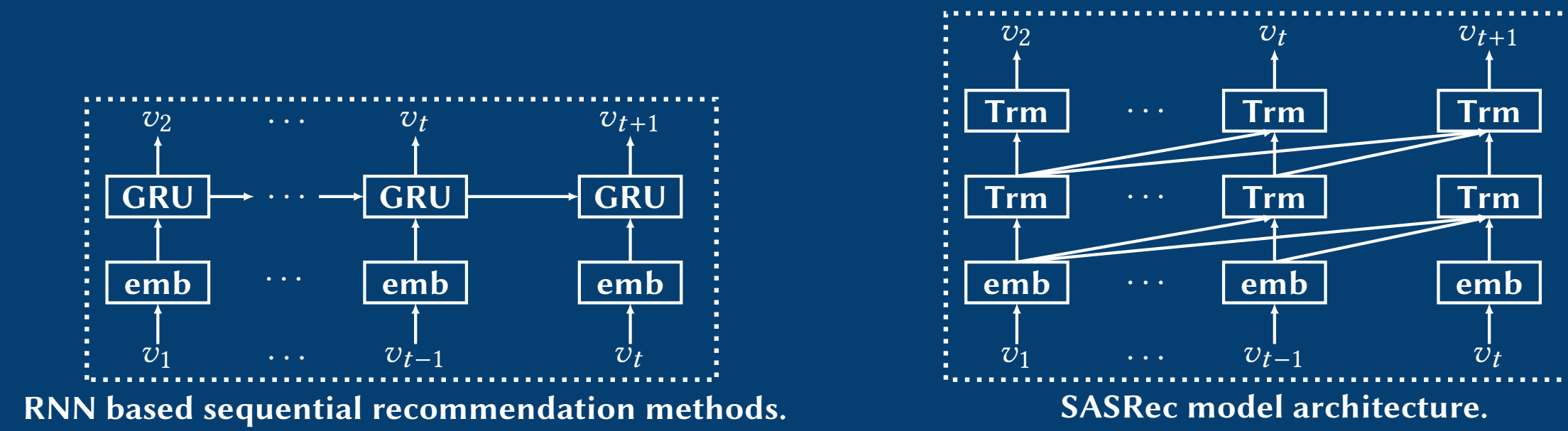
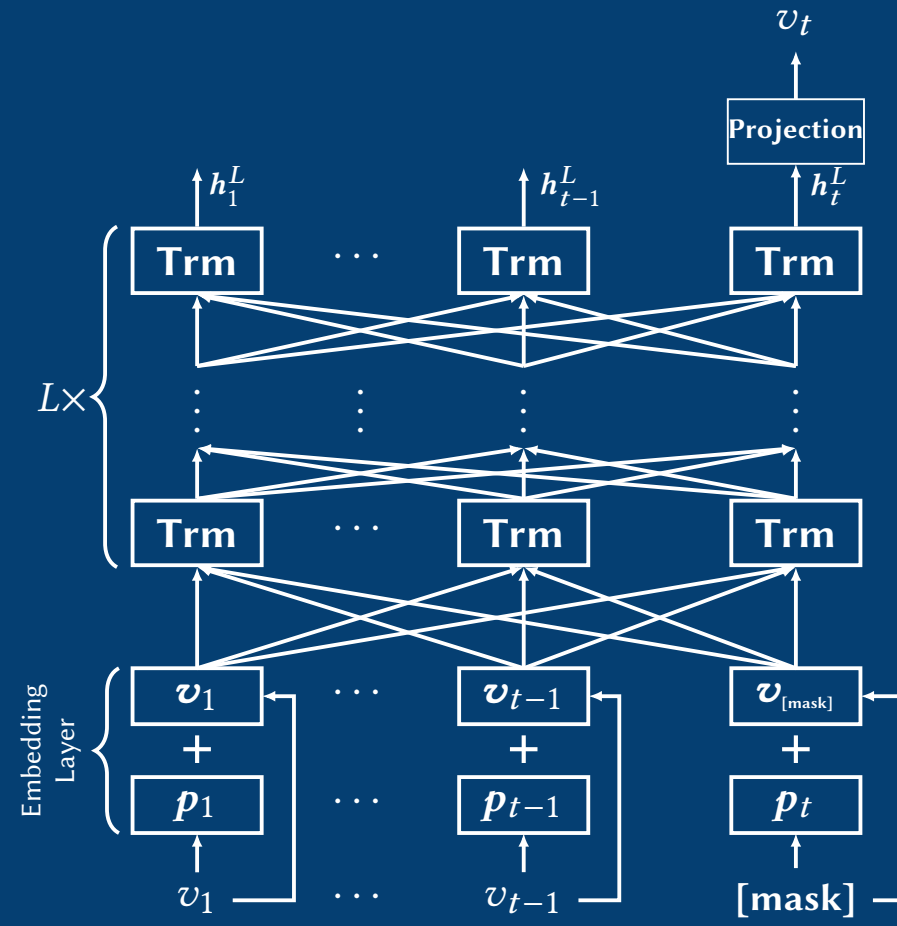


# Instead of left-to-right unidirectional model & predicting next



we use **bidirectional self-attention model** & **cloze task**



for representation learning in **sequential recommendation**

## BERT4Rec: Sequential Recommendation with Bidirectional Encoder Representations from Transformer

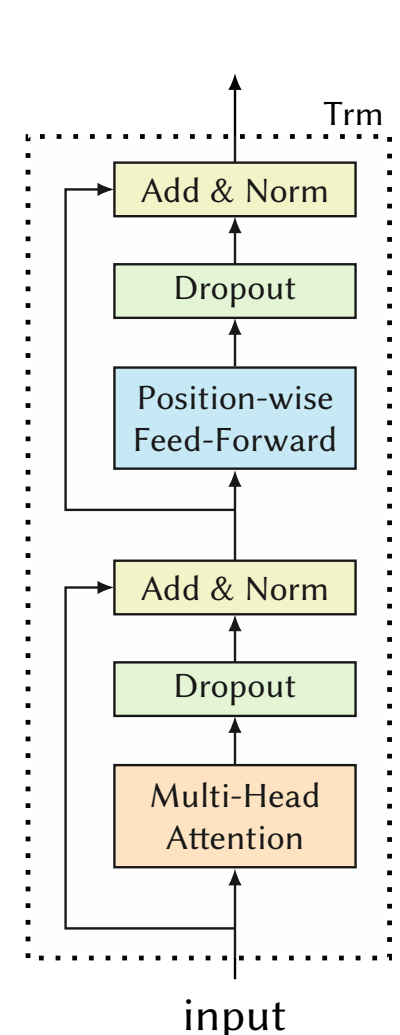
Fei Sun, Jun Liu, Jian Wu, Changhua Pei, Xiao Lin, Wenwu Ou, and Peng Jiang

### Motivation

$$p(v_{n_u+1} = v | \mathcal{S}_u), \mathcal{S}_u = [v_1^{(u)}, \dots, v_t^{(u)}, \dots, v_{n_u}^{(u)}]$$

- Unidirectional models restrict the power of hidden representations for items in the historical sequences.
- Rigid order assumption in unidirectional sequential models is not always right in user behavior sequences.
- User behaviors are often noisy due to a variety of unobservable external factors. They usually are roughly chronological, but not rigidly ordered.

### Model



$$\text{Att}(Q, K, V) = \text{softmax}\left(\frac{QK^T}{\sqrt{d/h}}\right)V$$
$$\text{head}_i = \text{Att}(\mathbf{H}^l \mathbf{W}_i^Q, \mathbf{H}^l \mathbf{W}_i^K, \mathbf{H}^l \mathbf{W}_i^V)$$
$$\text{MH}(\mathbf{H}^l) = [\text{head}_1; \text{head}_2; \dots; \text{head}_h] \mathbf{W}^O$$
$$\text{FFN}(\mathbf{x}) = \text{GELU}(\mathbf{x} \mathbf{W}^{(1)} + \mathbf{b}^{(1)}) \mathbf{W}^{(2)} + \mathbf{b}^{(2)}$$
$$\text{GELU}(x) = x \Phi(x), \Phi(x) \text{ is CDF of standard normal}$$
$$\text{PFFN}(\mathbf{H}^l) = [\text{FFN}(\mathbf{h}_1^l)^T; \dots; \text{FFN}(\mathbf{h}_h^l)^T]^T$$
$$\mathbf{A}^{l-1} = \text{LN}(\mathbf{H}^{l-1} + \text{Dropout}(\text{MH}(\mathbf{H}^{l-1})))$$
$$\text{Trm}(\mathbf{H}^{l-1}) = \text{LN}(\mathbf{A}^{l-1} + \text{Dropout}(\text{PFFN}(\mathbf{A}^{l-1})))$$
$$\mathbf{H}^l = \text{Trm}(\mathbf{H}^{l-1}), \quad \forall i \in [1, \dots, L]$$
$$\mathbf{h}_i^0 = \mathbf{v}_i + \mathbf{p}_i$$

**Learning** *cloze task*/random item mask

**Samples:** Input:  $[v_1, v_2, v_3, v_4, v_5]$  randomly mask  $\rightarrow [v_1, [\text{mask}]_1, v_3, [\text{mask}]_2, v_5]$   
Labels:  $[\text{mask}]_1 = v_2, [\text{mask}]_2 = v_4$

**Training:**  $\mathcal{L} = \frac{1}{|\mathcal{S}_u^m|} \sum_{v_m \in \mathcal{S}_u^m} -\log P(v_m = v_m^* | \mathcal{S}_u^l)$

$$P(v) = \text{softmax}(\text{GELU}(\mathbf{h}_t^L \mathbf{W}^P + \mathbf{b}^P) \mathbf{E}^T + \mathbf{b}^O)$$

**Predicting:**  $\mathcal{S}_u.\text{append}([\text{mask}])$  & predict

### Experiments

#### Datasets

Datasets	#users	#items	#actions	Avg. length	Density
Beauty	40,226	54,542	0.35m	8.8	0.02%
Steam	281,428	13,044	3.5m	12.4	0.10%
ML-1m	6040	3416	1.0m	163.5	4.79%
ML-20m	138,493	26,744	20m	144.4	0.54%

#### Task Settings & Evaluation Metrics & Baselines

Protocol	Metrics	Baselines
leave-one-out evaluation: $\mathcal{S}_u = [v_1^{(u)}, \dots, v_{n_u}^{(u)}]$	$\text{HR@}k = \frac{1}{ \mathcal{U} } \sum_{u \in \mathcal{U}} \mathbb{1}(\text{R}_{u, \mathcal{G}_u} \leq k)$	GRU4Rec
train: $[v_1^{(u)}, \dots, v_{n_u-2}^{(u)}]$ , val: $v_{n_u-1}^{(u)}$ , test: $v_{n_u}^{(u)}$	$\text{NDCG@}k = \frac{1}{ \mathcal{U} } \sum_{u \in \mathcal{U}} \frac{2^{\mathbb{1}(\text{R}_{u, \mathcal{G}_u} \leq k)} - 1}{\log_2(\text{R}_{u, \mathcal{G}_u} + 1)}$	Caser
Picking up the ground truth item from 100 randomly sampled (by spopularity) negative items	$\text{MRR} = \frac{1}{ \mathcal{U} } \sum_{u \in \mathcal{U}} \frac{1}{\text{R}_{u, \mathcal{G}_u}}$	GRU4Rec <sup>+</sup> SASRec

#### Overall Performances

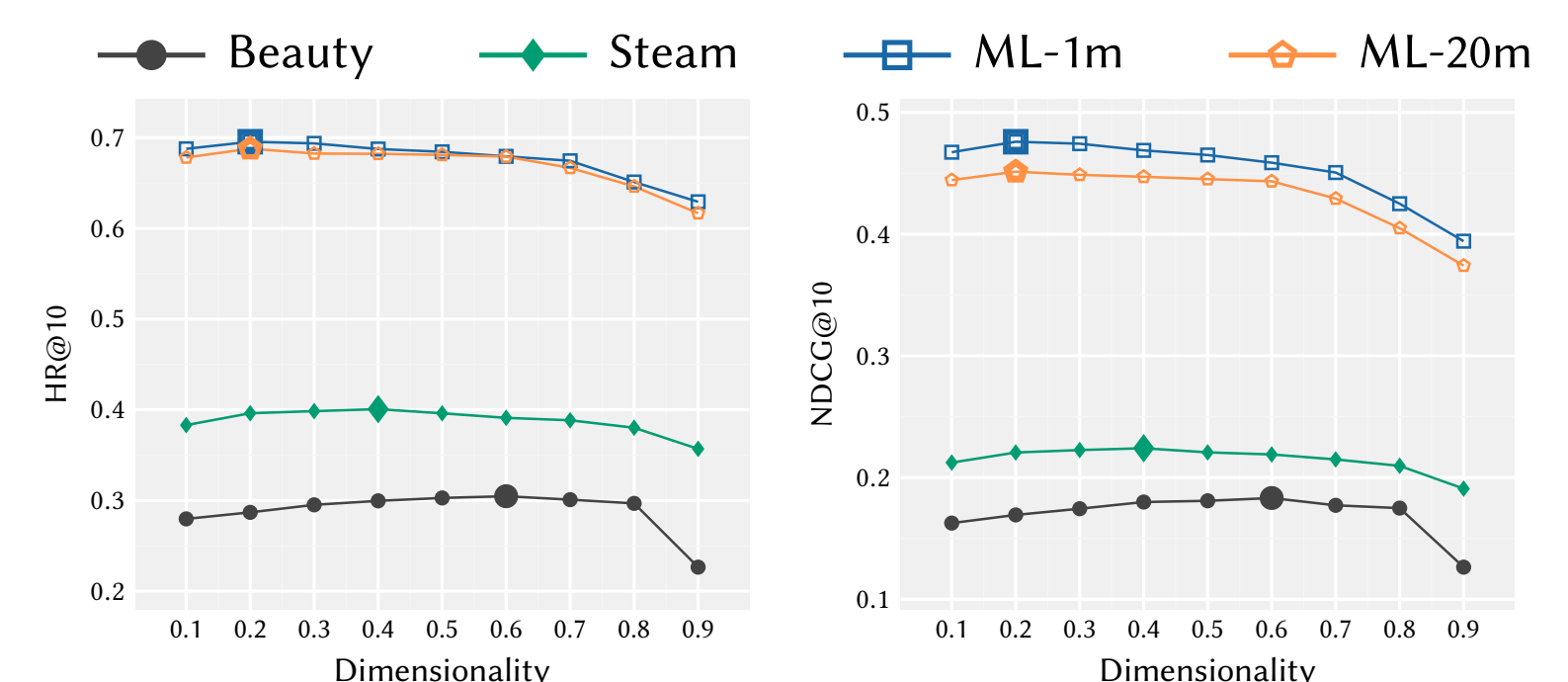
Datasets	Metric	GRU4Rec	GRU4Rec <sup>+</sup>	Caser	SASRec	BERT4Rec	Improv.
Beauty	HR@5	0.1315	0.1781	0.1625	0.1934	<b>0.2207</b>	14.12%
	HR@10	0.2343	0.2654	0.2590	0.2653	<b>0.3025</b>	14.02%
	NDCG@5	0.0812	0.1172	0.1050	0.1436	<b>0.1599</b>	11.35%
	NDCG@10	0.1074	0.1453	0.1360	0.1633	<b>0.1862</b>	14.02%
	MRR	0.1023	0.1299	0.1205	0.1536	<b>0.1701</b>	10.74%
Steam	HR@5	0.2171	0.2391	0.1766	0.2559	<b>0.2710</b>	5.90%
	HR@10	0.3313	0.3594	0.2870	0.3783	<b>0.4013</b>	6.08%
	NDCG@5	0.1370	0.1613	0.1131	0.1727	<b>0.1842</b>	6.66%
	NDCG@10	0.1802	0.2053	0.1484	0.2147	<b>0.2261</b>	5.31%
	MRR	0.1420	0.1757	0.1305	0.1874	<b>0.1949</b>	4.00%
ML-1m	HR@5	0.4673	0.5103	0.5353	0.5434	<b>0.5876</b>	8.13%
	HR@10	0.6207	0.6351	0.6692	0.6629	<b>0.6970</b>	4.15%
	NDCG@5	0.3196	0.3705	0.3832	0.3980	<b>0.4454</b>	11.91%
	NDCG@10	0.3627	0.4064	0.4268	0.4368	<b>0.4818</b>	10.32%
	MRR	0.3041	0.3462	0.3648	0.3790	<b>0.4254</b>	12.24%
ML-20m	HR@5	0.4657	0.5118	0.3804	0.5727	<b>0.6323</b>	10.41%
	HR@10	0.5844	0.6524	0.5427	0.7136	<b>0.7473</b>	4.72%
	NDCG@5	0.3090	0.3630	0.2538	0.4208	<b>0.4967</b>	18.04%
	NDCG@10	0.3637	0.4087	0.3062	0.4665	<b>0.5340</b>	14.47%
	MRR	0.2967	0.3476	0.2529	0.4026	<b>0.4785</b>	18.85%

#### Effectiveness of bidirectional self-attention & Cloze objective

Model	Beauty			ML-1m		
	HR@10	NDCG@10	MRR	HR@10	NDCG@10	MRR
SASRec	0.2653	0.1633	0.1536	0.6629	0.4368	0.3790
BERT4Rec (1 mask)	0.2940	0.1769	0.1618	0.6869	0.4696	0.4127
BERT4Rec	0.3025	0.1862	0.1701	0.6970	0.4818	0.4254

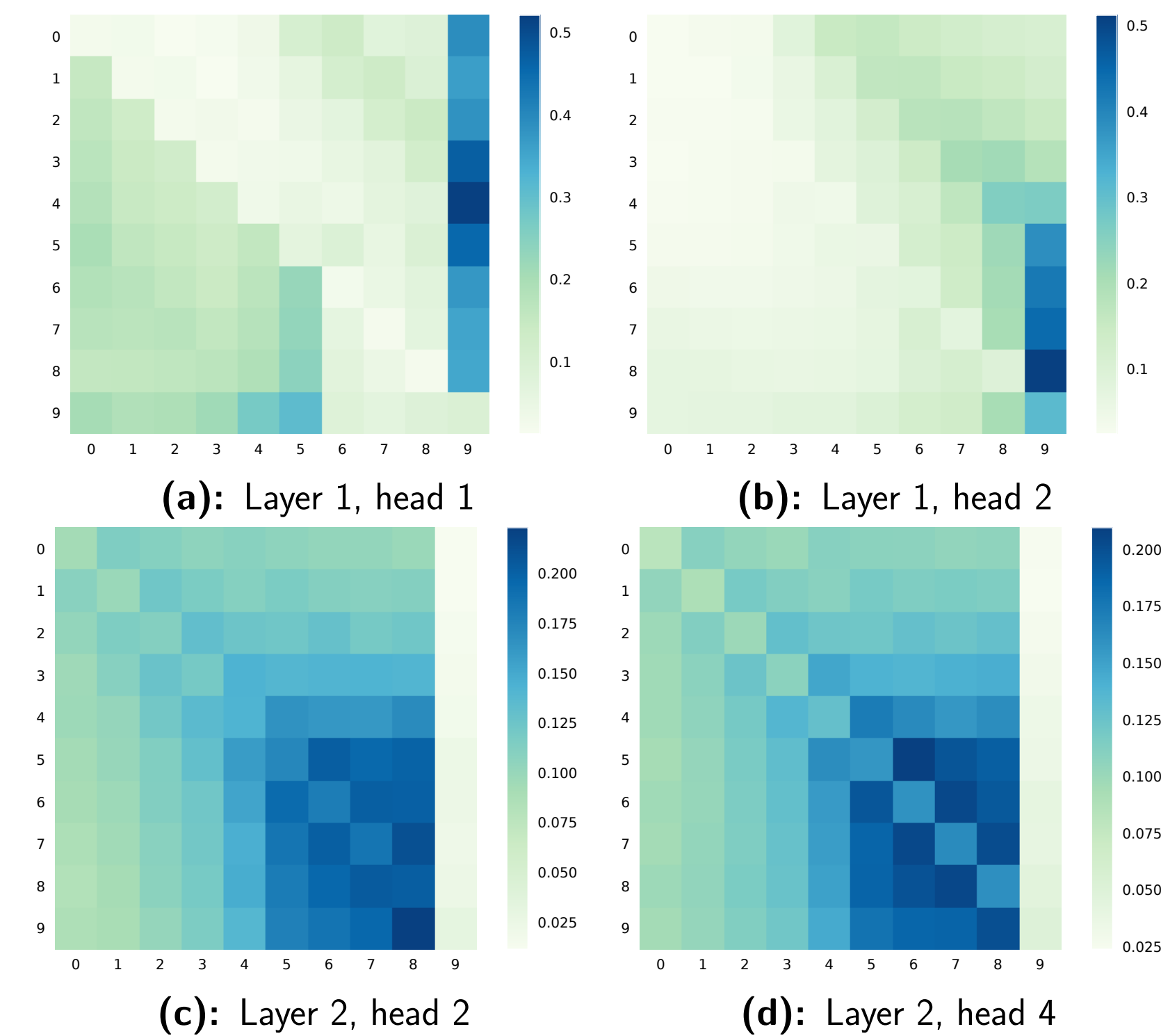
## Extra Results

### Impact of Mask Proportion $\rho$



Performance with different mask proportion  $\rho$  on  $d = 64$ . Bold symbols denote the best scores in each line.

### Attention Visualization



Heat-maps of average attention weights on Beauty, the last position "9" denotes "[mask]" (best viewed in color).

- Attention varies across different heads.
- Attention varies across different layers.
- Items tend to attend on the items at both sides.

### Ablation Study

Architecture	Dataset			
	Beauty	Steam	ML-1m	ML-20m
$L = 2, h = 2$	0.1832	0.2241	0.4759	0.4513
w/o PE	0.1741	0.2060	0.2155↓	0.2867↓
w/o PFFN	0.1803	0.2137	0.4544	0.4296
w/o LN	0.1642↓	0.2058	0.4334	0.4186
w/o RC	0.1619↓	0.2193	0.4643	0.4483
w/o Dropout	0.1658	0.2185	0.4553	0.4471
1 layer ( $L = 1$ )	0.1782	0.2122	0.4412	0.4238
3 layers ( $L = 3$ )	<b>0.1859</b>	<b>0.2262</b>	<b>0.4864</b>	<b>0.4661</b>
4 layers ( $L = 4$ )	<b>0.1834</b>	<b>0.2279</b>	<b>0.4898</b>	<b>0.4732</b>
1 head ( $h = 1$ )	<b>0.1853</b>	0.2187	0.4568	0.4402
4 heads ( $h = 4$ )	0.1830	<b>0.2245</b>	<b>0.4770</b>	<b>0.4520</b>
8 heads ( $h = 8$ )	0.1823	<b>0.2248</b>	0.4743	<b>0.4550</b>



CIKM2019

Alibaba Group  
阿里巴巴集团