Dataset	#User	#Item	#Interaction
Industrial	>40M	>6M	>1B
Spotify	0.16M	3.7M	1.2M
30Music	0.02M	4.5M	16M
Amazon-Book	0.54M	0.37M	1.09M

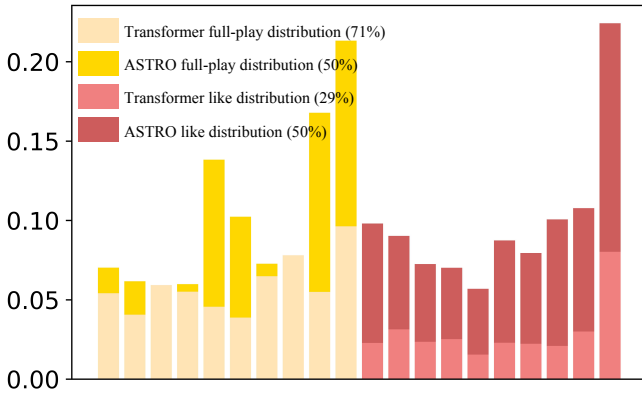


Figure 7: Attention Distribution of Different Models on Different Behaviors

special processing to the Industrial dataset. As a result, the data volume shown in the table is lower than the actual amount. However, it is clear that the scale of our industrial dataset still significantly

exceeds that of the other datasets. This substantial data volume provides a robust foundation for conducting scaling experiments.

B.2 Compared Methods

This study assesses the performance of various models, including DLRM, DIN, TWIN, Transformer, HSTU, and ASTRO variants. The subsequent section will introduce the detailed experimental setting in Industrial Dataset associated with each of these models.

- **DLRM:** The DLRM is selected as the baseline model after extensive online iterations, which has demonstrated remarkable effectiveness in both offline and online metrics. Its feature architecture incorporates statistical features and user behavior sequences. In terms of modeling methodology, we primarily employ a two-stage TWIN approach to extract user sequence representations. These representations are subsequently combined with statistical features and processed through MLP layers for feature interaction. The sequence length is set to 2000.
- **DIN:** DIN is a widely-used model that captures user interests through a target attention mechanism, focusing on the interaction between user historical behaviors and the target item. In terms of the experimental setup, we do not employ statistical features and only utilize behavior sequences with 1000 length.
- **TWIN:** A model that aligns the computation methods of GSU and ESU to enhance consistency and accuracy in modeling long-term user behavior. In the experimental setup, we also do not employ statistical features. The input sequence length for GSU is set to 2000, while for ESU, it is set to 1000.
- **Transformer:** A foundational model leveraging the Transformer architecture for sequence modeling. We employ a one-stage

Transformer encoder to directly process behavior sequences with a length of 2000.

- HSTU: The feature-level sequence undergoes temporal reorganization by chronological ordering of item-action pairs. This restructured sequence is then processed through the HSTU model at the structural level to predict target item-specific user action. The sequence length is set to 2000.
- ASTRO(-ATL,-BGF): In the previous method, the sequence length of 2000 is confined to a fixed time window and user behavior is not filtered. After introducing multi-scale sequence partitioning, the time period of the behavior sequence is extended to the entire user lifecycle (> 10 years), and extraction strategies are established based on the business logic. Ultimately, this process retains a sequence length of 200.
- ASTRO(-BGF): Based on multi-scale sequence partitioning, all traditional Transformer layers are replaced with Adaptive Transformer Layers. The attention distribution of each block is adaptively adjusted using the recommendation scenario and extraction strategy.
- ASTRO: This model integrates multi-scale sequence partitioning, adaptive transformer layer, and bit-wise gating fusion. Meanwhile, the layer number in the model is set to 2, which is consistent with the layer setting in the previous method.
- ASTRO-large: Compared to the ASTRO model, we have increased 4x sequence length and 6x layer number. As a result, the final sequence length is 800, and the layer number is 12.

C Attention Distribution Analysis

Sequence modeling in recommendation systems is essential for identifying items in a user's historical sequence that are similar to the target item, with this similarity quantified through attention scores. Before the adoption of attention mechanisms, sequence processing relied on average pooling, which assigned uniform attention scores to all items. The attention mechanism in DIN has been shown to enhance recommendation performance by assigning distinct attention scores to different items. The Transformer-based model emphasizes attention mechanisms. Thus this section explores the impact of multi-scale sequence partitioning on the attention distribution. The attention score will be introduced for two distinct behaviors: full-play and like.

Multi-scale sequence partitioning not only reduces computational complexity but also directly impacts attention distribution. Specifically, the proportion of "full-play" behaviors is higher than that of "like" behaviors. To address this imbalance, we distinguish between these two behaviors and extend the time window for "like" behaviors, thereby making the number of "like" and "full-play" behaviors approximately equal. We conduct comparative experiments between the Transformer and ASTRO. In the Transformer experiment, sequences contain equal proportions of "like" and "full-play" behaviors for computation. In contrast, the ASTRO model separates these behaviors and trains them in a block-wise manner, with attention distributions illustrated in Figure 7. In the Transformer structure, even with a similar number of "like" and "full-play" behaviors, the attention ratio for "full-play" behaviors reaches 71%. This is primarily because "full-play" samples have a higher proportion and

naturally occupy more gradients. Consequently, "full-play" behaviors receive higher attention scores due to the nature of the task, resulting in less attention to "like" behaviors and poorer performance on "like" tasks. In the ASTRO model, each behavior is assigned to a separate block with its own attention distribution. Through multi-scale sequence partitioning, we can not only increase the number of sparse yet important behaviors but also improve their attention distribution, thereby enhancing model effectiveness. Performance comparisons between the Transformer and ASTRO (-ATL, -BGF) across different datasets are shown in Table 1, which indicates that significant improvement can be achieved solely through multi-scale partitioning.