



# Recommender Transformers with Behavior Pathways

Zhiyu Yao\*

School of Software, BNRist, Tsinghua  
University  
Beijing, China  
yaozy19@mails.tsinghua.edu.cn

Xinyang Chen\*

Harbin Institute of Technology  
Shenzhen, China  
chenxinyang95@gmail.com

Sinan Wang\*

Alibaba Group  
China  
thusinan@foxmail.com

Qinyan Dai

School of Economics and Finance,  
Tsinghua University  
Shenzhen, China  
dqy22@mails.tsinghua.edu.cn

Yumeng Li

Alibaba Group  
China  
lym174806@alibaba-inc.com

Tanchao Zhu

Alibaba Group  
China  
tanchao.zhutc@alibaba-inc.com

Mingsheng Long<sup>†</sup>

School of Software, BNRist, Tsinghua  
University  
Beijing, China  
mingsheng@tsinghua.edu.cn

## ABSTRACT

Sequential recommendation requires the recommender to capture the evolving behavior characteristics from logged user behavior data for accurate recommendations. Nevertheless, user behavior sequences are viewed as a script with multiple ongoing threads intertwined. We find that only a small set of pivotal behaviors can be evolved into the user’s future action. As a result, the future behavior of the user is hard to predict. We conclude this characteristic for sequential behaviors of each user as the *behavior pathway*. Different users have their unique behavior pathways. Among existing sequential models, transformers have shown great capacity in capturing global-dependent characteristics. However, these models mainly provide a dense distribution over all previous behaviors using the self-attention mechanism, making the final predictions overwhelmed by the trivial behaviors not adjusted to each user. In this paper, we build the Recommender Transformer (RETR) with a novel Pathway Attention mechanism. RETR can dynamically plan the behavior pathway specified for each user, and sparingly activate the network through this behavior pathway to effectively capture evolving patterns useful for recommendation. The key design is a learned binary route to prevent the behavior pathway from being overwhelmed by trivial behaviors. Pathway attention is model-agnostic and can be applied to a series of transformer-based models for sequential recommendation. We empirically evaluate RETR on

seven intra-domain benchmarks and RETR yields state-of-the-art performance. On another five cross-domain benchmarks, RETR can capture more domain-invariant representations for sequential recommendation.

## CCS CONCEPTS

• Information systems → Information systems applications.

## KEYWORDS

Sequential Recommendation, Transformer Model

### ACM Reference Format:

Zhiyu Yao, Xinyang Chen\*, Sinan Wang\*, Qinyan Dai, Yumeng Li, Tanchao Zhu, and Mingsheng Long. 2024. Recommender Transformers with Behavior Pathways. In *Proceedings of the ACM Web Conference 2024 (WWW '24)*, May 13–17, 2024, Singapore, Singapore. ACM, New York, NY, USA, 12 pages. <https://doi.org/10.1145/3589334.3645528>

## 1 INTRODUCTION

Recommender systems [19, 30, 48] have been widely adopted in real-world industrial applications such as E-commerce and social media. Benefiting from the increase in computing power and model capacity, some recent efforts formulate recommendation as a time-series forecasting problem, known as *sequential recommendation* [6, 22]. The core idea of this field is to infer upcoming actions based on user’s historical behaviors, which are reorganized as time-ordered sequences. This intuitive modeling of recommendation is proved time-sensitive and context-aware to make precise predictions.

Recent advanced sequential recommendation models, such as SASRec [22], Bert4Rec [38] and SMRec [6], develop the transformer architecture to learn the sequential patterns. The transformer architecture brings these model powerful capacity to capture the characteristic of how users’ future behaviors are interacted with all previous behaviors. However, the user behavior may be casual or only associated with relevant subset behaviors. Interacting with all

\*Equal Contribution

<sup>†</sup>Corresponding Author

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

WWW '24, May 13–17, 2024, Singapore, Singapore

© 2024 Copyright held by the owner/author(s). Publication rights licensed to ACM.

ACM ISBN 979-8-4007-0171-9/24/05...\$15.00

<https://doi.org/10.1145/3589334.3645528>



**Figure 1: Typical examples of the behavior pathway for different users: correlated, casual, and drifted. The behavior pathway is outlined by the red boxes.**

previous behaviours will bring redundant information and make the most relevant behaviors overwhelmed by the trivial behaviors.

Recent advanced sequential recommendation models, such as SASRec [22], Bert4Rec [38] and S3-Rec [50], have achieved significant improvements. Transformers enable these models to recognize global-range sequential patterns, and to model how future behaviors are anchored in historical ones. The self-attention mechanism does make it possible to explore all previous behaviors of each user, with the whole neural network activated. However, misuse of all user information, regardless of whether it is informative or not, floods models with trivial ones, makes models dense in neuron connections and inefficient in computation, and results in key behaviors losing voice. And this clearly contradicts with the way our brain works.

The human being has many different parts of the brain specialized for various tasks, yet the brain only calls upon the relevant pieces for a given situation [47]. To some extent, user behavior sequences can be viewed as a script with multiple ongoing threads intertwined. And only key clues suggest what will happen next. In sequential recommendation, we find that only a small part of pivotal behaviors can be evolved into the user’s future action. And we conclude this characteristics of sequential behaviors as the *behavior pathway*.

Different users have their unique behavior pathways, and we have provided three typical examples: (a) *Correlated behavior pathway*: A user’s behavior pathway is closely associated with behaviors in a certain period. As shown in the first line of Figure 1, the mouse is clicked many times recently, leading to the final decision to buy a mouse. (b) *Casual behavior pathway*: A user’s behavior pathway is interested in a specific item at casual times. In the second line of Figure 1, the backpack is randomly clicked sequentially in a multi-hop manner. (c) *Drifted behavior pathway*: A user’s behavior pathway in a particular brand might drift over time. In the third line of Figure 1, the user was initially interested in a keyboard, but

suddenly became interested in buying a phone at last. It’s challenging to capture these potential behaviors dynamically for each user to make precise recommendations.

Motivated by the Pathways [10], a new way of thinking about AI, which builds a single model that is sparsely activated for all tasks with small pathways through the network called into action as needed, we propose a novel Recommender Transformer (RETR) with a Pathway Attention mechanism. RETR dynamically explores behavior pathways for different users and then captures evolving patterns through these pathways. The user-dependent pathway attention, which incorporates a pathway router, determines whether or not a behavior token will be maintained in the behavior pathway. The pathway router generates a customized binary route for each token based on their information redundancy. RETR has a stacked structure, and successive pathway routers constitute a hierarchical evolution of user behaviors. To enable the pathway router to be end-to-end optimized, we propose an adaptive Gumbel-Softmax sampling strategy to overcome the non-differentiable problem of sampling from a Bernoulli distribution.

To effectively capture the evolving patterns via the behavior pathway, our pathway attention mechanism makes RETR mainly attend to the obtained pathway. We force the model to focus on the most informative behaviors by using the query routed through the behavior pathway. We cut off the interaction from the off-pathway behaviors of the query. Compared with using all previous behaviors, our pathway attention mechanism is obviously more effective and can avoid the most informative tokens being overwhelmed by trivial behaviors. Besides, our pathway attention mechanism is model-agnostic and can be easily applied to the existing transformer-based models. To validate the effectiveness of our approach, we conduct experiments on seven intra-domain competitive datasets for sequential recommendations and RETR achieves state-of-the-art performance; Furthermore, our RETR also achieves consistent performance improvements under the cross-domain setting, indicating RETR can capture more domain-universal representation for sequential recommendation.

Our main contributions can be summarized as follows:

- We first propose the concept of behavior pathway for sequential recommendation, and find the key to the recommender is to dynamically capture the behavior pathway for each user.
- We propose the novel Recommender Transformer (RETR) with a novel pathway attention mechanism, which can generate the behavior pathway hierarchically and capture the evolving patterns dynamically through the pathway.
- We validate the effectiveness of RETR on 7 intra-domain benchmarks and 5 cross-domain benchmarks, both achieving state-of-the-art performance. RETR can capture more domain-invariant representations and our pathway attention can be applied together with a rich family of transformer-based models to yield consistent improvements.

## 2 RELATED WORK

**Traditional recommendation approaches.** Capturing evolving behavior characteristics is crucial for many online applications, such as advertising, social media and E-commerce, and it is the key challenge for sequential recommendation [1, 5, 9, 13, 22, 29, 33, 46, 51]. Traditional recommendation approach, such as the collaborative filtering (CF) [18] based on matrix approximation [24, 25], always assumes that the user’s behavior is static. However, in practice, user behaviors often change over time due to various reasons, making the CF deteriorate in a real-world application.

**Sequential recommendation approaches.** To overcome this challenge, some methods, such as FPMC [16] and HRM [41], use Markov chains to capture sequential patterns by learning user-specific transition matrices. Higher-order Markov Chains assume the next action is related to several previous actions. Benefit from this strong inductive bias, MC-based methods [15, 16] show superior performance in capturing short-term patterns. At the same time, there is a potential state space explosion problem when these approaches are faced with different possible sequences [42]. In recent years, many works have been using the deep neural network for sequential recommendation. The GRU4Rec [19] and the RepeatNet [34] adopt the recurrent network to capture dynamic patterns from the user behaviors dependent on sequence positions. The RNN-based models achieve competitive performance in capturing short-term behavior patterns but cannot capture long-term behavior patterns effectively. The CNN-based model, such as Caser [40], applies convolutional operations to extract transitions while tending to overlook the intrinsic relationship across user behaviors. The GNN-based methods, such as SRGNN [44], GCSAN [45], Jodie [26] and TGN [35] model behavior sequences as graph-structured data and incorporate an attention mechanism for a session-based recommendation. In addition, DIN [49] uses the gate mechanism to weight different user behaviors. However, concatenating all behaviors makes these models overlook the sequential characteristics. Recently, the MLP-based model like FMLP-Rec [51] uses the MLP as the backbone for sequential recommendation. However, these methods are still overwhelmed by the trivial behaviors.

**Transformer-based models for Sequential Recommendation.** SASRec [22], BertRec [38], S3-Rec [50], TGSRec [12], LightSANs [11] and SSE-PT [43] introduce the transformer architecture into sequential recommendation, which might lead to the over-parameterized architecture of Transformer-based methods. These models capture the evolving patterns by the self-attention mechanism, interacting with all previous behaviors. However, dense interactions will make the model not adapt to different users and overwhelm behavior pathways. Some methods like Locker [17] and Rec-denoiser [7] propose the sparse attention mechanism with learned mask, while they may overlook the ability of capturing the behavior pathway in the token level. To tackle this challenge, our paper builds the Recommender Transformer (RETR) with a new Pathway Attention mechanism that is dynamically activated for the behavior pathway of all users. Distinct from the previous routing architecture like Switch Transformer [14] using the MoE [37] structure for natural language tasks or TRAR [52] using the learned sparse attention for visual question answering, our RETR is designed explicitly for sequential recommendation. Our RETR uses the pathway router to adaptively route the sequential behavior of each user rather than routing the experts of feed-forward networks in switch transformer.

## 3 METHOD

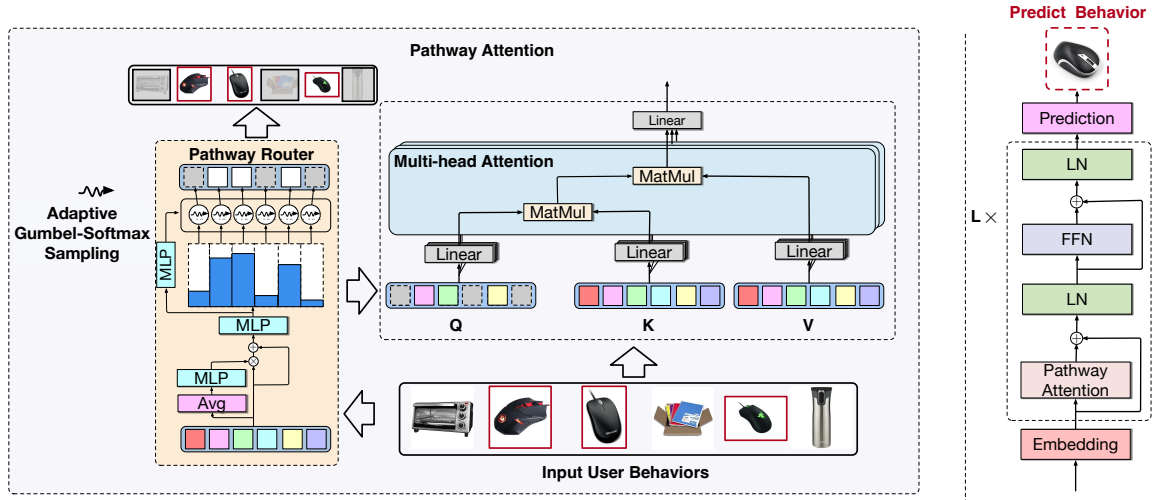
Suppose that we have a set of users and items, denoted by  $\mathcal{U}$  and  $\mathcal{I}$  respectively. In the task of sequential recommendation, chronologically-ordered behaviors of a user  $u \in \mathcal{U}$  could be represented by a user-interacted item sequence:  $\{i_1, \dots, i_n\}$ . Formally, given a user  $u$  with her or his behavior sequence  $\{i_1, \dots, i_n\}$ , the goal of sequential recommendation is to predict the next item the user  $u$  would interact with at the  $(n + 1)$ -th step, denoted as  $p(i_{n+1} | i_{1:n})$ .

As aforementioned, we highlight the key to sequential recommendation as the exploration of user-tailored behavior pathways, through which evolving characteristics could be learned. Motivated by this, we propose a novel *Recommender Transformer* (RETR) with a new *Pathway Attention*, the core subassembly of which is a pathway router. Besides the modification of architecture, we additionally introduce a hierarchical update strategy for the behavior pathway in the feed-forward procedure.

### 3.1 Recommender Transformer

Considering the limitation of overwhelming attention in Transformers [4] for sequential recommendation, we renovate the vanilla architecture to the Recommender Transformer (RETR) with a Pathway Attention mechanism, as shown in Figure 2.

**Model inputs.** To obtain the model inputs, we follow the sliding window practice and transform the user’s behavior sequence into a fixed-length- $N$  sequence  $s = (s_1, s_2, \dots, s_N)$ . Then we produce an item embedding matrix  $\mathbf{E}_I \in \mathbb{R}^{|\mathcal{I}| \times d}$ , where  $d$  is the embedding dimensionality. We perform a look-up operation from  $\mathbf{E}_I$  to retrieve the input embedding matrix  $\mathbf{E}_s \in \mathbb{R}^{N \times d}$  for sequence  $s$ . Besides, we also add a learnable position embedding  $\mathbf{P}_s \in \mathbb{R}^{N \times d}$  for sequence  $s$ . Finally, we can generate the input embedding of each behavior sequence  $s$  as  $\mathbf{X}_s = \mathbf{E}_s + \mathbf{P}_s \in \mathbb{R}^{N \times d}$ .



**Figure 2: The architecture of Recommender Transformer (RETR) on the right subfigure. Pathway Attention (left) explores the behavior pathway by the pathway router (orange module) and captures the evolving sequential characteristics of the use behaviors by the multi-head attention.**

**Overall architecture.** Recommender Transformer is characterized by stacking the Pathway Attention blocks and feed-forward layers alternately, containing  $L$  blocks. This stacking structure is conducive to learning behavior representations hierarchically. The overall equations of block  $l$  are formalized as

$$\begin{aligned} \widehat{\mathbf{Z}}^l, \mathbf{R}^l &= \text{Path-MSA}(\mathbf{Z}^{l-1}, \mathbf{R}^{l-1}) \\ \widehat{\mathbf{Z}}^l &= \text{LN}(\widehat{\mathbf{Z}}^l + \mathbf{Z}^{l-1}) \\ \mathbf{Z}^l &= \text{LN}(\text{FFN}(\widehat{\mathbf{Z}}^l) + \widehat{\mathbf{Z}}^l), \end{aligned} \quad (1)$$

where  $\mathbf{Z}^l \in \mathbb{R}^{N \times d}$ ,  $l \in \{1, \dots, L\}$  denotes the output of the  $l$ -th block. The initial input  $\mathbf{Z}^0 = \mathbf{X}_s \in \mathbb{R}^{N \times d}$  represents the raw behavior embedding.  $\mathbf{R}^{l-1} \in \mathbb{R}^{N \times 1}$  is the previous route from the  $(l-1)$ -th block and we initialize all elements in the route  $\mathbf{R}^0$  to 1. Path-MSA( $\cdot$ ) is to conduct the pathway multi-head self-attention. LN( $\cdot$ ) is to conduct layer normalization [3] and FFN represents the point-wise feed-forward network [4].

**3.1.1 Pathway Attention.** Note that the single-branch self-attention mechanism [4] in vanilla transformer cannot model the behavior pathway dynamically, resulting in key behaviors being overwhelmed by those non-pivotal or trivial ones. To solve this problem, we propose the Pathway Attention mechanism, as shown in Figure 2, which can dynamically attend to the behavior pathway of pivotal behavior tokens.

**Pathway router.** The pathway attention employs a sequence-adaptive pathway router to custom-tailor behavior pathway routes for users. The router generates a binary route  $\mathbf{R}^l \in \{0, 1\}^N$  to determine whether a behavior token would be part of the behavior pathway or not. Each router takes the pre-order route  $\mathbf{R}^{l-1}$  and user behavior tokens  $\mathbf{Z}^{l-1} \in \mathbb{R}^{N \times d}$  of the  $(l-1)$ -th block as its inputs. All elements in the route are initialized by 1 and are updated progressively in training.

Foremost, to suppress the potential disturbance to the model caused by the local drifted interest (Figure 1), it is crucial to incorporate the *global* information in the route generation. We apply the average pooling to all the preserved behavior tokens routed by  $\mathbf{R}^{l-1}$ , and produce the global sequential representation via a multilayer perceptron (MLP) module. Then, we combine this global representation with the inputs and employ a residual connection to maintain the original input information. Finally, we feed them to another MLP layer to predict the probabilities of keeping or dropping the behavior tokens. All MLP layers are column-wise and operate on the embedding dimensionality. The above procedure can be formulated as follows:

$$\begin{aligned} \mathbf{Z}_{\text{emb}}^l &= \mathbf{Z}^{l-1} + \mathbf{Z}^{l-1} \odot \text{MLP}\left(\frac{\sum_{i=1}^N \mathbf{R}_i^{l-1} \mathbf{Z}_i^{l-1}}{\sum_{i=1}^N \mathbf{R}_i^{l-1}}\right) \\ \boldsymbol{\pi} &= \text{Softmax}(\text{MLP}(\mathbf{Z}_{\text{emb}}^l)) \in \mathbb{R}^{N \times 2}, \end{aligned} \quad (2)$$

where  $\odot$  is the Hadamard product. For  $t \in \{1, 2, \dots, N\}$ , we let  $\boldsymbol{\pi}_t = [1 - \alpha_t, \alpha_t]$ , where the logit  $\alpha_t$  denotes the probability that the  $t$ -th behavior token is kept alive for the behavior pathway.

**Adaptive Gumbel-Softmax sampling from  $\boldsymbol{\pi}$  for router.** Our goal is to generate the binary route from  $\boldsymbol{\pi}$ . However, sampling from  $\boldsymbol{\pi}$  directly is non-differentiable, and it will impede the gradient-based training. Gumbel-Softmax [21] is an effective way to approximate the original non-differentiable sample from a discrete distribution with a differentiable sample from a Gumbel-Softmax distribution. Thus, we adapt the Gumbel-Softmax technique to achieve such a sampling procedure. Instead of directly sampling a keep-or-drop decision  $\widehat{\mathbf{R}}_t^l$  for the  $t$ -th behavior token from the distribution  $\boldsymbol{\pi}_t$ , we generate it as:

$$\widehat{\mathbf{R}}_t^l = \arg \max_{j \in \{0,1\}} (\log \pi_t(j) + G_t(j)), \quad (3)$$

where  $G_t = -\log(-\log U_t)$  is a standard Gumbel distribution, and  $U_t$  is sampled i.i.d. from a uniform distribution  $\text{Uniform}(0, 1)$ . To



remove the non-differentiable argmax operation in (3), the standard Gumbel-Softmax uses the reparameterization trick [21] as a differentiable approximation to relax the one-hot  $\widehat{\mathbf{R}}_t^l \in \{0, 1\}$  to  $v_t \in \mathbb{R}^2$ :

$$v_t(j) = \frac{\exp((\log \pi_t(j) + G_t(j))/\tau)}{\sum_{i \in \{0,1\}} \exp((\log \pi_t(i) + G_t(i))/\tau)}, j \in \{0, 1\}, \quad (4)$$

where  $\tau$  is the temperature parameter of the Softmax. However, it remains a well-known challenge to tune the temperature in Gumbel-Softmax since a low temperature will cause a high variance in gradient magnitude and a high temperature will lead to an over-smoothing probability. Furthermore, a fixed temperature is not adaptive across different datasets or behaviors of each user, which incurs huge tweaking cost. Motivated by these difficulties, our propose an adaptive variant of Gumbel-Softmax that introduces the token-specific weight mechanism into the Gumbel-Softmax:

$$\begin{aligned} \omega &= \text{ReLU}\left(\text{MLP}\left(\mathbf{Z}_{\text{emb}}^l\right)\right) \in \mathbb{R}^{N \times 1}, \\ v_t(j) &= \frac{\exp(\omega_t \log \pi_t(j) + G_t(j))}{\sum_{i \in \{0,1\}} \exp(\omega_t \log \pi_t(i) + G_t(i))}, j \in \{0, 1\}, \end{aligned} \quad (5)$$

where  $\omega_t$  is the weight specific for each token  $t$  of each sequence. For different user behaviors, we use an MLP module to dynamically introduce the weight from the inputs  $\mathbf{Z}_{\text{emb}}^l$ , avoiding the high variance in the gradient and mitigating the over-smoothing phenomenon. Our adaptive Gumbel-Softmax can make RETR dynamically adapt to diverse datasets and user behaviors without tuning the temperature.

**Hierarchical update strategy for router.** The preliminary route  $\widehat{\mathbf{R}}^l$ , sampled from  $\boldsymbol{\pi}$ , is not a final decision. In our design, once a token fails to be routed in a certain block, it would permanently lose the privilege to be part of the behavior pathway in the following feed-forward procedure. This constitutes a more efficient hierarchical pathway router strategy. Thus finally we formulate the route  $\mathbf{R}^l \in \mathbb{R}^{N \times 1}$  as the Hadamard product of  $\widehat{\mathbf{R}}^l$  and the pre-order route  $\mathbf{R}^{l-1}$  in the  $(l-1)$ -th block:

$$\mathbf{R}^l = \widehat{\mathbf{R}}^l \odot \mathbf{R}^{l-1}. \quad (6)$$

**Multi-head pathway attention.** The standard multi-head self-attention mechanism retrieves sequential characteristics by exploiting all behavior tokens, making the behavior pathway overwhelmed by the trivial behaviors. In the proposed pathway attention, the pathway router would be firstly applied to the input behavior tokens to route information. The pathway router would not pare down the number of tokens, but only the interactions between the off-pathway and on-pathway tokens, as these off-pathway tokens may also convey contextual information.

Specifically, for the query  $\mathbf{Q}$ , key  $\mathbf{K}$ , and value  $\mathbf{V}$  in the pathway attention: the query is routed by the pathway router through token-wise multiplication between  $\mathbf{R}_t^l$  and  $\mathbf{Z}_t^{l-1}$ , to prevent the pathway from being overwhelmed and to force the pathway attention to attend to the behavior pathway. The key and value are the original input behavior tokens, to ensure that the contextual information

from off-pathway behavior tokens can be captured as well:

$$\begin{aligned} \mathbf{Q}_m, \mathbf{K}_m, \mathbf{V}_m &= (\mathbf{R}^l \mathbf{Z}^{l-1}) \mathbf{W}_{\mathbf{Q}_m}^l, \mathbf{Z}^{l-1} \mathbf{W}_{\mathbf{K}_m}^l, \mathbf{Z}^{l-1} \mathbf{W}_{\mathbf{V}_m}^l \\ \widehat{\mathbf{Z}}_m^l &= \text{Softmax}\left(\frac{\mathbf{Q}_m \mathbf{K}_m^T}{\sqrt{d/h}}\right) \mathbf{V}_m^l, \end{aligned} \quad (7)$$

where  $m \in \{1, 2, \dots, h\}$  is the head index in the multi-head self-attention;  $\mathbf{W}_{\mathbf{Q}_m}^l, \mathbf{W}_{\mathbf{K}_m}^l, \mathbf{W}_{\mathbf{V}_m}^l \in \mathbb{R}^{d \times \frac{d}{h}}$  are transformation matrices learned from data. Finally, the outputs  $\{\widehat{\mathbf{Z}}_m^l \in \mathbb{R}^{N \times \frac{d}{h}}\}_{1 \leq m \leq h}$  of multiple heads are concatenated into  $\widehat{\mathbf{Z}}^l \in \mathbb{R}^{N \times d}$ . We use  $\widehat{\mathbf{Z}}^l, \mathbf{R}^l = \text{Path-MSA}(\mathbf{Z}^{l-1}, \mathbf{R}^{l-1})$  to summarize the above pathway attention. Its output is further transformed by (1) to form the final output of the  $l$ -th block  $\mathbf{Z}^l \in \mathbb{R}^{N \times d}$ .

In the prediction of the  $(t+1)$ -th behavior, only the first  $t$  observable behaviors should be taken into account. To avoid a future information leak and ensure causality, we apply a causal pathway attention in that a look-ahead mask is employed and all links between  $\mathbf{Q}_j$  and  $\mathbf{K}_i$  ( $j > i$ ) are removed.

**Model-agnostic pathway mechanism.** It is worth noting that our pathway attention is a lightweight module readily pluggable into any transformer-based model by replacing the self-attention mechanism with our *pathway attention* while remaining the architecture unchanged. To verify the effectiveness of our model-agnostic pathway mechanism, we apply our pathway attention to mainstream sequential recommendation transformers: BERTRec [38], SASRec [22], SMRec [6], S3-Rec [50], TGSRec [12], and LightSANS [11], which further enhances the performance and generalization of the these models.

## 3.2 Prediction Layer and Training Objective

**Prediction layer.** In the final layer of RETR, we calculate the user's preference score for the item  $k$  in the step  $(t+1)$  in the context of user behavior history as  $p(i_{t+1} = k | i_{1:t}) = e_k \cdot \mathbf{Z}_t^L$ , where  $e_k$  is the representation of item  $k$  from item embedding matrix  $\mathbf{E}_I$ , and  $\mathbf{Z}_t^L$  is the output of the  $L$ -th block in RETR at step  $t$ , with  $L$  being the number of RETR blocks.

**Training objective.** We adopt the pairwise ranking loss to optimize the RETR model parameters as:

$$\mathcal{L} = - \sum_{u \in \mathcal{U}} \sum_{t=1}^n \log \sigma(p(i_{t+1} | i_{1:t}) - p(i_{t+1}^- | i_{1:t})), \quad (8)$$

where we pair each ground-truth item  $i_{t+1}$  with a randomly sampled negative item  $i_{t+1}^-$ . In each epoch, we randomly generate one negative item for each time step in each sequence. This pairwise ranking loss is widely adopted in previous literature of sequential recommendation [22, 51].

## 4 EXPERIMENTS

We extensively evaluate the proposed Recommender Transformer (RETR) on seven intra-domain real-world benchmarks and five cross-domain benchmarks. Due to page limitation, we also include further ablation study results and visualization examples in Appendix A.1 and Appendix A.3 respectively.

**Table 1: Performance of state-of-the-art models under the intra-domain setting. Rec-denoiser (2022) uses BertRec as the backbone.**

Datasets	Metric	BERT4Rec	SASRec	SMRec	S3-Rec	SINE	TGSRec	LightSAN	Locker	Rec-denoiser	Jodie	TGN	RETR
Netflix	HR@10	0.4792	0.4622	0.4848	0.4917	0.4902	0.4887	0.4852	0.4897	0.4913	0.4813	0.4802	<b>0.5184</b>
	NDCG@10	0.3330	0.3202	0.3492	0.3571	0.3601	0.3512	0.3441	0.3557	0.3582	0.3368	0.3318	<b>0.3795</b>
	MRR	0.2652	0.2519	0.2725	0.2819	0.2796	0.2778	0.2785	0.2873	0.2886	0.2687	0.2612	<b>0.3175</b>
MSD	HR@10	0.4819	0.4766	0.5083	0.5315	0.5264	0.5137	0.4994	0.5495	0.5581	0.4825	0.4782	<b>0.5963</b>
	NDCG@10	0.4891	0.4831	0.5112	0.5381	0.5304	0.5279	0.5163	0.5495	0.5471	0.4872	0.4832	<b>0.6012</b>
	MRR	0.3120	0.3079	0.3302	0.3494	0.3667	0.3612	0.3451	0.3686	0.3723	0.3224	0.3102	<b>0.4025</b>
Taobao	HR@10	0.1261	0.1182	0.1272	0.1336	0.1580	0.1537	0.1590	0.1584	0.1603	0.1447	0.1421	<b>0.1803</b>
	NDCG@10	0.0425	0.0391	0.0531	0.0627	0.0873	0.0745	0.0794	0.0922	0.0952	0.0582	0.0571	<b>0.1218</b>
	MRR	0.0489	0.0436	0.0721	0.0788	0.0934	0.0802	0.0841	0.0928	0.0967	0.0628	0.0603	<b>0.1149</b>
Yelp	HR@10	0.7597	0.7373	0.7548	0.7597	0.7564	0.7533	0.7552	0.7503	0.7520	0.7492	0.7473	<b>0.7775</b>
	NDCG@10	0.4778	0.4642	0.4789	0.4937	0.4902	0.4887	0.4863	0.4935	0.4973	0.4792	0.4784	<b>0.5169</b>
	MRR	0.4026	0.3927	0.4023	0.4107	0.4093	0.4072	0.4086	0.4189	0.4214	0.3997	0.3985	<b>0.4378</b>
MovieLens	HR@10	0.8269	0.8233	0.8302	0.8352	0.8311	0.8303	0.8294	0.8349	0.8368	0.8277	0.8259	<b>0.8513</b>
	NDCG@10	0.5965	0.5936	0.6079	0.6172	0.6134	0.6081	0.6119	0.6003	0.6128	0.6009	0.5998	<b>0.6397</b>
	MRR	0.5614	0.5573	0.5703	0.5812	0.5801	0.5734	0.5791	0.5787	0.5812	0.5651	0.5627	<b>0.5978</b>
Tmall	HR@10	0.6196	0.6275	0.6476	0.6687	0.6512	0.6506	0.6399	0.6703	0.6729	0.6384	0.6362	<b>0.7214</b>
	NDCG@10	0.5025	0.5049	0.5192	0.5423	0.5411	0.5372	0.5415	0.5792	0.5830	0.5307	0.5198	<b>0.6197</b>
	MRR	0.4026	0.4804	0.4934	0.5194	0.5147	0.5121	0.5119	0.5373	0.5426	0.5003	0.4997	<b>0.5903</b>
Steam	HR@10	0.8656	0.8729	0.8792	0.8813	0.8765	0.8773	0.8832	0.8831	0.8892	0.8780	0.8731	<b>0.9079</b>
	NDCG@10	0.6283	0.6306	0.6408	0.6573	0.6502	0.6491	0.6519	0.6497	0.6523	0.6451	0.6399	<b>0.6835</b>
	MRR	0.5883	0.5925	0.6011	0.6135	0.5972	0.6003	0.6104	0.6114	0.6159	0.5873	0.5798	<b>0.6383</b>

**Intra-domain setting.** We evaluate RETR on seven intra-domain datasets: Netflix, MSD, Taobao, Yelp, Tmall, Steam, and MovieLens1M. All methods are trained from scratch on these datasets. The statistics of the seven datasets are summarized in Table 12 of Appendix A.2 and the description for these datasets can be found therein. All datasets are widely used for sequential recommendation task. It is notable that Netflix, MSD, Taobao and Steam are large-scale datasets.

**Cross-domain setting.** We evaluate the ability of RETR to capture the *domain-invariant* representation for sequential recommendation under the cross-domain setting. Following the training strategy in UniSRec [20], we pre-train our RETR on multiple datasets “Grocery and Gourmet Food”, “Home and Kitchen”, “CDs and Vinyl”, “Kindle Store” and “Movies and TV”, and then fine-tune the pre-trained RETR respectively on different target datasets “Prime Pantry”, “Industrial and Scientific”, “Musical Instruments”, “Arts, Crafts and Sewing” and “Office Products”. These are all sub-categories in the Amazon datasets [31]. The detailed descriptions of source and target datasets are described in Appendix A.2.

**Evaluation metrics.** We apply top-k Hit Ratio (HR@k), top-k Normalized Discounted Cumulative Gain (NDCG@k) and Mean Reciprocal Rank (MRR) for evaluation, reporting HR@10, NDCG@10 and MRR of the results. Besides, following the standard strategy in SASRec [22], we pair the ground-truth item with 100 randomly sampled negative items that the user has not interacted with. All

metrics are calculated according to the ranking of the items and we report the average score.

**Baseline methods.** We compare our RETR with several state-of-the-art sequence recommendation models. Specifically, we compare our RETR with state-of-the-art transformer-based sequential recommendation models: Rec-denoiser [7], Locker [17], SASRec [22], BertRec [38], SMRec [6], S3-Rec [50], TGSRec [12] and LightSANs [11]. These methods adopt the attention mechanism to make precise recommendations. Note that Rec-denoiser [7] and Locker [17] are novel transformer-based models with learnable sparse attention. Besides, we also compare our RETR with state-of-the-art graph-based sequential recommendation methods: Jodie [26] and TGN [35]. We further compare our approach with some cross-domain recommendation models RecGURU [27] and UniSRec [20]. All baseline methods are configured using default parameters of the original paper or optimal parameters which can produce their best results through a grid search.

**Implementation details.** Our model is supervised by the pairwise rank loss in (8), using the ADAM [23] optimizer with an initial learning rate of 0.001. Batch size is set to 512. The maximum number of training epochs for all methods is set to 300. All hyperparameters are tuned on the validation set. The training process is early stopped within 10 epochs. Our RETR has  $L = 2$  layers, and each layer has  $h = 4$  heads (the ablation study of multi-head attention can be found in Appendix ??) and dimension  $d$  is set to be 256. The maximum sequence length  $N$  is set to 200 for MovieLens1M and

**Table 2: Performance comparison of different competitive methods under the cross-domain setting. X-\* indicates the model pre-trained on multiple datasets and finetuned on a target dataset (“X” stands for “cross-domain”) following the training procedure of UniSRec (2022).**

Target Datasets	Metric	X-SMRec	X-LightSAN	X-Locker	X-Rec-denoiser	UniSRec	RecGURU	X-RETR
Scientific	HR@10	0.1304	0.1315	0.1312	0.1329	0.1235	0.1023	<b>0.1459</b>
	NDCG@10	0.0706	0.0725	0.0744	0.0752	0.0634	0.0572	<b>0.0865</b>
Pantry	HR@10	0.0713	0.0725	0.0741	0.0734	0.0693	0.0469	<b>0.0847</b>
	NDCG@10	0.0327	0.0321	0.0346	0.0338	0.0311	0.0209	<b>0.0425</b>
Instruments	HR@10	0.1293	0.1278	0.1316	0.1301	0.1267	0.1113	<b>0.1353</b>
	NDCG@10	0.0792	0.0778	0.0830	0.0813	0.0748	0.0681	<b>0.0893</b>
Arts	HR@10	0.1281	0.1275	0.1283	0.1293	0.1239	0.1084	<b>0.1378</b>
	NDCG@10	0.0749	0.0738	0.0749	0.0763	0.0712	0.0651	<b>0.0853</b>
Office	HR@10	0.1319	0.1321	0.1336	0.1343	0.1280	0.1145	<b>0.1426</b>
	NDCG@10	0.0842	0.0858	0.0857	0.0869	0.0831	0.0768	<b>0.0998</b>

100 for the other intra-domain and cross-domain datasets. In the cross-domain setting, we further pre-train the proposed approach and baseline methods from multiple datasets following the training strategy in UniSRec [20] for 300 epochs using the default parameters from UniSRec [20]. All models are pretrained without using the item ID embeddings. In the phase of finetuning on the target domain, we use the item ID embeddings and fix the backbone for all competing pretrained models. All experiments are repeated three times, implemented in PyTorch [32], and conducted on a single NVIDIA 3090 GPU.

#### 4.1 Intra-domain Results

The results of different methods on seven intra-domain datasets are shown in Table 1. We can easily find that transformer-based models, SASRec [22], BertRec [38], SMRec [6], S3-Rec [50], TGSRec [12] and LightSANs [11], achieve competitive performance on most datasets, indicating that the transformer-based models have a better capacity to capture sequential behaviors of complex characteristics. These models can capture the interaction information between all previous user behaviors via the attention mechanism. Besides, the graph-based models like Jodie [26] and TGN [35] also achieve competitive performance. Rec-denoiser [7] and Locker [17] introduce novel sparse attention mechanisms, thereby achieving better performance compared with other baselines. This validates the effectiveness of the learned mask attention. While these baselines are strong competitors, our RETR can achieve state-of-the-art performance by a large margin on most datasets compared with the Rec-denoiser and Locker.

**Intra-domain setting.** We evaluate RETR on seven intra-domain datasets: Netflix, MSD, Taobao, Yelp, Tmall, Steam, and MovieLens1M. All methods are trained from scratch on these datasets. The statistics of the seven datasets are summarized in Table 12 of Appendix A.2 and the description for these datasets can be found therein. All datasets are widely used for sequential recommendation task. It is notable that Netflix, MSD, Taobao and Steam are large-scale datasets.

**Results on Yelp, MovieLens1M and Tmall.** Our RETR achieves competitive performance on Yelp and Tmall. These datasets are sparse, containing less action information. Thus they have lots of noisy logged information. By effectively capturing the behavior pathway, RETR is not affected by this trivial behavior information and captures the most informative behavior representation to achieve better performance. Note that under the Tmall benchmark, RETR gains 7% HR@10, 12% NDCG@10 and 14% MRR against the strongest baseline SMRec [6]. Besides, for the MoveLens1M benchmark, RETR also achieves the best performance among all competing baselines.

**Results on large-scale datasets.** Our RETR can consistently achieve state-of-the-art results on large-scale datasets (Netflix, MSD, Taobao, and Steam). These datasets are challenging and difficult to capture pivotal behavior pathway useful for precise recommendation from the rich but noisy user’s behaviors. Especially for the Taobao dataset, RETR gains relative improvements of 12% HR@10, 37% NDCG@10 and 20% MRR against the strongest baseline SINE [39]. It provides evidence that RETR can achieve competitive performance in both small- and large-scale datasets. The substantial performance gains of our RETR indicate that focusing more on the behavior pathway enables RETR to capture sequential characteristics more efficiently and effectively than the vanilla self-attention mechanism, which considers all previous user behaviors and is easily overwhelmed.

#### 4.2 Cross-domain Results

To verify the ability of RETR to capture domain-invariant representations, we evaluate pre-trained RETR on 5 target datasets under the cross-domain setting. The multi-domain pre-training version of RETR, denoted as X-RETR, can be effectively transferred to new domains. We also provide the results for multi-domain pre-training version of Rec-denoiser, SASRec, SMRec, and LightSAN, denoted as X-Rec-denoiser, UniSRec, X-SMRec, and X-LightSAN respectively. Technically, we follow the pretraining strategy of UniSRec [20] to train all models. As shown in Table 2, the X-RETR already achieves competitive cross-domain performance, outperforming

**Table 3: Ablation study of model-agnostic pathway attention on MovieLens. Results in each column are obtained without/with pathway attention. (↑: positive improvement using pathway attention.)**

Metric	BERT4Rec	SMRec	S3-Rec	TGSRec	LightSAN	RETR
HR@10	0.8269 / 0.8506 ↑	0.8302 / 0.8585 ↑	0.8352 / <b>0.8594</b> ↑	0.8303 / 0.8497 ↑	0.8294 / 0.8529 ↑	0.8513
NDCG@10	0.5965 / 0.6389 ↑	0.6079 / 0.6486 ↑	0.6172 / <b>0.6487</b> ↑	0.6081 / 0.6325 ↑	0.6119 / 0.6475 ↑	0.6397
MRR	0.5614 / 0.5998 ↑	0.5703 / 0.6031 ↑	0.5812 / <b>0.6052</b> ↑	0.5734 / 0.5973 ↑	0.5791 / 0.5993 ↑	0.5978

the state-of-the-art cross-domain method UniSRec by a large margin on most target datasets. Compared with other multi-domain pre-trained backbones, X-RETR achieves the highest performance and empowers better transferability among different backbones. Specially, X-RETR gains **12%** HR@10 and **22.5%** NDCG@10 compared with X-SMRec, on the Scientific benchmark. These results indicate that RETR can extract domain-invariant representations for sequential recommendation, indicating that a stronger backbone is crucial in parallel with transfer-learning method for enhancing the transferability. RETR can be regarded as a general backbone to capture more domain-invariant representations.

### 4.3 Ablation Study

Here we provide an ablation study of RETR. For more ablation results, please refer to Appendix A.1, including the ablation study of different hyperparameters and effectiveness of each model component, and the explicit ablation study of adaptive Gumbel-Softmax. Further visual examples are provided in Appendix A.3

**Pathway attention towards different transformer-based models.** As described before, our RETR yields state-of-the-art performance on all datasets. We further apply our pathway mechanism towards different transformer-based models like BERTRec [38], SMRec [6], S3-Rec [50], TGSRec [12], and LightSANs [11]. In Table 3, we observe that our pathway attention can improve the performance of all baseline transformer-based models substantially. RETR can be further enhanced using advanced backbones alternative to the vanilla transformers and achieve the best results among all competing methods. These results provide strong evidences that our proposed pathway attention is model-agnostic to transformer-based methods and not limited to a particular architectural choice.

**Table 4: Ablation study of the effectiveness of each model component. Experiments are conducted on the Yelp Dataset.**

Model	MRR
<b>RETR</b>	<b>0.4378</b>
RETR w/o Pathway Router	0.3887
RETR w/o hierarchical update	0.4234
SASRec	0.3927

#### Effectiveness of each model component and number of blocks.

In Table 4, we analyze the efficacy of each component in RETR on the Yelp dataset and have the following observations. First, we remove the pathway router module and randomly choose whether it

can be maintained or dropped for each input behavior token. Removing the pathway router decreases the prediction performance a lot (MRR: 0.4354  $\rightarrow$  0.3887), showing the necessity of learning behavior pathway effectively based on a data-dependent module. Second, discarding the hierarchical update strategy for the behavior pathway also decreases the prediction performance, suggesting that this strategy is crucial for RETR to get a more accurate behavior pathway. In Table 5, we adjust the number of blocks for RETR on Yelp. We find that the performance first increases rapidly with the growth of the block number and achieves the best performance at  $L = 2$ . We perform a similar grid search on other datasets.

**Table 5: Ablation study of the number of blocks for each RETR block. Experiments are conducted on the Yelp Dataset.**

Model (# number of blocks)	MRR
RETR ( $L = 1$ )	0.4197
<b>RETR (<math>L = 2</math>)</b>	<b>0.4378</b>
RETR ( $L = 3$ )	0.4342
RETR ( $L = 4$ )	0.4340

## 5 CONCLUSION

A sequential recommender is designed to make accurate recommendations based on users' historical behaviors. However, the users' behaviors are dynamic and come in a continually evolving manner. A user's current decision may only call upon the interest from the certain relevant behaviors of the past. We conclude these sequential characteristics as the behavior pathway. We propose the Recommender Transformer (RETR) with a novel pathway attention mechanism to tackle these challenges. The pathway attention mechanism develops a pathway router to dynamically allocate the behavior pathway for each user and capture the evolving patterns. RETR can capture more domain-invariant representations and the pathway attention is model-agnostic and can be easily applied to a series of transformer-based methods. RETR achieves state-of-the-art performance on seven intra-domain datasets and five cross-domain benchmarks for sequential recommendation.

## ACKNOWLEDGMENTS

This work was supported by the National Key Research and Development Plan (2021YFC3000905), National Natural Science Foundation of China (62022050, U2342217 and 62306085), and Alibaba Group through Alibaba Innovative Research Program.

## REFERENCES

- [1] Gediminas Adomavicius and Alexander Tuzhilin. 2005. Toward the next generation of recommender systems: A survey of the state-of-the-art and possible extensions. *TKDE* (2005).
- [2] Nabiha Asghar. 2016. Yelp dataset challenge: Review rating prediction. *arXiv preprint arXiv:1605.05362* (2016).
- [3] Jimmy Lei Ba, Jamie Ryan Kiros, and Geoffrey E Hinton. 2016. Layer normalization. *arXiv preprint arXiv:1607.06450* (2016).
- [4] Prateep Bhattacharjee and Sukhendu Das. 2017. Temporal Coherency based Criteria for Predicting Video Frames using Deep Multi-stage Generative Adversarial Networks. In *NeurIPS*. 4268–4277.
- [5] Shuqing Bian, Wayne Xin Zhao, Kun Zhou, Jing Cai, Yancheng He, Cunxiang Yin, and Ji-Rong Wen. 2021. Contrastive Curriculum Learning for Sequential User Behavior Modeling via Data Augmentation. In *CIKM*.
- [6] Chao Chen, Haoyu Geng, Nianzu Yang, Junchi Yan, Daiyue Xue, Jianping Yu, and Xiaokang Yang. 2021. Learning Self-Modulating Attention in Continuous Time Space with Applications to Sequential Recommendation. In *ICML*.
- [7] Huiyuan Chen, Yusan Lin, Menghai Pan, Lan Wang, Chin-Chia Michael Yeh, Xiaoting Li, Yan Zheng, Fei Wang, and Hao Yang. 2022. Denoising Self-Attentive Sequential Recommendation. In *Proceedings of the 16th ACM Conference on Recommender Systems*.
- [8] Rewon Child, Scott Gray, Alec Radford, and Ilya Sutskever. 2019. Generating long sequences with sparse transformers. *arXiv preprint* (2019).
- [9] Qiang Cui, Shu Wu, Qiang Liu, Wen Zhong, and Liang Wang. 2018. MV-RNN: A multi-view recurrent neural network for sequential recommendation. *TKDE* (2018).
- [10] Jeff Dean. 2021. Introducing pathways: A nextgeneration ai architecture. *Google Blog* (2021).
- [11] Xinyan Fan, Zheng Liu, Jianxun Lian, Wayne Xin Zhao, Xing Xie, and Ji-Rong Wen. 2021. Lighter and better: low-rank decomposed self-attention networks for next-item recommendation. In *SIGIR*.
- [12] Ziwei Fan, Zhiwei Liu, Jiawei Zhang, Yun Xiong, Lei Zheng, and Philip S Yu. 2021. Continuous-time sequential recommendation with temporal graph collaborative transformer. In *CIKM*.
- [13] Hui Fang, Danning Zhang, Yiheng Shu, and Guibing Guo. 2020. Deep learning for sequential recommendation: Algorithms, influential factors, and evaluations. *TOIS* (2020).
- [14] William Fedus, Barret Zoph, and Noam Shazeer. 2021. Switch transformers: Scaling to trillion parameter models with simple and efficient sparsity. *arXiv preprint arXiv:2101.03961* (2021).
- [15] Ruining He, Chen Fang, Zhaowen Wang, and Julian McAuley. 2016. Vista: A visually, socially, and temporally-aware model for artistic recommendation. In *RecSys*.
- [16] Ruining He and Julian McAuley. 2016. Fusing similarity models with markov chains for sparse sequential recommendation. In *ICDM*.
- [17] Zhankui He, Handong Zhao, Zhe Lin, Zhaowen Wang, Ajinkya Kale, and Julian McAuley. 2021. Locker: Locally Constrained Self-Attentive Sequential Recommendation. In *CIKM*.
- [18] Jonathan L Herlocker, Joseph A Konstan, Al Borchers, and John Riedl. 1999. An algorithmic framework for performing collaborative filtering. In *SIGIR*.
- [19] Balázs Hidasi, Alexandros Karatzoglou, Linas Baltrunas, and Domonkos Tikk. 2015. Session-based recommendations with recurrent neural networks. *arXiv preprint arXiv:1511.06939* (2015).
- [20] 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 *KDD*.
- [21] Eric Jang, Shixiang Gu, and Ben Poole. 2016. Categorical reparameterization with gumbel-softmax. *arXiv preprint arXiv:1611.01144* (2016).
- [22] Wang-Cheng Kang and Julian McAuley. 2018. Self-attentive sequential recommendation. In *ICDM*.
- [23] Diederik Kingma and Jimmy Ba. 2015. Adam: A method for stochastic optimization. In *ICLR*.
- [24] Yehuda Koren. 2008. Factorization meets the neighborhood: a multifaceted collaborative filtering model. In *KDD*.
- [25] Yehuda Koren, Robert Bell, and Chris Volinsky. 2009. Matrix factorization techniques for recommender systems. *Computer* (2009).
- [26] Srijan Kumar, Xikun Zhang, and Jure Leskovec. 2019. Predicting dynamic embedding trajectory in temporal interaction networks. In *KDD*.
- [27] Chenglin Li, Mingjun Zhao, Huanming Zhang, Chenyun Yu, Lei Cheng, Guoqiang Shu, Beibei Kong, and Di Niu. 2022. RecGURU: Adversarial Learning of Generalized User Representations for Cross-Domain Recommendation. In *WSDM*.
- [28] Shiyang Li, Xiaoyong Jin, Yao Xuan, Xiyu Zhou, Wenhui Chen, Yu-Xiang Wang, and Xifeng Yan. 2019. Enhancing the locality and breaking the memory bottleneck of transformer on time series forecasting. *NeurIPS* (2019).
- [29] Yuli Liu, Christian Walder, and Lexing Xie. 2022. Determinantal Point Process Likelihoods for Sequential Recommendation. *arXiv preprint arXiv:2204.11562* (2022).
- [30] Jie Lu, Dianshuang Wu, Mingsong Mao, Wei Wang, and Guangquan Zhang. 2015. Recommender system application developments: a survey. *Decision Support Systems* (2015).
- [31] Jianmo Ni, Jiacheng Li, and Julian McAuley. 2019. Justifying recommendations using distantly-labeled reviews and fine-grained aspects. In *EMNLP-IJCNLP*.
- [32] Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan, Trevor Killeen, Zeming Lin, Natalia Gimelshein, Luca Antiga, et al. 2019. Pytorch: An imperative style, high-performance deep learning library. *NeurIPS* (2019).
- [33] 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 *CIKM*.
- [34] Pengjie Ren, Zhumin Chen, Jing Li, Zhaochun Ren, Jun Ma, and Maarten De Rijke. 2019. Repeatnet: A repeat aware neural recommendation machine for session-based recommendation. In *AAAI*.
- [35] Emanuele Rossi, Ben Chamberlain, Fabrizio Frasca, Davide Eynard, Federico Monti, and Michael Bronstein. 2020. Temporal graph networks for deep learning on dynamic graphs. *arXiv preprint arXiv:2006.10637* (2020).
- [36] Ramprasaath R Selvaraju, Michael Cogswell, Abhishek Das, Ramakrishna Vedantam, Devi Parikh, and Dhruv Batra. 2017. Grad-cam: Visual explanations from deep networks via gradient-based localization. In *ICCV*.
- [37] Noam Shazeer, Azalia Mirhoseini, Krzysztof Maziarz, Andy Davis, Quoc Le, Geoffrey Hinton, and Jeff Dean. 2017. Outrageously large neural networks: The sparsely-gated mixture-of-experts layer. *arXiv preprint arXiv:1701.06538* (2017).
- [38] 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 *CIKM*.
- [39] Qiaoyu Tan, Jianwei Zhang, Jiangchao Yao, Ninghao Liu, Jingren Zhou, Hongxia Yang, and Xia Hu. 2021. Sparse-interest network for sequential recommendation. In *WSDM*.
- [40] Jiayi Tang and Ke Wang. 2018. Personalized top-n sequential recommendation via convolutional sequence embedding. In *WSDM*.
- [41] Pengfei Wang, Jiafeng Guo, Yanyan Lan, Jun Xu, Shengxian Wan, and Xueqi Cheng. 2015. Learning hierarchical representation model for nextbasket recommendation. In *SIGIR*.
- [42] Chao-Yuan Wu, Amr Ahmed, Alex Beutel, Alexander J Smola, and How Jing. 2017. Recurrent recommender networks. In *WSDM*.
- [43] Liwei Wu, Shuqing Li, Cho-Jui Hsieh, and James Sharpnack. 2020. SSE-PT: Sequential recommendation via personalized transformer. In *Fourteenth ACM Conference on Recommender Systems*. 328–337.
- [44] Shu Wu, Yuyuan Tang, Yanqiao Zhu, Liang Wang, Xing Xie, and Tieniu Tan. 2019. Session-based recommendation with graph neural networks. In *AAAI*.
- [45] Chengfeng Xu, Pengpeng Zhao, Yanchi Liu, Victor S Sheng, Jiajie Xu, Fuzhen Zhuang, Junhua Fang, and Xiaofang Zhou. 2019. Graph Contextualized Self-Attention Network for Session-based Recommendation. In *IJCAI*.
- [46] An Yan, Shuo Cheng, Wang-Cheng Kang, Mengting Wan, and Julian McAuley. 2019. CosRec: 2D convolutional neural networks for sequential recommendation. In *CIKM*.
- [47] Zeenat F Zaidi. 2010. Gender differences in human brain: a review. *The open anatomy journal* (2010).
- [48] Wayne Xin Zhao, Shanlei Mu, Yupeng Hou, Zihan Lin, Yushuo Chen, Xingyu Pan, Kaiyuan Li, Yujie Lu, Hui Wang, Changxin Tian, et al. 2021. Rebole: Towards a unified, comprehensive and efficient framework for recommendation algorithms. In *CIKM*.
- [49] 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 *KDD*.
- [50] Kun Zhou, Hui Wang, Wayne Xin Zhao, Yutao Zhu, Sirui Wang, Fuzheng Zhang, Zhongyuan Wang, and Ji-Rong Wen. 2020. S3-rec: Self-supervised learning for sequential recommendation with mutual information maximization. In *CIKM*.
- [51] Kun Zhou, Hui Yu, Wayne Xin Zhao, and Ji-Rong Wen. 2022. Filter-enhanced MLP is All You Need for Sequential Recommendation. In *WWW*.
- [52] Yiyi Zhou, Tianhe Ren, Chaoyang Zhu, Xiaoshuai Sun, Jianzhuang Liu, Xinghao Ding, Mingliang Xu, and Rongrong Ji. 2021. Trar: Routing the attention spans in transformer for visual question answering. In *ICCV*.



## A APPENDIX

### A.1 Further Ablation Study

**Replace the proposed pathway-based method with other sparse attention methods.** We replace the proposed pathway-based method with two sparse attention methods: LogSparse [28] and sparse attention [8]. As shown in Table 6, our RETR using the pathway attention remarkably outperforms other competing methods with two sparse attention methods on Tmall. These experiment results show that sparse attention methods cannot capture the exact behavior pathway and show worse performance than RETR.

**Table 6: Ablation study of sparse attention methods on the Tmall dataset.**

Model	NDCG@10	HR@10	MRR
RETR w/ LogSparse	0.4923	0.6015	0.4735
RETR w/ Sparse-Att	0.4871	0.5873	0.4620
RETR	<b>0.6197</b>	<b>0.7214</b>	<b>0.5903</b>

**Number of heads and maximum sequence length.** In Table 7, we adjust the number of heads for RETR on Yelp. We find that the performance first increases with the growth of the head number and achieves the best performance at  $h = 4$ . We perform a similar grid search on other datasets. In Table 8, we adjust the maximum sequence length  $N$  for RETR on Yelp. As shown in Table 8, we find that the performance of our RETR first increases rapidly with the growth of the block number and achieves the best performance at  $N = 100$ . We perform a similar grid search on other datasets.

**Table 7: Ablation of the head number for RETR on the Yelp Dataset.**

Model (# $h$ )	MRR
RETR ( $h = 1$ )	0.4317
RETR ( $h = 2$ )	0.4345
<b>RETR (<math>h = 4</math>)</b>	<b>0.4378</b>
RETR ( $h = 8$ )	0.4369

**Table 8: Ablation of the maximum sequence length for RETR on the Yelp Dataset.**

Model (# maximum sequence length )	MRR
RETR ( $N = 25$ )	0.4237
RETR ( $N = 50$ )	0.4294
<b>RETR (<math>N = 100</math>)</b>	<b>0.4378</b>
RETR ( $N = 200$ )	0.4362

**Quantitative results on whether the proposed model effectively captures a useful pathway.** We give quantitative results to validate that our RETR can effectively capture various behavior pathways. Specifically, we evaluate our RETR using a subset of sequences derived from the obtained behavior pathway on Tmall.

Technically, we first train RETR on Tmall. For each user, we take the captured behavior pathway from our RETR as the inputs to retrain a RETR rather than using the whole user’s behaviors. As shown in Table 9, we find that using the behavior pathway as the inputs can achieve comparable results as the original RETR which uses complete user behaviors. It provides the evidence that our RETR can aptly capture the useful pathway for each user.

**Table 9: Quantitative results on the Tmall dataset.**

Model	NDCG@10	HR@10	MRR
RETR w/ pathway	0.6112	0.7142	0.5831
RETR	0.6197	0.7214	0.5903
SASRec	0.5049	0.6275	0.4804
SASRec w/ pathway	0.5778	0.6812	0.5425
SASRec w/ off-pathway	0.4383	0.5697	0.4215

**Table 10: Comparison Parameters and GFLOPs on the Yelp Dataset.**

Model	Param (M)	GFLOPs	MRR
RETR	5.0	9.6	<b>0.4378</b>
SASRec [22]	5.0	9.6	0.3927
SINE [39]	5.1	9.7	0.4011
SMRec [6]	5.2	9.9	0.4023

**Why use the cross-attention mechanism?** We choose the cross-attention mechanism for three main reasons: (1) The cross-attention mechanism can force the pathway attention to attend to the behavior pathway; (2) It can ensure that the contextual information from off-pathway behavior tokens can be captured, using the original input behavior tokens as the key and value; (3) Our pathway cross-attention mechanism avoids the trivial interaction between the off-pathway tokens, while the previous self-attention mechanism for sequential models can be overwhelmed by the trivial information in the off-pathway behavior tokens. To verify our explanation, we further conduct evaluation experiments on Tmall. Specifically, we train RETR on Tmall, and then use the trained RETR on Tmall to capture the behavior pathway for each user in Tmall. We use the pathway behaviors and off-pathway behaviors as the inputs to train SASRec respectively. As shown in Table 9, we can see that SASRec achieves better performance using the behavior pathway as the inputs compared with the original SASRec using the whole user’s behavior as the inputs. On the contrary, the off-pathway inputs hurt SASRec’s performance seriously. Finally, our RETR achieves the best performance, indicating that the pathway-offpathway cross-attention is more effective than the pathway self-attention.

**Evaluation on efficiency.** The efficiency is compared between SASRec [22], SINE [39] and SMRec [6] on the Yelp dataset. The computation cost is measured with gigabit floating-point operations (GFLOPs) on the self-attention module with position encoding. Meanwhile, the model scale measured with parameters is also presented. As shown in Table 10, our RETR has almost the same

**Table 11: Ablation study of RETR with quantitative results on the Yelp and Tmall datasets. “-” indicates failure case of model training.**

Model	Yelp			Tmall		
	NDCG@10	HR@10	MRR	NDCG@10	HR@10	MRR
RETR w/ standard Gumbel-Softmax ( $\tau = 2$ )	0.5113	0.7695	0.4328	0.6049	0.7084	0.5787
RETR w/ standard Gumbel-Softmax ( $\tau = 1$ )	0.5124	0.7719	0.4345	0.6084	0.7105	0.5803
RETR w/ standard Gumbel-Softmax ( $\tau = 0.8$ )	0.5136	0.7730	0.4354	0.6103	0.7138	0.5822
RETR w/ standard Gumbel-Softmax ( $\tau = \mathbf{0.6}$ )	0.5152	0.7749	0.4360	0.6095	0.7129	0.5817
RETR w/ standard Gumbel-Softmax ( $\tau = 0.4$ )	-	-	-	-	-	-
RETR w/ standard Gumbel-Softmax ( $\tau = 0.2$ )	-	-	-	-	-	-
RETR w/ Gumbel-Softmax (temperature annealing)	0.5134	0.7727	0.4350	0.6093	0.7133	0.5814
RETR w/ adaptive Gumbel-Softmax	<b>0.5169</b>	<b>0.7775</b>	<b>0.4378</b>	<b>0.6197</b>	<b>0.7214</b>	<b>0.5903</b>

number of parameters or GFLOPs, compared with SASRec, indicating that our pathway router is a light-weight module. Our pathway attention does not bring more costs. It’s worth noticing that the parameter scales and GLOPs of other competing transformers (apart from SASRec) are larger than RETR, but our RETR achieves higher performance. This result shows that our RETR is more efficient and effective than other competing attention-based models.

**Ablation study for adaptive Gumbel-Softmax.** The temperature parameter  $\tau$  is a crucial hyperparameter for the standard Gumbel-Softmax. A fixed temperature cannot be adaptive across different datasets or users. It is widely-known to be uneasy to tune the temperature parameter, in that a lower value may lead to high variances in gradients and a higher value may lead to over-smoothing probabilities. To mitigate these technical issues, we propose a novel adaptive Gumbel-Softmax mechanism to eliminate the need of temperature tuning, which can produce token-specific weights automatically adjusted to varying behaviors of each user.

As shown in Table 11, we find that RETR with the standard Gumbel-Softmax achieves the highest performance in different datasets at different temperatures ( $\tau = 0.8$  for Tmall and  $\tau = 0.6$  for Yelp). These results show that a fixed temperature is not adaptive across diverse datasets. However, lower temperatures ( $\tau = 0.2, 0.4$ ) cause the failure of model training for RETR due to the high variance in gradients. Higher temperatures ( $\tau = 1, 2$ ) may perform worse because of the over-smoothing probabilities in Gumbel-Softmax. It proves difficult to tune the temperature for the standard Gumbel-Softmax. Previous work like TRAR [52] develops a schedule that starts with a high temperature and gradually anneals it to a small but non-zero value. This schedule can make the training more stable, but it cannot achieve the best performance adaptively across different datasets.

In contrast, RETR with the proposed adaptive Gumbel-Softmax featured by the token-specific weight mechanism achieves the best performance compared with the standard Gumbel-Softmax under different temperatures. These results indicate that the adaptive Gumbel-Softmax is much more effective in capturing the behavior pathway and can be dynamically adapted to different datasets without tuning the temperature. This adaptive mechanism can simultaneously overcome the over-smoothing phenomenon and dynamically avoid the high variances in gradients.

## A.2 Descriptions of the Datasets

**Intra-domain datasets.** Here are descriptions of the seven datasets: (1) **Netflix**: Netflix dataset is a large-scale movie rating dataset released by Netflix. (2) **MSD**: The Million Song Dataset (MSD) is a large-scale, metadata-rich and open-source dataset on Kaggle. (3) **Taobao**: Taobao dataset [39] contains user behaviors in Taobao’s recommender system. In experiments, we only use the click behaviors. (4) **Yelp** [2]: Yelp is a dataset for business recommendation. We only use the transaction records after January 1st, 2019. (5) **Tmall**: Tmall contains users’ shopping logs on Tmall online shopping platform, which is from the IJCAI-15 competition. (6) **Steam** [22]: Steam dataset is collected from a large online video game distribution platform. This dataset includes 2,567,538 users, 15,474 games and 7,793,069 English reviews from October 2010 to January 2018. (7) **MovieLens1M**: this is a widely used benchmark dataset for evaluating collaborative filtering algorithms. The version we use is MovieLens-1M, which includes 1 million user ratings.

**Table 12: Statistics of the intra-domain datasets.**

Dataset	Users	Items	Actions
Netflix	463,435	17,769	57,000,000
MSD	571,355	41,140	34,000,000
Taobao	987,994	4,162,024	100,150,807
Yelp	30,431	20,033	316,354
MovieLens1M	6,040	3,416	1,000,000
Tmall	66,909	37,367	427,797
Steam	334,730	13,047	3,700,000

**Cross-domain datasets.** We choose five categories from Amazon review datasets [20]: Grocery and Gourmet Food, Home and Kitchen, CDs and Vinyl, Kindle Store, and Movies and TV, as the source domain datasets for pre-training. For the target datasets, we choose another five categories from Amazon review datasets [20]: Prime Pantry, Industrial and Scientific, Musical Instruments, Arts, Crafts and Sewing, and Office Products, as target domain datasets to evaluate the proposed approach under the cross-domain setting. The detail statistics are shown in Table 13.

Following previous works [22], we group the interaction records by users or sessions for all datasets and sort them by the timestamps

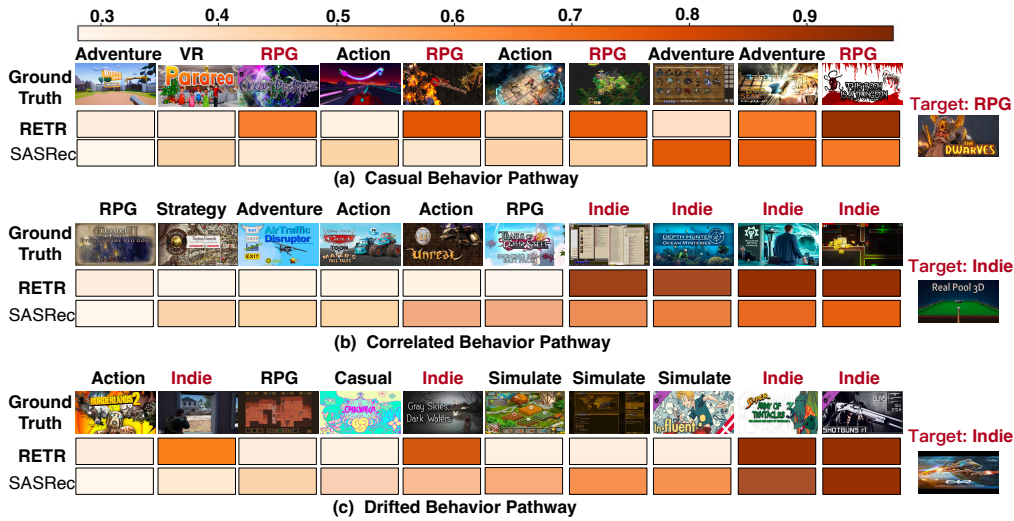


Figure 3: Visualizations of behavior heatmaps for RETR and SASRec of three random users in Steam dataset. They are corresponding to casual, correlated and drifted behavior pathways respectively.

Table 13: Statistics of the cross-domain datasets.

Source Dataset	Users	Items	Actions	Target Dataset	Users	Items	Actions
Food	115,349	39,670	1,027,4137	Scientific	8,442	4,385	59,427
CDs	94,010	64,439	1,118,563	Pantry	13,101	4,898	126,9626
Kindle	138,436	98,111	2,204,596	Instruments	24,962	9,964	208,926
Movies	281,700	59.203	3,226,731	Arts	45,486	21,019	395,150
Home	731,913	185,552	6,451,926	Office	87,436	25,986	684,837

in ascending order. We split the historical sequence for each user into three parts: (1) the most recent behavior for testing, (2) the second most recent behavior for validation, and (3) all remaining behaviors for training. During testing, the input sequences contain training behaviors and validation behaviors. We filter less popular items and inactive users with fewer than five interaction records.

### A.3 Visual Examples

**Setups.** We also provide qualitative visualizations for our RETR, and SASRec [22]. Technically, we use the GradCAM [36] to generate behavior heatmaps of the output of the last layer in each model. Three random examples of users' historical behaviors in the Steam dataset are shown in sequential order through subplots (a)–(c) in Figure 3. We provide attention heatmaps of each example at the last ten time steps. We can observe three main behavior pathway characteristics corresponding to three behavior sequences respectively: (a) *Casual behavior pathway*: RPG games are randomly clicked by the user, while the user has a continuing interest in RPGs. (b) *Correlated behavior pathway*: The user has recently been interested in indie games. (c) *Drifted behavior pathway*: The user has recently been interested in simulation games but chooses an indie game at last.

**Visualization results.** We elaborate the three representative categories of behavior pathway in recommender systems with model-learned attention heatmaps. (1) *Casual behavior pathway*: As shown in Figure 3(a), the RGB game is randomly clicked at casual times. Our RETR can capture all the RPG casual behavior pathways, while the SASRec focuses on the incorrect recent adventure games. The SASRec cannot capture the early clicked RPG game. This phenomenon proves that our RETR can deal with the casual behavior pathway effectively. (2) *Correlated behavior pathway*: For the correlated behavior pathway, we also provide an example which is shown in Figure 3(b). The indie game is clicked many times recently, leading to the final decision to an indie game. Our RETR can effectively capture the correlated behavior pathway. However, the SASRec provides higher attention scores on the recent RPG games. On the contrary, our RETR pays no attention to these wrong results, showing that it has a greater ability to cope with the correlated behavior pathway. (3) *Drifted behavior pathway*: As shown in Figure 3(c). The user was initially interested in the indie game, but suddenly became interested in simulation games recently and chose an indie game at last. Our RETR captures the drifted behavior pathway for the indie game and has not concentrated on the old drifted pathway – simulation games, while the SASRec is affected by the trivial behaviors of simulation games. These visualization results strongly show that our RETR can capture various behavior pathways dynamically for each user.