

Modeling Sequences as Star Graphs to Address Over-smoothing in Self-attentive Sequential Recommendation

Bo Peng

The Ohio State University
Columbus, OH, US
peng.707@buckeyemail.osu.edu

Srinivasan Parthasarathy

The Ohio State University
Columbus, OH, US
srini@cse.ohio-state.edu

Ziqi Chen

The Ohio State University
Columbus, OH, US
chen.8484@buckeyemail.osu.edu

Xia Ning

The Ohio State University
Columbus, OH, US
ning.104@osu.edu

ABSTRACT

Self-attention (SA) mechanisms have been widely used in developing sequential recommendation (SR) methods, and demonstrated state-of-the-art performance. However, in this paper, we show that self-attentive SR methods substantially suffer from the over-smoothing issue that item embeddings within a sequence become increasingly similar across attention blocks. As widely demonstrated in the literature, this issue could lead to a loss of information in individual items, and significantly degrade models' scalability and performance. To address the over-smoothing issue, in this paper, we view items within a sequence constituting a star graph and develop a method, denoted as MSSG, for SR. Different from existing self-attentive methods, MSSG introduces an additional internal node to specifically capture the global information within the sequence, and does not require information propagation among items. This design fundamentally addresses the over-smoothing issue and enables MSSG a linear time complexity with respect to the sequence length. We compare MSSG with ten state-of-the-art baseline methods on six public benchmark datasets. Our experimental results demonstrate that MSSG significantly outperforms the baseline methods, with an improvement of as much as 10.10%. Our analysis shows the superior scalability of MSSG over the state-of-the-art self-attentive methods. Our complexity analysis and run-time performance comparison together show that MSSG is both theoretically and practically more efficient than self-attentive methods. Our analysis of the attention weights learned in SA-based methods indicates that on sparse recommendation data, modeling dependencies in all item pairs using the SA mechanism yields limited information gain, and thus, might not benefit the recommendation performance¹.

¹We will release the source code and processed datasets upon the acceptance of this paper.

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 ACM must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from permissions@acm.org.

Conference acronym 'XX, June 03–05, 2018, Woodstock, NY

© 2024 Association for Computing Machinery.
ACM ISBN 978-1-4503-XXXX-X/18/06...\$15.00
<https://doi.org/XXXXXXXX.XXXXXXX>

CCS CONCEPTS

• Information systems → Recommender systems.

KEYWORDS

Sequential Recommendation, Over-smoothing, Self-attention

ACM Reference Format:

Bo Peng, Ziqi Chen, Srinivasan Parthasarathy, and Xia Ning. 2024. Modeling Sequences as Star Graphs to Address Over-smoothing in Self-attentive Sequential Recommendation. In *Proceedings of Make sure to enter the correct conference title from your rights confirmation email (Conference acronym 'XX)*. ACM, New York, NY, USA, 12 pages. <https://doi.org/XXXXXXXX.XXXXXXX>

1 INTRODUCTION

SR aims to identify and recommend the next item of users' interest based on their historical interactions. It has been widely employed in applications such as online retail [9] and video streaming [1], and has been drawing increasing attention from the research community. With the prosperity of deep learning, deep neural networks, especially recurrent neural networks (RNNs) [12] have emerged as the prominent architecture of SR methods. From a graph perspective, RNNs represent each interaction sequence using a sequential graph as shown in Figure 1a, and recurrently integrate information from items within the sequence. This design, however, may not effectively capture the long-range dependencies within sequences as it is challenging to propagate information across long distances [31].

To better capture the long-range dependencies, SA mechanisms [31] have recently been utilized for SR [14, 37], and demonstrated state-of-the-art performance [14]. Given a sequence, SA mechanisms learn attention weights for all pairs of items within sequences, and aggregates items using these attention weights. From a graph perspective, this is essentially equivalent to viewing items within a sequence as nodes in a fully connected graph² [4] as illustrated in Figure 1b, where v_s is the t -th item in the sequence, and the edge weights are determined by the attention weights. The fully connected graph allows each item node to propagate information to all others in a single step, thereby enabling SA-based methods to effectively capture long-range dependencies within a sequence.

Though promising, modeling sequences as fully connected graphs leads to two main issues. First, as will be demonstrated in our analysis (Section 6.3), SA-based SR methods substantially suffer from

²For simplicity, we drop the self-loops in the graph.

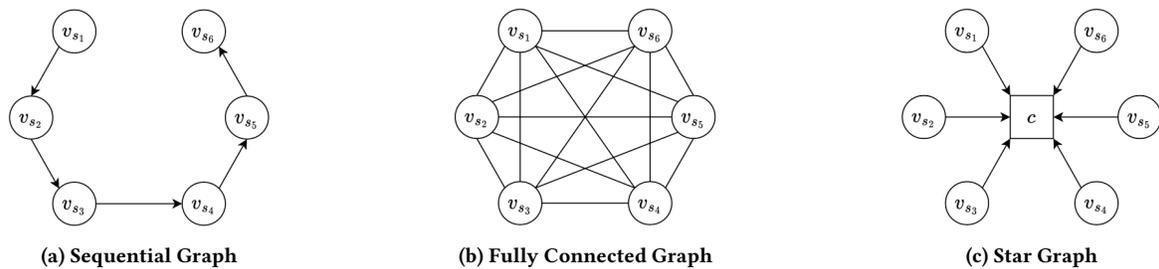


Figure 1: Illustration of a sequence represented by different graphs, in which v_{s_t} is the t -th item in the sequence and c is the internal node in the star graph.

the over-smoothing issue [3] that the embeddings of items within a sequence could become more and more similar across attention blocks. Learning similar embeddings for all items within a sequence leads to the loss of information on individual items, and thus, could substantially degrade the model performance [3, 13, 27, 32]. In addition, as shown in Shi *et al.* [27], the over-smoothing issue also deteriorates the learning of deep models, thereby degrading the model scalability. Second, SA-based SR methods generally suffer from quadratic time complexities, which limit their utilities in large-scale recommendation applications.

Addressing the above issues while still effectively capturing the long-range dependencies within sequences represents a critical research challenge in developing SR methods. To tackle this challenge, in this paper, we develop MSSG, a method that models user interaction sequences using star graphs for SR. Specifically, MSSG stacks multiple blocks to recursively capture users’ intent from their interaction sequence [14]. MSSG has an attention layer and a feed-forward layer in each block. In the attention layer, as illustrated in Figure 1c, MSSG models each sequence as a star graph, where edge weights are determined by the attention weights, and introduces an additional internal node c to integrate information from all item nodes v_{s_t} , thereby capturing the user’s intent. In the feed-forward layer, MSSG updates the embedding of the internal node via non-linear networks to more accurately estimate the user’s intent. Previous work [3] shows that requiring information propagation among all item nodes could be the essential reason for the over-smoothing issue. Thus, MSSG does not propagate information among items. Instead, it utilizes an internal node for information aggregation. This design allows MSSG to fundamentally address the over-smoothing issue (Section 6.3). Beyond addressing the over-smoothing issue, by modeling sequences as star graphs, MSSG also achieves a linear time complexity with respect to the sequence length as will be proven in Section 4.5. In addition, similarly to that in SA mechanisms, MSSG could effectively capture long-range dependencies within sequences as each item node could propagate information to the internal node within one step.

We extensively compare MSSG with ten state-of-the-art baseline methods on six benchmark datasets. Our experimental results demonstrate that MSSG could remarkably outperform the state-of-the-art baseline methods on the six datasets, with an improvement of up to 10.10% at Recall@10 (Section 5.3.2). Our experimental results also show that MSSG consistently outperforms the state-of-the-art SA-based SR method on users with different activity levels (Section 6.2). In addition, our analysis reveals that existing

SA-based SR methods suffer from the over-smoothing issue (Section 6.3), which limits their scalability in learning deep models (Section 6.4). Our complexity analysis (Section 4.5) and run-time performance comparison (Section 6.5) together show that MSSG is both theoretically and practically more efficient than the state-of-the-art SA-based SR methods. Our analysis also suggests that on sparse recommendation data, leveraging SA mechanisms to capture item dependencies could yield limited information gain.

We summarize our major contributions as follows: 1) To the best of our knowledge, this is the first work demonstrating the over-smoothing issue in SA-based SR methods; 2) We develop MSSG for SR in which we model sequences using star graphs to address the over-smoothing issue and enable linear time complexity; 3) MSSG outperforms ten state-of-the-art baseline methods on six benchmark datasets; 4) Our analysis suggests that MSSG could achieve both superior scalability and superior run-time performance compared to SA-based SR methods; 5) Our analysis also indicates that modeling item dependencies using the SA mechanism on sparse recommendation data could yield limited information gain, and thus, might not improve the recommendation performance.

2 RELATED WORK

2.1 Sequential Recommendation

Numerous SR methods have been developed, particularly using Markov Chains (MCs) and neural networks. Specifically, Rendle *et al.* [25] developed FPMC, a method in which MCs are utilized to model the transitions among items. In recent years, neural networks such as RNNs and convolutional neural networks (CNNs) have been widely adapted for SR. For instance, Li *et al.* [15] developed NARM, which incorporates attention mechanisms into RNNs to better capture users’ long-term preferences. Tang *et al.* [29] developed a CNNs-based model Caser that employs vertical and horizontal convolutional filters to capture the synergies among items for better recommendation. Besides RNNs and CNNs-based methods, SA-based methods are also widely developed for SR. Kang *et al.* [14] developed a SA-based method SASRec, which stacks multiple SA blocks to recursively learn users’ preferences. Sun *et al.* [28] integrated the cloze objective and SA mechanisms, and developed a bidirectional sequence modeling method, denoted as Bert4Rec, for comprehensive user intent modeling. Zhou *et al.* [37] developed FMLP for SR, in which they replaced the attention map within each SA block with a predefined transformation matrix for more efficient item aggregation. He *et al.* [10] constrained SA blocks to

attend to local items, thereby enhancing the modeling of users' short-term preferences. Recently, simple shallow methods have been introduced for SR, and shown promising performance. For example, Ma *et al.* [18] developed HGN, which uses a single gating layer to adaptively aggregate items and capture users' preferences. Peng *et al.* [22] developed HAM in which a single pooling layer is employed to learn the associations among items.

We notice that SGNN-HN developed in Pan *et al.* [21] also leverages star graphs for SR. However, MSSG differs significantly from SGNN-HN. Specifically, MSSG is developed to address the over-smoothing issue in SA-based SR methods. To this end, MSSG utilizes the internal node to aggregate information from individual items, and does not allow information propagation among items. In contrast, SGNN-HN aims to improve graph neural networks (GNNs) in capturing long-range dependencies in sequences, and leverages the internal node as an anchor node to better propagate information among items. MSSG and SGNN-HN are developed to address different issues, and have substantially different architectures. Thus, MSSG is not an extension of SGNN-HN.

2.2 Over-smoothing

The over-smoothing issue was firstly identified in GNNs [6] by Li *et al.* [16]. They observed that as GNNs propagate and mix information between neighboring nodes across layers, all nodes could eventually have identical embeddings. As a result, GNNs lose the information on individual nodes, which could substantially deteriorate the model expressiveness [3]. Recently, Shi *et al.* [27] has proven that theoretically, SA-based methods could also suffer from over-smoothing and result in sub-optimal performance. Numerous approaches have been introduced to mitigate this issue. For example, Chen *et al.* [3] optimized the graph topology based on the model prediction to mitigate over-smoothing in GNNs. Shi *et al.* [27] developed hierarchical fusion strategies (ConcatFusion and MaxFusion) to alleviate over-smoothing by fusing embeddings from both shallow layers and deep layers as final output.

3 DEFINITION AND NOTATIONS

In this paper, we tackle the SR problem that given the historical interactions of users, we recommend the next item of users' interest. In this paper, $\mathbb{U} = \{u_1, u_2, \dots\}$ is the set of all the users, where u_i represents the i -th user, and $|\mathbb{U}|$ is the total number of users. Similarly, we represent the set of all the items as $\mathbb{V} = \{v_1, v_2, \dots\}$. $|\mathbb{V}|$ is the total number of items. The historical interactions of u_i is represented as a sequence $S_i = \{v_{s_1(i)}, v_{s_2(i)}, \dots\}$, where $v_{s_t(i)}$ is the t -th interacted item in S_i and $|S_i|$ is the length of the sequence. Given S_i , the next item that u_i will interact with is referred to as the ground-truth next item, denoted as $v_g(i)$. The goal of MSSG is to correctly recommend $v_g(i)$ for u_i . When no ambiguity arises, we will eliminate i in S_i , $v_{s_t(i)}$ and $v_g(i)$. We use uppercase letters to denote matrices, lower-case bold letters to denote row vectors and lower-case non-bold letters to represent scalars. Table A1 (Appendix A) shows the key notations used in this paper.

Over-smoothing in SA-based methods: the embeddings of items within a sequence become increasingly similar across attention blocks, thereby deteriorating the model scalability and performance.

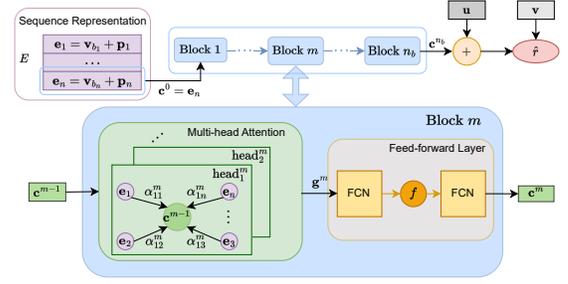


Figure 2: The overall architecture of MSSG. MSSG models sequences using star graphs and utilizes the internal node c to integrate information from all item nodes.

4 METHOD

Figure 2 presents the overall architecture of MSSG. MSSG treats sequences as star graphs to address the over-smoothing issue. MSSG also explicitly captures users' long-term preferences to enable comprehensive user intent modeling. We present the detail of MSSG in the following sections.

4.1 Sequence Representation

Following SASRec, we focus on the most recent n items in users' interaction histories to generate recommendations. Specifically, we transform each interaction sequence S to a fixed-length sequence $B = \{v_{b_1}, v_{b_2}, \dots\}$, which contains the last n items in S (i.e., $v_{b_t} = v_{s_{|S_i| - n + t}}$). If $|S_i|$ is shorter than n , we pad empty items, denoted as v_0 , at the beginning of B . In MSSG, following the literature [14, 31], we represent items and positions in each sequence using learnable embeddings. We learn an item embedding matrix $V \in \mathbb{R}^{|\mathbb{V}| \times d}$ in which the j -th row \mathbf{v}_j is the embedding of item v_j , and d is the dimension of embeddings. We use a constant zero vector as the embedding of padding items. Similarly, we learn a position embedding matrix $P \in \mathbb{R}^{n \times d}$, and \mathbf{p}_t is the embedding of the t -th position. Given V and P , we represent the sequence B using a matrix E :

$$E = [\mathbf{e}_1; \mathbf{e}_2; \dots; \mathbf{e}_n] = [\mathbf{v}_{b_1} + \mathbf{p}_1; \mathbf{v}_{b_2} + \mathbf{p}_2; \dots; \mathbf{v}_{b_n} + \mathbf{p}_n], \quad (1)$$

where \mathbf{v}_{b_t} is the embedding of v_{b_t} .

4.2 Modeling Sequences as Star Graphs

MSSG employs a multi-block architecture to recursively model users' intent. Distinct from existing SA-based SR methods [5, 14], in each block, MSSG treats each item within B as a node in a star graph, and introduces an additional internal node to estimate the user's intent by aggregating information from all item nodes using multi-head attention [31]. Across blocks, MSSG updates only the embedding of the internal node to enhance the estimation of the user's intent, while maintaining fixed embeddings on item nodes to preserve their individual information. In MSSG, we have n_b blocks, and each block contains an attention layer and a feed-forward layer. In what follows, we present the attention layer and the feed-forward layer in each block in detail.

4.2.1 Attention layer. MSSG employs multi-head attention to adaptively aggregate information from item nodes to the central internal

node in each block as follows:

$$\alpha_k^m = \text{softmax} \left(\frac{(\mathbf{c}^{m-1} Q_k^m) (EZ_k^m)^\top}{\sqrt{d}} \right), \quad \text{head}_k^m = \alpha_k^m (EA_k^m), \quad (2)$$

$$\mathbf{g}^m = [\text{head}_1^m, \text{head}_2^m, \dots, \text{head}_{n_h}^m] O^m,$$

where \mathbf{c}^{m-1} is the embedding of the internal node from the $(m-1)$ -th block, and \mathbf{g}^m is the output of the attention layer in the m -th block. Each attention layer has n_h attention heads and head_k^m is the output from the k -th attention head. α_k^m is the attention weights learned in the k -th attention head. Q_k^m, Z_k^m and $A_k^m \in \mathbb{R}^{d \times \frac{d}{n_h}}$ are learnable parameters in the k -th attention head. $O^m \in \mathbb{R}^{d \times d}$ is a learnable parameter to integrate the attention heads.

Unlike SA-based methods, which require every item to aggregate information from all others to capture global information within the sequence, MSSG introduces an additional internal node to specifically capture the global information and estimate the user's intent. This approach results in a linear time complexity as will be shown in Section 4.5, while could still effectively capture long-range dependencies within the sequence.

4.2.2 Feed-forward layer. After each attention layer, we include a feed-forward layer to endow MSSG with non-linearity, and enable more expressive models [31]. Particularly, given \mathbf{g}^m , we stack two fully-connected networks (FCNs) as the feed-forward layer:

$$\mathbf{c}^m = (f(\mathbf{g}^m W_1^m + \mathbf{b}_1^m)) W_2^m + \mathbf{b}_2^m, \quad (3)$$

where $f(\cdot)$ is an activation function such as ReLU [19] and GELU [11], $W_1^m, W_2^m \in \mathbb{R}^{d \times d}$ and $\mathbf{b}_1^m, \mathbf{b}_2^m$ are learnable parameters. It is worth noting that different from SA-based methods, MSSG updates only the embedding of the internal node across blocks, while fixing the item node embeddings. As a result, MSSG could preserve the information on individual items and will not be affected by the over-smoothing issue, even when learning deep models (Section 6.3). Following SASRec [14], we connect each layer with its previous layer using residual connections [8] to mitigate the degradation issue [26] during training. Recent studies [17, 22] show that the most recent interacted item could be a strong indicator of the next item of users' interest. Thus, for each sequence, MSSG employs \mathbf{e}_n , the embedding of the last item in E (Equation 1), as the initial embedding of the internal node in the first block (i.e., $\mathbf{c}^0 = \mathbf{e}_n$).

4.3 Comprehensive User Intent Modeling

As demonstrated in the literature [18, 22], both users' short-term preferences and long-term preferences play important roles in generating accurate recommendations. By considering \mathbf{e}_n , the embedding of the most recent interacted item, as the initial embedding of the internal node, the internal node embedding should be able to effectively capture the user's short-term preferences across blocks. However, it may not also fully capture the user's long-term preferences as the long-term preferences of the user could be different from her short-term preferences. Thus, to enable comprehensive user intent modeling, MSSG learns user embeddings to specifically capture users' long-term preferences. In particular, MSSG learns an

embedding matrix $U \in \mathbb{R}^{|\mathbb{U}| \times d}$ to represent the long-term preferences of all the users. The i -th row \mathbf{u}_i in U represents the long-term preferences of user u_i .

4.4 Recommendation Scores and Network Training

To generate recommendation scores for user u_i , MSSG integrates the embedding of the internal node from the last block (i.e., $\mathbf{c}_i^{n_b}$) and the learned user embedding \mathbf{u}_i , as follows:

$$\hat{r}_{ij} = (\mathbf{c}_i^{n_b} + \mathbf{u}_i) \mathbf{v}_j^\top, \quad (4)$$

where \hat{r}_{ij} is the recommendation score of user u_i on item v_j , and n_b is the total number of blocks in MSSG. For each user, we recommend items with the top- k highest recommendation scores. Following SASRec, we employ the binary cross-entropy loss to minimize the negative log-likelihood of correctly recommending the ground truth next item as follows:

$$\min_{\Theta} - \sum_{S_i \in \mathbb{T}, v_j \notin S_i} [\log(\sigma(\hat{r}_{ij})) + \log(1 - \sigma(\hat{r}_{ij}))], \quad (5)$$

where \mathbb{T} is the set of all the training sequences; Θ is the set of learnable parameters (e.g., V and P); and \hat{r}_{ij} is the recommendation score of user u_i on her ground-truth next item $v_j(i)$. For each training sequence, we randomly sample one negative item v_j for optimization. All the learnable parameters are randomly initialized, and are optimized in an end-to-end manner.

4.5 Complexity Analysis

Previous work [30? ?] shows that the time complexity of each SA block is $O(n^2 d + nd^2)$, where n is the length of the transformed sequence (Equation 1) and d is the dimension of embeddings. The quadratic time complexity limits the utility of SA-based methods in latency-sensitive recommendation applications. In contrast, MSSG models sequences as star graphs and could achieve a time complexity of $O(nd^2)$, which is linear with respect to the sequence length n . This allows MSSG to be theoretically more efficient than SA-based methods. Particularly, each SA block requires $6nd^2 + 2n^2 d + 2nd$ operations to compute, whereas each block in MSSG requires only $2nd^2 + 4d^2 + 2nd + 2d$ operations. The difference amounts to $4d^2(n-1) + 2d(n^2-1)$, which is quadratic with respect to both d and n . As will be shown in Section 6.5, the better time complexity could translate into superior run-time performance of MSSG over SA-based methods on modern GPUs.

5 MATERIALS

5.1 Baseline Methods

We compare MSSG with ten state-of-the-art baseline methods. Specifically, we compare MSSG with the MC-based method FPMC [25] and the CNNs-based method Caser [29]. We also compare MSSG with the RNNs-based method NARM [15] and two shallow methods HAM [22] and HGN [18]. We compare MSSG with three state-of-the-art SA-based methods SASRec [14], Bert4Rec [28] and FMLP [37]. We further compare MSSG with SASRec equipped with two hierarchical fusion strategies: 1) SASRec with ConcatFusion (SASRec-Cat); and 2) SASRec with MaxFusion (SASRec-Max). These two fusion strategies are developed in Shi *et al.* [27] to alleviate over-smoothing. Note that, we do not consider baseline methods that the implementation

is not publicly available (e.g., SGNN-HN) to enable a fair comparison. HGN and SASRec have been compared with a comprehensive set of other methods including GRU4REC [12] and NextItRec [35], and have outperformed those methods. Therefore, we compare MSSG with HGN and SASRec instead of the methods that they outperform.

To ensure a fair comparison, we exhaustively tune the hyper-parameters of all the baseline methods. Thus, the performance of baseline methods reported in our experiments is generally better than that reported in the literature [18, 23]. For example, on the Beauty dataset (Section 5.2), the Recall@10 (Section 5.3.2) of SASRec reported in our experiments is 20.3% higher than that reported in a recent work [23]. We report the implementation detail of MSSG and baseline methods, and the search range for each hyper-parameter in Appendix B.

5.2 Datasets

We evaluate MSSG and baseline methods on six public benchmark datasets: 1) Amazon-Beauty (Beauty) and Amazon-Toys (Toys) are from Amazon reviews [20]; 2) Goodreads-Children (Children) and Goodreads-Comics (Comics) are from the Goodreads website [? ?]; and 3) MovieLens-1M (ML-1M) and MovieLens-20M (ML-20M) are from the MovieLens website [7]. We discuss the dataset pre-processing and statistics in Appendix C.

5.3 Experimental Protocol

5.3.1 Training, validation and testing sets. Following SASRec, on all the datasets, given the historical interactions S of each user, we use the last item (i.e., $v_{s|s}$) in the history for testing, the second last item (i.e., $v_{s|s-1}$) for validation, and all the previous items for training. We tune hyper-parameters on the validation sets for all the methods using grid search, and use the best-performing hyper-parameters in terms of Recall@10 (Section 5.3.2) for testing.

5.3.2 Evaluation metrics. Following the literature [14, 18, 22, 29], we use Recall@ k and NDCG@ k to evaluate MSSG and the baseline methods. We refer the audience of interest to Peng *et al.* [22] for the detailed definitions of both Recall@ k and NDCG@ k .

6 EXPERIMENTAL RESULTS

6.1 Overall Performance

Table 1 presents the overall performance of MSSG, its variant MSSG-u and the state-of-the-art baseline methods at Recall@ k and NDCG@ k on six benchmark datasets. In MSSG-u, we remove user embeddings (i.e., u_i) when calculating recommendation scores (Equation 4). We introduce MSSG-u to evaluate the effectiveness of the learnable user embeddings in MSSG. For NARM, on ML-20M, with the implementation provided by the authors, we get the out-of-memory (OOM) issue on NVIDIA Volta V100 GPUs with 16 GB memory.

As shown in Table 1, overall, MSSG is the best-performing method on the six datasets. In terms of Recall@10 and Recall@20, MSSG significantly outperforms all the baseline methods on five out of the six datasets except ML-20M. We observe a similar trend at both NDCG@10 and NDCG@20. We notice that on ML-20M, MSSG considerably underperforms SASRec. However, without user embeddings, MSSG-u could still significantly outperform SASRec on

ML-20M. These results demonstrate the strong performance of MSSG and its variant MSSG-u on different recommendation scenarios.

6.1.1 Comparing MSSG to FPMC and Caser. As presented in Table 1, MSSG outperforms the MC-based method FPMC and CNNs-based method Caser by a significant margin on all the six benchmark datasets. FPMC views the transitions among items as an MC and leverages only the most recent interacted items to predict the next item of users' interest. Consequently, FPMC may not be able to fully utilize the information in early interacted items, and result in sub-optimal performance. Caser employs CNNs to integrate items within sequences. Though CNNs have been shown effective in capturing local structures, they may struggle to capture global information within sequences [33]. In contrast, by modeling sequences as star graphs, MSSG could effectively aggregate information from all items to capture the global information, and thus, enable superior recommendation performance over FPMC and Caser.

6.1.2 Comparison between MSSG and NARM. Table 1 shows that compared to the RNNs-based method NARM, MSSG demonstrates superior performance on all the six datasets at all the evaluation metrics. NARM primarily uses RNNs to recurrently learn users' preferences. As demonstrated in the literature [31], due to the recurrent nature, RNNs may struggle to model long-range dependencies within sequences, which could limit their ability to capture global information. In MSSG, as illustrated in Figure 1c, the central internal node could aggregate information from every item node within one step. Thus, compared to NARM, MSSG could better capture long-range dependencies within sequences, and achieve substantial performance improvement on all the datasets.

6.1.3 Comparison between MSSG and HAM. Table 1 also presents that MSSG significantly outperforms the best-performing shallow method HAM on all the six datasets. Compared to HAM, MSSG achieves a remarkable average improvement of 7.8% and 6.0% at Recall@10 and NDCG@10, respectively, across the six datasets. HAM employs the simple mean pooling to aggregate items and estimate users' intent. Though efficient, this simple approach may not be sufficiently expressive to accurately estimate users' intent from their diverse interactions. Different from HAM, MSSG stacks multiple non-linear blocks to enable expressive models, which leads to a more accurate user intent modeling as compared to HAM.

6.1.4 Comparison between MSSG and SASRec. Table 1 presents that both MSSG and MSSG-u substantially outperform the best-performing SA-based method SASRec. Similarly to MSSG, SASRec stacks multiple SA blocks to recursively capture users' intent, and has been demonstrated state-of-the-art performance [14, 22] in SR. However, as shown in Table 1, MSSG consistently outperforms SASRec on five out of six datasets (i.e., Beauty, Toys, Children, Comics and ML-1M). Particularly, compared to SASRec, in terms of Recall@10, MSSG achieves a significant average improvement of 14.2% over the five datasets. In terms of NDCG@10, MSSG also outperforms SASRec with a significant average improvement of 15.1% over the five datasets. Similarly, MSSG-u also remarkably outperforms SASRec across the six datasets with an average improvement of 7.0% and 7.2% at Recall@10 and NDCG@10, respectively. From a graph perspective, SASRec models sequences as fully connected graphs, which allow each item to aggregate information from other items

Table 1: Overall Performance

Dataset	Metric	FPMC	Caser	NARM	HGN	HAM	SASRec	SASRec-Cat	SASRec-Max	Bert4Rec	FMLP	MSSG-u	MSSG	imprv (%)
Beauty	Recall@10	0.0750	0.0647	0.0554	0.0782	<u>0.0865</u>	0.0728	0.0704	0.0741	0.0549	0.0622	0.0874	0.0923	6.71*
	Recall@20	0.1086	0.0869	0.0835	0.1121	<u>0.1221</u>	0.1083	0.1024	0.1106	0.0776	0.0952	0.1249	0.1286	5.32*
	NDCG@10	0.0430	0.0371	0.0291	0.0434	<u>0.0493</u>	0.0398	0.0390	0.0397	0.0298	0.0330	0.0486	0.0516	4.67*
	NDCG@20	0.0514	0.0427	0.0362	0.0520	<u>0.0582</u>	0.0487	0.0470	0.0489	0.0355	0.0413	0.0581	0.0607	4.30*
Toys	Recall@10	0.0869	0.0675	0.0557	0.0926	<u>0.1010</u>	0.0909	0.0843	0.0903	0.0539	0.0787	0.1031	0.1055	4.46*
	Recall@20	0.1173	0.0895	0.0816	0.1240	<u>0.1303</u>	0.1278	0.1123	0.1240	0.0756	0.1072	0.1392	0.1434	10.10*
	NDCG@10	0.0517	0.0396	0.0308	0.0540	<u>0.0620</u>	0.0521	0.0496	0.0515	0.0290	0.0444	0.0604	0.0625	0.81
	NDCG@20	0.0594	0.0452	0.0373	0.0619	<u>0.0694</u>	0.0614	0.0566	0.0600	0.0344	0.0516	0.0696	0.0720	3.75*
Children	Recall@10	0.1298	0.1303	0.1126	0.1536	0.1729	<u>0.1752</u>	0.1471	0.1706	0.1435	0.1675	0.1709	0.1924	9.82*
	Recall@20	0.1843	0.1832	0.1735	0.2118	<u>0.2366</u>	0.2337	0.2064	0.2320	0.1930	0.2249	0.2334	0.2589	9.43*
	NDCG@10	0.0763	0.0764	0.0585	0.0917	<u>0.1038</u>	0.1086	0.0861	0.1037	0.0884	0.1040	0.1038	0.1176	8.29*
	NDCG@20	0.0900	0.0897	0.0738	0.1063	0.1199	<u>0.1234</u>	0.1010	0.1191	0.1009	0.1184	0.1195	0.1343	8.83*
Comics	Recall@10	0.2517	0.2320	0.1304	0.2857	0.3055	<u>0.3196</u>	0.2975	0.3144	0.2363	0.3163	0.3193	0.3317	3.79*
	Recall@20	0.2986	0.2791	0.1896	0.3322	0.3513	<u>0.3650</u>	0.3456	0.3649	0.2846	0.3631	0.3680	0.3833	5.01*
	NDCG@10	0.1875	0.1642	0.0720	0.2061	0.2319	<u>0.2430</u>	0.2219	0.2327	0.1478	0.2377	0.2435	0.2485	2.26*
	NDCG@20	0.1993	0.1760	0.0869	0.2178	0.2435	<u>0.2545</u>	0.2340	0.2454	0.1600	0.2496	0.2558	0.2615	2.75*
ML-1M	Recall@10	0.1614	<u>0.2874</u>	0.2349	0.2428	0.2745	<u>0.2623</u>	0.2397	0.2598	0.1611	0.2409	0.2866	0.3000	4.38*
	Recall@20	0.2298	<u>0.3955</u>	0.3422	0.3439	0.3707	<u>0.3709</u>	0.3455	0.3737	0.2291	0.3571	0.3993	0.4132	4.48*
	NDCG@10	0.0841	<u>0.1619</u>	0.1252	0.1374	0.1576	0.1463	0.1296	0.1349	0.0832	0.1261	0.1582	0.1684	4.01*
	NDCG@20	0.1013	<u>0.1892</u>	0.1524	0.1629	0.1818	0.1681	0.1561	0.1636	0.1003	0.1553	0.1867	0.1970	4.12*
ML-20M	Recall@10	0.0992	0.1739	OOM	0.1588	0.1673	<u>0.1922</u>	0.1859	0.1915	0.1106	0.1707	0.1957	0.1786	1.82*
	Recall@20	0.1622	0.2649	OOM	0.2390	0.2493	<u>0.2919</u>	0.2844	0.2909	0.1866	0.2672	0.2959	0.2753	1.37*
	NDCG@10	0.0478	0.0907	OOM	0.0845	0.0895	<u>0.1004</u>	0.0967	0.1002	0.0515	0.0862	0.1015	0.0924	1.10*
	NDCG@20	0.0636	0.1136	OOM	0.1047	0.1102	<u>0.1255</u>	0.1214	0.1251	0.0706	0.1105	0.1267	0.1167	0.96*

For each dataset, the best performance in MSSG-u and MSSG is in **bold**, and the best performance among the baseline methods is underlined. The column "imprv" presents the percentage improvement of MSSG or MSSG-u over the best-performing baseline methods. The * indicates that the improvement is statistically significant at 95% confidence level.

within the sequence. Though promising, this design results in quadratic time complexities and leads to the over-smoothing issue (Section 6.3). As shown in the literature [2, 34], this issue could substantially limit the model scalability and degrade the recommendation performance. In contrast, by modeling sequences as star graphs, MSSG could fundamentally address the over-smoothing issue while still being able to capture the long-range dependencies within sequences. Consequently, as presented in Table 1, MSSG could enable significant improvement over SASRec on different recommendation scenarios.

6.1.5 Comparison between MSSG and SASRec-Max. Table 1 also shows that overall, both MSSG and MSSG-u achieve superior performance over SASRec-Max on the six datasets. For example, at Recall@10, MSSG and MSSG-u achieves an average improvement of 11.4% and 7.7%, respectively, over the six datasets compared to SASRec-Max. SASRec-Max mitigates over-smoothing by fusing embeddings from both shallow layers and deep layers using max pooling. In contrast, MSSG fundamentally addresses over-smoothing by modeling sequences using star graphs, thereby avoiding information propagation among items. The substantial improvement of MSSG over SASRec-Max shows that our approach is more effective than hierarchical fusion [27] in both addressing over-smoothing, and improving recommendation performance.

6.1.6 Performance Summary across Datasets. Table 2 shows the average improvement of MSSG-u and MSSG over the baseline methods Caser, HAM and SASRec across the six datasets. We focus on these baseline methods since they achieve the best performance in terms of at least one evaluation metric on the datasets. As shown in Table 2, both MSSG-u and MSSG significantly outperform the baseline methods across the six datasets. For example, MSSG-u and MSSG

Table 2: Performance Improvement of MSSG-u and MSSG (%)

Method	Metric	Caser	HAM	SASRec
MSSG-u	Recall@10	28.2*	4.6*	7.0*
	Recall@20	28.5*	6.5*	5.7*
	NDCG@10	29.6*	2.5	7.2*
	NDCG@20	29.8*	3.8*	7.0*
MSSG	Recall@10	32.8*	7.8*	10.6*
	Recall@20	32.5*	9.3*	8.7*
	NDCG@10	34.7*	6.0*	11.2*
	NDCG@20	34.4*	7.0*	10.6*

In this table, the column Caser, HAM and SASRec represents the percentage improvement of MSSG-u and MSSG over the corresponding method. The * indicates that the improvement is statistically significant at 85% confidence level.

achieves a significant average improvement of 7.0% and 10.6%, respectively, over SASRec at Recall@10 across the six datasets.

6.2 Comparison on Users with Different Activity Levels

We also compare the performance of MSSG and MSSG-u with the best-performing SA-based SR method SASRec on users with different activity levels. Specifically, we quantify users' activity levels using the number of interactions in their history, and bin users into different buckets based on their activity levels. In this analysis, we use ten buckets in total: five buckets correspond to the top-10%, ..., top 40-50% most active users; the other five buckets are for the bottom 40-50%, ..., bottom-10% most active users. Figure 3 presents the performance of MSSG, MSSG-u and SASRec on Beauty, Toys and Children. We present the results on Comics, ML-1M and ML-20M in Appendix D.1 (Figure A1). The results on Comics, ML-1M and ML-20M have a similar trend with that shown in Figure 3.

As shown in Figure 3, MSSG outperforms SASRec by a remarkable margin on both active and less active users across the three datasets. Similarly, MSSG-u substantially outperforms SASRec on

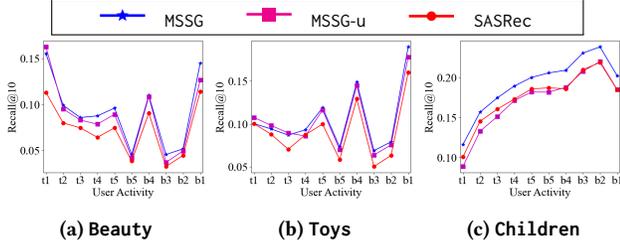


Figure 3: Performance on users of different activity levels

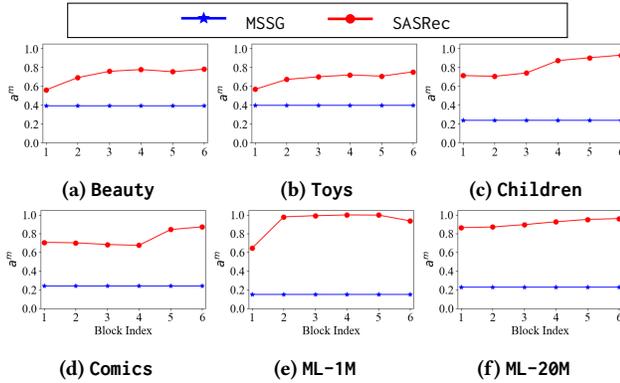


Figure 4: The average similarity a^m in different blocks of MSSG and SASRec

users of different activity levels on Beauty and Toys, and achieves highly competitive performance with SASRec on Children. These results signify that by modeling sequences using star graphs, MSSG and MSSG-u could more accurately estimate the intent of users with different activity levels compared to SASRec. Note that users’ activity level is not the only factor in determining the difficulty of estimating their intent. As a result, we do not expect a strictly increasing or decreasing performance curve as activity levels change.

6.3 Analysis on Over-smoothing

We analyze if the state-of-the-art SA-based SR method SASRec suffers from the over-smoothing issue. Specifically, for each dataset, we train a SASRec model consisting of six SA blocks. Apart from the number of blocks, we use the best-performing hyper-parameters of SASRec on each dataset to train the model. Subsequently, we calculate the average similarity between the embeddings of items within a sequence for each SA block. We denote the average similarity calculated from the output of the m -th SA block as a^m . We discuss the details on the calculation of a^m in Appendix D.2.

We use a^m to assess whether the SASRec model suffers from the over-smoothing issue. In particular, the increase of a^m over m implies that generally, the embeddings of items within a user’s interaction sequence become more and more similar across blocks. This serves as strong evidence to demonstrate that the SASRec model suffers from the over-smoothing issue. We also evaluate if MSSG suffers from over-smoothing by calculating a^m from the MSSG model trained using the same hyper-parameters as in SASRec. We present a^m calculated from MSSG and SASRec in Figure 4. Note that, we observe a similar trend to that in Figure 4 when using a smaller number of blocks. Due to the space limit, we present results for using six blocks only.

As illustrated in Figure 4, in SASRec, a^m increases across blocks on all the six datasets. For example, on Beauty, the average similarity rises from 0.56 (i.e., a^1) in the first block to 0.78 (i.e., a^6) in the last block. These results demonstrate that SASRec substantially suffers from the over-smoothing issue. In contrast, as shown in Figure 4, a^m from MSSG remains constant across blocks, indicating that MSSG is not affected by over-smoothing. The key distinction between MSSG and SASRec lies in their sequence modeling approaches. While SASRec utilizes fully connected graphs to model sequences, MSSG instead models sequences using star graphs. As a result, MSSG does not require information propagation among item nodes, and allows each item to maintain its individual information across blocks. Consequently, MSSG could fundamentally address the over-smoothing issue. As shown in the literature [27, 32], the over-smoothing issue could substantially deteriorate the model scalability, and thus, degrade the model performance in the task of interest. Therefore, by addressing this issue, MSSG could achieve superior performance over SASRec on benchmark datasets as shown in Table 1.

6.4 Analysis on Scalability

The huge success of Transformer [31] highlights the importance of a method’s scalability in determining its utility for real-world applications. In this analysis, we evaluate the scalability of MSSG and SASRec with respect to the number of blocks and the embedding dimensions. Particularly, we assess the ability of MSSG and SASRec in learning deep (i.e., with a large number of blocks) and wide (i.e., with a large embedding dimension) models on benchmark datasets. To enable a fair comparison, in this analysis, we apply the best-performing hyper-parameters of SASRec on each dataset for SASRec and MSSG.

Figure 5 presents the performance of SASRec and MSSG at Recall@10 over different numbers of blocks on the six benchmark datasets. We limit the maximum number of blocks (n_b) to six, as we got the out-of-memory issue with SASRec when $n_b > 6$.

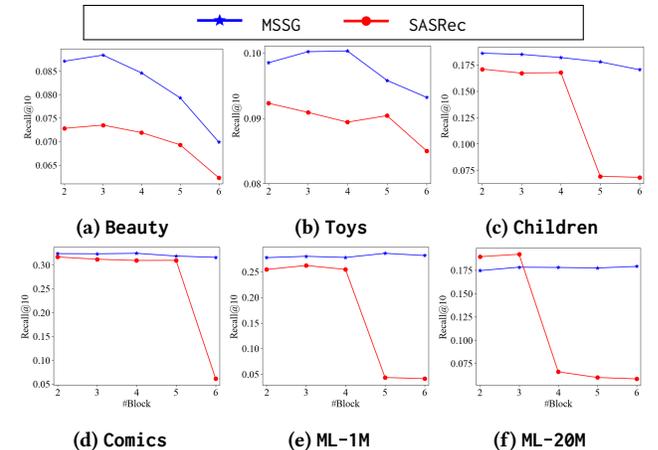


Figure 5: Performance over different numbers of blocks

Figure 5 demonstrates that MSSG outperforms SASRec at the scalability with respect to n_b . In particular, MSSG maintains reasonable performance on all datasets, while SASRec fails on four out of the six datasets (Children, Comics, ML-1M, and ML-20M) when $n_b=6$. As discussed in Section 6.3, the state-of-the-art

SA-based SR method SASRec substantially suffers from the over-smoothing issue. Consequently, SASRec struggles to learn deep models, as the item embeddings become increasingly similar across SA blocks [27, 32, 36]. In contrast, as shown in Section 6.3, MSSG is not affected by over-smoothing, and thus, could enable superior scalability with respect to n_b , compared to SASRec. It is worth noting that, on ML-1M and ML-20M, the best performing MSSG-u and MSSG models have at least 4 blocks (Appendix B), which reveals that the better scalability of MSSG at n_b could translate into superior recommendation performance on real datasets.

Due to the space limit, we present the scalability comparison between MSSG and SASRec in terms of embedding dimensions (i.e., d) in Appendix D.3. As shown in Appendix D.3 (Figure A2), MSSG also substantially outperforms SASRec at the scalability over d . We refer the audience to Appendix D.3 for a more detailed discussion.

6.5 Comparison on Run-time Performance

We conduct an analysis to evaluate the run-time performance of MSSG and SASRec during testing. To enable a fair comparison, similarly to that in Section 6.4, we apply the best-performing hyper-parameters on SASRec for MSSG and SASRec, and compare their run-time performance during testing when using different numbers of blocks. In addition, we perform the evaluation for both MSSG and SASRec using NVIDIA Volta V100 GPUs, and report the average computation time per user over five runs in Figure 6. Moreover, we use the same evaluation script for both MSSG and SASRec to avoid any run-time performance differences that might arise from differences in the implementation. We focus on the run-time performance during testing due to the fact that it could signify the models' latency in real-time recommendation, which could significantly affect the user experience and thus revenue.

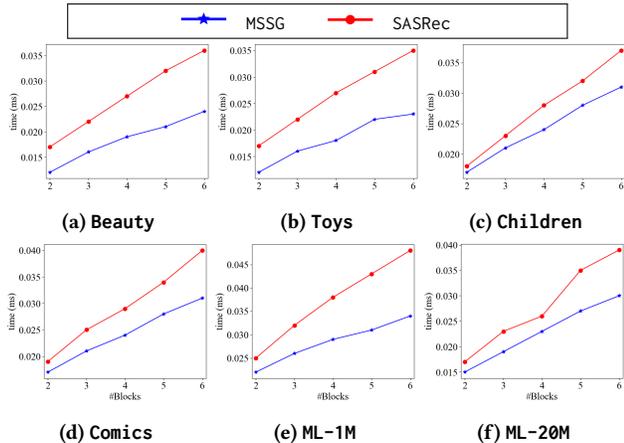


Figure 6: Run-time performance over the number of blocks

As shown in Figure 6, on all the datasets, MSSG achieves stronger run-time performance compared to SASRec over different numbers of blocks, and the improvement increases as the number of blocks increases. Particularly, when using the best-performing n_b of SASRec on each dataset (e.g., $n_b = 2$ on Beauty as in Appendix B), MSSG achieves an average speedup of 16.6% compared to SASRec over the six datasets. Note that, as shown in Figure 5, with the same hyper-parameters, MSSG also outperforms SASRec in terms of the

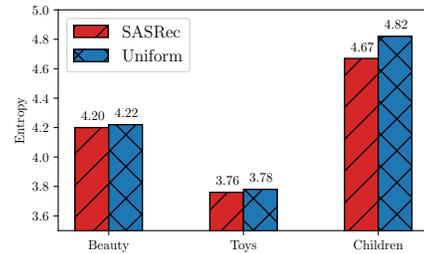


Figure 7: The average entropy of attention weight distribution in SASRec and uniform distribution

recommendation performance on all the datasets except ML-20M. These results demonstrate that compared to SASRec, MSSG could generate more accurate recommendations in lower latency, thus significantly improving the user experience. Note that on GPUs, all the computations are performed in parallel. However, when calculating the time complexity (Section 4.5), we assume the computations are serial. Therefore, in terms of the run-time performance on GPUs, the improvement of MSSG over SASRec may not be as significant as that suggested by the time complexity. However, in many applications, the recommendation model could be deployed on edge devices with limited computing resources. In these applications, as suggested by the time complexity comparison (Section 4.5), MSSG could achieve a more substantial speedup over SASRec.

6.6 Analysis on Attention Weights in SASRec

We conduct an analysis to investigate the information gain derived from the attention weights learned in SASRec. Particularly, we measure the information gain by comparing the average entropy from the attention weight distributions and that from uniform distributions (i.e., weigh items equally). We present more details on the calculation of the average entropy in Appendix D.4. A higher difference between the average entropies indicates a larger information gain derived from the learned attention weights. Figure 7 shows the average entropy from attention weight distributions and uniform distributions on Beauty, Toys and Children. Following SASRec, we focus on the attention weights learned in the first SA block in this analysis.

As shown in Figure 7, on all the three datasets, the difference between the average entropy from attention weight distributions and uniform distributions is highly marginal, indicating that SASRec achieves limited information gain from the learned attention weights. An important difference between MSSG and SASRec is that SASRec captures the dependencies in each pair of items by learning attention weights, while MSSG does not explicitly model these dependencies. As shown in Figure 7, on the notoriously sparse recommendation data [22], limited information gain could be achieved by capturing the dependency in each item pair. As a result, as shown in Table 1 and Table 2, without explicitly modeling dependencies in item pairs, MSSG and MSSG-u could still achieve significant improvement over SASRec on benchmark datasets.

7 CONCLUSION

In this paper, we identify the over-smoothing issue in SA-based SR methods. To address this issue, we model sequences using star graphs and develop MSSG for SR. Different from SA-based methods

in which each item could aggregate information from all others, MSSG introduces an internal node for information aggregation and does not propagate information among item nodes. Consequently, MSSG could fundamentally address the over-smoothing issue and achieve a linear time complexity with respect to the sequence length. We extensively evaluate MSSG against ten state-of-the-art baseline methods on six benchmark datasets. Our experimental results show that overall, MSSG outperforms baseline methods on all the datasets except for ML-20M, with an improvement of up to 10.10%. On ML-20M, a variant of MSSG could still significantly outperform all the baseline methods. Our analysis shows that MSSG achieves superior scalability with respect to both the number of blocks and the embedding dimensions compared to the state-of-the-art SA-based SR method SASRec. This improved scalability could lead to enhanced recommendation performance in benchmark datasets. In addition, our complexity analysis and run-time performance comparison together demonstrate that MSSG is both theoretically and practically more efficient than SASRec. Thus, MSSG could be particularly desirable in applications with limited computing resources. Our analysis also suggests that SASRec achieves limited information gain by explicitly modeling dependencies between items using the SA mechanism. Therefore, without modeling the dependencies, MSSG could still outperform SASRec in benchmark datasets.

REFERENCES

- [1] Francois Belletti, Minmin Chen, and Ed H Chi. 2019. Quantifying long range dependence in language and user behavior to improve rnnns. In *Proceedings of the 25th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining*. 1317–1327.
- [2] Chen Cai and Yusu Wang. 2020. A note on over-smoothing for graph neural networks. *arXiv* (2020).
- [3] Deli Chen, Yankai Lin, Wei Li, Peng Li, Jie Zhou, and Xu Sun. 2020. Measuring and relieving the over-smoothing problem for graph neural networks from the topological view. In *Proceedings of the AAAI conference on artificial intelligence*, Vol. 34. 3438–3445.
- [4] Vijay Prakash Dwivedi and Xavier Bresson. 2020. A generalization of transformer networks to graphs. *arXiv* (2020).
- [5] Ziwei Fan, Zhiwei Liu, Yu Wang, Alice Wang, Zahra Nazari, Lei Zheng, Hao Peng, and Philip S Yu. 2022. Sequential Recommendation via Stochastic Self-Attention. *arXiv* (2022).
- [6] Will Hamilton, Zhitao Ying, and Jure Leskovec. 2017. Inductive representation learning on large graphs. *Advances in neural information processing systems* 30 (2017).
- [7] F. Maxwell Harper and Joseph A. Konstan. 2015. The MovieLens Datasets: History and Context. *ACM Trans. Interact. Intell. Syst.* 5, 4, Article 19 (dec 2015), 19 pages. <https://doi.org/10.1145/2827872>
- [8] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. 2016. Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*. 770–778.
- [9] Ruining He and Julian McAuley. 2016. Ups and downs: Modeling the visual evolution of fashion trends with one-class collaborative filtering. In *proceedings of the 25th international conference on world wide web*. 507–517.
- [10] Zhankui He, Handong Zhao, Zhe Lin, Zhaowen Wang, Ajinkya Kale, and Julian McAuley. 2021. Locker: Locally constrained self-attentive sequential recommendation. In *Proceedings of the 30th ACM International Conference on Information & Knowledge Management*. 3088–3092.
- [11] Dan Hendrycks and Kevin Gimpel. 2016. Gaussian error linear units (gelus). *arXiv* (2016).
- [12] Balázs Hidasi and Alexandros Karatzoglou. 2018. Recurrent neural networks with top-k gains for session-based recommendations. In *Proceedings of the 27th ACM international conference on information and knowledge management*. 843–852.
- [13] Wenbing Huang, Yu Rong, Tingyang Xu, Fuchun Sun, and Junzhou Huang. 2020. Tackling over-smoothing for general graph convolutional networks. *arXiv* (2020).
- [14] Wang-Cheng Kang and Julian McAuley. 2018. Self-attentive sequential recommendation. In *International Conference on Data Mining*. 197–206.
- [15] Jing Li, Pengjie Ren, Zhumin Chen, Zhaochun Ren, Tao Lian, and Jun Ma. 2017. Neural attentive session-based recommendation. In *Proceedings of the Conference on Information and Knowledge Management*. 1419–1428.
- [16] Qimai Li, Zhichao Han, and Xiao-Ming Wu. 2018. Deeper insights into graph convolutional networks for semi-supervised learning. In *Proceedings of the AAAI conference on artificial intelligence*, Vol. 32.
- [17] Qiao Liu, Yifu Zeng, Refuoc Mokhosi, and Haibin Zhang. 2018. STAMP: Short-Term Attention/Memory Priority Model for Session-Based Recommendation. In *Proceedings of the 24th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining (London, United Kingdom) (KDD '18)*. Association for Computing Machinery, New York, NY, USA, 1831–1839. <https://doi.org/10.1145/3219819.3219950>
- [18] Chen Ma, Peng Kang, and Xue Liu. 2019. Hierarchical gating networks for sequential recommendation. In *Proceedings of the ACM SIGKDD international conference on knowledge discovery & data mining*. 825–833.
- [19] Vinod Nair and Geoffrey E Hinton. 2010. Rectified linear units improve restricted boltzmann machines. In *Icml*.
- [20] Jianmo Ni, Jiacheng Li, and Julian McAuley. 2019. Justifying Recommendations using Distantly-Labeled Reviews and Fine-Grained Aspects. In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*. Association for Computational Linguistics, Hong Kong, China, 188–197. <https://doi.org/10.18653/v1/D19-1018>
- [21] Zhiqiang Pan, Fei Cai, Wanyu Chen, Honghui Chen, and Maarten De Rijke. 2020. Star graph neural networks for session-based recommendation. In *Proceedings of the 29th ACM international conference on information & knowledge management*. 1195–1204.
- [22] Bo Peng, Zhiyun Ren, Srinivasan Parthasarathy, and Xia Ning. 2021. HAM: hybrid associations models for sequential recommendation. *IEEE Transactions on Knowledge and Data Engineering* (2021).
- [23] Shashank Rajput, Nikhil Mehta, Anima Singh, Raghunandan H Keshavan, Trung Vu, Lukasz Heldt, Lichan Hong, Yi Tay, Vinh Q Tran, Jonah Samost, et al. 2023. Recommender Systems with Generative Retrieval. *arXiv preprint arXiv:2305.05065* (2023).
- [24] Steffen Rendle, Christoph Freudenthaler, Zeno Gantner, and Lars Schmidt-Thieme. 2012. BPR: Bayesian personalized ranking from implicit feedback. *arXiv* (2012).
- [25] Steffen Rendle, Christoph Freudenthaler, and Lars Schmidt-Thieme. 2010. Factorizing personalized markov chains for next-basket recommendation. In *Proceedings of the international conference on World wide web*. 811–820.
- [26] Prasun Roy, Subhankar Ghosh, Saumik Bhattacharya, and Umapada Pal. 2018. Effects of degradations on deep neural network architectures. *arXiv* (2018).
- [27] Han Shi, Jiahui Gao, Hang Xu, Xiaodan Liang, Zhenguo Li, Lingpeng Kong, Stephen Lee, and James T Kwok. 2022. Revisiting Over-smoothing in BERT from the Perspective of Graph. *arXiv* (2022).
- [28] Fei Sun, Jun Liu, Jian Wu, Changhua Pei, Xiao Lin, Wenwu Ou, and Peng Jiang. 2019. BERT4Rec: Sequential recommendation with bidirectional encoder representations from transformer. In *Proceedings of the 28th ACM international conference on information and knowledge management*. 1441–1450.
- [29] Jiaxi Tang and Ke Wang. 2018. Personalized top-n sequential recommendation via convolutional sequence embedding. In *Proceedings of the international conference on web search and data mining*. 565–573.
- [30] Yi Tay, Mostafa Dehghani, Dara Bahri, and Donald Metzler. 2020. Efficient transformers: A survey. *arXiv* (2020).
- [31] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. 2017. Attention is all you need. *Advances in neural information processing systems* 30 (2017).
- [32] Keyulu Xu, Chengtao Li, Yonglong Tian, Tomohiro Sonobe, Ken-ichi Kawarabayashi, and Stefanie Jegelka. 2018. Representation learning on graphs with jumping knowledge networks. In *International Conference on Machine Learning*. 5453–5462.
- [33] Rikiya Yamashita, Mizuho Nishio, Richard Kinh Gian Do, and Kaori Togashi. 2018. Convolutional neural networks: an overview and application in radiology. *Insights into imaging* 9 (2018), 611–629.
- [34] Chaoqi Yang, Ruijie Wang, Shuochao Yao, Shengzhong Liu, and Tarek Abdelzaher. 2020. Revisiting over-smoothing in deep GCNs. *arXiv* (2020).
- [35] Fajie Yuan, Alexandros Karatzoglou, Ioannis Arapakis, Joemon M Jose, and Xi-anngnan He. 2019. A Simple Convolutional Generative Network for Next Item Recommendation. In *Proceedings of the Twelfth ACM International Conference on Web Search and Data Mining*. ACM.
- [36] Lingxiao Zhao and Leman Akoglu. 2019. Pairnorm: Tackling oversmoothing in gnns. *arXiv* (2019).
- [37] Kun Zhou, Hui Yu, Wayne Xin Zhao, and Ji-Rong Wen. 2022. Filter-enhanced MLP is all you need for sequential recommendation. In *Proceedings of the ACM Web Conference 2022*. 2388–2399.

A NOTATIONS

Table A1: Key Notations

notations	meanings
\mathbb{U}	the set of users
\mathbb{V}	the set of items
S_i	the historical interaction sequence of user u_i
$v_{s_t}(i)$	the t -th item in S_i
B_i	the fixed-length sequence converted from S_i
$v_{b_t}(i)$	the t -th item in B_i
$v_g(i)$	the ground-truth next item that user u_i will interact with

Table A1 summarizes the key notations used in this paper.

B REPRODUCIBILITY

We implement MSSG and MSSG-u in Python 3.9.13 with PyTorch 1.10.2. We use Adam optimizer with learning rate 1e-3 for MSSG and MSSG-u on all the datasets. For all the baseline methods except FPMC and Bert4Rec, we use the implementation provided by the authors in GitHub. For FPMC and Bert4Rec, we use the implementation in RecBole³, a widely used library to benchmark recommendation methods. For SASRec, FMLP, MSSG-u and MSSG, we search the embedding dimension d in {64, 128, 256, 512}, the length of the fixed-length sequences n in {50, 75, 100, 125, 150, 175, 200}, the number of heads n_h in {1, 2, 4, 8, 16}, and the number of blocks n_b in {1, 2, 3, 4, 5}. For Bert4Rec, we search d in {64, 128, 256, 512}, n_h in {1, 2, 4, 8, 16} and n_b in {1, 2, 3, 4, 5}. For SASRec-Cat and SASRec-Max, the search range of d , n and n_h is the same with that in SASRec. We search n_b in SASRec-Cat and SASRec-Max in {2, 3, 4, 5} as SASRec-Cat and SASRec-Max are equivalent to SASRec when $n_b=1$. We use GELU [11] as the activation function in the feed-forward layer for SASRec, SASRec-Cat, SASRec-Max, MSSG-u and MSSG. We use the PyTorch implementation in GitHub⁴ for SASRec. For FPMC, we search d in {64, 128, 256, 512}. For Caser we search d in {64, 128, 256, 512}, the length of the subsequences n_s in {4, 5, 6}, the number of negative items during training n_p in {1, 2}, the number of vertical filters in CNNs n_v in {1, 2, 4}, and the number of horizontal filters in CNNs n_f in {4, 8, 16}. For NARM we search d in {64, 128, 256, 512} and the learning rate l_r in {1e-2, 1e-3, 1e-4}. For HGN we search d in {64, 128, 256, 512}, n_s in {3, 4, 5}, n_p in {1, 2}, and the regularization factor λ in {0, 1e-3, 1e-4}. For HAM we search d in {64, 128, 256, 512}, n_s in {3, 4, 5}, n_p in {1, 2}, λ in {0, 1e-3, 1e-4}, the number of items in low order n_l in {1, 2, 3}, and the order of item synergies n_o in {1, 2, 3}. Following the instruction in RecBole, we use the Bayesian personalized ranking loss [24] and cross-entropy loss for FPMC and Bert4Rec, respectively. We report the best-performing hyper-parameters of MSSG, MSSG-u and baseline methods in Table A2.

C DATASET PRE-PROCESSING AND STATISTICS

Following SASRec, for Beauty, Toys and ML-1M, we only keep the users and items with at least 5 ratings. For Children, Comics and ML-20M, which are not used in SASRec, following HAM, we keep users with at least 10 ratings, and items with at least 5 ratings. Following

³<https://recbole.io/>

⁴<https://github.com/pmixer/SASRec.pytorch>

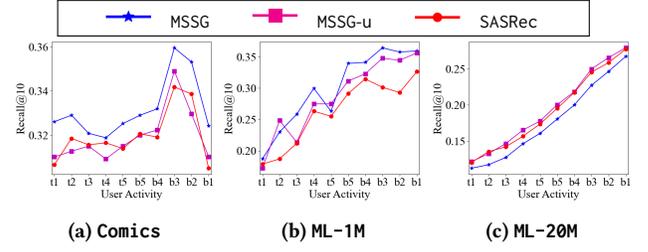


Figure A1: Performance on users of different activity levels

the literature [14, 22], we consider the ratings as users’ implicit feedback, and convert ratings into binary values. Particularly, for ratings with a range from 1 to 5, we convert ratings 4 and 5 to binary value 1, or 0 otherwise. Table A3 presents the statistics of the processed datasets.

D MORE EXPERIMENTAL RESULTS AND DETAILS

D.1 Comparison on Users with Different Activity Levels (Cont.)

Figure A1 shows the performance of MSSG, MSSG-u and SASRec in users of different activity levels on the Comics, ML-1M and ML-20M datasets. As shown in Figure A1, the performance on Comics, ML-1M and ML-20M has a similar trend with that on Beauty, Toys and Children (Figure 3 in Section 6.2). Particularly, MSSG substantially outperforms SASRec on Comics and ML-1M in both active and less active users. Similarly, in most activity levels, MSSG-u considerably outperforms SASRec on ML-1M and ML-20M. On Comics, MSSG-u also achieves highly comparable performance with SASRec.

D.2 Calculating Average Similarities Between Item embeddings

Given the transformed sequence of user u_i : $B_i = \{v_{b_1}(i), v_{b_2}(i), \dots\}$ (Section 4.1), we calculate the average similarity of the embeddings of items within a sequence for each SA block as follows:

$$a_i^m = \frac{1}{|B_i|^2} \sum_{j=1}^{|B_i|} \sum_{k=1}^{|B_i|} \cos(\mathbf{e}_j^m, \mathbf{e}_k^m), \quad (6)$$

where $|B_i|$ is the length of B_i ; \mathbf{e}_j^m is the embedding of the j -th item in B_i after the m -th SA block; $\cos(\cdot)$ is the cosine similarity; a_i^m is the average similarity between the embeddings of items in B_i after the m -th SA block. Note that as illustrated in Figure 1b, SASRec updates the embedding of each item within the sequence in each SA block. Thus, we expect different embeddings for the same item in different blocks. We eliminate i in $\mathbf{e}_j^m(i)$ in Equation 6 for simplicity.

Given a_i^m , we average it over all the users in the dataset:

$$a^m = \frac{1}{|\mathbb{U}|} \sum_{u_i \in \mathbb{U}} a_i^m, \quad (7)$$

where a^m is the average similarity over all the users and \mathbb{U} is the set of all the users.

Table A2: Best-performing Hyper-parameters in MSSG, MSSG-u and Baseline Methods

Dataset	FPMC		Caser				NARM		HGN				HAM				SASRec					
	d	d	n_s	n_p	n_o	n_f	d	l_r	d	n_s	n_p	λ	d	n_s	n_p	λ	n_l	n_o	d	n	n_h	n_b
Beauty	512	512	4	1	2	8	512	1e-3	512	3	2	1e-3	512	3	2	1e-3	1	1	128	75	4	2
Toys	512	512	4	1	1	4	512	1e-3	512	3	1	1e-4	512	3	1	1e-3	1	1	128	50	2	3
Children	128	512	4	1	2	4	256	1e-4	256	3	1	0	256	4	2	1e-4	1	2	256	175	2	1
Comics	256	512	4	1	1	4	256	1e-3	512	3	1	0	512	3	1	0	1	3	256	200	1	1
ML-1M	256	128	6	1	4	16	512	1e-4	128	4	1	1e-3	512	5	1	1e-3	2	1	256	150	4	3
ML-20M	128	256	6	1	2	8	OOM	OOM	128	4	2	1e-3	256	5	2	1e-3	3	3	256	150	2	3

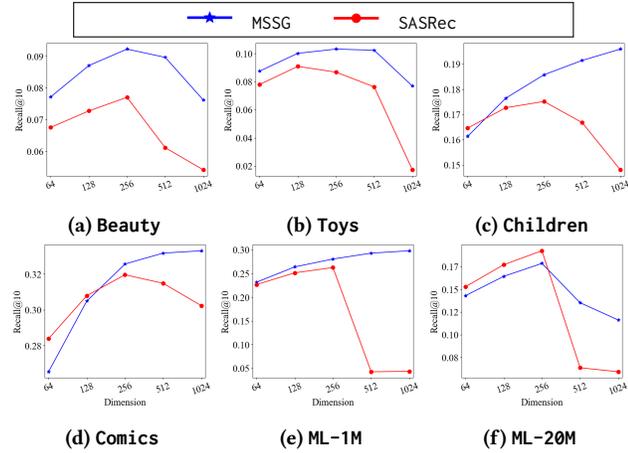
Dataset	SASRec-Cat			SASRec-Max				Bert4Rec			FMLP				MSSG-u				MSSG				
	d	n	n_h	n_b	d	n	n_h	n_b	d	n_h	n_b	d	n	n_h	n_b	d	n	n_h	n_b	d	n	n_h	n_b
Beauty	256	50	4	3	128	50	2	3	512	4	1	128	200	1	3	256	75	4	2	256	75	16	3
Toys	256	75	8	5	128	100	16	3	512	4	1	128	50	1	3	256	50	4	4	256	50	8	4
Children	256	100	2	2	256	200	2	2	256	4	2	256	150	1	1	256	200	1	2	512	100	1	1
Comics	256	200	1	2	256	100	1	3	128	4	1	512	200	1	1	512	200	1	1	512	200	1	1
ML-1M	256	75	4	3	512	200	2	4	256	4	3	512	175	1	1	512	200	2	5	512	200	4	4
ML-20M	256	200	2	4	512	150	2	2	128	4	3	512	150	1	2	512	150	8	5	256	100	4	4

In the table, "OOM" represents the out-of-memory issue.

Table A3: Dataset Statistics

dataset	#users	#items	#intrns	#intrns/u	#u/i
Beauty	22,363	12,101	198,502	8.9	16.4
Toys	19,412	11,924	167,597	8.6	14.1
Children	48,296	32,871	2,784,423	57.6	84.7
Comics	34,445	33,121	2,411,314	70.0	72.8
ML-1M	6,040	3,952	1,000,209	165.6	253.1
ML-20M	129,780	13,663	9,926,480	76.5	726.5

In this table, "#users", "#items" and "#intrns" represents the number of users, items and user-item interactions, respectively. The column "#intrns/u" has the average number of interactions of each user. The column "#u/i" has the average number of interactions on each item.

**Figure A2: Performance over different embedding dimensions**

D.3 Scalability over Embedding Dimensions

Figure A2 presents the performance of SASRec and MSSG at Recall@10 over different embedding dimensions on the six datasets. We observe that when the embedding dimension $d=2048$, both SASRec and MSSG were unable to complete training within 24 hours. Thus, we limit the maximum embedding dimension to 1024 in this analysis. As shown in Figure A2, MSSG is also more scalable than SASRec with respect to d . On all the datasets, MSSG could perform reasonably well when $d=1024$. However, SASRec fails on three out of the six datasets

(Toys, ML-1M and ML-20M) when $d=1024$. These results demonstrate the superior scalability of MSSG over SASRec with respect to d . We also notice that as illustrated in Figure A2, the recommendation performance of MSSG improves as d increases on both Children and Comics. These results suggest that better scalability in d also leads to improved recommendation performance in real-world datasets.

D.4 Calculating Entropy from Attention Weight Distributions

Given the attention map A from the first SA block in SASRec, we calculate the average entropy of the attention weight distributions over all the items as follows:

$$\frac{-1}{\sum_{u_i \in \mathcal{U}} \sum_{j=1}^n \mathbb{1}(v_{b_j}(i) \neq v_0)} \sum_{u_i \in \mathcal{U}} \sum_{j=1}^n \mathbb{1}(v_{b_j}(i) \neq v_0) \sum_{k=1}^j (A_{ijk} \log(A_{ijk})), \quad (8)$$

where $\mathbb{1}(x)$ is an indicator function (i.e., $\mathbb{1}(x) = 1$ if x is true, otherwise 0); n is the length of the transformed sequence (Section 4.1); $v_{b_j}(i) \neq v_0$ excludes the padding items; and $-\sum_{k=1}^j (A_{ijk} \log(A_{ijk}))$ is the entropy calculated from the attention weights of $v_{b_j}(i)$.

Received 20 February 2007; revised 12 March 2009; accepted 5 June 2009