

第二十七届全国信息检索学术会议

CCIR 2021

CCIR

中国·大连

2021年10月29-31日

Contrastive Learning: A Recommendation Perspective

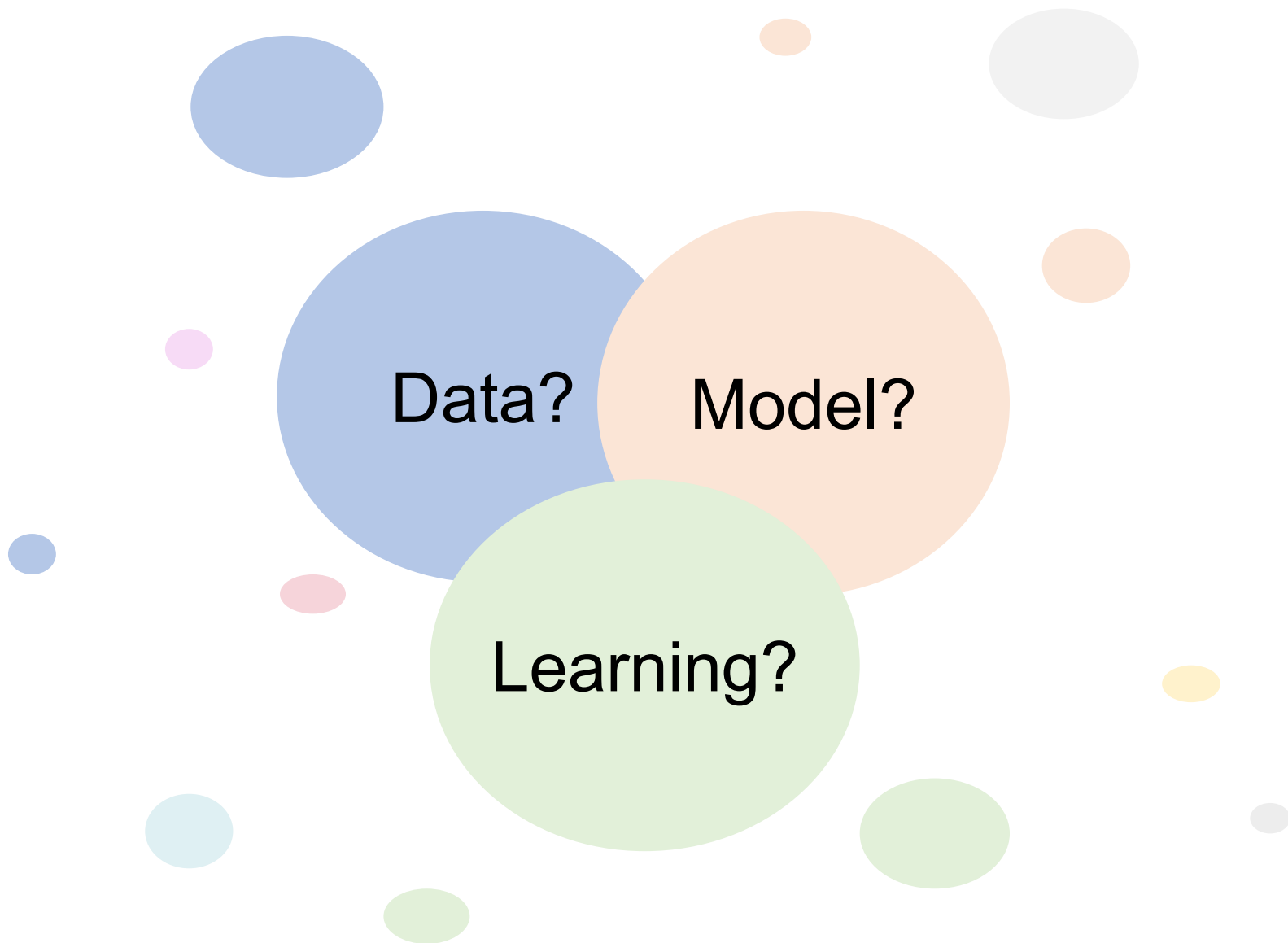
Pengjie Ren (任鹏杰)
IRLab, Shandong University

renpengjie@sdu.edu.cn



山东大学信息检索实验室
Information Retrieval Lab

What matters for DL?



What matters for DL? Data?



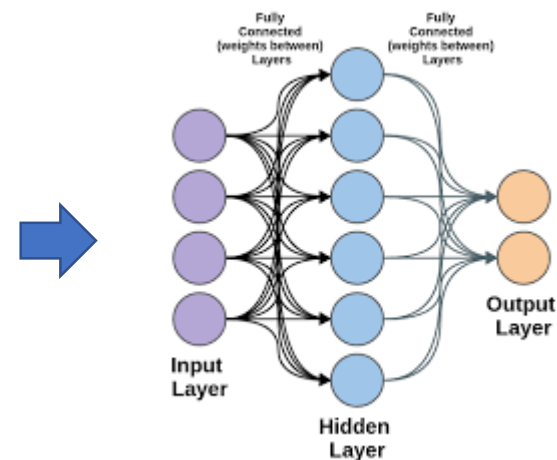
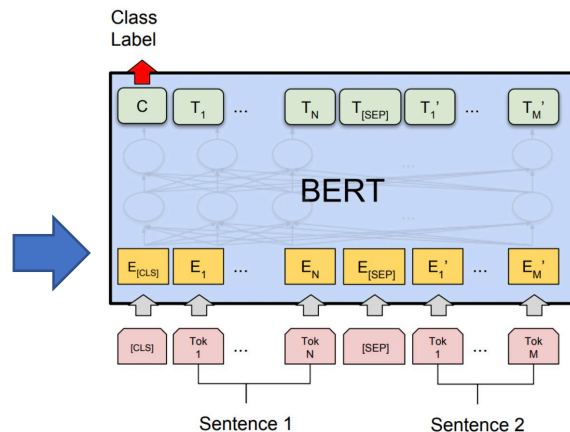
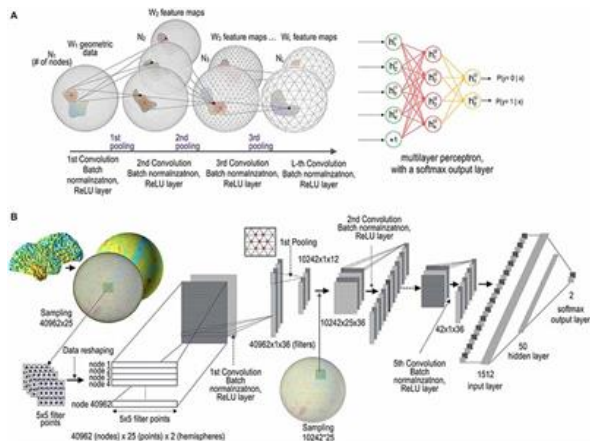
What matters for DL? Model?

Model is getting simpler.

Everything is you need.

Attention is all you need.

MLP is all you need.



Ashish Vaswani et al. Attention Is All You Need. NeurIPS 2017.

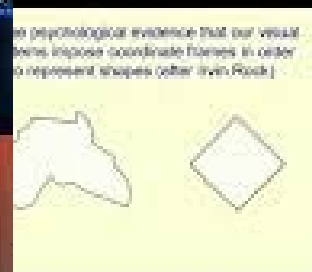
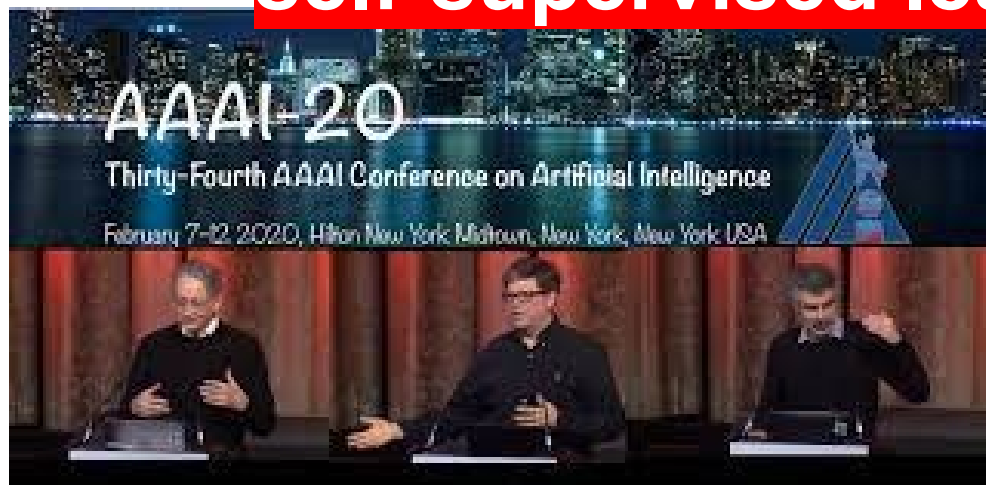
Ilya Tolstikhin et al. MLP-Mixer: An all-MLP Architecture for Vision. NeurIPS 2021.

Luke Melas-Kyriazi. Do You Even Need Attention? A Stack of Feed-Forward Layers Does Surprisingly Well on ImageNet. arXiv 2021.

What matters for DL? Learning!



self-supervised learning (SSL)



Do we have evidence?

MS MARCO Document Ranking

date	description
2021/04/25	🏆 PROP_step400K base + doc2query top1000(ensemble v0.2)
2021/04/25	Knowledge Retrieval

SQuAD2.0 tests the ability of a system to answer questions, but also abstain when presented with questions based on the provided paragraph.

Rank	Model	Date
1	Human Performance (Stanford University) (Rajpurbaney et al., 2018)	Feb 21, 2021
2	IE-Net (ensemble) (RICOH_SOURCE)	Feb 24, 2021
3	SA-Net on Alibaba QIANXUN	Apr 06, 2020
4	SA-Net-V2 (ensemble) (QIANXUN)	May 05, 2020

Rank	Model	F1	HEQQ	HEQD
	Human Performance (Choi et al. EMNLP '18)	81.1	100	100
1	RoR (Single model) Anonymous	74.9	72.2	16.4
2	EL-QA (Single model) JD AI Research	74.6	71.6	16.3

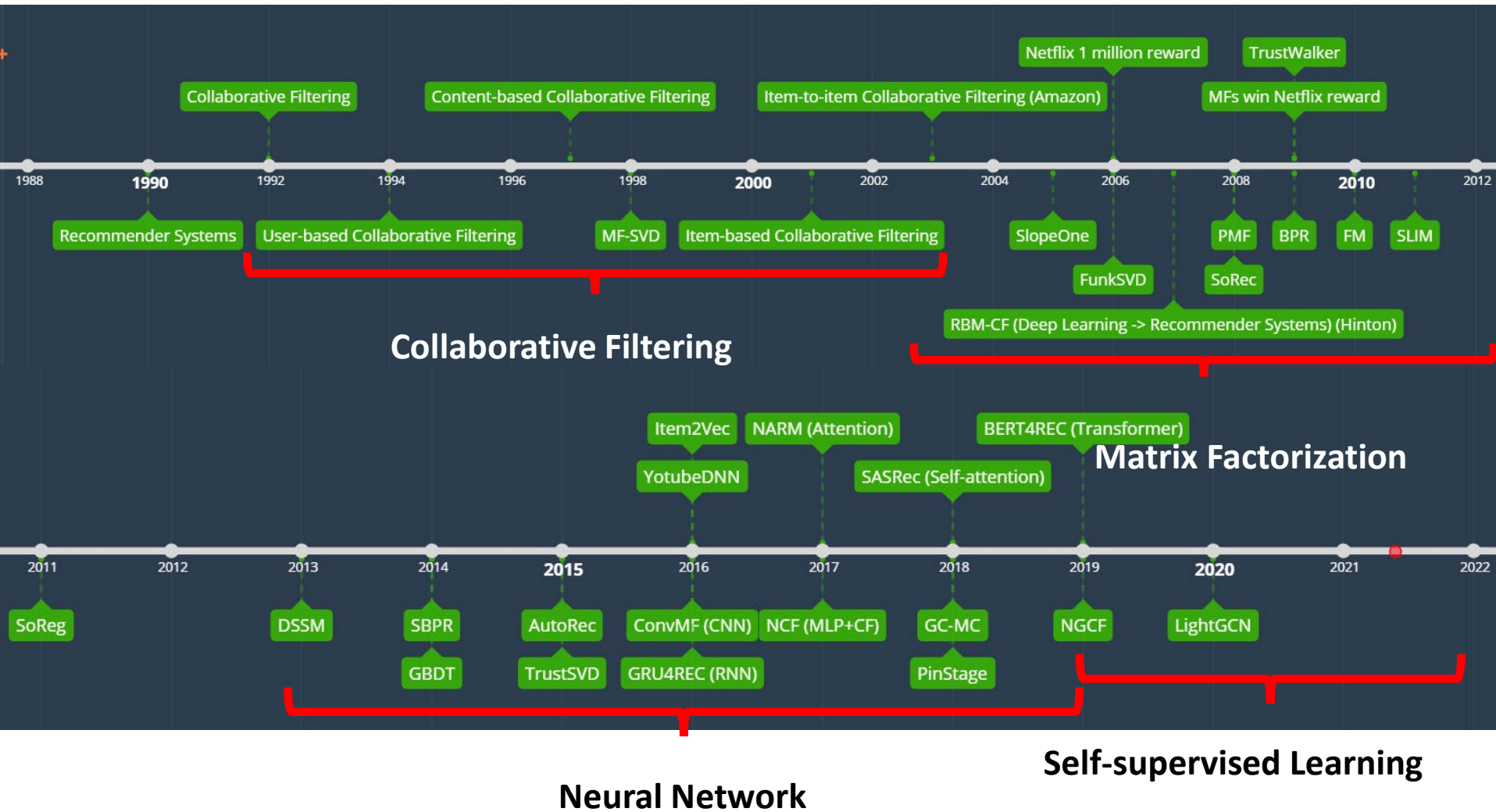
WMT 2014 EN-DE

Models are evaluated on the English-German dataset of the Ninth Workshop on Statistical Machine Translation (WMT 2014) based on BLEU.

Model	BLEU	Paper / Source
Transformer Big + BT (Edunovic et al., 2018)	35.0	Understanding Back-Translation at Scale
DeepL	33.3	DeepL Press release
Admin (Liu et al., 2020)	30.1	Very Deep Transformers for Neural Machine Translation
MUSE (Zhao et al., 2019)	29.9	MUSE: Parallel Multi-Scale Attention for Sequence to Sequence Learning
DynamicConv (Wu et al., 2019)	29.7	Pay Less Attention With Lightweight and Dynamic Convolutions

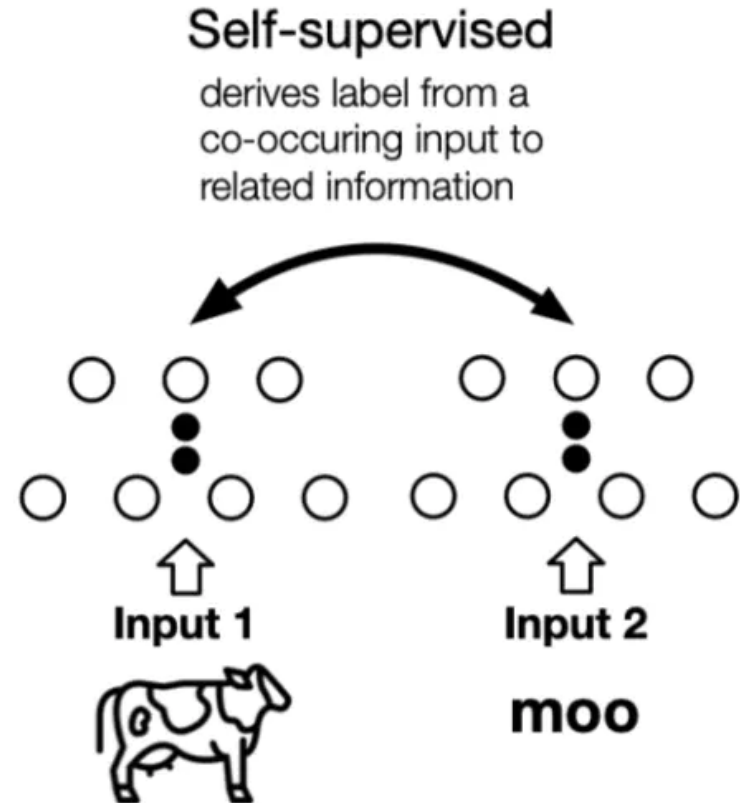
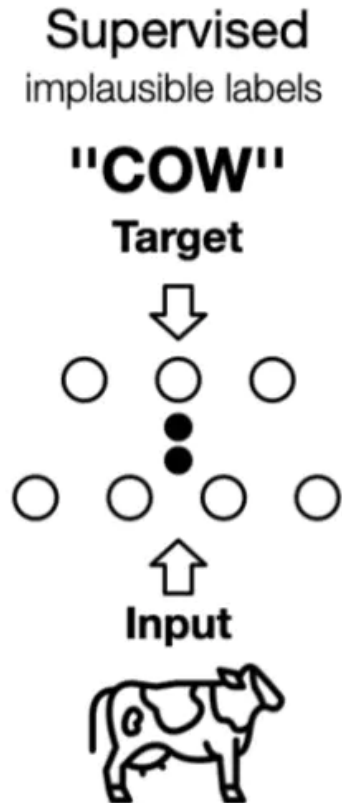
So where are we with SSL for recommendation?

A brief history of recommender systems

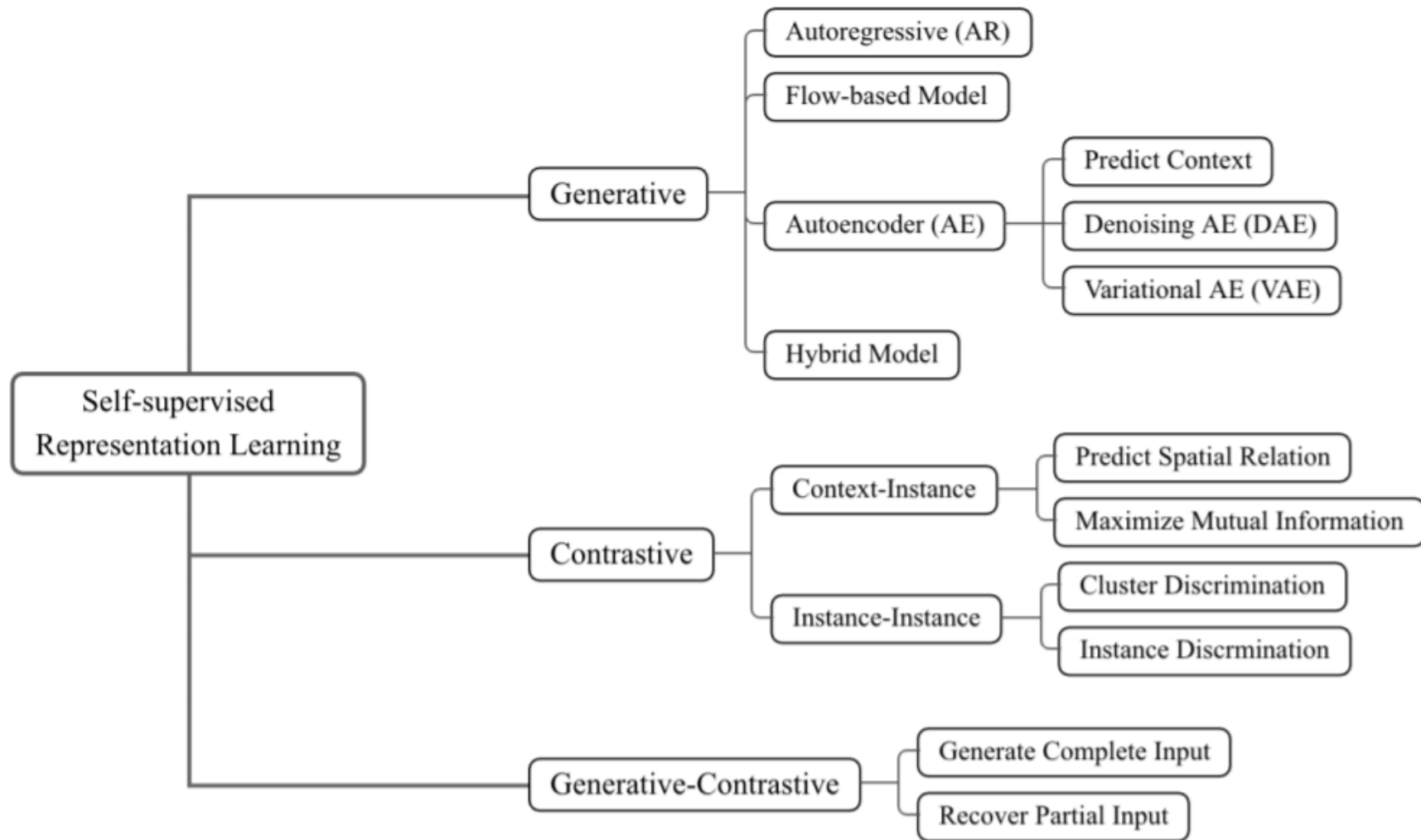


Then let's talk about SSL a bit...

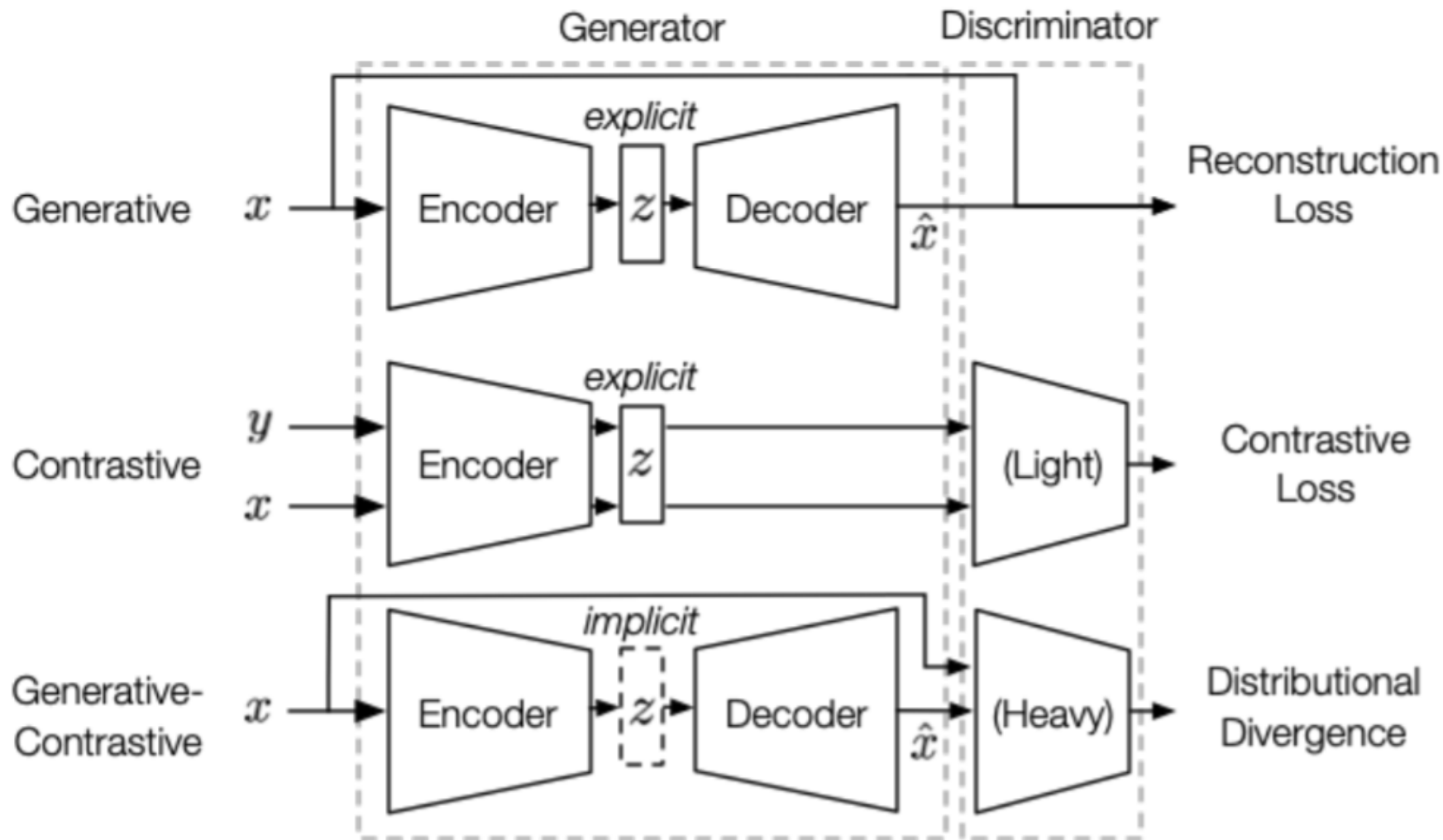
Supervised vs. unsupervised vs. self-supervised



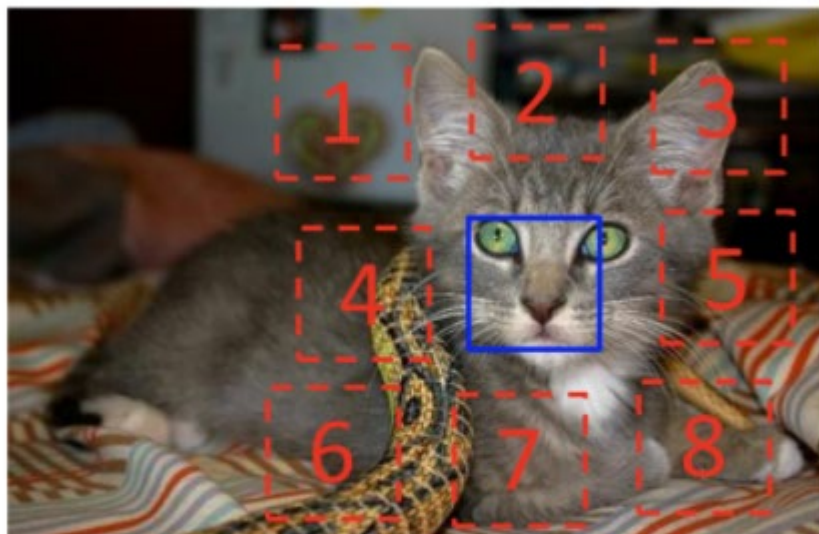
Self-supervised learning in general



Self-supervised learning in general

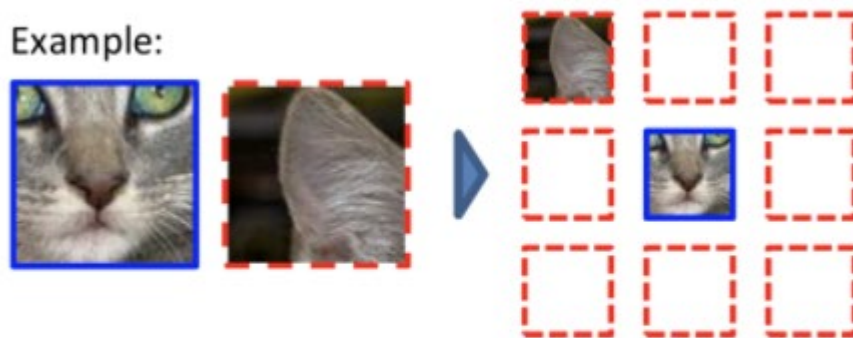


Contrastive self-supervised learning

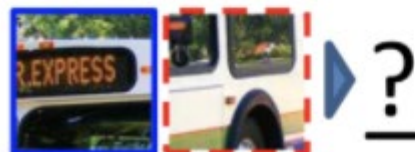


$$X = (\text{cat face}, \text{cat ear}); Y = 3$$

Example:



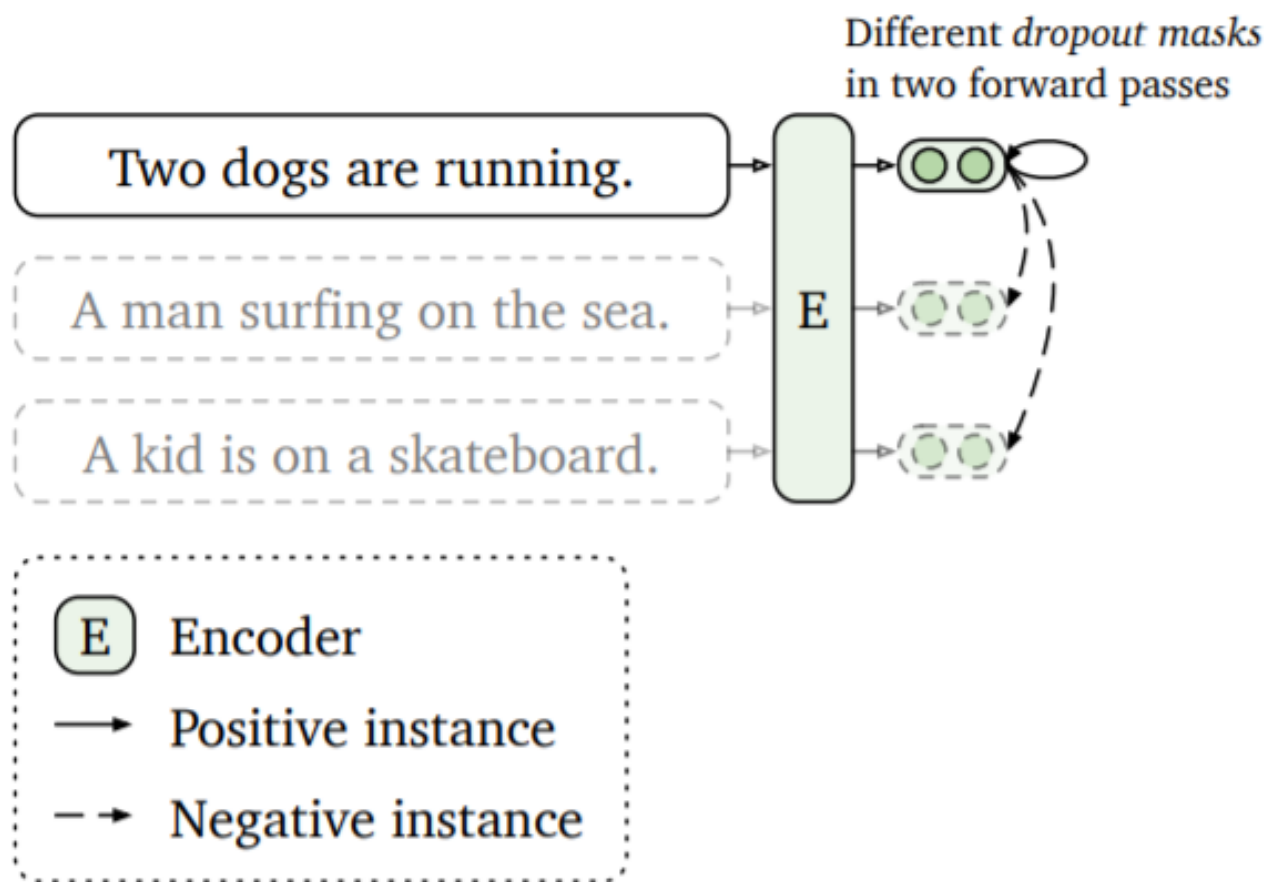
Question 1:



Question 2:

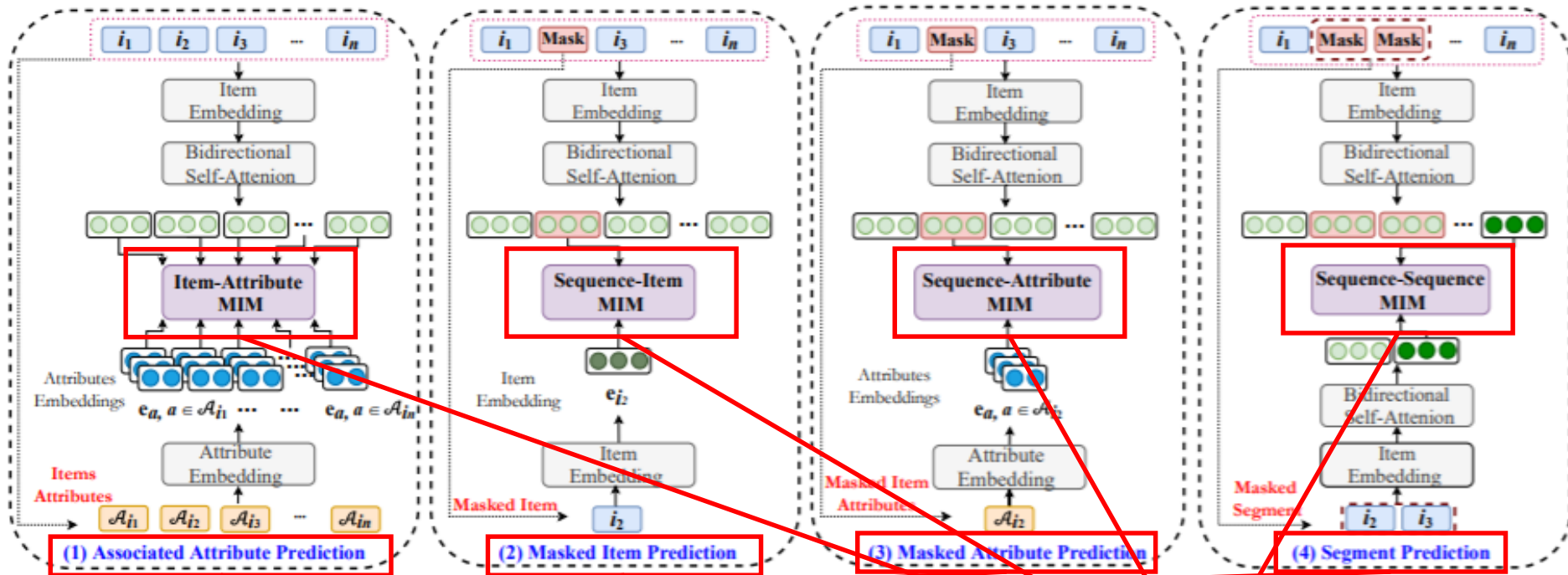


Contrastive self-supervised learning



S³-Rec

Next Item Prediction overemphasizes the final performance, the association or fusion between context data and sequence data has not been well captured and utilized for sequential recommendation.

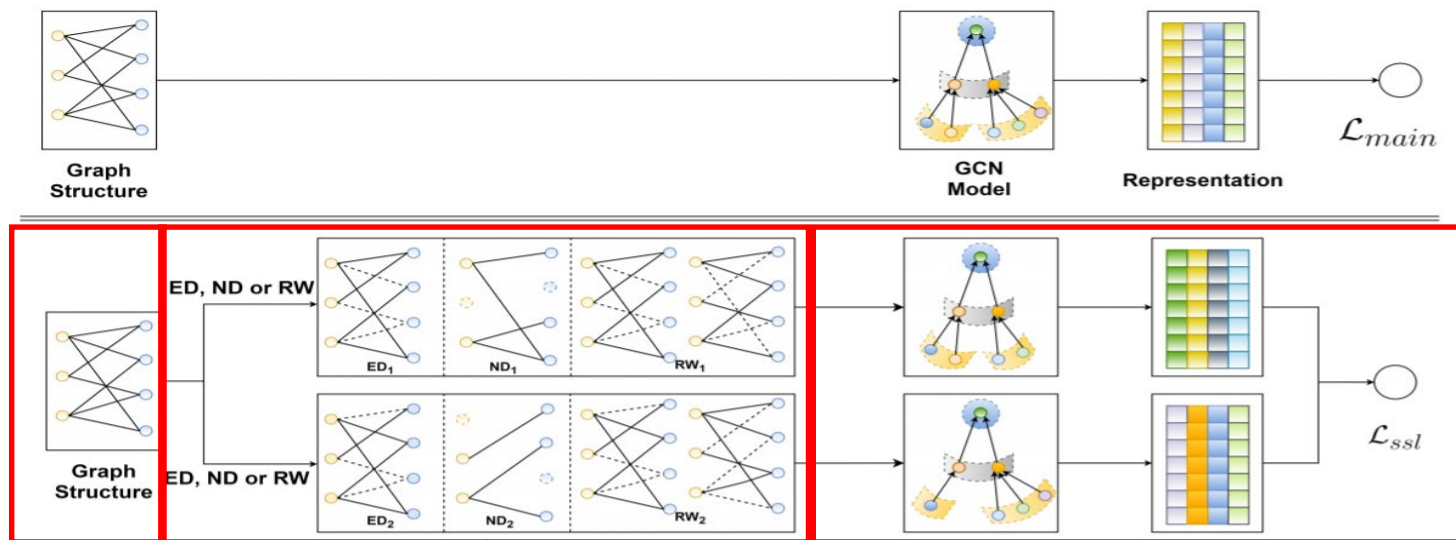


- (1) Associated Attribute Prediction (AAP)
- (2) Masked Item Prediction (MIP)
- (3) Masked Attribute Prediction (MAP)
- (4) Segment Prediction (SP)

Mutual Information Maximization

SGL

- (1) High-degree nodes exert larger impact on the representation learning.
- (2) Representations are vulnerable to noisy interactions.



Interaction between
users and items
→ Bipartite graph
→ Encoder: GCN

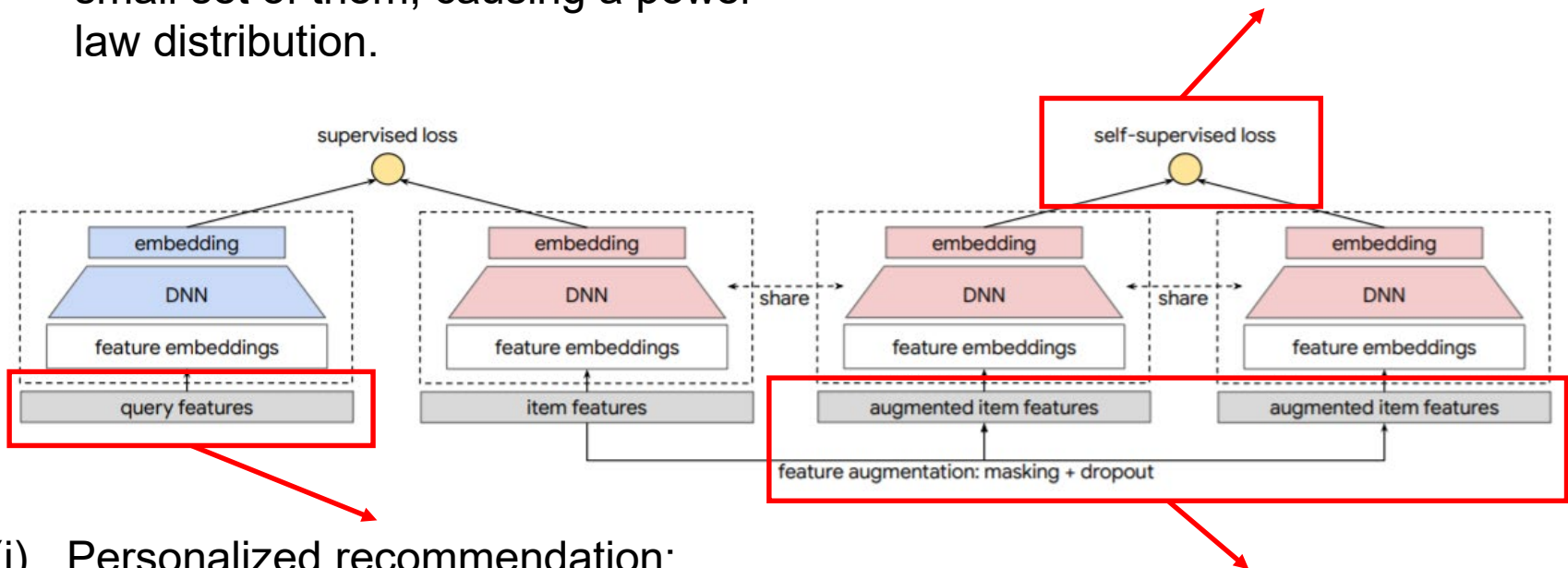
Augmentation:
1. Edge Dropout (ED)
2. Node Dropout (ND)
3. Random Walk (RW)

InfoNCE: Maximize the
agreement of positive
pairs and minimize that
of negative pairs

MSSL

Data sparsity: With millions to billions of items in the corpus, users tend to provide feedback for a very small set of them, causing a power-law distribution.

InfoNCE: Maximize the agreement of positive pairs and minimize that of negative pairs.



- (i) Personalized recommendation: when the query is a user;
- (ii) Item to item recommendation: when the query is also an item; and
- (iii) Search: when the query is a piece of free text.

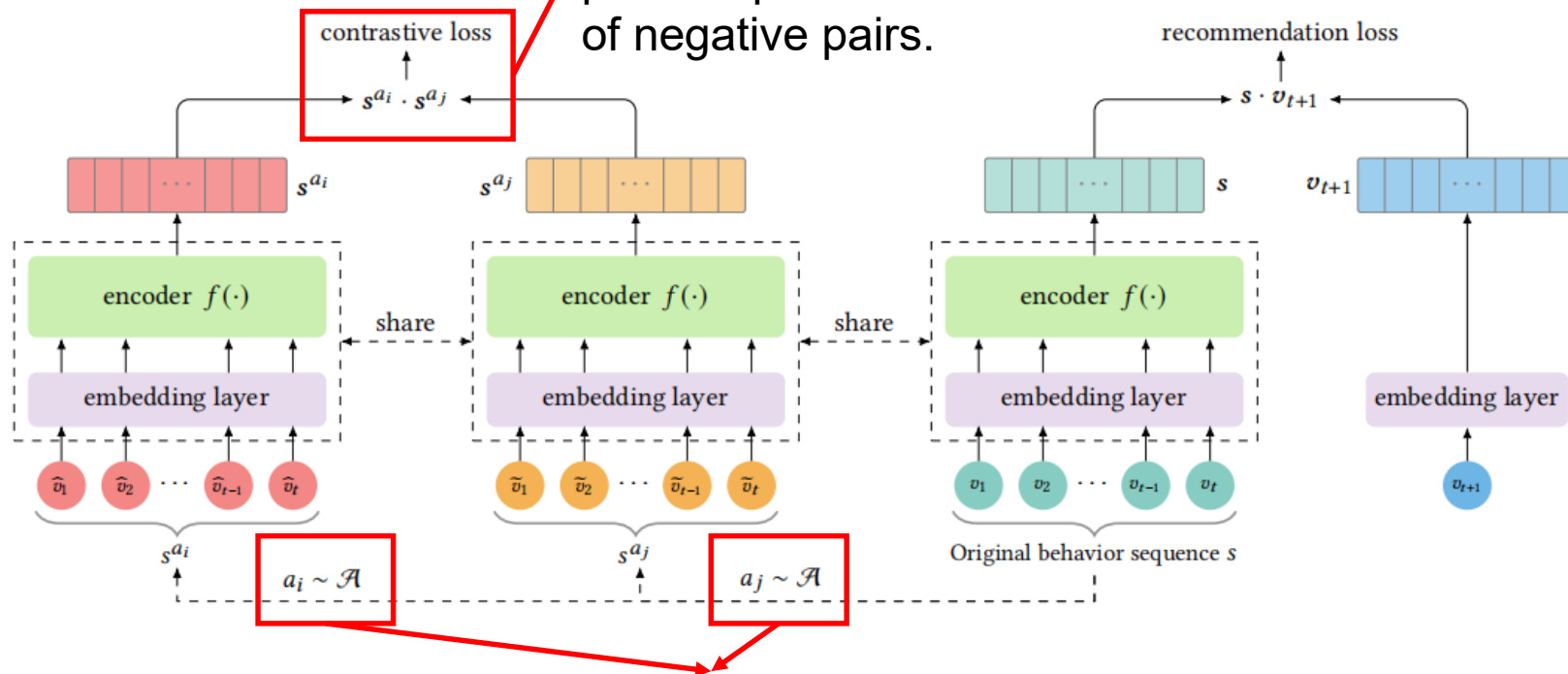
Augmentation:

1. Masking: Applying a masking pattern on the set of item features.
2. Dropout: Randomly dropouting categorical feature values.

CL4SRec

The data sparsity of SR makes it difficult to get high-quality user representations.

CE: Maximize the agreement of positive pairs and minimize that of negative pairs.



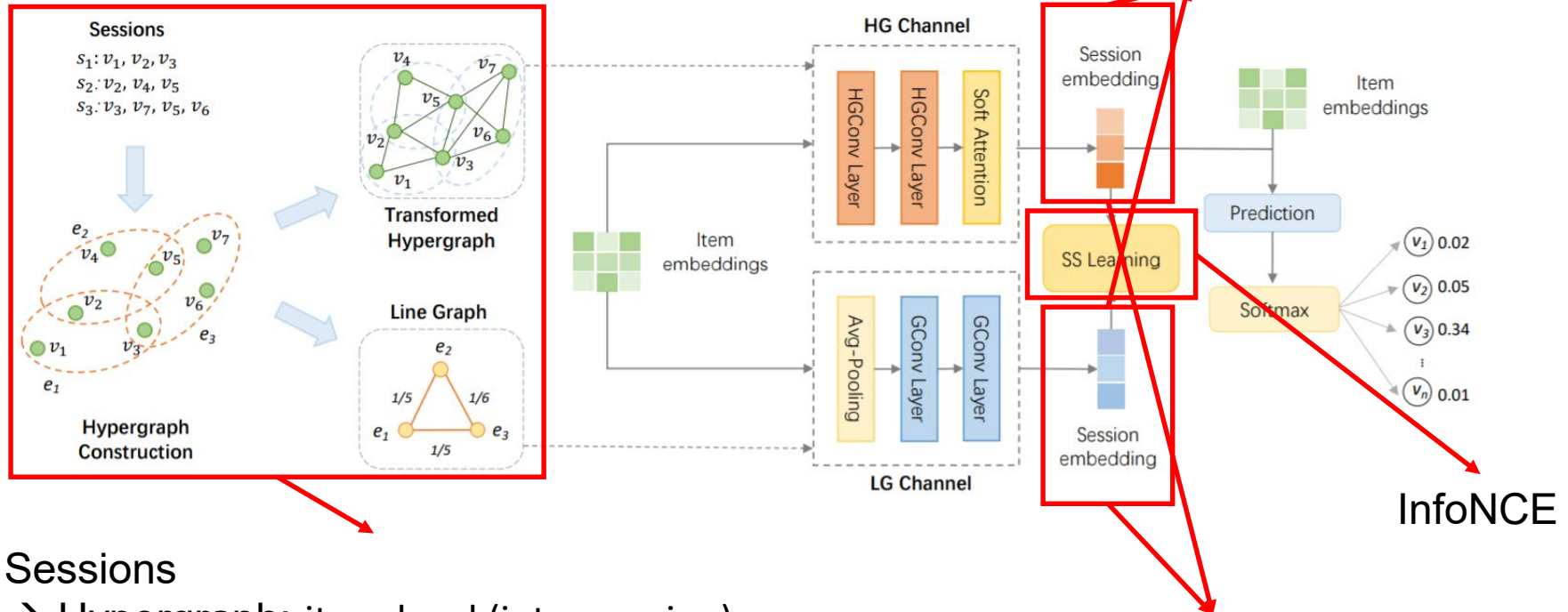
Augmentation:

1. Item crop: Randomly selecting a continuous sub-sequence.
2. Item mask: Randomly masking a proportion of items.
3. Item reorder: Randomly shuffling a proportion of items.

DHCN

Regarding the two channels as two views characterizing different aspects of sessions.

The two groups of embeddings know little about each other but can mutually complement.



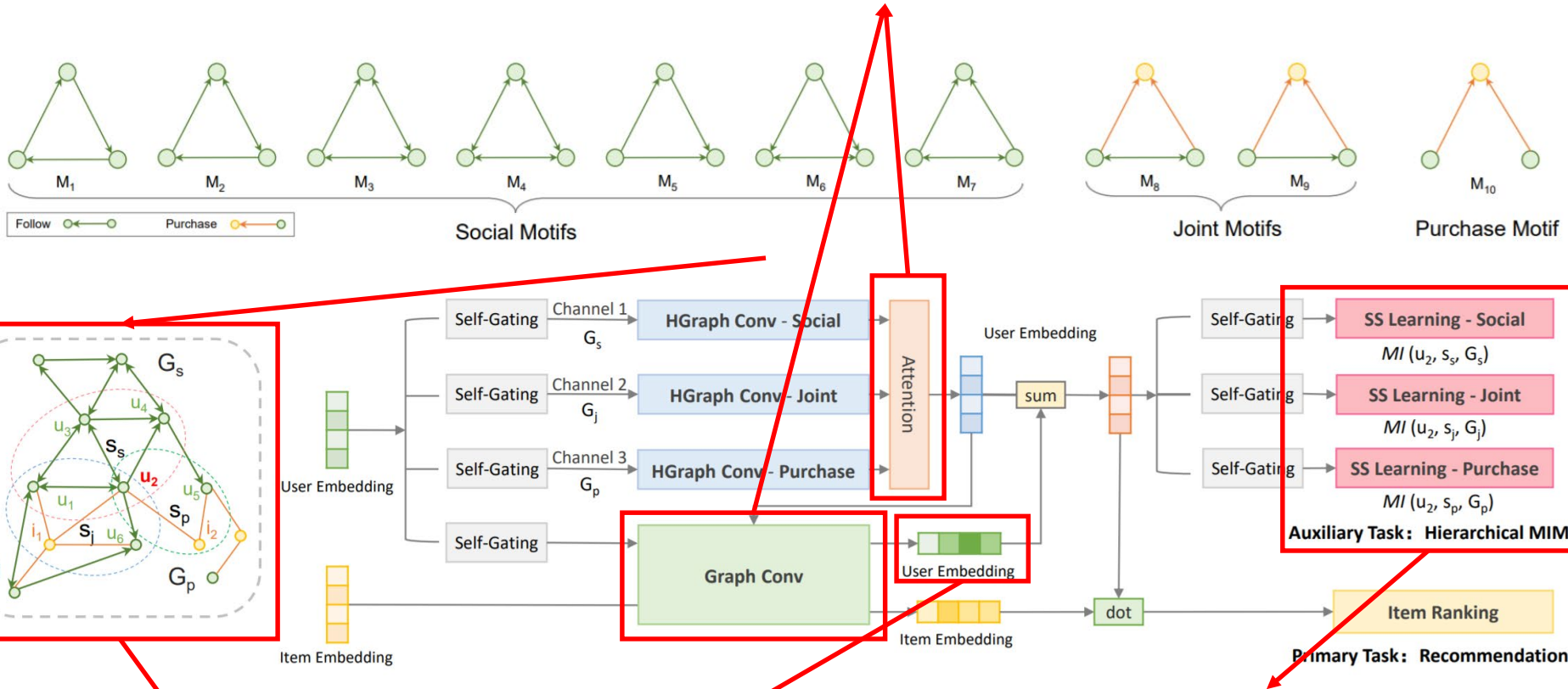
Sessions

- Hypergraph: item-level (intra-session) structural information
- Line graph: session-level (inter-session) structural information
- Encoder: GNN

Augmentation:
Corrupted session embedding, aka.
row-wise and column-wise shuffling.

MHCN

The aggregation operations might lead to a loss of high-order information.



social network and
user-item interaction
→ Motifs → Hypergraph
→ Encoder: GNN

Augmentation:
Corrupted user embedding,
aka. row-wise and column-
wise shuffling.

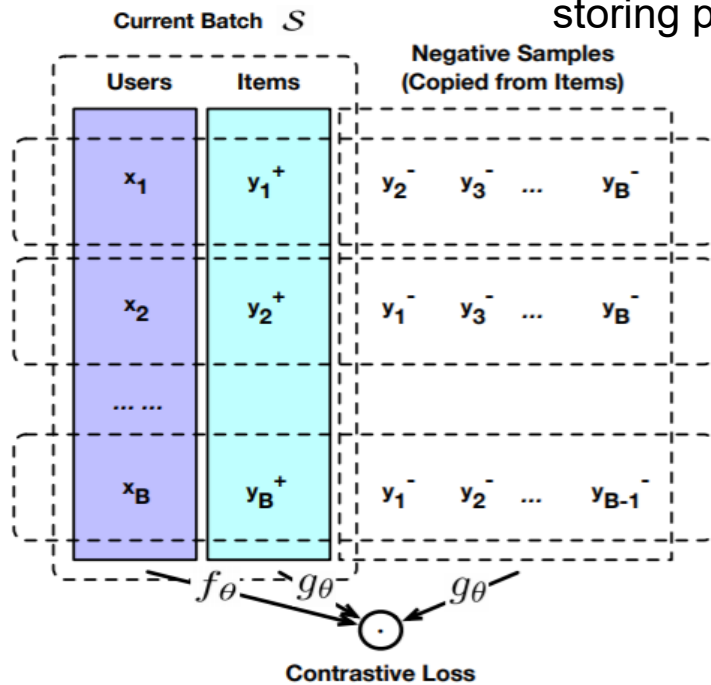
InfoNCE: Hierarchically
maximizing the mutual
information between
representations of the user, the
user-centered sub-hypergraph,
and the hypergraph.

CLRec

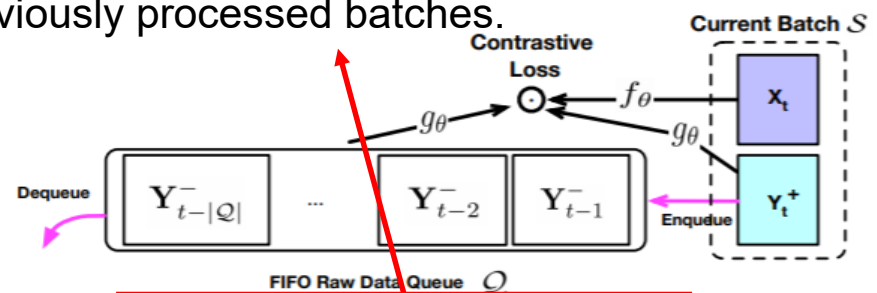
Live recommender systems face severe exposure bias.

Reducing the exposure bias by adjusting sampling strategies.

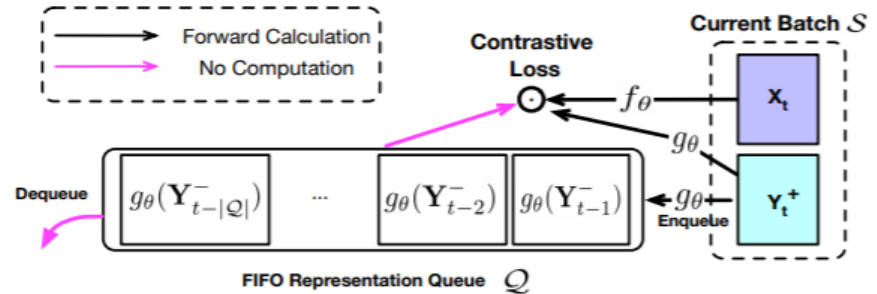
(b) Sampling instances from a fixed-size FIFO queue storing previously processed batches.



(a) Intra-Batch Contrastive Learning



(b) Inter-Batch Contrastive Learning



(c) Cached-Batch Contrastive Learning

(a) Sampling instances in the present batch.

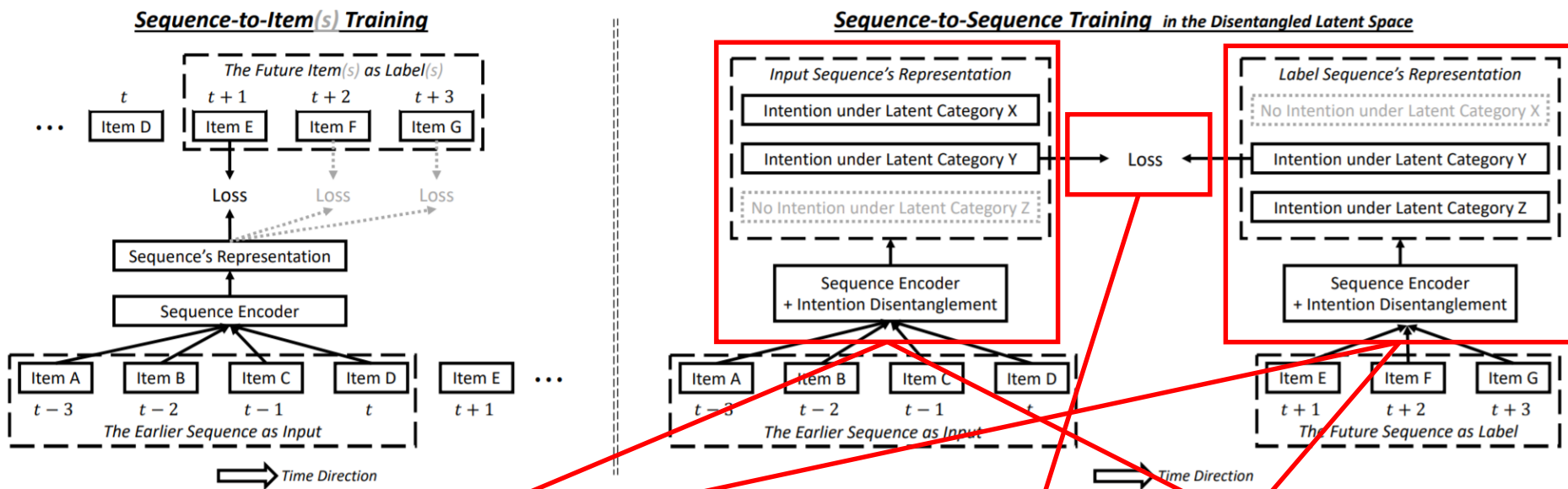
(c) Differing from variant (b) in that the queue caches the representations.

A popular choice of contrastive loss is equivalent to reducing the exposure bias.

Seq2seq for recommendation

Seq2item training is myopic and can easily lead to non-diverse recommendation lists.

Mining extra signals for supervision by looking at the longer-term future.



SASRec+Attention
 → Intention clustering
 → Intention weighting
 → Intention aggregation

$$\mathcal{L}_{s2s}(\theta, u, t, k) = -\ln p_{\theta}(\phi_{\theta}^{(k)}(\mathbf{x}_{T_u:t+1}^{(u)}) | \phi_{\theta}^{(k)}(\mathbf{x}_{1:t}^{(u)}))$$

MrTransformer

The relation between different sequences is under explored.

Discriminating the common and unique preference representations between a pair of sequences.



CE: Next item prediction based on the recombined representation.

Transformer+Coverage

- Preference editing
- Preference identification
- Preference separation
- Preference recombination

MSE: The common representation is close enough to each other.

MSE: The recombined representation is close enough to the original preference representation.

How pretraining helps recommendation?

Learning from Imperfection

Make full use of available data.

Make use of more data.



Low-resource

Sparse

.....

Noisy



Heterogeneous

Cold-start

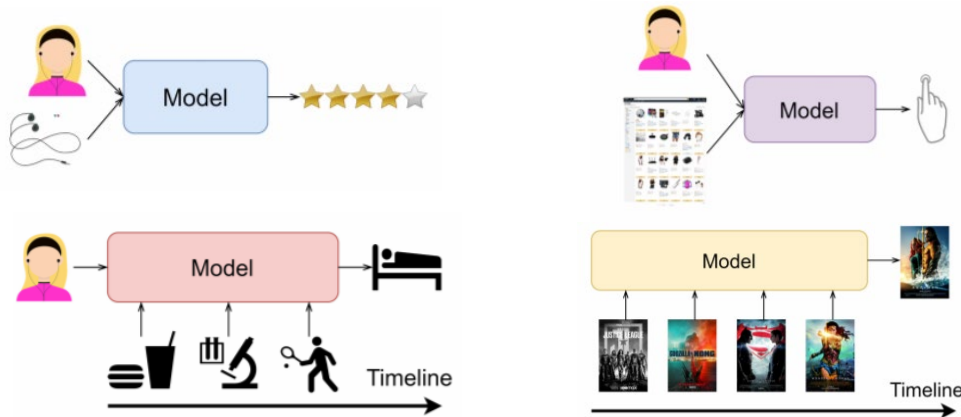
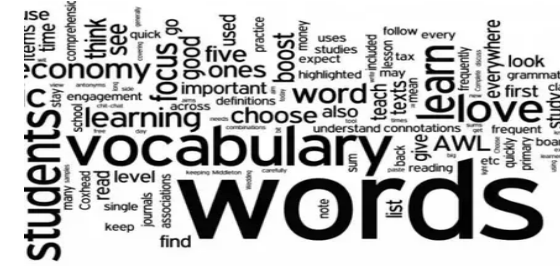
Computational efficiency

Multimodal pretraining

Model architecture

Theoretical foundations

What is unique for recommendation pretraining?

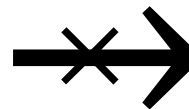


Universal Recommender

The Big Idea



Item



User



Thank you for your attention!



Pengjie Ren

renpengjie@sdu.edu.cn

<https://pengjieren.github.io/>

<https://ir.sdu.edu.cn/>



山东大学
SHANDONG UNIVERSITY

山东大学信息检索实验室

Information Retrieval Lab