



# Ensemble Modeling with Contrastive Knowledge Distillation for Sequential Recommendation

Hanwen Du  
hwdu@stu.suda.edu.cn  
Soochow University  
Suzhou, Jiangsu, China

Huanhuan Yuan  
hhyuan@stu.suda.edu.cn  
Soochow University  
Suzhou, Jiangsu, China

Pengpeng Zhao<sup>\*</sup>  
ppzhao@suda.edu.cn  
Soochow University  
Suzhou, Jiangsu, China

Fuzhen Zhuang  
zhuangfuzhen@buaa.edu.cn  
Beihang University  
Beijing, China

Guanfeng Liu  
guanfeng.liu@mq.edu.au  
Macquarie University  
Sydney, Australia

Lei Zhao  
zhaol@suda.edu.cn  
Soochow University  
Suzhou, Jiangsu, China

Yanchi Liu  
yanchi.liu@rutgers.edu  
Rutgers University  
New Brunswick, New Jersey, USA

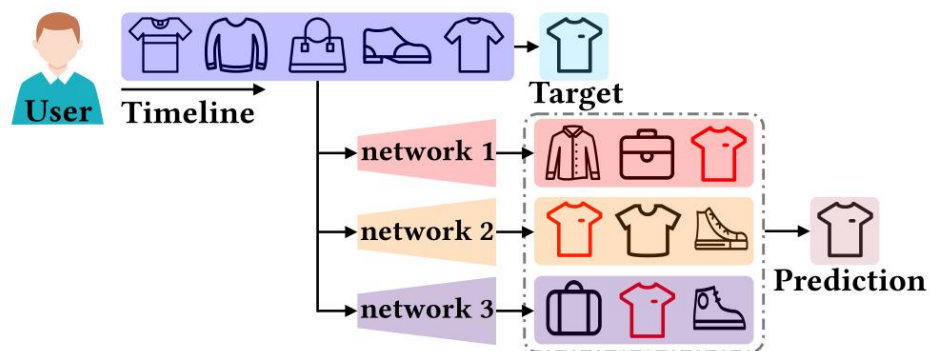
Victor S. Sheng  
victor.sheng@ttu.edu  
Texas Tech University  
Lubbock, Texas, USA

code: <https://github.com/hw-du/EMKD>.

SIGIR 2023



# Introduction



**Figure 1: An illustration of ensemble modeling for sequential recommendation. Three parallel networks make different predictions based on users' historical interactions. Although each individual network is unable to make an accurate prediction, combining the predictions of these networks together will get the correct result.**

We propose a novel framework called Ensemble Modeling with Contrastive Knowledge Distillation for sequential recommendation (EMKD). To the best of our knowledge, this is the first work to apply the ensemble modeling to sequential recommendation.

We propose a novel contrastive knowledge distillation approach that facilitates knowledge transfer and distills knowledge from both the representation level and the logits level

## Method

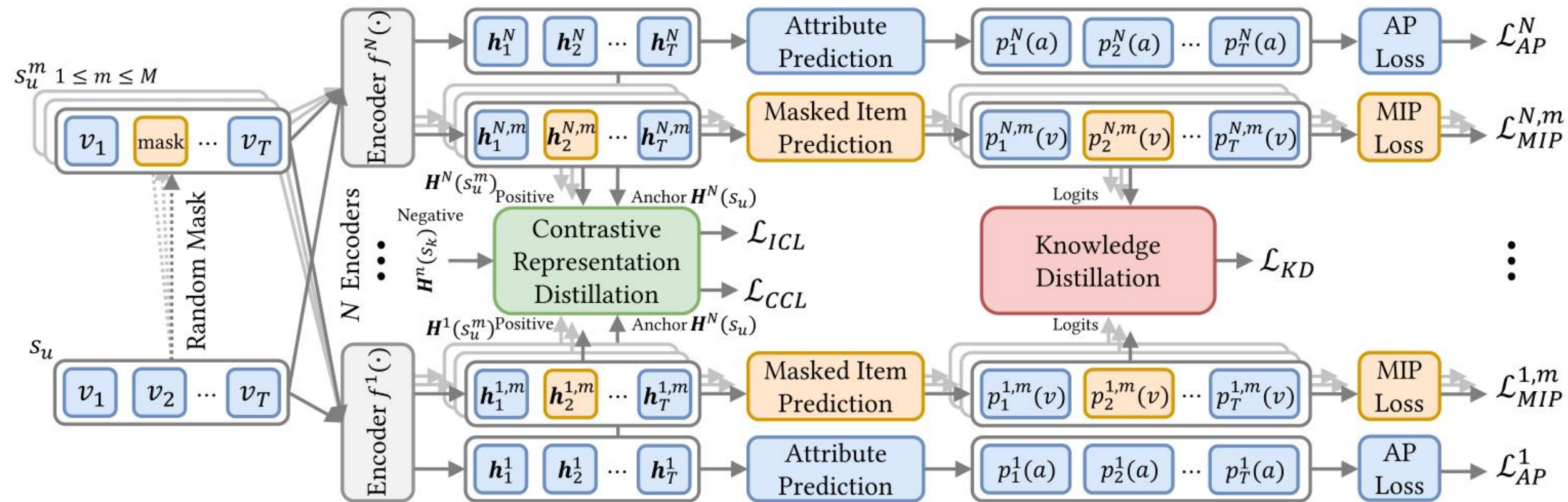
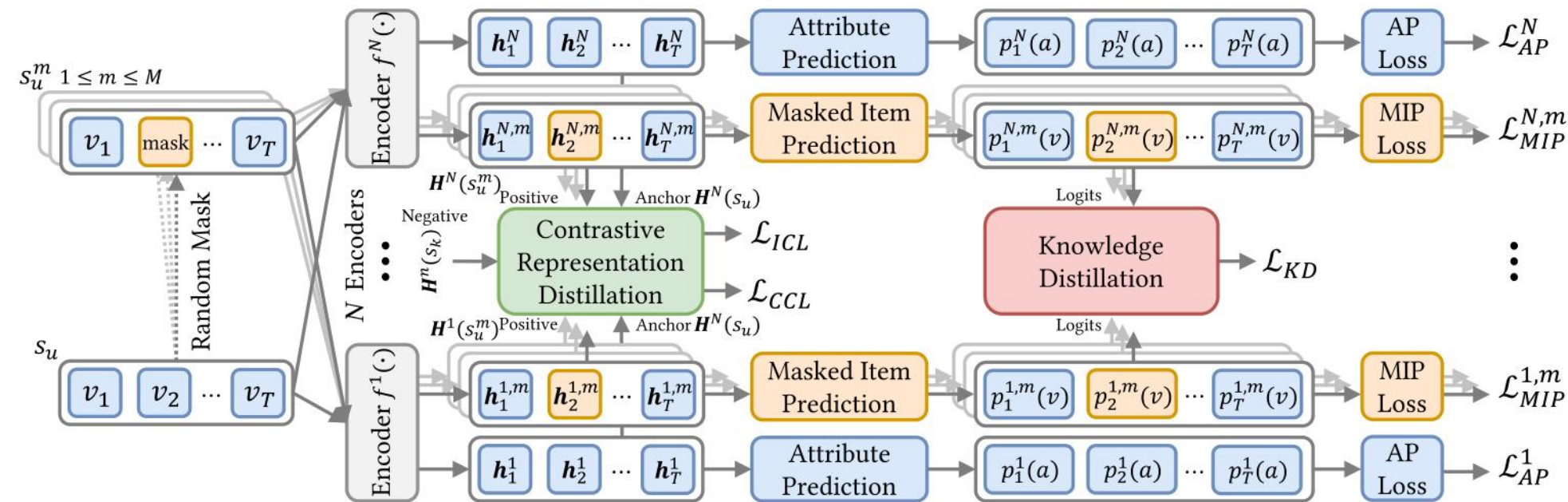


Figure 2: An overview of EMKD with  $N$  parallel networks  $f^1(\cdot), \dots, f^N(\cdot)$ . For each original sequence  $s_u$ , we generate  $M$  different masked sequences. The hidden representations of the original sequence  $H^1(s_u), \dots, H^N(s_u)$  serve as the anchor for contrastive representation distillation and are used for the attribute prediction task, while the hidden representations of the masked sequences  $H^1(s_u^m), \dots, H^N(s_u^m)$  serve as positive samples for contrastive representation distillation and are used for the masked item prediction task. Negative samples  $H^n(s_k)$  ( $1 \leq n \leq N$ ) for contrastive representation distillation are collected from the same batch. We compute the Kullback-Leibler divergence on the logits of the masked item prediction task between different networks for knowledge distillation.

## Method



$$E(s_u) = [v_1 + p_1, v_2 + p_2, \dots, v_T + p_T] \quad (1)$$

$$H(s_u) = [h_1, h_2, \dots, h_T] = f(s_u) = Trm(E(s_u)) \quad (2)$$

$$s_u^m = [\hat{v}_1, \hat{v}_2, \dots, \hat{v}_T], 1 \leq m \leq M \quad (3)$$

$$\hat{v}_t = \begin{cases} v_t, & t \notin \mathcal{I}_u^m \\ [\text{mask}], & t \in \mathcal{I}_u^m \end{cases} \quad (3)$$

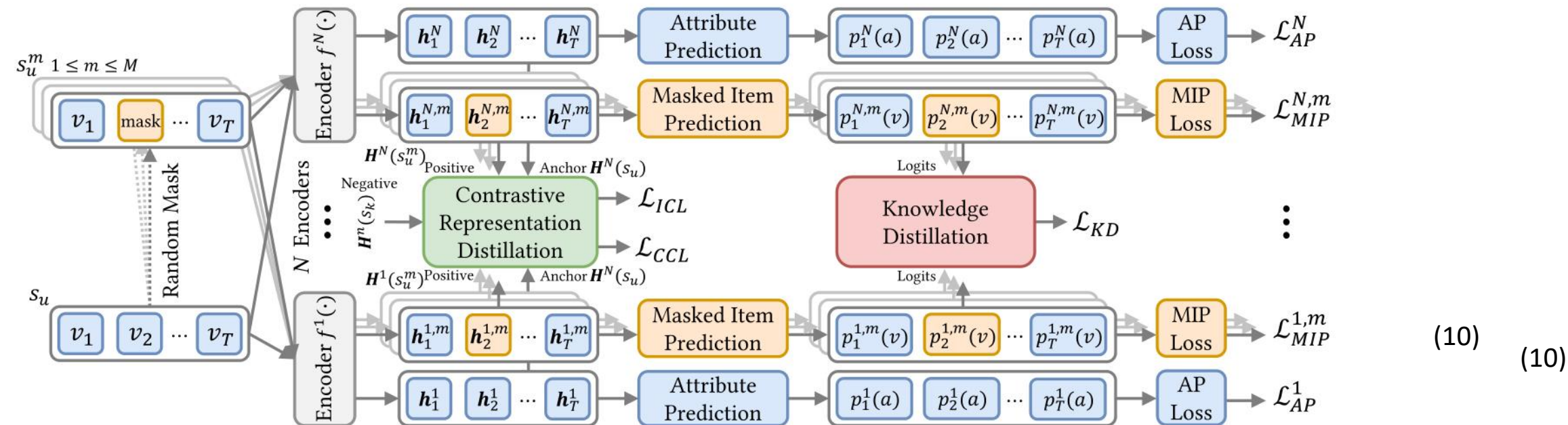
$$H^n(s_u^m) = [h_1^{n,m}, h_2^{n,m}, \dots, h_T^{n,m}] = f^n(s_u^m), 1 \leq n \leq N, 1 \leq m \leq M \quad (4)$$

$$p_t^{n,m}(v) = h_t^{n,m} \mathbf{W} + \mathbf{b}, 1 \leq t \leq T \quad (5)$$

$$q_t^{n,m}(v_i) = \frac{\exp(p_t^{n,m}(v_i))}{\sum_{j=1}^{|\mathcal{V}|} \exp(p_t^{n,m}(v_j))} \quad (6)$$

$$\mathcal{L}_{MIP}^{n,m} = - \sum_{t \in \mathcal{I}_u^m} \sum_{i=1}^{|\mathcal{V}|} y_i \log q_t^{n,m}(v_i)$$

## Method



$$\mathcal{L}_{MIP} = \sum_{n=1}^N \sum_{m=1}^M \mathcal{L}_{MIP}^{n,m} \quad (7)$$

$$\mathcal{L}_{AP} = \sum_{n=1}^N \mathcal{L}_{AP}^n \quad (11)$$

$$\mathbf{H}^n(s_u) = [\mathbf{h}_1^n, \mathbf{h}_2^n, \dots, \mathbf{h}_T^n] = f^n(s_u), 1 \leq n \leq N \quad (8)$$

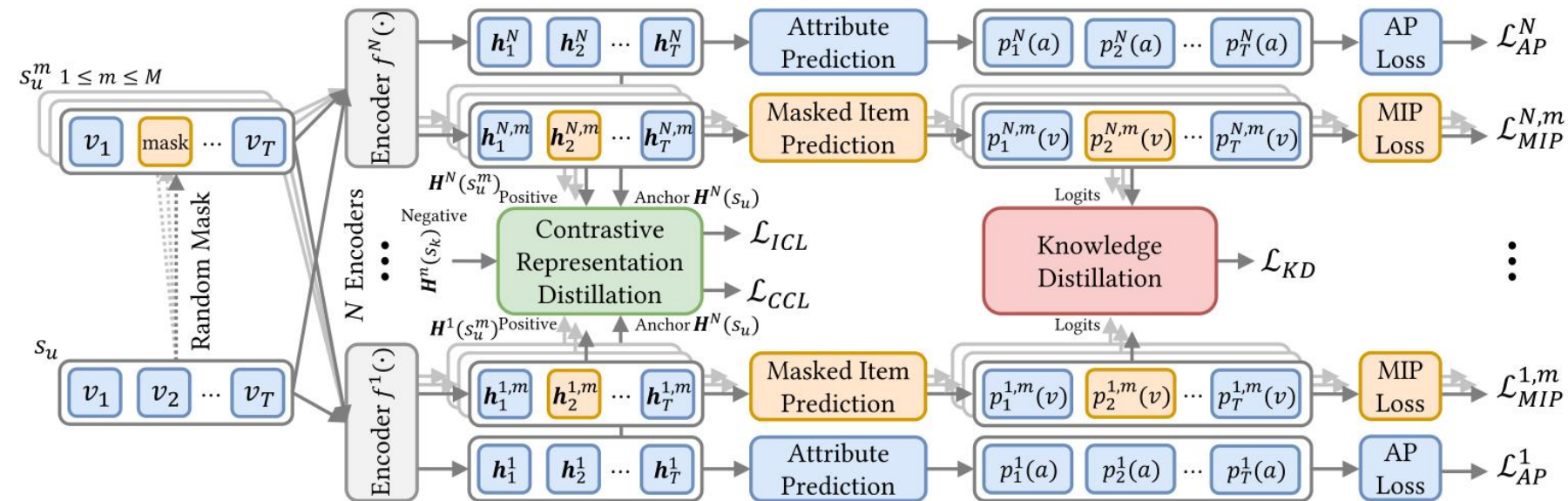
$$\mathcal{L}_{ICL}^n = - \sum_{m=1}^M \log \frac{\exp(g(\mathbf{H}^n(s_u), \mathbf{H}^n(s_u^m))/\tau)}{\sum_{k=1, k \neq u}^{\mathcal{B}} \exp(g(\mathbf{H}^n(s_u), \mathbf{H}^n(s_k))/\tau)} \quad (12)$$

$$p_t^n(a) = \sigma(\mathbf{h}_t^n \mathbf{W}' + \mathbf{b}'), 1 \leq t \leq T \quad (9)$$

$$\mathcal{L}_{AP}^n = - \sum_{t=1}^T \sum_{i=1}^{|\mathcal{A}|} [y_i \log p_t^n(a_i) + (1 - y_i) \log(1 - p_t^n(a_i))] \quad (10)$$

$$\mathcal{L}_{ICL} = \sum_{n=1}^N \mathcal{L}_{ICL}^n \quad (13)$$

## Method



$$\mathcal{L}_{CCL}^{x,y} = - \sum_{m=1}^M \log \frac{\exp(g(\mathbf{H}^x(s_u), \mathbf{H}^y(s_u^m))/\tau)}{\sum_{k=1, k \neq u}^B \exp(g(\mathbf{H}^x(s_u), \mathbf{H}^y(s_k))/\tau)} \quad (14)$$

$$\mathcal{L}_{CCL} = \sum_{x=1}^N \sum_{y=1, y \neq x}^N \mathcal{L}_{CCL}^{x,y} \quad (15)$$

$$z_t^{x,m}(v_i) = \frac{\exp(p_t^{x,m}(v_i)/\tau)}{\sum_{j=1}^{|\mathcal{V}|} \exp(p_t^{x,m}(v_j)/\tau)}, z_t^{y,m}(v_i) = \frac{\exp(p_t^{y,m}(v_i)/\tau)}{\sum_{j=1}^{|\mathcal{V}|} \exp(p_t^{y,m}(v_j)/\tau)} \quad (16)$$

$$\mathcal{L}_{KD}^{x,y} = \sum_{m=1}^M \sum_{t \in \mathcal{I}_u^m} \sum_{i=1}^{|\mathcal{V}|} z_t^{x,m}(v_i) \log \frac{z_t^{x,m}(v_i)}{z_t^{y,m}(v_i)} \quad (17)$$

$$\mathcal{L}_{KD} = \sum_{x=1}^N \sum_{y=1, y \neq x}^N \mathcal{L}_{KD}^{x,y} \quad (18)$$

$$\mathcal{L}_{EMKD} = \mathcal{L}_{MIP} + \mathcal{L}_{AP} + \lambda(\mathcal{L}_{ICL} + \mathcal{L}_{CCL}) + \mu \mathcal{L}_{KD} \quad (19)$$

$$\mathbf{H}^n(s_{u'}) = [h_2^n, h_3^n, \dots, h_T^n, h_{\text{mask}}^n] = f^n(s_{u'})$$

$$p(v) = \frac{1}{N} \sum_{n=1}^N (\mathbf{h}_{\text{mask}}^n \mathbf{W} + \mathbf{b}) \quad (20)$$



# Experiments

**Table 1: Performance comparison (NDCG@10) between the original model and the ensemble models. We independently train two parallel networks initialized with different random seeds and compare the result with the original model.**

Model	GRU4Rec		Caser		SASRec	
	Original	Ensemble(2×)	Original	Ensemble(2×)	Original	Ensemble(2×)
Beauty	0.0175	0.0199	0.0212	0.0247	0.0284	0.0365
Toys	0.0097	0.0102	0.0168	0.0193	0.0320	0.0378
ML-1M	0.0649	0.0720	0.0734	0.0786	0.0918	0.1032



# Experiments

**Table 2: Dataset statistics after preprocessing.**

Datasets	Beauty	Toys	ML-1M
#users	22,363	19,412	6,040
#items	12,101	11,924	3,953
#actions	198,502	167,597	1,000,209
avg. actions/user	8.9	8.6	163.5
avg. actions/item	16.4	14.1	253.0
sparsity	99.93%	99.93%	95.81%
#attributes	1,221	1,027	18
avg. attributes/item	5.1	4.3	1.7





# Experiments

**Table 3: Overall performance of different methods for sequential recommendation. The best score and the second-best score in each row are bolded and underlined, respectively. The last column indicates improvements over the best baseline method.**

Dataset	Metric	GRU4Rec	Caser	SASRec	BERT4Rec	FDSA	S <sup>3</sup> -Rec	MMInfoRec	CL4SRec	DuoRec	EMKD	Improv.
Beauty	HR@5	0.0206	0.0254	0.0371	0.0364	0.0317	0.0382	0.0527	0.0396	<u>0.0559</u>	<b>0.0702</b>	25.58%
	HR@10	0.0332	0.0436	0.0592	0.0583	0.0496	0.0634	0.0739	0.0630	<u>0.0867</u>	<b>0.0995</b>	14.76%
	NDCG@5	0.0139	0.0154	0.0233	0.0228	0.0184	0.0244	<u>0.0378</u>	0.0232	0.0331	<b>0.0500</b>	32.28%
	NDCG@10	0.0175	0.0212	0.0284	0.0307	0.0268	0.0335	<u>0.0445</u>	0.0307	0.0430	<b>0.0594</b>	33.48%
Toys	HR@5	0.0121	0.0205	0.0429	0.0371	0.0269	0.0440	<u>0.0579</u>	0.0503	0.0539	<b>0.0745</b>	28.67%
	HR@10	0.0184	0.0333	0.0652	0.0524	0.0483	0.0705	<u>0.0818</u>	0.0736	0.0744	<b>0.1016</b>	24.21%
	NDCG@5	0.0077	0.0125	0.0248	0.0259	0.0227	0.0286	<u>0.0408</u>	0.0264	0.0340	<b>0.0534</b>	30.88%
	NDCG@10	0.0097	0.0168	0.0320	0.0309	0.0281	0.0369	<u>0.0484</u>	0.0339	0.0406	<b>0.0622</b>	28.51%
ML-1M	HR@5	0.0806	0.0912	0.1078	0.1308	0.0953	0.1128	0.1454	0.1142	<u>0.1930</u>	<b>0.2315</b>	19.95%
	HR@10	0.1344	0.1442	0.1810	0.2219	0.1645	0.1969	0.2248	0.1815	<u>0.2865</u>	<b>0.3239</b>	13.05%
	NDCG@5	0.0475	0.0565	0.0681	0.0804	0.0597	0.0668	0.0856	0.0705	<u>0.1327</u>	<b>0.1616</b>	21.78%
	NDCG@10	0.0649	0.0734	0.0918	0.1097	0.0864	0.0950	0.1203	0.0920	<u>0.1586</u>	<b>0.1915</b>	20.74%



# Experiments

**Table 4: Ablation study (NDCG@10) on three datasets. Bold score indicates the performance under the default setting.  $\uparrow$  indicates the performance better than the default setting.**

Architecture	Dataset		
	Beauty	Toys	ML-1M
(1) EMKD( $\times 3$ )	<b>0.0594</b>	<b>0.0622</b>	<b>0.1915</b>
(2) Remove ICL	0.0529	0.0545	0.1679
(3) Remove CCL	0.0552	0.0560	0.1807
(4) Remove KD	0.0537	0.0571	0.1758
(5) Independent Training	0.0452	0.0484	0.1476
(6) Single Encoder	0.0363	0.0375	0.1183
(7) EMKD( $\times 2$ )	0.0536	0.0568	0.1792
(8) EMKD( $\times 4$ )	0.0591	0.0629 $\uparrow$	0.1930 $\uparrow$
(9) Remove AP	0.0578	0.0609	0.1831

**Table 5: Performance comparison (NDCG@10) of models with different parameter sizes on three datasets. \* indicates the default setting for each model.**

Architecture	Beauty		Toys		ML-1M	
	Params.	NDCG@10	Params.	NDCG@10	Params.	NDCG@10
SASRec-2 Layers*	4.69M	0.0284	4.65M	0.0320	2.51M	0.0918
SASRec-4 Layers	6.27M	0.0301	6.23M	0.0313	4.09M	0.0896
SASRec-6 Layers	7.85M	0.0298	7.80M	0.0332	5.67M	0.0857
SASRec-8 Layers	9.43M	0.0279	9.38M	0.0305	7.24M	0.0932
SASRec-10 Layers	11.01M	0.0282	10.96M	0.0310	8.82M	0.0881
BERT4Rec-2 Layers*	7.80M	0.0307	7.71M	0.0309	3.53M	0.1097
BERT4Rec-4 Layers	9.38M	0.0328	9.29M	0.0312	5.11M	0.1113
BERT4Rec-6 Layers	10.96M	0.0332	10.87M	0.0306	6.69M	0.1100
BERT4Rec-8 Layers	12.54M	0.0310	12.45M	0.0298	8.27M	0.1093
BERT4Rec-10 Layers	14.12M	0.0319	14.03M	0.0293	9.85M	0.1099
EMKD( $\times 2$ )	9.36M	0.0536	9.28M	0.0568	5.08M	0.1792
EMKD( $\times 3$ )*	14.05M	0.0594	13.91M	0.0622	7.62M	0.1915

# Experiments

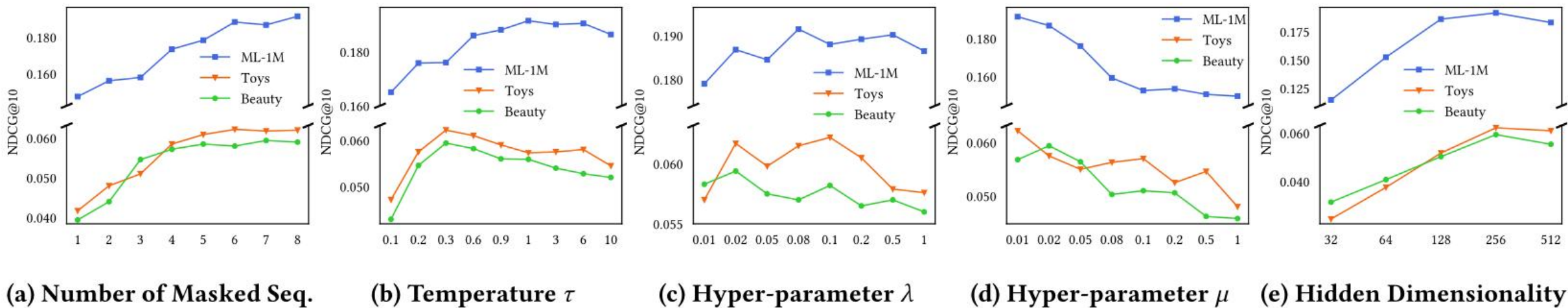
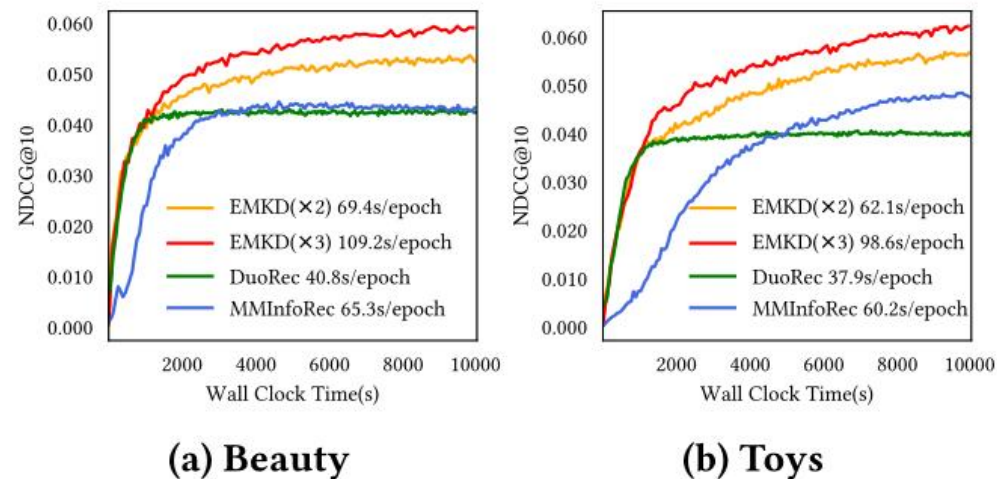
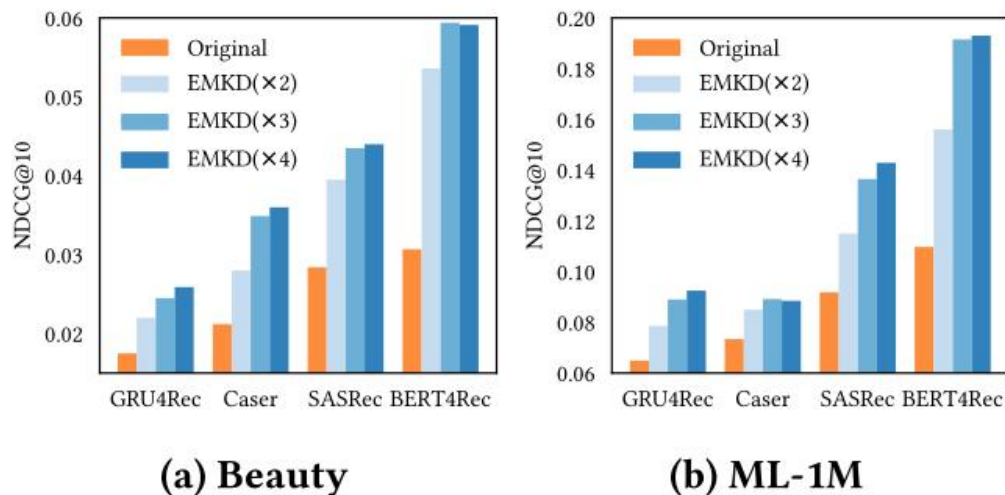


Figure 3: Performance (NDCG@10) comparison w.r.t different hyper-parameters on three datasets.

# Experiments



**Figure 4: Performance comparison (NDCG@10) of different models enhanced by EMKD on Beauty and ML-1M datasets. We design three variants for each group of base sequence encoder with 2,3,4 parallel networks respectively.**

**Figure 5: Training efficiency (NDCG@10) on Beauty and Toys datasets. The training speed of EMKD is slightly lower than MMInfoRec, while the convergence speed of EMKD is comparable with DuoRec.**



**Thanks**