

# Recommender Forest for Efficient Retrieval

Chao Feng<sup>1\*</sup>, Wuchao Li<sup>1\*</sup>, Defu Lian<sup>1</sup>, Zheng Liu<sup>2</sup>, Enhong Chen<sup>1</sup>

<sup>1</sup>School of Computer Science and Technology  
University of Science and Technology of China, Hefei, China

<sup>2</sup>Microsoft Research Asia, Beijing, China

{chaofeng, liwuchao}@mail.ustc.edu.cn

{liandefu, cheneh}@ustc.edu.cn, zhengliu@microsoft.com



中国科学技术大学

University of Science and Technology of China



# Indexes in recommender system

- Recommender system needs to select the **top-n** items for users from a **massive-scale** item set.
- For the sake of efficient recommendation, RS usually calls for the collaboration of representation learning and Approximate Nearest Neighbour search (ANNs) index.
- Two kinds of indexes: **Independent training** & **Joint training**



# Independent training index

- Process
  1. Users and items are represented by embeddings in the same latent space.
  2. The item embeddings are organized with a specific ANNs index, like SCANN and HNSW.
- Limitation
  - The representation model is independently learned and can be incompatible with the ANNs index



# Joint training index: TDM & JTM

- The item set is organized with binary tree structure:
  - internal node: cluster center.
  - leaf node: item.
- A preference model is learned to route from the root to the leaf nodes for the top-n recommendation results.
- Achieve empirical gains over the conventional two-stage methods.

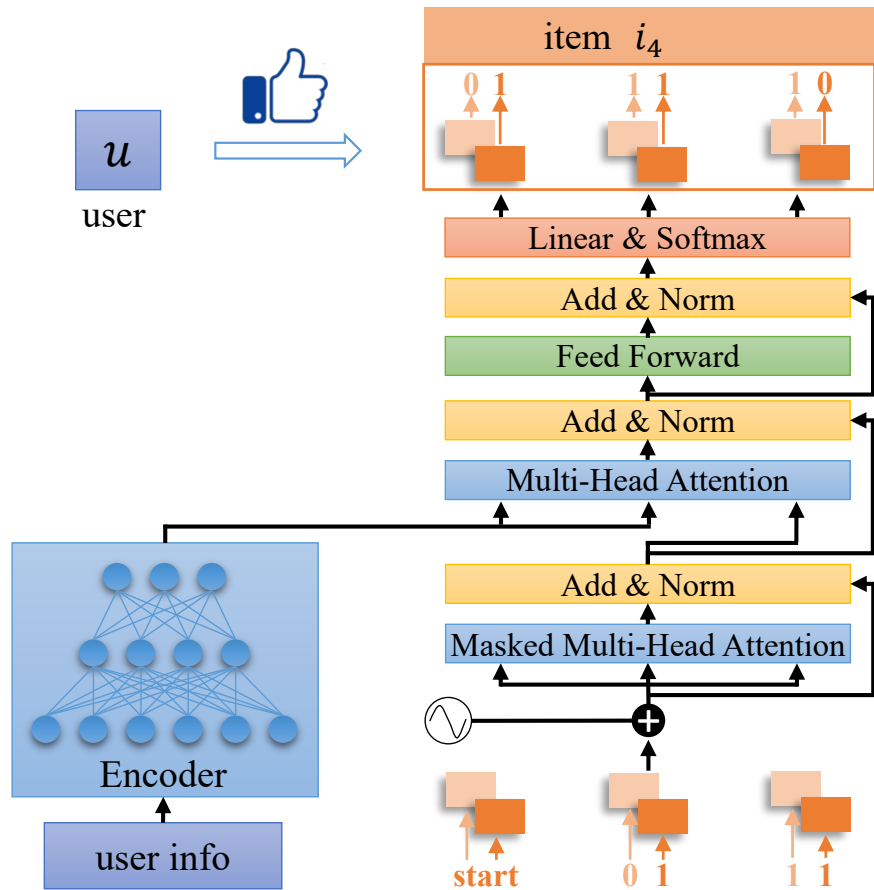


# Joint training index: TDM & JTM

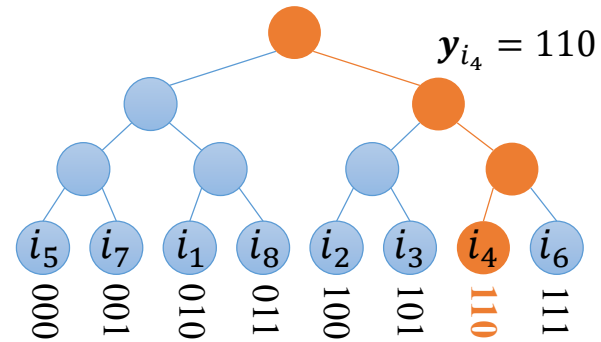
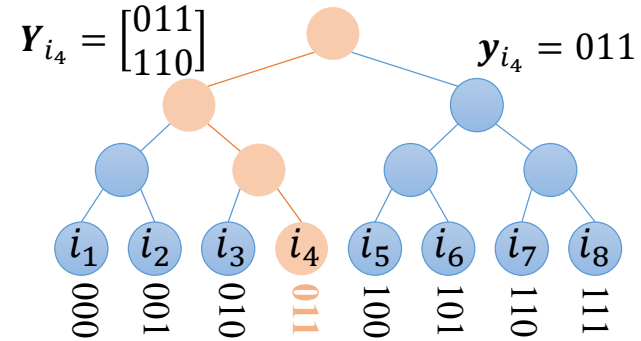
- Limitation
  - It is challenging to route to items located around partition boundaries.
  - A routing decision is made without consideration of the routing trajectory.
  - Memory-consuming, given that the number of internal nodes is at the same magnitude as the leaf nodes.



# Introducing RecForest



Each item is represented by a sequence in each tree



The sequence length is  $\lceil \log_k N \rceil$

# Construction of Tree and Forest

- Tree construction (Balanced-Kmeans)
  - Pre-train item embeddings: DIN or BPR
  - Randomly sample  $K^H - N$  items ( $H$  is the height of the tree) from the entire item set.
  - For each cluster, the included items are evenly partitioned via Kmeans w.r.t. their embedding similarity.
- Forest construction
  - Different trees are naturally diversified due to the inherent randomness.



# Experiments



# Comparison with Baselines

	NDCG@20	NDCG@40	Memory	Time	NDCG@20	NDCG@40	Memory	Time
Method	Movie				Amazon			
DIN	<b>0.5440</b>	0.5473	-	193.87	<b>0.2766</b>	<b>0.3039</b>	-	492.64
YoutubeDNN	0.5329	<b>0.5484</b>	-	<b>29.38</b>	0.2195	0.2491	-	<b>120.91</b>
JTM	0.5149	0.5075	10.80	12.05	0.1533	0.1683	75.99	6.64
TDM	0.4684	0.4651	10.80	9.33	0.0856	0.0949	75.99	6.61
SCANN	0.4665	0.4695	3.64	18.64	0.1529	0.1780	14.66	4.48
IPNSW	0.5330	0.5486	10.08	15.52	0.2255	0.2548	66.46	10.28
RecForest	<b>0.5580</b>	<b>0.5682</b>	<b>3.21</b>	<b>8.33</b>	<b>0.2339</b>	<b>0.2576</b>	<b>7.32</b>	<b>3.79</b>
Method	Gowalla				Tmall			
DIN	<b>0.2798</b>	<b>0.3095</b>	-	186.41	<b>0.2275</b>	<b>0.2491</b>	-	4057.69
YoutubeDNN	0.2312	0.2637	-	<b>53.55</b>	0.1736	0.1975	-	<b>1086.75</b>
JTM	0.2595	0.2484	77.56	2.64	0.0749	0.0849	151.19	30.11
TDM	0.1723	0.1775	77.56	2.55	0.0257	0.0272	151.19	29.42
SCANN	0.1839	0.2083	15.48	1.86	0.1105	0.1226	28.10	20.88
IPNSW	0.2464	0.2805	70.39	4.73	0.1696	0.1902	132.72	52.90
RecForest	<b>0.3783</b>	<b>0.3963</b>	<b>7.39</b>	<b>1.82</b>	<b>0.2059</b>	<b>0.2261</b>	<b>9.29</b>	<b>18.88</b>

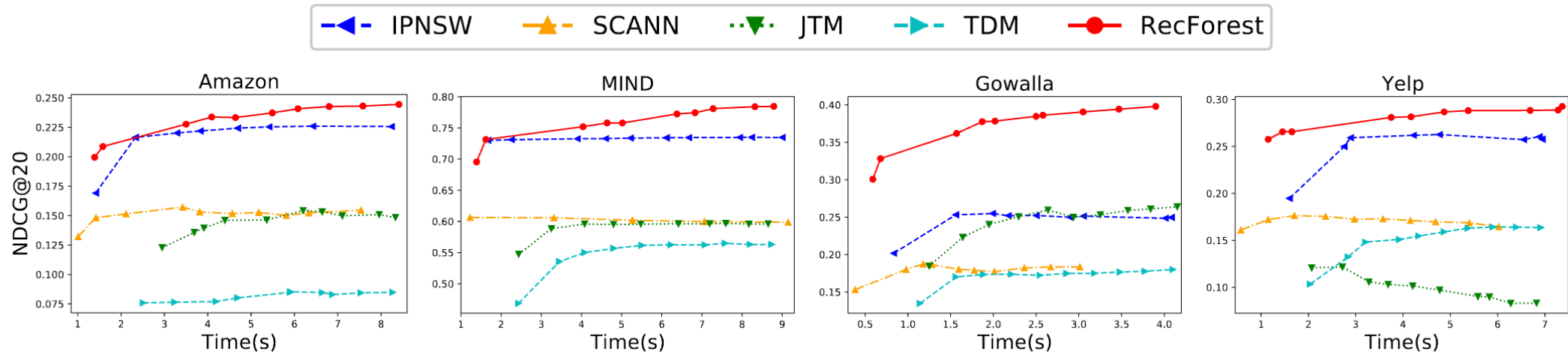
More accurate

Lighter

Faster



# Comparison with Baselines



- RecForest strikes the best balance between query time and retrieval accuracy.
- With the increase of beam size, the accuracy of RecForest can improve more significantly than baselines.

# Why is RecForest better ?

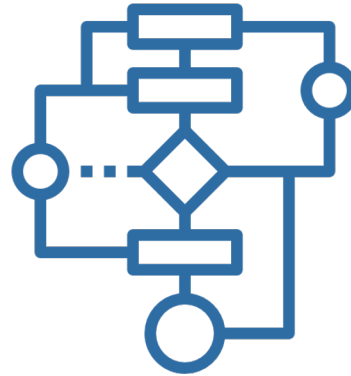
- The retrieval of the near-boundary can be challenging
  - multiple K-ary trees
- Not considering the trajectory history
  - using Transformer Decoder
- Memory-consuming:
  - There are mere K vectors (corresponding to the K different branches) in each K-ary tree.



# More content in the paper



Complexity analysis



Algorithm flow



More experiments

Thanks