



# Attention Calibration for Transformer-based Sequential Recommendation

Peilin Zhou\*

The Hong Kong University  
of Science and Technology  
(Guangzhou)  
zhoupalin@gmail.com

Qichen Ye\*

Peking University  
yeeeqichen@pku.edu.cn

Yueqi Xie

The Hong Kong University  
of Science and Technology  
yxieay@connect.ust.hk

Jingqi Gao

Upstage  
mrgao.ary@gmail.com

Shoujin Wang

University of Technology  
Sydney  
shoujin.wang@uts.edu.au

Jae Boum Kim

The Hong Kong University  
of Science and Technology  
jbkim@cse.ust.hk

Chenyu You

Yale University  
chenyu.you@yale.edu

Sunghun Kim<sup>†</sup>

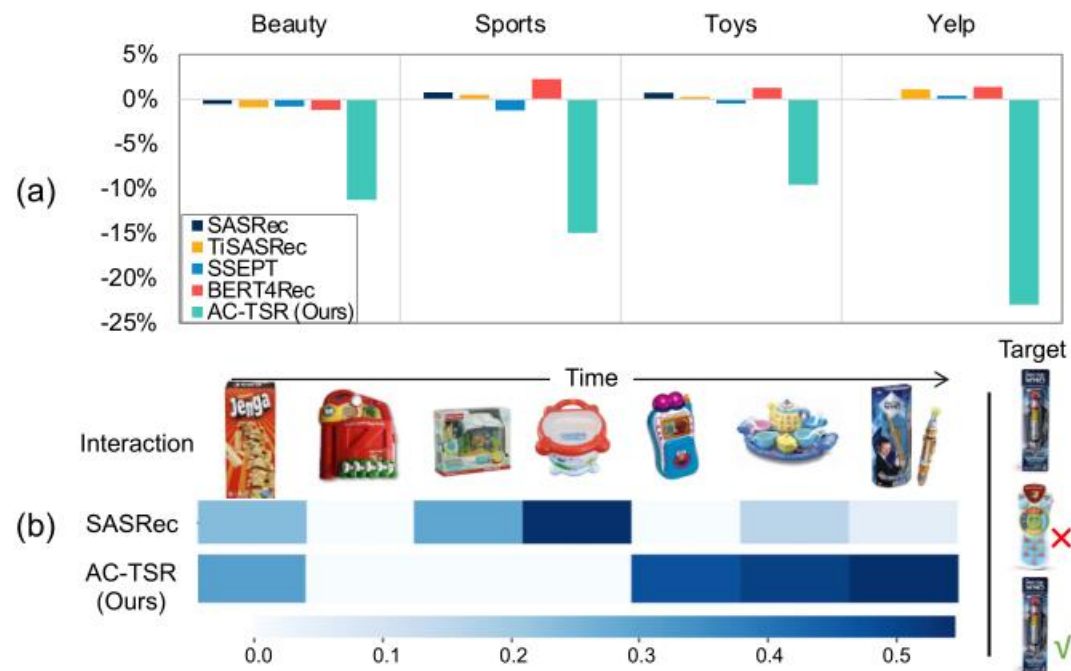
The Hong Kong University  
of Science and Technology  
(Guangzhou)  
hunkim@cse.ust.hk

code: <https://github.com/AIM-SE/AC-TSR>.

CIKM 2023



# Introduction



After careful and in-depth analysis, we found the aforementioned unreliable or inaccurate assignment of attention weights could be mainly attributed to the following two factors:

- (1) Sub-optimal position encoding.
- (2) Noisy input.

**Figure 1: (a) Removing the highest attention weight from transformer-based SRS does not lead to a significant decrease in model performance and even improves performance in some cases; (b) Visualization of the attention weights from SASRec and our proposed AC-TSR.**

# Method

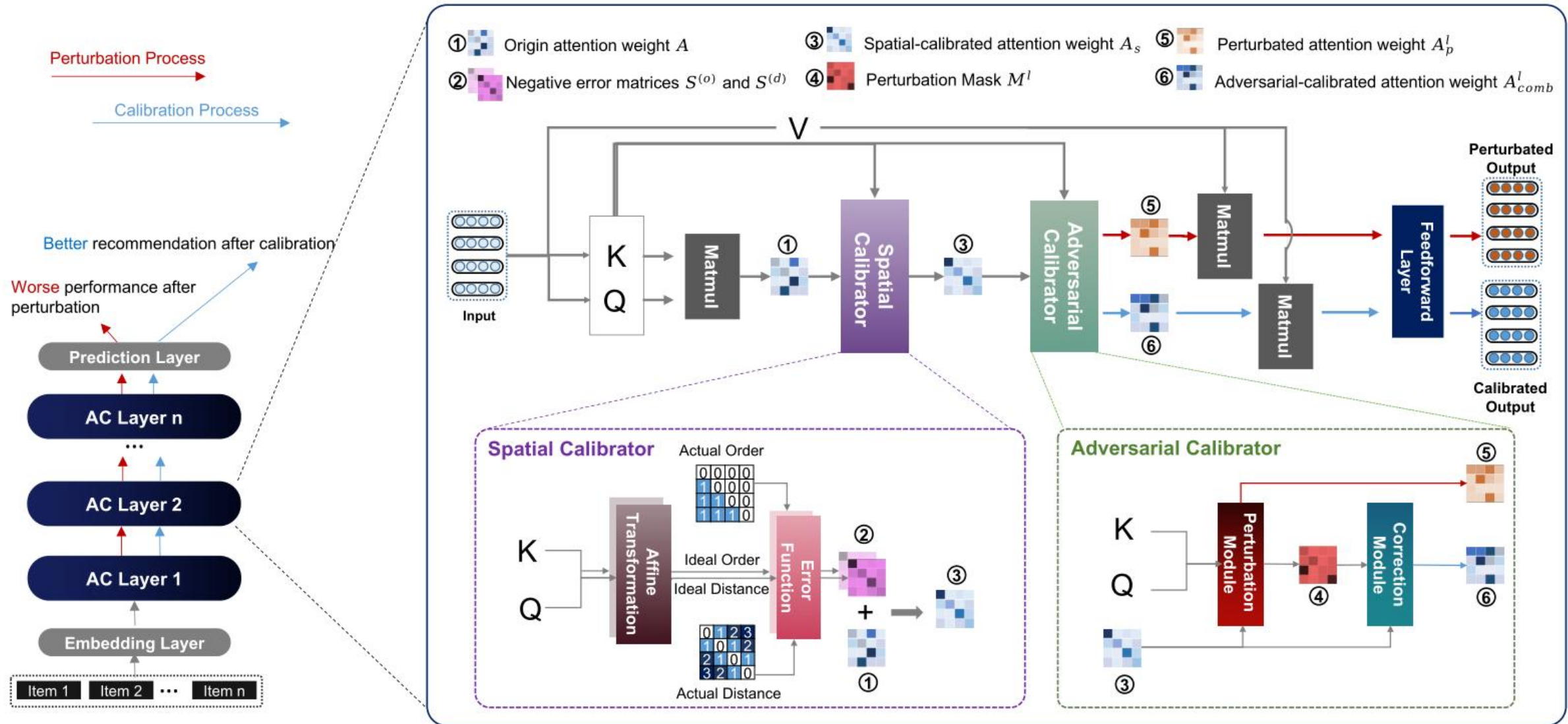
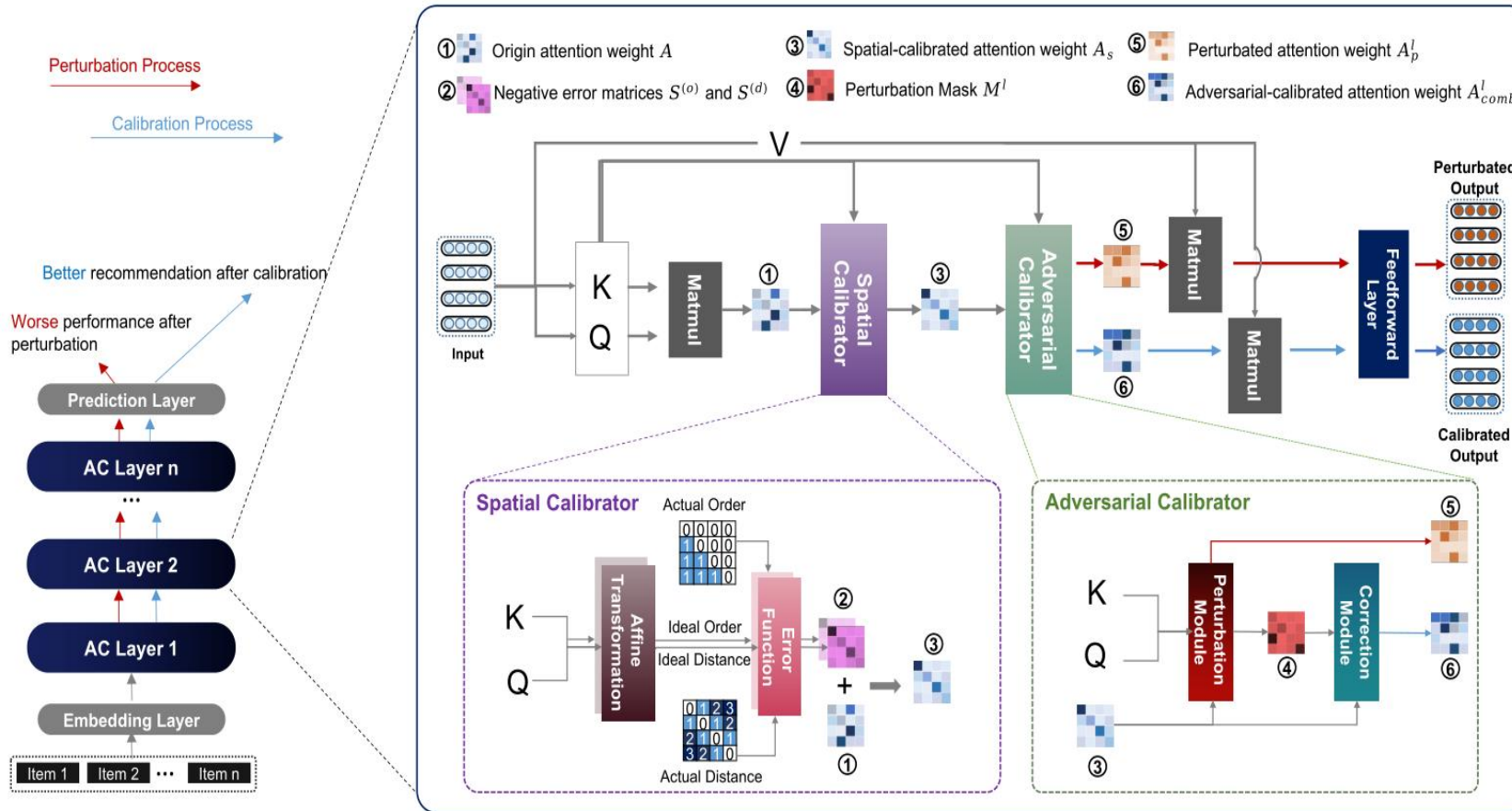


Figure 2: Overview of the proposed AC-TSR framework.

# Method



$$H = \text{Self-Attention}(Q, K, V) = \text{softmax}\left(\frac{QK^T}{\sqrt{d}}\right)V \quad (1)$$

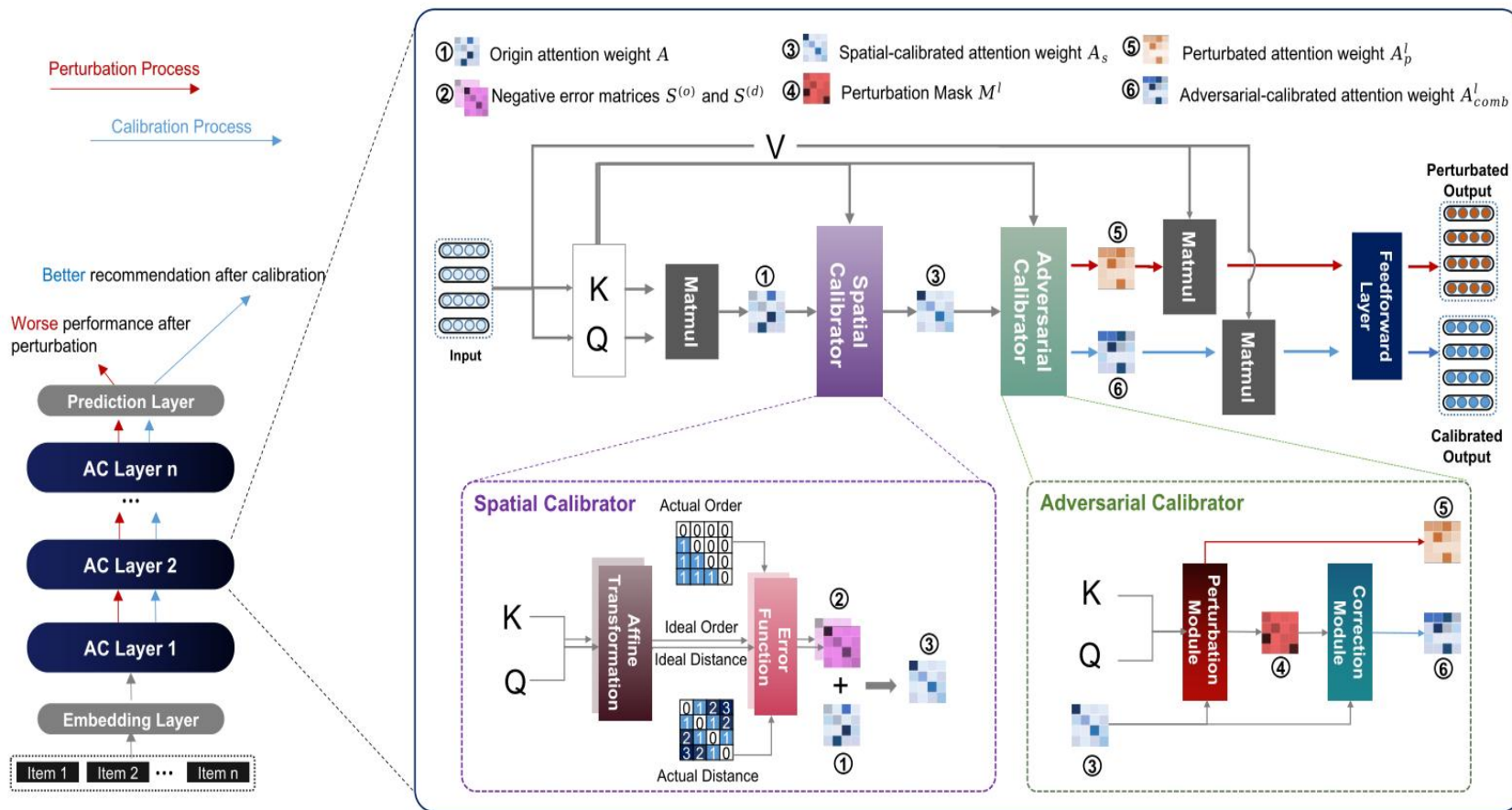
$$F_i = \text{FFN}(H_i) = \text{ReLU}(H_i W_1 + b_1) W_2 + b_2 \quad (2)$$

$$F_i^L = \text{FFN}(H_i^L) \quad (3)$$

$$\hat{y} = \text{softmax}\left(\text{TF}_n^{L^T}\right) \quad (4)$$

$$\mathcal{L} = - \sum_{i=1}^{|\mathcal{I}|} y_i \log(\hat{y}_i) \quad (5)$$

# Method



$$o_{ij} = \mathbb{I}(i < j) = \begin{cases} 1, & i < j \\ 0, & \text{otherwise} \end{cases} \quad (6)$$

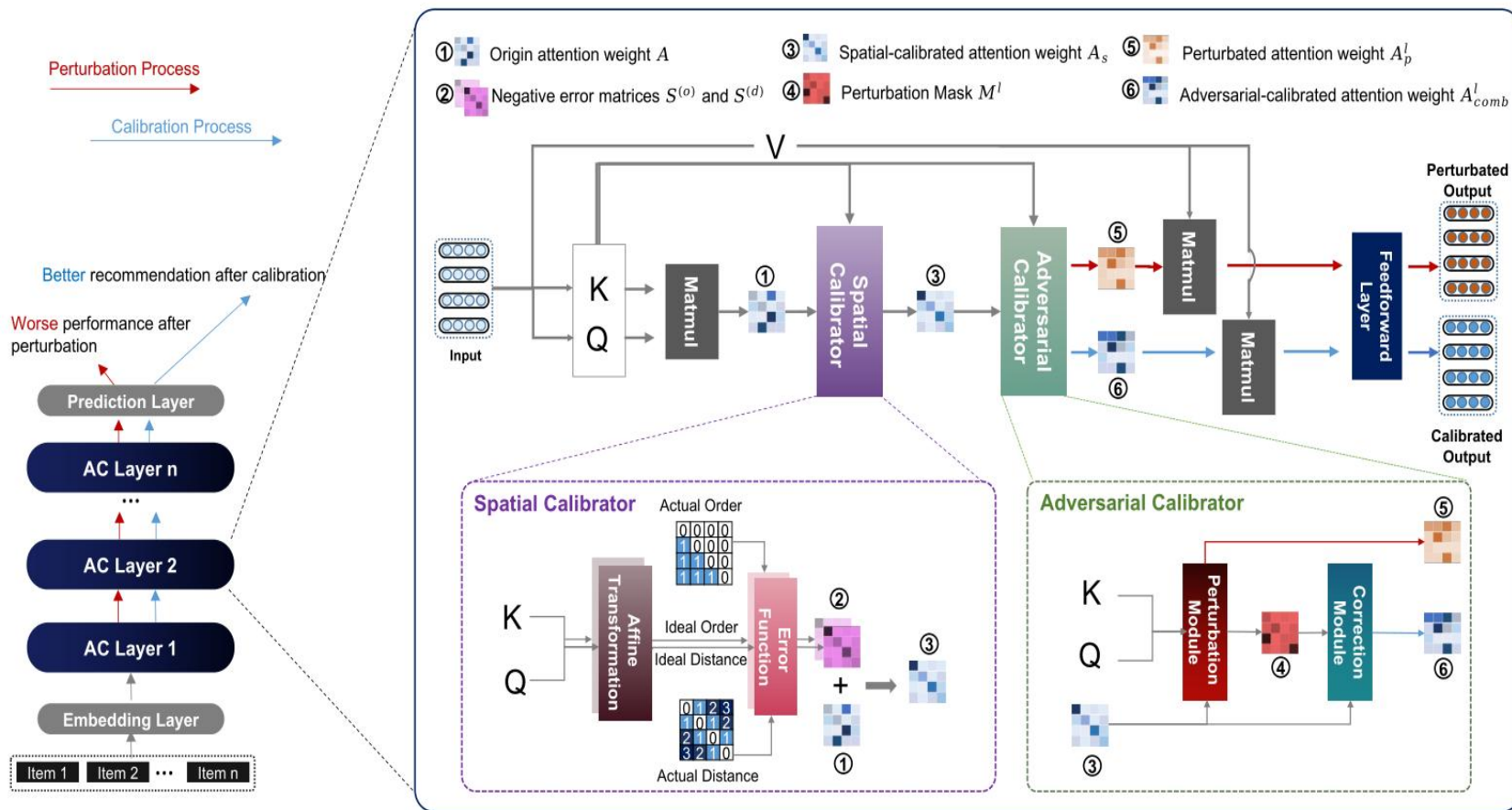
$$d_{ij} = \ln(1 + |i - j|) \quad (7)$$

$$\hat{o}_{ij} = \text{sigmoid} \left( \text{affine}^{(o)} \left( \begin{bmatrix} \mathbf{q}_i^l; \mathbf{k}_j^l \end{bmatrix} \right) \right) \quad (8)$$

$$\hat{d}_{ij} = \text{affine}^{(d)} \left( \begin{bmatrix} \mathbf{q}_i^l; \mathbf{k}_j^l \end{bmatrix} \right) \quad (9)$$

$$s_{ij}^{(o)} = o_{ij} \ln(\hat{o}_{ij}) + (1 - o_{ij})(1 - \ln(\hat{o}_{ij})) \quad (10)$$

## Method



$$s_{ij}^{(d)} = -\frac{\theta^2 (d_{ij} - \hat{d}_{ij})^2}{2} \quad (11)$$

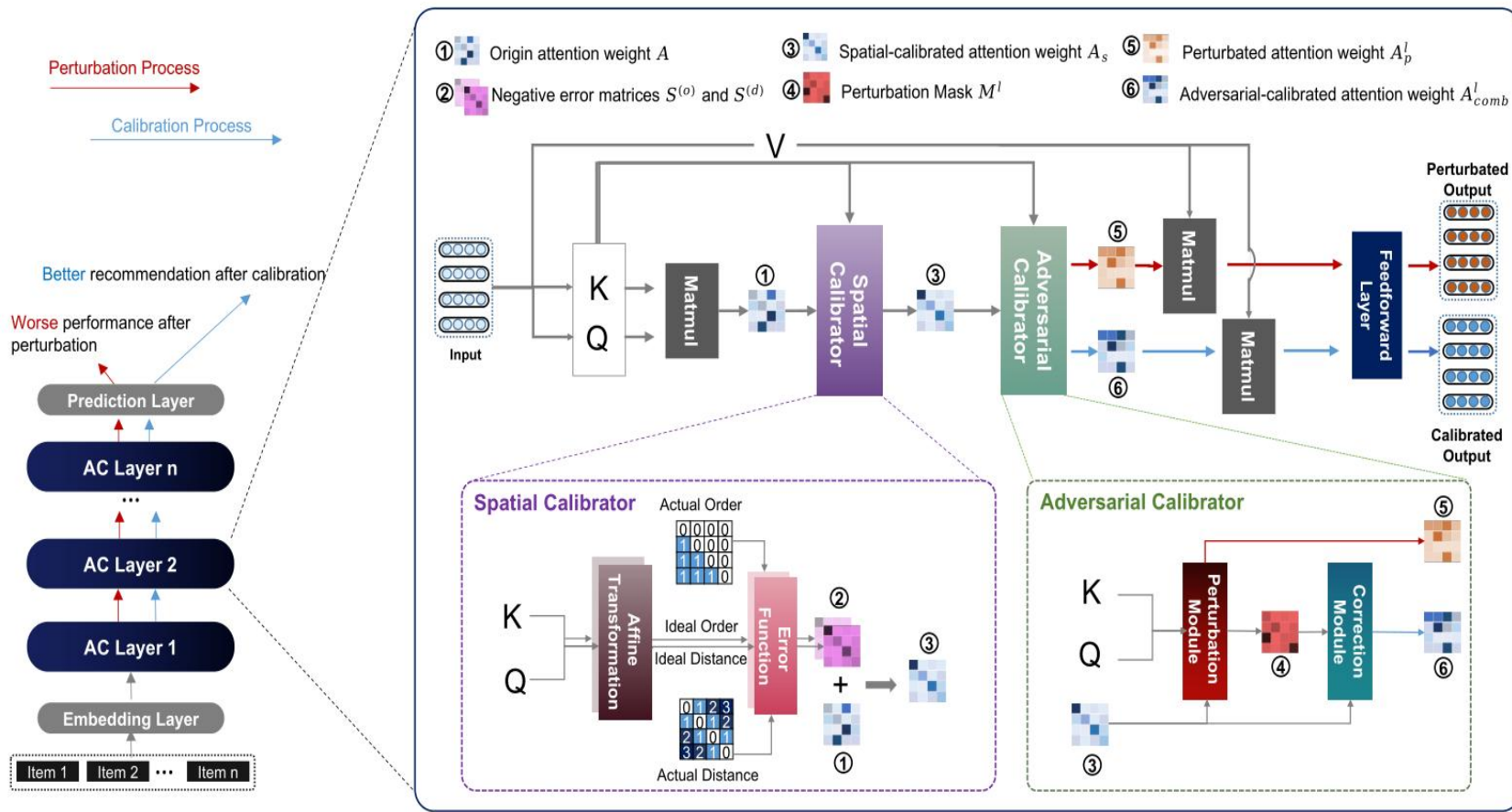
$$A_s = \text{softmax} \left( \frac{QK^T}{\sqrt{d}} + s^{(o)} + s^{(d)} \right) \quad (12)$$

$$A_p^l = M^l \odot A_s^l + (1 - M^l) \odot \mu \quad (13)$$

$$M^l = \text{sigmoid} \left( \frac{Q^l W_{Q_p}^l (K^l W_{K_p}^l)^T}{\sqrt{d}} \right) \quad (14)$$

$$A_c^l = A_s^l \odot e^{1-M^l} \quad (15)$$

# Method



$$A_{comb}^l = g * A_s^l + (1 - g) * A_c^l \quad (16)$$

$$g = \sigma(Q^l W_g^l + b_g^l) \quad (17)$$

$$\mathcal{L}_P = - \sum_{i=1}^{|I|} y_i \log(\hat{y}_i^P) \quad (18)$$

$$\mathcal{L}_C = - \sum_{i=1}^{|I|} y_i \log(\hat{y}_i^C) \quad (19)$$

$$\mathcal{L}_{P_{final}}(\theta^P) = -\mathcal{L}_P(\theta) + \alpha \mathcal{L}_{norm}(\theta^P) \quad (20)$$

$$\mathcal{L}_{norm}(\theta^P) = \sum_{l=0}^L \|1 - \mathbf{m}^l\|_2 \quad (21)$$

$$\mathcal{L}_{final} = \mathcal{L}_{P_{final}} + \mathcal{L}_C \quad (22)$$



# Experiments

**Table 1: Overall performance. The highest results are denoted in bold, while the runner-up results are underscored. "\*" denotes the statistical significance for  $p < 0.01$  compared to the best baseline methods with paired  $t$ -test.**

SR Model	Beauty				Sports				Toys				Yelp			
	Recall		NDCG		Recall		NDCG		Recall		NDCG		Recall		NDCG	
	@10	@20	@10	@20	@10	@20	@10	@20	@10	@20	@10	@20	@10	@20	@10	@20
PopRec	0.0157	0.0242	0.0076	0.0097	0.0146	0.0244	0.0078	0.0103	0.0105	0.0172	0.0060	0.0077	0.0099	0.0161	0.0051	0.0067
BPR	0.0375	0.0590	0.0168	0.0222	0.0302	0.0480	0.0144	0.0188	0.0344	0.0560	0.0151	0.0205	0.0589	0.0830	0.0324	0.0384
GRU4Rec	0.0654	0.1002	0.0322	0.0410	0.0386	0.0609	0.0195	0.0251	0.0449	0.0708	0.0221	0.0287	0.0418	0.0679	0.0206	0.0271
Caser	0.0474	0.0731	0.0239	0.0304	0.0227	0.0364	0.0118	0.0153	0.0361	0.0566	0.0186	0.0238	0.0380	0.0608	0.0197	0.0255
LightSANS	0.0770	0.1177	0.0358	0.0461	0.0509	0.0781	0.0226	0.0294	0.0768	0.1116	0.0354	0.0442	0.0630	0.0904	0.0385	0.0453
Locker	0.0802	0.1197	0.0365	0.0464	0.0508	0.0753	0.0225	0.0286	0.0755	0.1094	0.0345	0.0430	0.0603	0.0869	0.0380	0.0446
SASRec	0.0779	0.1152	0.0353	0.0447	0.0504	0.0760	0.0224	0.0289	0.0776	0.1100	0.0352	0.0434	0.0618	0.0879	0.0387	0.0453
w/ AC	<u>0.0817*</u>	<u>0.1218*</u>	<b>0.0375*</b>	<u>0.0454*</u>	<u>0.0532*</u>	<u>0.0817*</u>	<u>0.0235*</u>	<u>0.0307*</u>	<u>0.0825*</u>	<u>0.1166*</u>	<u>0.0371*</u>	<u>0.0456*</u>	<b>0.0664*</b>	<b>0.0955*</b>	<b>0.0407*</b>	<b>0.0480*</b>
Improve.	4.88%	5.73%	6.23%	1.57%	5.56%	7.50%	4.91%	6.23%	6.31%	6.00%	5.40%	5.07%	7.44%	8.65%	5.17%	5.96%
BERT4Rec	0.0557	0.0868	0.0279	0.0358	0.0313	0.0502	0.0155	0.0202	0.0489	0.0769	0.0253	0.0324	0.0467	0.0710	0.0264	0.0325
w/ AC	<u>0.0628*</u>	<u>0.0929*</u>	<u>0.0318*</u>	<u>0.0394*</u>	<u>0.0381*</u>	<u>0.0607*</u>	<u>0.0196*</u>	<u>0.0253*</u>	<u>0.0643*</u>	<u>0.0924*</u>	<u>0.0339*</u>	<u>0.0410*</u>	<u>0.0481*</u>	<u>0.0769*</u>	<u>0.0265*</u>	<u>0.0337*</u>
Improve.	12.73%	7.03%	13.98%	10.06%	21.73%	20.92%	26.45%	25.25%	31.49%	20.16%	33.99%	26.54%	3.00%	8.31%	0.38%	3.69%
SSE-PT	0.0587	0.0936	0.0278	0.0366	0.0363	0.0580	0.0184	0.0239	0.0560	0.0837	0.0255	0.0325	0.0556	0.0779	0.0323	0.0379
w/ AC	<u>0.0629*</u>	<u>0.1001*</u>	<u>0.0293*</u>	<u>0.0387*</u>	<u>0.0379*</u>	<u>0.0589*</u>	<u>0.0191*</u>	<u>0.0244*</u>	<u>0.0614*</u>	<u>0.0896*</u>	<u>0.0282*</u>	<u>0.0353*</u>	<u>0.0565*</u>	<u>0.0821*</u>	<u>0.0330*</u>	<u>0.0394*</u>
Improve.	7.16%	6.94%	5.40%	5.74%	4.41%	1.55%	3.80%	2.09%	9.64%	7.05%	10.59%	8.62%	1.62%	5.39%	2.17%	3.96%
TiSASRec	0.0794	0.1208	0.0356	0.0461	0.0523	0.0799	0.0230	0.0300	0.0819	0.1171	0.0367	0.0456	0.0618	0.0909	0.0387	0.0460
w/ AC	<b>0.0823*</b>	<b>0.1227*</b>	<u>0.0373*</u>	<b>0.0474*</b>	<b>0.0548*</b>	<b>0.0837*</b>	<b>0.0241*</b>	<b>0.0313*</b>	<b>0.0831*</b>	<b>0.1208*</b>	<b>0.0375*</b>	<b>0.0470*</b>	<u>0.0654*</u>	<u>0.0939*</u>	<u>0.0401*</u>	<u>0.0473*</u>
Improve.	3.65%	1.57%	4.78%	2.82%	4.78%	4.76%	4.78%	4.33%	1.47%	3.16%	2.18%	3.07%	5.83%	3.30%	3.62%	2.83%





# Experiments

**Table 2: Model Complexity.**

Model	# Parameters	Inference speed	Recall@20			
			Beauty	Sports	Toys	Yelp
SASRec	0.87M	2482.33/s	0.1152	0.0760	0.1100	0.0879
AC-SASRec	0.90M	917.54/s	0.1218	0.0817	0.1166	0.0955
AC-SASRec-lite	0.87M	2482.33/s	0.1164	0.0768	0.1150	0.0913



# Experiments

**Table 3: Ablation study of AC-TSR on Beauty dataset.**

Settings	Spatial		Adv.	Recall		NDCG	
	order	distance		@10	@20	@10	@20
(A)	✓	✓	✓	0.0817	0.1218	0.0375	0.0476
(B)	✓	✗	✓	0.0791	0.1210	0.0364	0.0470
(C)	✗	✓	✓	0.0792	0.1201	0.0372	0.0475
(D)	✗	✗	✓	0.0800	0.1202	0.0367	0.0469
(E)	✓	✓	✗	0.0802	0.1197	0.0365	0.0464
(F)	✓	✗	✗	0.0776	0.1168	0.0353	0.0452
(G)	✗	✓	✗	0.0806	0.1188	0.0367	0.0463
(H)	✗	✗	✗	0.0779	0.1152	0.0353	0.0447

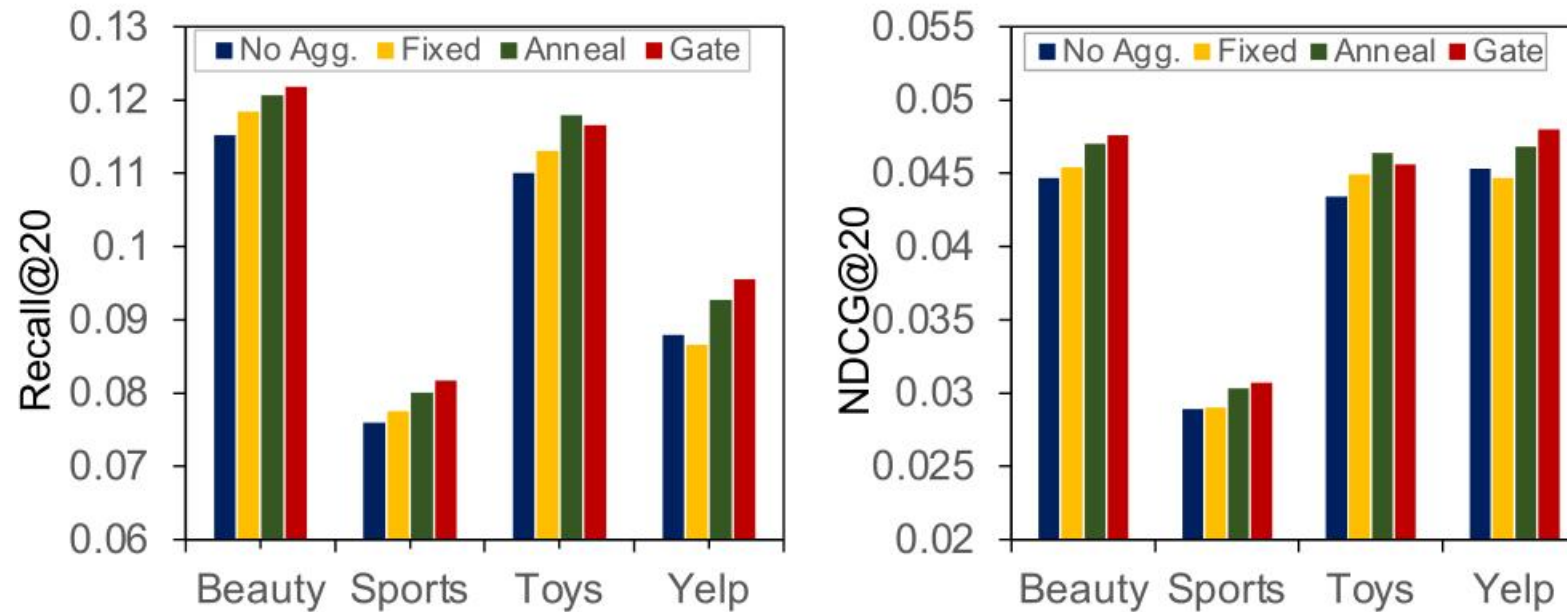


# Experiments

**Table 4: Impact of different positional encoding strategies. The SASRec is chosen as the backbone.**

Position Encoding Strategy	Sports		Toys	
	Recall@20	NDCG@20	Recall@20	NDCG@20
Remove Position	0.0775	0.0294	0.1170	0.0456
Absolute Position	0.0760	0.0289	0.1100	0.0434
Relative Position	0.0753	0.0285	0.1172	0.0461
Decoupled Position	0.0769	0.0295	0.1153	0.0449
Spatial Calibrator (Ours)	<b>0.0785</b>	<b>0.0298</b>	<b>0.1193</b>	<b>0.0462</b>

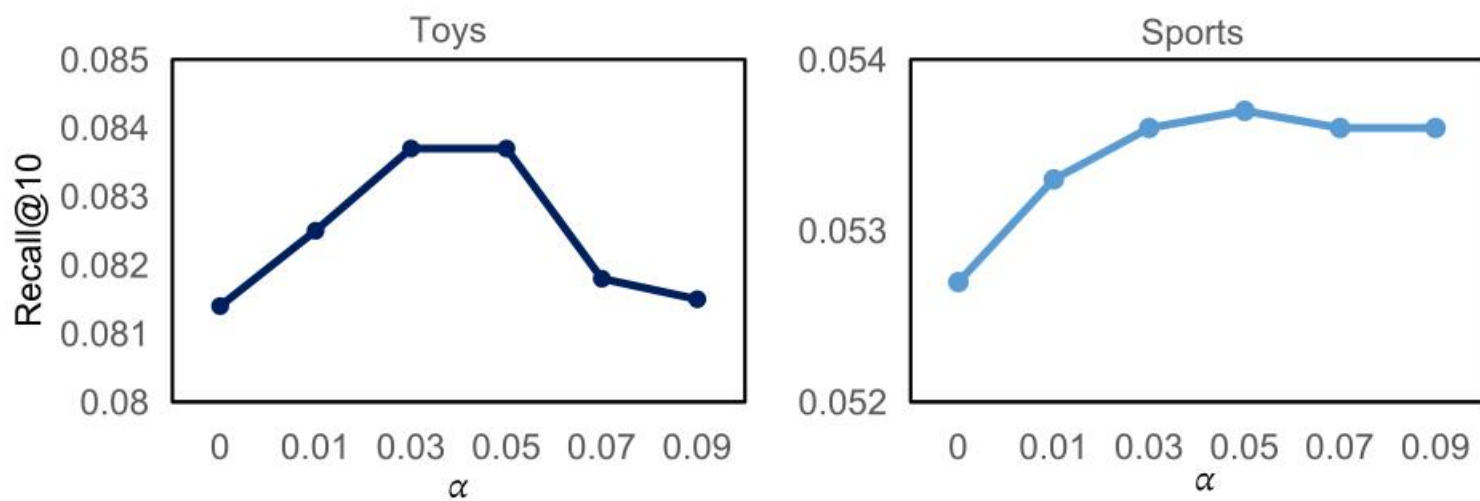
# Experiments



**Figure 3: Impact of different aggregation strategies in Correction Module.**

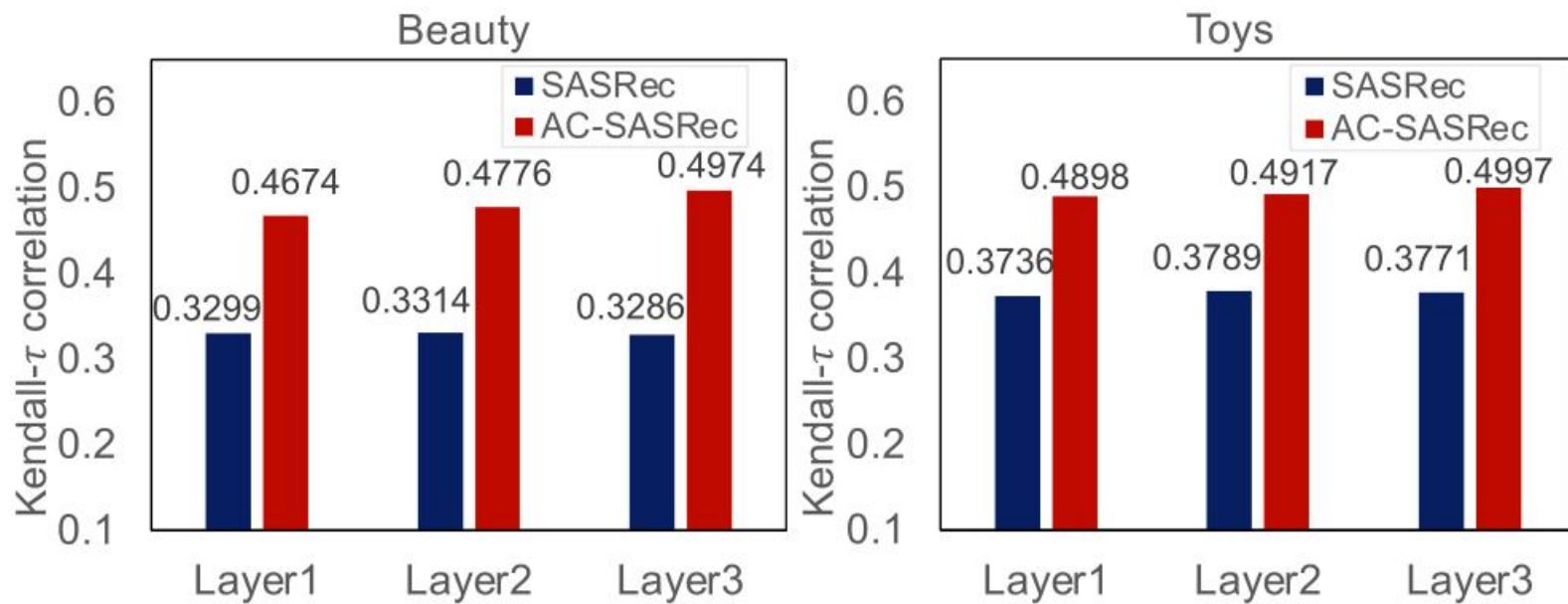


# Experiments



**Figure 4: Effect of balance parameter  $\alpha$ .**

# Experiments



**Figure 5: Comparison of the mean Kendall- $\tau$  correlation between attention weights and gradient importance measures. The results verify that our AC method can improve Kendall- $\tau$  correlation by a large margin.**

# Experiments

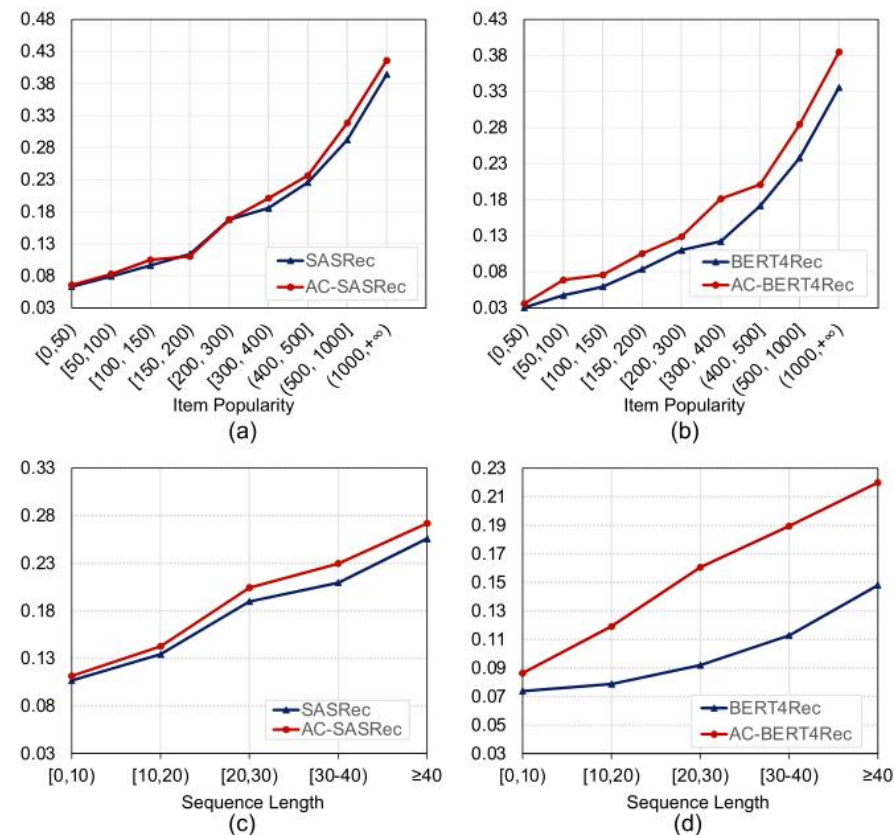


Figure 6: Performance comparison (Recall@20) between AC-TSR and TSR under different sequence lengths (*i.e.*, number of training interactions of users) and item popularity (*i.e.*, number of training interactions of items) on Amazon Beauty.



**Thanks**