

self-attention layer and lack inductive bias¹. ii) The second limitation pertains to the low-pass filtering nature of self-attention. By focusing on the entire range of data, self-attention may unintentionally smoothen out important and detailed patterns in embedding, resulting in the *oversmoothing* problem. The oversmoothing issue poses a significant challenge in the SR domain, as it may hinder the ability of model to capture crucial temporal dynamics and provide accurate predictions. In Table 1, most Transformer-based models are limited to low-pass filters. These models do not consider high-pass filters. Note that FMLPRec (Zhou et al. 2022) attempts to learn a filter, but it tends to gravitate towards the low-pass filter (cf. Fig. 2 (b)). As shown in Fig. 1, the low-pass filter only captures the ongoing preferences of the user, i.e., an Apple fanatic, and it may be difficult to capture preferences based on new interests or trends (e.g., snorkel mask to buy for vacation). When recommending items for the next time ($|\mathcal{S}^{u_1}| + 1$), it is undemanding to recommend long-term interests, but recommending short-term interests is a challenging task.

In this paper, we address these two limitations and present **Beyond Self-Attention for Sequential Recommendation** (BSARec), a novel model that uses inductive bias via Fourier transform with self-attention. By using the Fourier transform, BSARec gains access to the inductive bias of frequency information, enabling the capture of essential patterns and periodicity that may be overlooked by self-attention alone. This enhances inductive bias and has the potential to improve recommendation performance.

To tackle the oversmoothing issue, we introduce our own designed frequency rescaler to apply high-pass filters into BSARec’s architecture. Our frequency rescaler can capture high-frequency behavioral patterns, such as interests driven by short-term trends, as well as low-frequency patterns, such as long-term interests, in a user’s behavioral patterns (cf. Fig. 1). Additionally, our method provides a perspective to improve the performance of SR models and solve the problem of oversmoothing.

To evaluate the efficacy of BSARec, we conduct extensive experiments on 6 benchmark datasets. Our experimental results demonstrate that BSARec consistently outperforms 7 baseline methods regarding recommendation performance. Additionally, we conduct a series of experiments that underscore the necessity of our approach and verify its effectiveness in mitigating the oversmoothing problem, leading to improved recommendation accuracy and enhanced generalization capabilities. The contributions of this work are as follows:

- We unveil the low-pass filtering nature of the self-attention of Transformer-based SR models, resulting in the problem of oversmoothing.
- We propose a novel model, **Beyond Self-Attention for Sequential Recommendation** (BSARec), that leverages the Fourier transform to balance between our inductive

¹We mean by inductive bias a pre-determined attention structure that is not trained but injected by us when designing our model. Therefore, we call it as *attentive inductive bias*.

bias and self-attention. Further, we design the rescaler for high-pass filters to mitigate the oversmoothing issue.

- Extensive evaluation on 6 benchmark datasets demonstrates BSARec’s outperformance over 7 baseline methods, validating its effectiveness in improving recommendation performance.

Preliminaries

Problem Formulation

The goal of SR is to predict the user’s next interaction with an item given their historical interaction sequences. Given a set of users \mathcal{U} and items \mathcal{V} , we can sort the interacted items of each user $u \in \mathcal{U}$ chronologically in a sequence as $\mathcal{S}^u = [v_1^u, v_2^u, \dots, v_{|\mathcal{S}^u|}^u]$, where v_i^u denotes the i -th interacted item in the sequence. The aim is to recommend a Top- k list of items as potential next items in a sequence. Formally, we predict $p(v_{|\mathcal{S}^u|+1}^u = v | \mathcal{S}^u)$.

Self-Attention for Sequential Recommendation

The basic idea behind the self-attention mechanism is that elements within sequences are correlated but hold varying levels of significance concerning their positions in the sequence. Self-attention uses dot-products between items in the sequence to infer their correlations defined as:

$$\mathbf{A} = \text{softmax} \left(\frac{\mathbf{Q}\mathbf{K}^T}{\sqrt{d}} \right), \quad (1)$$

where $\mathbf{Q} = \mathbf{E}_{\mathcal{S}^u} \mathbf{W}_Q$, $\mathbf{K} = \mathbf{E}_{\mathcal{S}^u} \mathbf{W}_K$, and d is the scale factor. The scaled dot-product component learns the latent correlation between items. Other components in Transformer are utilized in SASRec, including the point-wise feed-forward network, residual connection, and layer normalization. Our method uses this self-attention matrix and adds an inductive bias to find the trade-off between the two methods.

Discrete vs. Graph Fourier Transform

This subsection introduces the concept of the frequency domain and the Fourier transform, providing a cohesive foundation for the proposed method.

The Discrete Fourier Transform (DFT) is a linchpin in digital signal processing (DSP), projecting a sequence of values into the frequency domain (or the Fourier domain). We typically use $\mathcal{F} : \mathbb{R}^N \rightarrow \mathbb{C}^N$ to denote DFT with the Inverse DFT (IDFT) $\mathcal{F}^{-1} : \mathbb{C}^N \rightarrow \mathbb{R}^N$. Applying \mathcal{F} to a signal is equal to multiplying it from the left by a DFT matrix. The rows of this matrix consist of the Fourier basis $\mathbf{f}_j = [e^{2\pi i(j-1) \cdot 0} \dots e^{2\pi i(j-1)(N-1)}]^T / \sqrt{N} \in \mathbb{R}^N$, where i is the imaginary unit and j denotes the j -th row. For the spectrum of \mathbf{x} , let it be represented as $\bar{\mathbf{x}} = \mathcal{F}\mathbf{x}$. We can define $\bar{\mathbf{x}}_{\text{lfc}} \in \mathbb{C}^c$ containing the c lowest elements of $\bar{\mathbf{x}}$, and $\bar{\mathbf{x}}_{\text{hfc}} \in \mathbb{C}^{N-c}$ as the vector containing the remaining elements. The low-frequency components (LFC) of the sequence signal \mathbf{x} are defined as:

$$\text{LFC}[\mathbf{x}] = [\mathbf{f}_1, \mathbf{f}_2, \dots, \mathbf{f}_c] \bar{\mathbf{x}}_{\text{lfc}} \in \mathbb{R}^N. \quad (2)$$

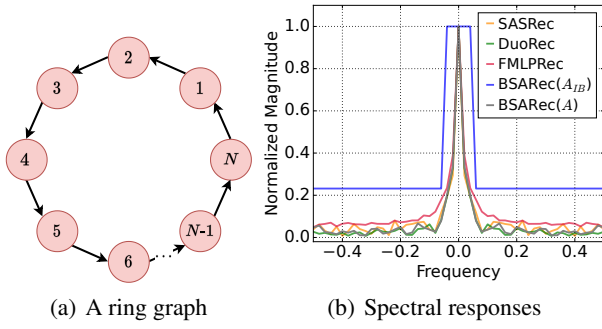


Figure 2: (a) A ring graph with N nodes, and (b) visualization of the filter of the self-attentions in LastFM.

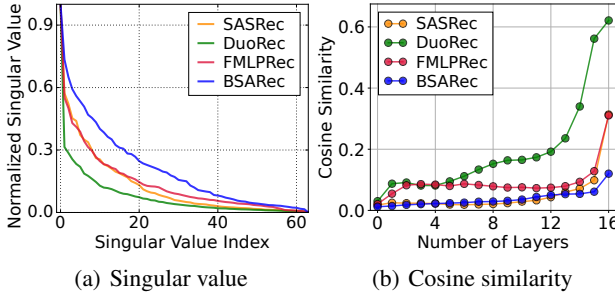


Figure 3: Visualization of oversmoothing in LastFM. The singular values and cosine similarity of user sequence output embedding.

Conversely, the high-frequency components (HFC) are:

$$HFC[\mathbf{x}] = [\mathbf{f}_{c+1}, \mathbf{f}_{c+2}, \dots, \mathbf{f}_N] \bar{\mathbf{x}}_{hfc} \in \mathbb{R}^N. \quad (3)$$

Note that we use real-valued DFT and multiplying with the Fourier bases in Eqs. 2 and 3 means IDFT. For more descriptions, interested readers should refer to Appendix (Shin et al. 2023).

The Graph Fourier Transform (GFT) can be considered as a generalization of DFT toward graphs. In other words, DFT is a special case of GFT, where a ring graph of N nodes is used (see Fig. 2 (a)) (Sandryhaila and Moura 2014). In fact, DFT is a method to project a sequence of values onto the eigenspace of the Laplacian matrix of the ring graph (which is the same as the Fourier domain).

The frequency concept can also be described with the ring graph. The number of neighboring nodes with different signs on their signals corresponds to the frequency. Therefore, low-frequency information means a series of signals over N nodes whose signs do not change often. In the case of the SR in our work, where N nodes mean N item embeddings, such low-frequency information means a long-standing interest of a user (see Fig. 1).

Motivation

In this section, we show that self-attention in the spectral domain is a low-pass filter that continuously erases high-frequency information. We visualize the spectrum of

self-attention of the Transformer-based sequential model as shown in Fig. 2 (b). It shows that the spectrum is concentrated in the low frequency region, and it reveals that self-attention is a low-pass filter. We further make theoretical justifications for the low-pass filter of self-attention.

Theorem 1 (Self-Attention is a low-pass filter). *Let $\mathbf{A} = \text{softmax}(\mathbf{Q}\mathbf{K}^T/\sqrt{d})$. Then \mathbf{A} inherently acts as a low-pass filter. For all $\mathbf{x} \in \mathbb{R}^N$, in other words, $\lim_{t \rightarrow \infty} \|HFC[\mathbf{A}^t(\mathbf{x})]\|_2 / \|LFC[\mathbf{A}^t(\mathbf{x})]\|_2 = 0$.*

Theorem 1 is ensured by the Perron-Frobenius theorem (Meyer and Stewart 2023; He and Wai 2021), revealing that the attention matrix is a low-pass filter independent of the input key and query matrices. A proof of Theorem 1 and the formal definition of the low-pass filter are provided in Appendix (Shin et al. 2023). If the self-attention matrix is applied successively, the final output loses all feature expressiveness as the number of layers increases to infinity.

Therefore, the self-attention causes the oversmoothing problem that Transformer-based SR models lose feature representation in deep layers (see Fig. 3). As can be seen from the empirical analysis of Fig. 3, as the number of layers of these models increases, the cosine similarity increases, and the singular value tends to decay rapidly² (Fan et al. 2023). This inevitably causes the model to fail to capture the user’s detailed preferences, and performance degradation is a natural result.

We not only alleviate oversmoothing using a high-pass filter as motivation against this background, but also try to capture short-term preferences of user behavior patterns through inductive bias.

Proposed Method

Here, we introduce the overview of BSARec, the method behind our BSARec, and its relation to previous models.

Embedding Layer

Given a user’s action sequence \mathcal{S}^u and the maximum sequence length N , the sequence is first truncated by removing earliest item if $|\mathcal{S}^u| > N$ or padded with 0s to get a fixed length sequence $\mathbf{s} = (s_1, s_2, \dots, s_N)$. With an item embedding matrix $\mathbf{M} \in \mathbb{R}^{|\mathcal{V}| \times D}$, we define the embedding representation of the sequence \mathbf{E}^u , where D is the latent dimension size and $\mathbf{E}_i^u = \mathbf{M}_{s_i}$. To make our model sensitive to the positions of items, we adopt positional embedding to inject additional positional information while maintaining the same embedding dimensions of the item embedding. A trainable positional embedding $\mathbf{P} \in \mathbb{R}^{N \times D}$ is added to the sequentially ordered item embedding matrix \mathbf{E}^u . Moreover, dropout and layer normalization are also implemented:

$$\mathbf{E}^u = \text{Dropout}(\text{LayerNorm}(\mathbf{E}^u + \mathbf{P})). \quad (4)$$

Beyond Self-Attention Encoder

We develop item encoders by stacking beyond self-attention (BSA) blocks based on the embedding layer. It consists of

²This indicates that the largest singular value predominates and the other outliers are much smaller, and there is a potential risk of losing embedding rank.

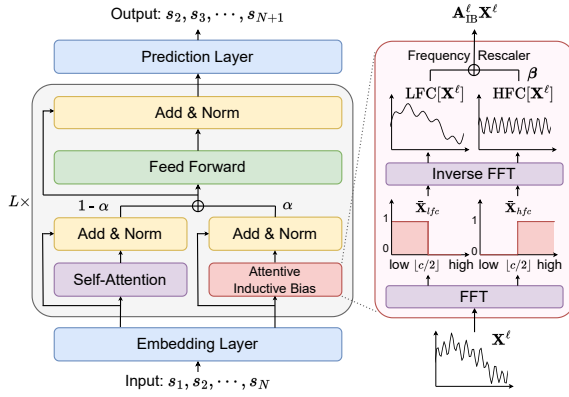


Figure 4: Architecture of our proposed BSARec. We propose a BSA encoder that uses both an inductive bias with a frequency rescaler and original self-attention.

3 modules (see Fig. 4): BSA layer, attentive inductive bias with frequency rescaler, and feed forward network.

Beyond Self-Attention Layer Let $\tilde{\mathbf{A}}^\ell$ be a beyond self-attention (BSA), $\mathbf{A}_{\text{IB}}^\ell$ be a rescaled filter matrix for the l -th layer, and \mathbf{X}^ℓ is the input for the l -th layer. When $l = 0$, we set $\mathbf{X}^0 = \mathbf{E}^u$. We use the following BSA layer:

$$\mathbf{S}^\ell = \tilde{\mathbf{A}}^\ell \mathbf{X}^\ell = \alpha \mathbf{A}_{\text{IB}}^\ell \mathbf{X}^\ell + (1 - \alpha) \mathbf{A}^\ell \mathbf{X}^\ell, \quad (5)$$

where the first term corresponds to DSP, where the discrete Fourier transform is utilized, $\alpha \leq 1$ is a coefficient to (de-)emphasize the inductive bias. Therefore, our main design point is to trade off between the verified inductive bias and the trainable self-attention.

For the multi-head version used in BSARec, the multi-head self-attention (MSA) is defined as:

$$\hat{\mathbf{X}}^\ell = \text{MSA}(\mathbf{X}^\ell) = [\mathbf{S}^1, \mathbf{S}^2, \dots, \mathbf{S}^h] \mathbf{W}^O, \quad (6)$$

where h is the number of heads and the projection matrix $\mathbf{W}^O \in \mathbb{R}^{D \times D}$ is the learnable parameter.

Attentive Inductive Bias with Frequency Rescaler We propose a filter that injects the attentive inductive bias and, at the same time, adjusts the scale of the frequency by dividing it into low and high frequency components:

$$\mathbf{A}_{\text{IB}}^\ell \mathbf{X}^\ell = \text{LFC}[\mathbf{X}^\ell] + \beta \text{HFC}[\mathbf{X}^\ell], \quad (7)$$

where β is a trainable parameter to scale the high-pass filter. In particular, β can be either a vector with D dimension or a scalar parameter.

Meaning of our Attentive Inductive Bias We note that DFT is used in Eq. (7), which assumes the ring graph in Fig. 2 (a) — in the perspective of self-attention, this inductive bias says that an item to purchase is influenced by its previous item. This attentive inductive bias does not need to be trained since we know that it presents universally for SR.

However, we do not stop at utilizing the inductive bias in a naïve way but extract its low and high-frequency information to learn how to optimally mix them in Eq. (7). To be

more specific, suppose a ring graph of N item embeddings. $\text{LFC}[\cdot]$ on them extracts their common signals that do not greatly change following the ring graph topology whereas $\text{HFC}[\cdot]$ extracts locally fluctuating signals (see Fig. 1). By selectively utilizing the high-pass information, we can prevent the oversmoothing problem (see Fig. 3). If relying on $\text{LFC}[\cdot]$ only, we cannot prevent the oversmoothing problem.

In addition, we also learn the self-attention matrix \mathbf{A}^ℓ in Eq. (5) and combine it with our attentive inductive bias $\mathbf{A}_{\text{IB}}^\ell$. By separating \mathbf{A}^ℓ from $\tilde{\mathbf{A}}^\ell$, the self-attention mechanism focuses on capturing non-obvious attentions in \mathbf{A}^ℓ .

Point-wise Feed-Forward Network and Layer Outputs

The multi-head attention function is primarily based on linear projection. A point-wise feed-forward network is applied to import nonlinearity to the self-attention block. The process is defined as follows:

$$\tilde{\mathbf{X}}^\ell = (\text{GELU}(\hat{\mathbf{X}}^\ell \mathbf{W}_1^\ell + \mathbf{b}_1^\ell)) \mathbf{W}_2^\ell + \mathbf{b}_2^\ell, \quad (8)$$

where $\mathbf{W}_1^\ell, \mathbf{W}_2^\ell \in \mathbb{R}^{D \times D}$ and $\mathbf{b}_1^\ell, \mathbf{b}_2^\ell \in \mathbb{R}^{D \times D}$ are learnable parameters. The dropout layer, residual connection, and layer normalization operations are applied as follows:

$$\mathbf{X}^{\ell+1} = \text{LayerNorm}(\mathbf{X}^\ell + \hat{\mathbf{X}}^\ell + \text{Dropout}(\tilde{\mathbf{X}}^\ell)). \quad (9)$$

Prediction Layer and Training

In the final layer of BSARec, we calculate the user’s preference score for the item i derived from user’s historical interactions. This score is given by:

$$\hat{y}_i = p(v_{|\mathcal{S}^u|+1}^u = v | \mathcal{S}^u) = \mathbf{e}_v^T \mathbf{X}_{|\mathcal{S}^u|}^L, \quad (10)$$

where \mathbf{e}_v is the representation of item v from \mathbf{M} , and $\mathbf{X}_{|\mathcal{S}^u|}^L$ is the output of the L -layer blocks at step $|\mathcal{S}^u|$. This dot product computes the similarity between these two vectors to give us the preference score \hat{y}_i .

The cross-entropy (CE) loss function is usually used in SR since the next item prediction task is treated as a classification task over the whole item set (Zhang et al. 2019; Qiu et al. 2022; Du et al. 2023). We adopt the CE loss to optimize the model parameter as:

$$\mathcal{L} = -\log \frac{\exp(\hat{y}_g)}{\sum_{i \in |\mathcal{V}|} \exp(\hat{y}_i)}, \quad (11)$$

where $g \in |\mathcal{V}|$ is the ground truth item.

Relation to Previous Models

Several Transformer-based SR models can be a special case of BSARec, and the comparison with existing models is as follows: i) When α is 0 in BSARec, our model is reduced to SASRec. This is because pure self-attention is used as it is. However, one difference is that their loss functions are different. BSARec uses the CE loss, while SASRec uses the BCE loss. Even in the case of DuoRec, which extends SASRec with contrastive learning, it can be a BSARec with $\alpha = 0$ except for contrastive learning. ii) In the case of the FMLPRec, it uses DFT only without self-attention. Nevertheless, the most significant difference is that the filter matrix itself in FMLPRec is a learnable matrix. Because of this,

Datasets	Metric	Caser	GRU4Rec	SASRec	BERT4Rec	FMLPRec	DuoRec	FEARec	BSARec	Improv.
Beauty	HR@5	0.0125	0.0169	0.0340	0.0469	0.0346	<u>0.0707</u>	0.0706	0.0736	4.10%
	HR@10	0.0225	0.0304	0.0531	0.0705	0.0559	0.0965	0.0982	0.1008	2.65%
	HR@20	0.0403	0.0527	0.0823	0.1073	0.0869	0.1313	<u>0.1352</u>	0.1373	1.55%
	NDCG@5	0.0076	0.0104	0.0221	0.0311	0.0222	0.0501	<u>0.0512</u>	0.0523	2.15%
	NDCG@10	0.0108	0.0147	0.0283	0.0387	0.0291	0.0584	<u>0.0601</u>	0.0611	1.66%
	NDCG@20	0.0153	0.0203	0.0356	0.0480	0.0369	0.0671	<u>0.0694</u>	0.0703	1.30%
Sports	HR@5	0.0091	0.0118	0.0188	0.0275	0.0220	0.0396	0.0411	0.0426	3.65%
	HR@10	0.0163	0.0187	0.0298	0.0428	0.0336	0.0569	<u>0.0589</u>	0.0612	3.90%
	HR@20	0.0260	0.0303	0.0459	0.0649	0.0525	0.0791	<u>0.0836</u>	0.0858	2.63%
	NDCG@5	0.0056	0.0079	0.0124	0.0180	0.0146	0.0276	<u>0.0286</u>	0.0300	4.90%
	NDCG@10	0.0080	0.0101	0.0159	0.0229	0.0183	0.0331	<u>0.0343</u>	0.0360	4.96%
	NDCG@20	0.0104	0.0131	0.0200	0.0284	0.0231	0.0387	<u>0.0405</u>	0.0422	4.20%
Toys	HR@5	0.0095	0.0121	0.0440	0.0412	0.0432	0.0770	0.0783	0.0805	2.81%
	HR@10	0.0161	0.0211	0.0652	0.0635	0.0671	0.1034	<u>0.1054</u>	0.1081	2.56%
	HR@20	0.0268	0.0348	0.0929	0.0939	0.0974	0.1369	<u>0.1397</u>	0.1435	2.72%
	NDCG@5	0.0058	0.0077	0.0297	0.0282	0.0288	0.0568	<u>0.0574</u>	0.0589	2.61%
	NDCG@10	0.0079	0.0106	0.0366	0.0353	0.0365	0.0653	<u>0.0661</u>	0.0679	2.72%
	NDCG@20	0.0106	0.0140	0.0435	0.0430	0.0441	0.0737	<u>0.0747</u>	0.0768	2.81%
Yelp	HR@5	0.0117	0.0130	0.0149	0.0256	0.0159	<u>0.0271</u>	0.0262	0.0275	1.48%
	HR@10	0.0197	0.0221	0.0249	0.0433	0.0287	<u>0.0442</u>	0.0437	0.0465	5.20%
	HR@20	0.0337	0.0383	0.0424	0.0717	0.0490	<u>0.0717</u>	0.0691	0.0746	4.04%
	NDCG@5	0.0070	0.0080	0.0091	0.0159	0.0100	0.0170	0.0165	0.0170	0.00%
	NDCG@10	0.0096	0.0109	0.0123	0.0216	0.0142	<u>0.0225</u>	0.0221	0.0231	2.67%
	NDCG@20	0.0131	0.0150	0.0167	0.0287	0.0192	<u>0.0294</u>	0.0285	0.0302	2.72%
LastFM	HR@5	0.0303	0.0312	0.0413	0.0294	0.0367	<u>0.0431</u>	0.0431	0.0523	21.35%
	HR@10	0.0431	0.0404	0.0633	0.0459	0.0560	0.0624	0.0587	0.0807	27.49%
	HR@20	0.0642	0.0541	0.0927	0.0596	0.0826	<u>0.0963</u>	0.0826	0.1174	21.91%
	NDCG@5	0.0227	0.0217	0.0284	0.0198	0.0243	<u>0.0300</u>	0.0304	0.0344	13.16%
	NDCG@10	0.0268	0.0245	0.0355	0.0252	0.0306	<u>0.0361</u>	0.0354	0.0435	20.50%
	NDCG@20	0.0321	0.0280	0.0429	0.0286	0.0372	<u>0.0446</u>	0.0414	0.0526	17.94%
ML-1M	HR@5	0.0927	0.1005	0.1374	0.1512	0.1316	0.1838	0.1834	0.1944	5.77%
	HR@10	0.1556	0.1657	0.2137	0.2346	0.2065	0.2704	0.2705	0.2757	1.92%
	HR@20	0.2488	0.2664	0.3245	0.3440	0.3137	<u>0.3738</u>	0.3714	0.3884	3.91%
	NDCG@5	0.0592	0.0619	0.0873	0.1021	0.0846	<u>0.1252</u>	0.1236	0.1306	4.31%
	NDCG@10	0.0795	0.0828	0.1116	0.1289	0.1087	<u>0.1530</u>	0.1516	0.1568	2.48%
	NDCG@20	0.1028	0.1081	0.1395	0.1564	0.1356	<u>0.1790</u>	0.1771	0.1851	3.41%

Table 2: Performance comparison of different methods on 6 datasets. The best results are in boldface and the second-best results are underlined. ‘Improv.’ indicates the relative improvement against the best baseline performance.

FMLPRec’s filter is inevitably learned as a low-pass filter, while BSARec uses a filter rescaler to simultaneously use a high-pass filter. iii) Similar to BSARec, FEARec separates low-frequency and high-frequency information in the frequency domain. However, FEARec allows the frequency domain to be learned separately before entering its Transformer’s encoder. BSARec adaptively uses low and high-frequency information by using a frequency rescaler in a step to inject an inductive bias. FEARec is designed with a complex model structure using contrastive learning and frequency normalization. However, our model shows better performance with a much simpler architecture.

Experiments

Experimental Setup

Datasets We evaluate our model on 6 SR datasets where the sparsity and domain varies: i,ii,iii) Amazon Beauty, Sports, Toys (McAuley et al. 2015), iv) Yelp, v) ML-1M (Harper and Konstan 2015), and vi) LastFM. We followed the data pre-processing procedure from Zhou et al. (2020, 2022), where all reviews and ratings are regarded

as implicit feedback. The detailed dataset statistics are presented in Appendix (Shin et al. 2023).

Baselines To verify the effectiveness of our model, we compare our method with well-known SR baselines with three categories:

- RNN or CNN-based sequential models: GRU4Rec (Hidasi et al. 2016) and Caser (Tang and Wang 2018).
- Transformer-based sequential models: SASRec (Kang and McAuley 2018), BERT4Rec (Sun et al. 2019), and FMLPRec (Zhou et al. 2022).
- Transformer-based sequential models with contrastive learning: DuoRec (Qiu et al. 2022) and FEARec (Du et al. 2023).

Implementation Details Our method is implemented in PyTorch on an NVIDIA RTX 3090 with 16 GB memory. We search the best hyperparameters for baselines based on their recommended hyperparameters. We conduct experiments under the following hyperparameters: the coefficient α is in $\{0.1, 0.3, 0.5, 0.7, 0.9\}$, and c is chosen from $\{1, 3, 5, 7, 9\}$. The number of BSA blocks L is set to 2, and the number of heads in Transformer h is in $\{1, 2, 4\}$. The dimension of

Methods	Beauty		Toys	
	HR@20	NDCG@20	HR@20	NDCG@20
BSARec	0.1373	0.0703	0.1435	0.0768
Only \mathbf{A}	0.1265	0.0657	0.1320	0.0720
Only \mathbf{A}_{IB}	0.1338	0.0677	0.1402	0.0744
Scalar β	0.1333	0.0685	0.1435	0.0756

Table 3: Ablation studies on $\tilde{\mathbf{A}}$ and β . More results in other datasets are in Appendix (Shin et al. 2023).

D is set to 64, and the maximum sequence length N is set to 50. For training, the Adam optimizer is optimized with a learning rate in $\{5 \times 10^{-4}, 1 \times 10^{-3}\}$, and the batch size is set to 256. The best hyperparameters are in Appendix (Shin et al. 2023) for reproducibility.

Metrics To measure the recommendation accuracy, we commonly use widely used Top- k metrics, HR@ k (Hit Rate) and NDCG@ k (Normalized Discounted Cumulative Gain) to evaluate the recommended list, where k is set to 5, 10, and 20. To ensure a fair and comprehensive comparison, we analyze the ranking results across the full item set without negative sampling (Krichene and Rendle 2020).

Experimental Results

Table 2 presents the detailed recommendation performance. Overall, our proposed method, BSARec, clearly marks the best accuracy. First, compared to existing RNN-based and CNN-based methods, Transformer-based methods show better performance in modeling interaction sequences in SR. Second, in Transformer-based methods, BERT4Rec and FMLPRec models outperform SASRec. In particular, FMLPRec redesigned the self-attention of the existing Transformer only with MLP, but it still does not perform well in all datasets. Third, there is no doubt that models using contrastive learning show higher results than models that do not. DuoRec and FEARec significantly outperform SASRec, BERT4Rec, and FMLPRec.

Surprisingly, however, BSARec records the best performance across all datasets and all metrics. The most surprising thing is that it can show better performance than DuoRec and FEARec without using contrastive learning. In LastFM, BSARec shows a performance improvement of 27.49% based on HR@10. Thus, our model leaves a message that it can show good performance without going to complex model design by adding contrastive learning.

Ablation, Sensitivity, and Additional Studies

Ablation Studies As ablation study models, we define the following models: i) the first ablation model has only the self-attention term i.e., \mathbf{A} , ii) the second ablation model has only the attentive inductive bias term, i.e., \mathbf{A}_{IB} in Eq. 5, and iii) the third ablation model uses β as a single parameter. For Beauty and Toys, the ablation study model with only \mathbf{A}_{IB} outperforms the case with only \mathbf{A} (e.g., HR@20 in Beauty by \mathbf{A}_{IB} of 0.1338 versus 0.1265 by \mathbf{A}). However, BSARec,

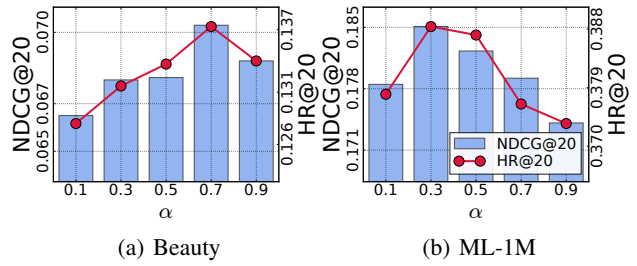


Figure 5: Sensitivity to α . More results in other datasets are in Appendix (Shin et al. 2023).

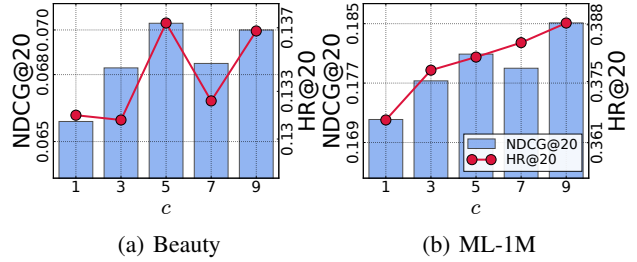


Figure 6: Sensitivity to c . More results in other datasets are in Appendix (Shin et al. 2023).

which utilizes them all, outperforms them. This shows that both are required to achieve the best accuracy.

Sensitivity to α Fig. 5 shows the NDCG@20 and HR@20 by varying the α . For Beauty, we find our BSARec, a larger value of α is preferred. For ML-1M, with $\alpha = 0.3$, we can achieve the best accuracy. The trade-off between the self-attention matrix and the inductive bias differs for each dataset from these results.

Sensitivity to c Fig. 6 shows the NDCG@20 and HR@20 by varying the c . For Beauty, the best accuracy is achieved when c is 5. For ML-1M, the larger the value of c , the better performance is reached.

Visualization of Learned β In Fig. 7 (a), we show learned β at each layer for all datasets. We can see that a higher weight is learned in the first layer than in the second layer, which confirms that putting more weight on high-frequency in the first layer is effective. In particular, LastFM and Beauty show higher β weights than other datasets.

Case Study We introduce a case study obtained from our experiment. In Fig. 7 (b), we analyze one of the heavy users in LastFM. The user u_{322} constantly listens to artists, mainly in the rock genre. In other models, u_{322} cannot capture sudden interaction changes in the next step. Only BSARec recommends an artist from the pop genre as the next artist u_{322} will listen to. This shows that BSARec can capture high-frequency signals that are abrupt changes in user preference.

Model Complexity and Runtime Analyses

To evaluate the overhead of BSARec, we evaluate the number of parameters and runtime per epoch during training.

Acknowledgements

Noseong Park is the corresponding author. This work was supported by an IITP grant funded by the Korean government (MSIT) (No.2020-0-01361, Artificial Intelligence Graduate School Program (Yonsei University); No.2022-0-00113, Developing a Sustainable Collaborative Multi-modal Lifelong Learning Framework).

References

- Chen, Q.; Zhao, H.; Li, W.; Huang, P.; and Ou, W. 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.
- Choi, J.; Hong, S.; Park, N.; and Cho, S.-B. 2023a. Blurring-Sharpener Process Models for Collaborative Filtering. In *SIGIR*.
- Choi, J.; Hong, S.; Park, N.; and Cho, S.-B. 2023b. GREED: Graph Neural Reaction-Diffusion Networks. In *ICML*.
- Choi, J.; Jeon, J.; and Park, N. 2021. LT-OCF: Learnable-Time ODE-based Collaborative Filtering. In *CIKM*.
- Choi, J.; Wi, H.; Kim, J.; Shin, Y.; Lee, K.; Trask, N.; and Park, N. 2023c. Graph Convolutions Enrich the Self-Attention in Transformers! *arXiv preprint arXiv:2312.04234*.
- Choi, J.; Wi, H.; Lee, C.; Cho, S.-B.; Lee, D.; and Park, N. 2023d. RDGCL: Reaction-Diffusion Graph Contrastive Learning for Recommendation. *arXiv preprint arXiv:2312.16563*.
- Dong, Y.; Cordonnier, J.-B.; and Loukas, A. 2021. Attention is not all you need: Pure attention loses rank doubly exponentially with depth. In *ICML*, 2793–2803. PMLR.
- Dosovitskiy, A.; Beyer, L.; Kolesnikov, A.; Weissenborn, D.; Zhai, X.; Unterthiner, T.; Dehghani, M.; Minderer, M.; Heigold, G.; Gelly, S.; et al. 2020. An image is worth 16x16 words: Transformers for image recognition at scale. *ICLR*.
- Du, X.; Yuan, H.; Zhao, P.; Qu, J.; Zhuang, F.; Liu, G.; Liu, Y.; and Sheng, V. S. 2023. Frequency Enhanced Hybrid Attention Network for Sequential Recommendation. In *SIGIR*, 78–88.
- Fan, Z.; Liu, Z.; Peng, H.; and Yu, P. S. 2023. Addressing the Rank Degeneration in Sequential Recommendation via Singular Spectrum Smoothing. *arXiv preprint arXiv:2306.11986*.
- Gao, C.; Zheng, Y.; Li, N.; Li, Y.; Qin, Y.; Piao, J.; Quan, Y.; Chang, J.; Jin, D.; He, X.; et al. 2023. A survey of graph neural networks for recommender systems: Challenges, methods, and directions. *ACM Transactions on Recommender Systems*, 1(1): 1–51.
- Gong, C.; Wang, D.; Li, M.; Chandra, V.; and Liu, Q. 2021. Vision transformers with patch diversification. *arXiv preprint arXiv:2104.12753*.
- Guo, X.; Wang, Y.; Du, T.; and Wang, Y. 2023. Contra-norm: A contrastive learning perspective on oversmoothing and beyond. In *ICLR*.
- Hansen, C.; Hansen, C.; Maystre, L.; Mehrotra, R.; Brost, B.; Tomasi, F.; and Lalmas, M. 2020. Contextual and sequential user embeddings for large-scale music recommendation. In *RecSys*, 53–62.
- Harper, F. M.; and Konstan, J. A. 2015. The movielens datasets: History and context. *Acm transactions on interactive intelligent systems (tiis)*, 5(4): 1–19.
- He, R.; and McAuley, J. 2016. Fusing similarity models with markov chains for sparse sequential recommendation. In *ICDM*, 191–200. IEEE.
- He, X.; Deng, K.; Wang, X.; Li, Y.; Zhang, Y.; and Wang, M. 2020. LightGCN: Simplifying and Powering Graph Convolution Network for Recommendation. In *SIGIR*.
- He, Y.; and Wai, H.-T. 2021. Identifying first-order lowpass graph signals using perron frobenius theorem. In *ICASSP 2021-2021 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, 5285–5289. IEEE.
- Hidasi, B.; Karatzoglou, A.; Baltrunas, L.; and Tikk, D. 2016. Session-based recommendations with recurrent neural networks. In *ICLR*.
- Hong, S.; Jo, M.; Kook, S.; Jung, J.; Wi, H.; Park, N.; and Cho, S.-B. 2022. TimeKit: A Time-series Forecasting-based Upgrade Kit for Collaborative Filtering. In *2022 IEEE International Conference on Big Data (Big Data)*, 565–574. IEEE.
- Huang, X.; Qian, S.; Fang, Q.; Sang, J.; and Xu, C. 2018. CSAN: Contextual self-attention network for user sequential recommendation. In *ACM MM*, 447–455.
- Jiang, J.; Zhang, P.; Luo, Y.; Li, C.; Kim, J. B.; Zhang, K.; Wang, S.; Xie, X.; and Kim, S. 2023. AdaMCT: adaptive mixture of CNN-transformer for sequential recommendation. In *CIKM*.
- Jiang, S.; Qian, X.; Mei, T.; and Fu, Y. 2016. Personalized travel sequence recommendation on multi-source big social media. *IEEE Transactions on Big Data*, 2(1): 43–56.
- Kang, W.-C.; and McAuley, J. 2018. Self-attentive sequential recommendation. In *ICDM*, 197–206. IEEE.
- Kong, T.; Kim, T.; Jeon, J.; Choi, J.; Lee, Y.-C.; Park, N.; and Kim, S.-W. 2022. Linear, or Non-Linear, That is the Question! In *WSDM*, 517–525.
- Krichene, W.; and Rendle, S. 2020. On sampled metrics for item recommendation. In *Proceedings of the 26th ACM SIGKDD international conference on knowledge discovery & data mining*, 1748–1757.
- Lee, Y.-C.; Kim, S.-W.; and Lee, D. 2018. gOCCF: Graph-theoretic one-class collaborative filtering based on uninteresting items. In *AAAI*, volume 32.
- Li, J.; Wang, Y.; and McAuley, J. 2020. Time interval aware self-attention for sequential recommendation. In *WSDM*, 322–330.
- Li, Q.; Han, Z.; and Wu, X.-M. 2018. Deeper insights into graph convolutional networks for semi-supervised learning. In *AAAI*.
- Lin, G.; Gao, C.; Zheng, Y.; Chang, J.; Niu, Y.; Song, Y.; Gai, K.; Li, Z.; Jin, D.; Li, Y.; et al. 2023. Mixed Attention Network for Cross-domain Sequential Recommendation. *arXiv preprint arXiv:2311.08272*.

- Liu, Q.; Yan, F.; Zhao, X.; Du, Z.; Guo, H.; Tang, R.; and Tian, F. 2023. Diffusion Augmentation for Sequential Recommendation. In *CIKM*, 1576–1586.
- McAuley, J.; Targett, C.; Shi, Q.; and Van Den Hengel, A. 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.
- Meyer, C. D.; and Stewart, I. 2023. *Matrix analysis and applied linear algebra*. SIAM.
- Qiu, R.; Huang, Z.; Yin, H.; and Wang, Z. 2022. Contrastive learning for representation degeneration problem in sequential recommendation. In *WSDM*, 813–823.
- Rendle, S.; Freudenthaler, C.; and Schmidt-Thieme, L. 2010. Factorizing personalized markov chains for next-basket recommendation. In *TheWebConf (former WWW)*, 811–820.
- Rusch, T. K.; Chamberlain, B.; Rowbottom, J.; Mishra, S.; and Bronstein, M. 2022. Graph-Coupled Oscillator Networks. In *ICML*, volume 162, 18888–18909.
- Sandryhaila, A.; and Moura, J. M. 2014. Discrete signal processing on graphs: Frequency analysis. *IEEE Transactions on Signal Processing*, 62(12): 3042–3054.
- Schedl, M.; Zamani, H.; Chen, C.-W.; Deldjoo, Y.; and Elahi, M. 2018. Current challenges and visions in music recommender systems research. *International Journal of Multimedia Information Retrieval*, 7: 95–116.
- Shin, Y.; Choi, J.; Wi, H.; and Park, N. 2023. An Attentive Inductive Bias for Sequential Recommendation Beyond the Self-Attention. *arXiv preprint arXiv:2312.10325*.
- Sun, F.; Liu, J.; Wu, J.; Pei, C.; Lin, X.; Ou, W.; and Jiang, P. 2019. BERT4Rec: Sequential recommendation with bidirectional encoder representations from transformer. In *CIKM*, 1441–1450.
- Tang, J.; and Wang, K. 2018. Personalized top-n sequential recommendation via convolutional sequence embedding. In *WSDM*, 565–573.
- Vaswani, A.; Shazeer, N.; Parmar, N.; Uszkoreit, J.; Jones, L.; Gomez, A. N.; Kaiser, Ł.; and Polosukhin, I. 2017. Attention is all you need. In *NeurIPS*.
- Wang, P.; Zheng, W.; Chen, T.; and Wang, Z. 2022. Anti-Oversmoothing in Deep Vision Transformers via the Fourier Domain Analysis: From Theory to Practice. In *ICLR*.
- Wu, J.; Cai, R.; and Wang, H. 2020. Déjà vu: A contextualized temporal attention mechanism for sequential recommendation. In *TheWebConf (former WWW)*, 2199–2209.
- Wu, L.; Li, S.; Hsieh, C.-J.; and Sharpnack, J. 2020. SSE-PT: Sequential recommendation via personalized transformer. In *RecSys*, 328–337.
- Wu, S.; Sun, F.; Zhang, W.; Xie, X.; and Cui, B. 2022. Graph neural networks in recommender systems: a survey. *ACM Computing Surveys*, 55(5): 1–37.
- Ying, R.; He, R.; Chen, K.; Eksombatchai, P.; Hamilton, W. L.; and Leskovec, J. 2018. Graph Convolutional Neural Networks for Web-Scale Recommender Systems. In *KDD*.
- Yue, Z.; Wang, Y.; He, Z.; Zeng, H.; McAuley, J.; and Wang, D. 2023. Linear Recurrent Units for Sequential Recommendation. *arXiv preprint arXiv:2310.02367*.
- Zhang, T.; Zhao, P.; Liu, Y.; Sheng, V. S.; Xu, J.; Wang, D.; Liu, G.; Zhou, X.; et al. 2019. Feature-level Deeper Self-Attention Network for Sequential Recommendation. In *IJCAI*, 4320–4326.
- Zhou, D.; Kang, B.; Jin, X.; Yang, L.; Lian, X.; Jiang, Z.; Hou, Q.; and Feng, J. 2021. Deepvit: Towards deeper vision transformer. *arXiv preprint arXiv:2103.11886*.
- Zhou, K.; Wang, H.; Zhao, W. X.; Zhu, Y.; Wang, S.; Zhang, F.; Wang, Z.; and Wen, J.-R. 2020. S3-rec: Self-supervised learning for sequential recommendation with mutual information maximization. In *CIKM*, 1893–1902.
- Zhou, K.; Yu, H.; Zhao, W. X.; and Wen, J.-R. 2022. Filter-enhanced MLP is all you need for sequential recommendation. In *TheWebConf (former WWW)*, 2388–2399.
- Zhou, P.; Ye, Q.; Xie, Y.; Gao, J.; Wang, S.; Kim, J. B.; You, C.; and Kim, S. 2023. Attention Calibration for Transformer-based Sequential Recommendation. In *CIKM*, 3595–3605.