







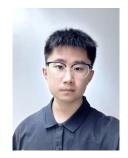
Joint Modeling in Deep Recommender Systems

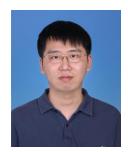




















Pengyue Jia¹ Jingtong Gao¹ Yejing Wang¹ Yuhao Wang¹ Xiaopeng Li¹ Qidong Liu¹ Yichao Wang² Bo Chen² Huifeng Guo² Ruiming Tang²

¹City University of Hong Kong, ²Huawei Noah's Ark Lab





Agenda













Joint Modeling in RS





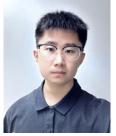
Multi-task Recommendation



Yuhao Wang

Multi-scenario Recommendation







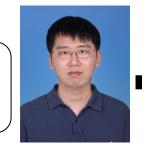
Pengyue Jia Xiaopeng Li

Multi-behavior Recommendation



Jingtong Gao

Multi-modal Recommendation



Qidong Liu

Conclusion

Future Work



Yichao Wang

Agenda













Joint Modeling in RS





Multi-task Recommendation



Recommendation





Yuhao Wang

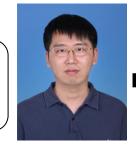
Pengyue Jia Xiaopeng Li





Jingtong Gao

Multi-modal Recommendation



Multi-scenario

Qidong Liu







Yichao Wang









Age of Information Explosion



Information overload



Items can be Products, News, Movies, Videos, Friends, etc.







- > Recommendation has been widely applied in online services
 - E-commerce, Content Sharing, Social Networking, etc.











Product Recommendation

Frequently bought together









- > Recommendation has been widely applied in online services
 - E-commerce, Content Sharing, Social Networking, etc.





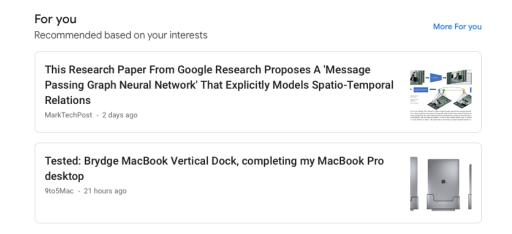


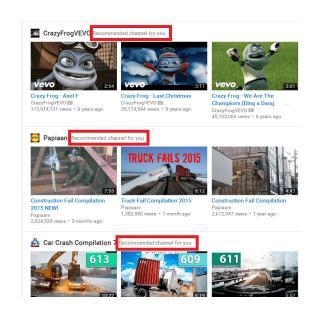






News/Video/Image Recommendation











- > Recommendation has been widely applied in online services
 - E-commerce, Content Sharing, Social Networking, etc.

facebook









Friend Recommendation



Deep Recommender Architecture









Advantages

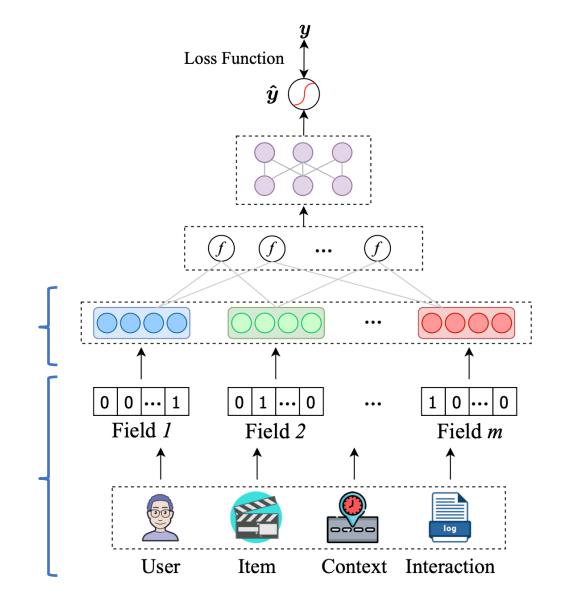
- **Feature representations of users** and items
- Non-linear relationships between users and items

Feature Embedding Layer

High/low-frequency features embedding sizes

Input Layer

Feature selection



Deep Recommender Architecture







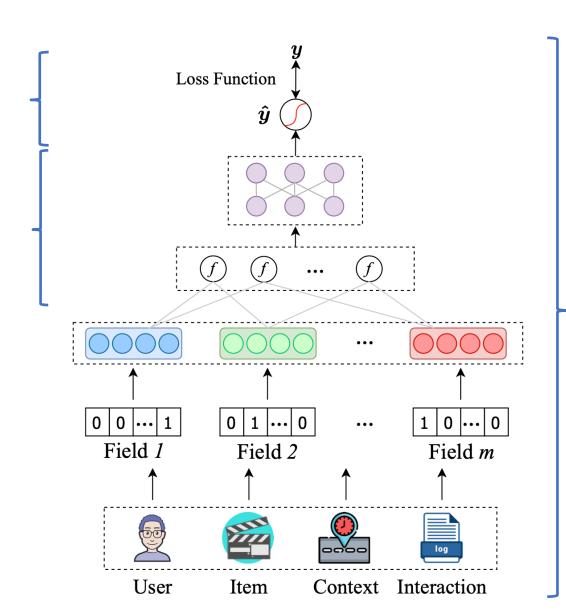


Output Layer

BCE, BPR, MSE

Feature Interaction Layer

Pooling, convolution, and the number of layers, inner product, outer product, convolution, etc.



System Design

Hardware infrastructure, data pipeline, information transfer, implementation, deployment, optimization, evaluation, etc.

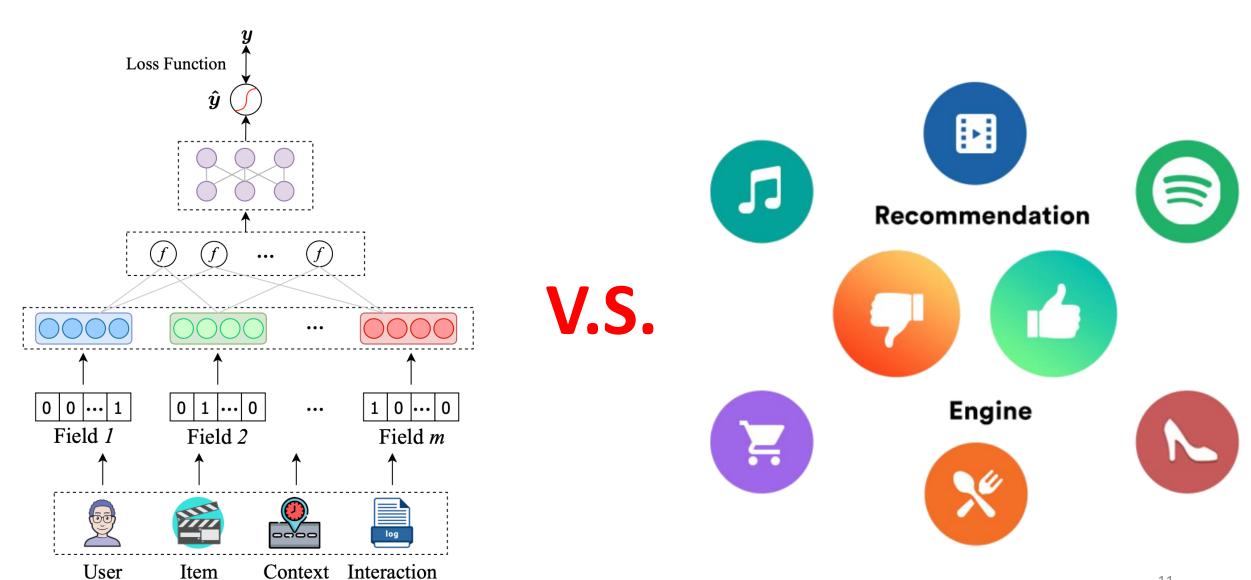
Deep Recommender Architecture











Agenda













Yejing Wang

Joint Modeling in RS





Multi-task Recommendation



Yuhao Wang

Multi-scenario Recommendation



Pengyue Jia Xiaopeng Li

Multi-behavior Recommendation



Jingtong Gao

Multi-modal Recommendation



Qidong Liu



Future Work



Joint Modeling in Recommendations

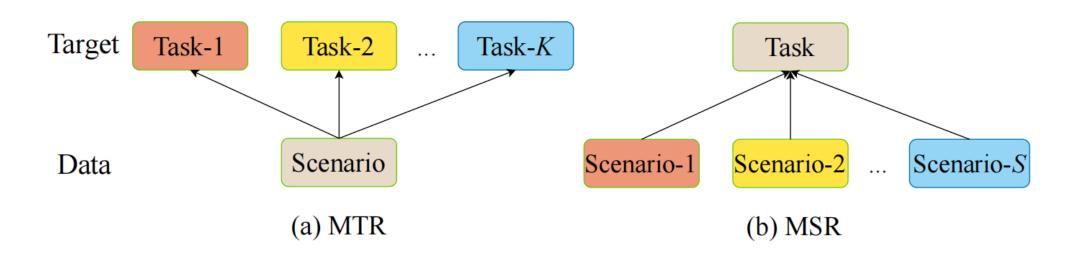








- > Handling the inter-dependency between users and items under more complex circumstances
- Advantages
 - One model for several situations
 - Performance improvement caused by information sharing in different situations
- > Two typical representatives:
 - Multi-task recommendation (MTR)
 - Multi-scenario recommendation (MSR)



Joint Modeling in Recommendations

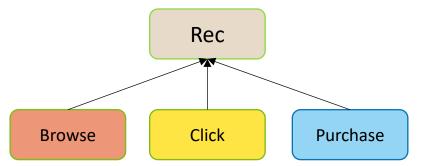




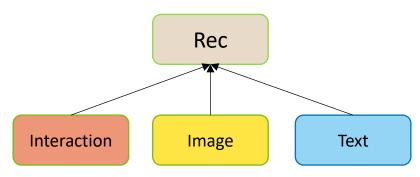




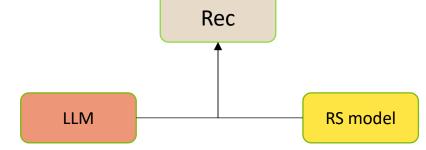
- More joint modeling methods:
 - Multi-behavior recommendation
 - Multi-modal recommendation
 - Large language model-based recommendation



Multi-behavior recommendation



Multi-modal recommendation



Large language model-based recommendation

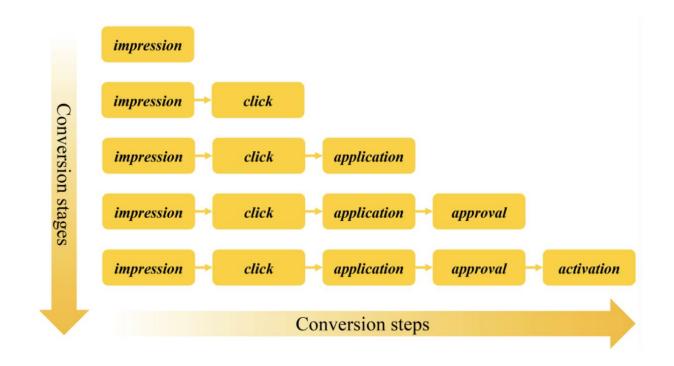






- ➤ Multi-Task Recommendation:
 - Independent tasks: Comments, repost, likes, bookmarks
 - Multi-stage conversion tasks: click, application, approval, activation ...





How to extract useful information from other tasks?

How to capture task dependences and resolve the sparsity issue?

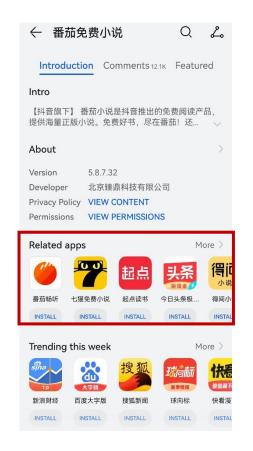


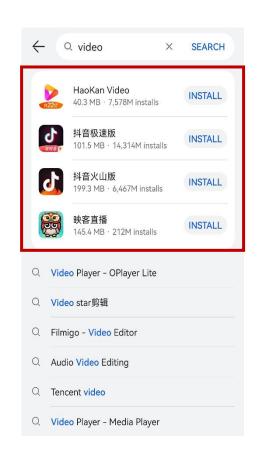






➤ Multi-Scenario Recommendation: construct multiple scenarios for user diverse requirements.

















➤ Multi-Behavior Modeling: click, download, like, buy





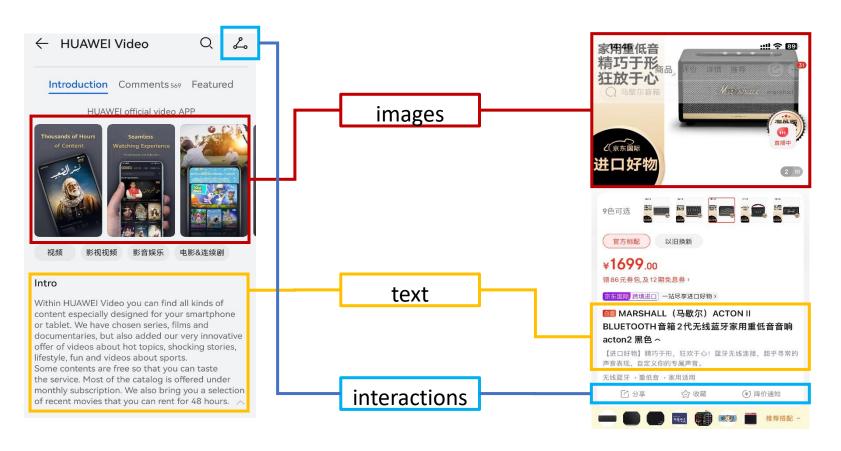
How to learn the relationship between different type of behaviors?







➤ Multi-Modal Modeling: user interactions, images, text ...



How to extract and align data from different modalities?

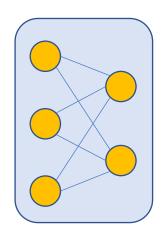








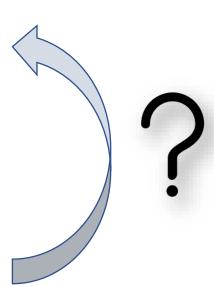
➤ Large Language Model-based Recommendation



Trained on labeled data with supervised learning

Collaborative signals

ID-based in-domain collaborative knowledge



DRS



LLM

Pre-trained on large-scale corpora with self-supervised learning

Semantic signals

Generalization, reasoning and open-world knowledge

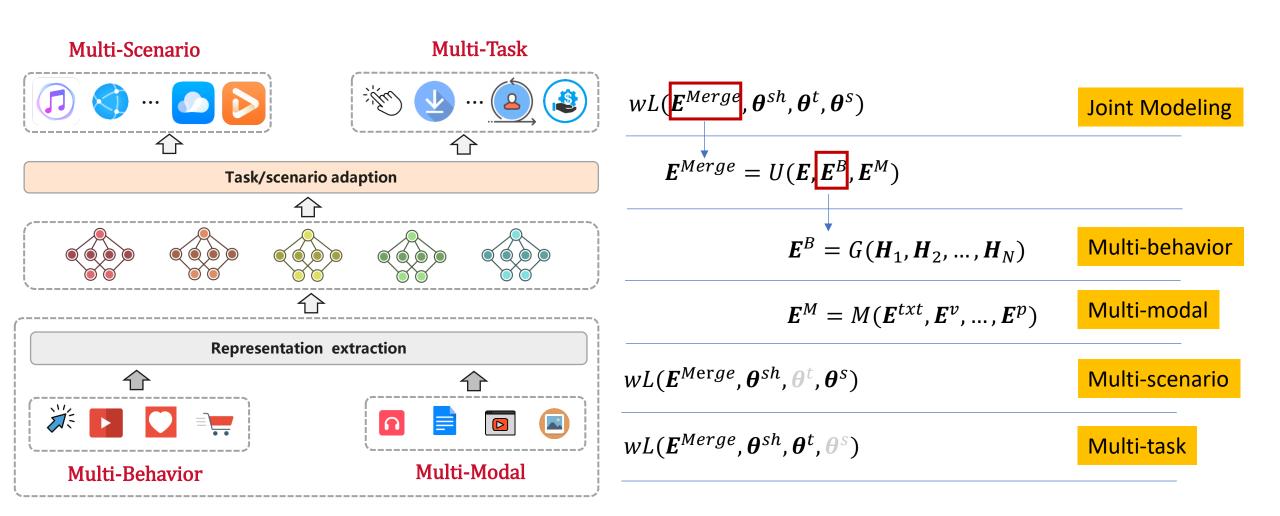
Relations and Formulations of Joint Modeling











Agenda













Yejing Wang

Joint Modeling in RS





Multi-task Recommendation

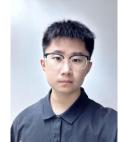


Yuhao Wang

Multi-scenario Recommendation



Pengyue Jia Xiaopeng Li

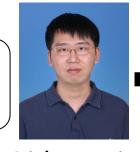


Multi-behavior Recommendation



Jingtong Gao

Multi-modal Recommendation



Qidong Liu



Future Work



Yichao Wang

MTR & MTDRS









Multi-Task Recommendation (MTR)

Multi-Task Deep Recommender Systems (MTDRS)

>How

Multi-Task Learning (MTL) + Deep Neural Networks

>Why

- Learning high-order feature interactions and
- Modeling complex user-item interaction behaviors

Benefits & Challenges







Benefits

- Mutual enhancement among tasks
- Higher efficiency of computation and storage

≻Challenges

- Effectively and efficiently capture useful information & relevance among tasks
- Data sparsity
- Unique sequential dependency

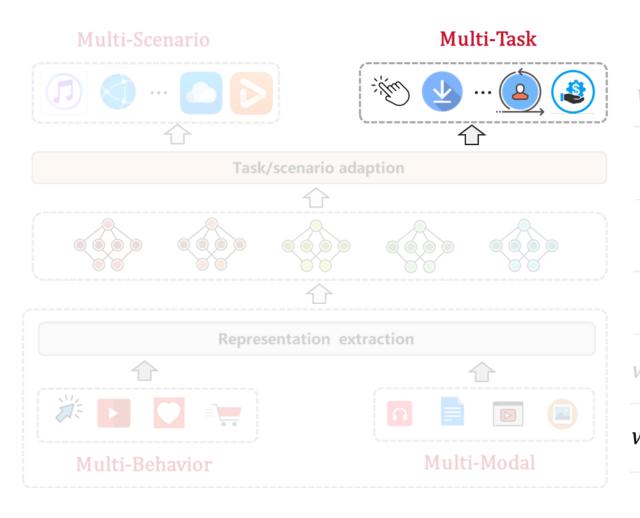
Multi-Task Modeling

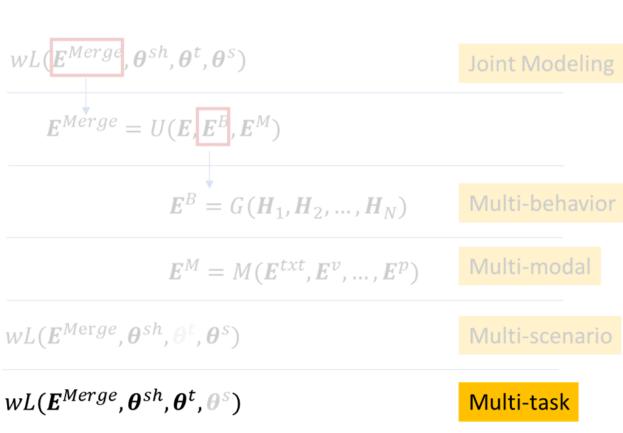












Formulation









> Problem:

- Learning MTL model with task-specific parameters $(\theta^1, ..., \theta^K)$ and shared parameter θ^s , which outputs the K task-wise predictions
- > Optimization problem:

$$\underset{\{\theta^{1},...,\theta^{K}\}}{\arg\min} \mathcal{L}\left(\theta^{s},\theta^{1},\cdots,\theta^{K}\right) = \underset{\{\theta^{1},...,\theta^{K}\}}{\arg\min} \sum_{k=1}^{K} \omega^{k} L^{k}\left(\theta^{s},\theta^{k}\right)$$

- $\mathcal{L}(\theta^s, \theta^k)$: loss function for k-th task with parameter θ^s, θ^k
- ω^k : loss weight for k-th task

BCE loss
$$L^k\left(heta^s, heta^k
ight) = -\sum_{n=1}^N \left[y_n^k \log\left(\hat{y}_n^k
ight) + \left(1 - y_n^k
ight) \log\left(1 - \hat{y}_n^k
ight)
ight]$$

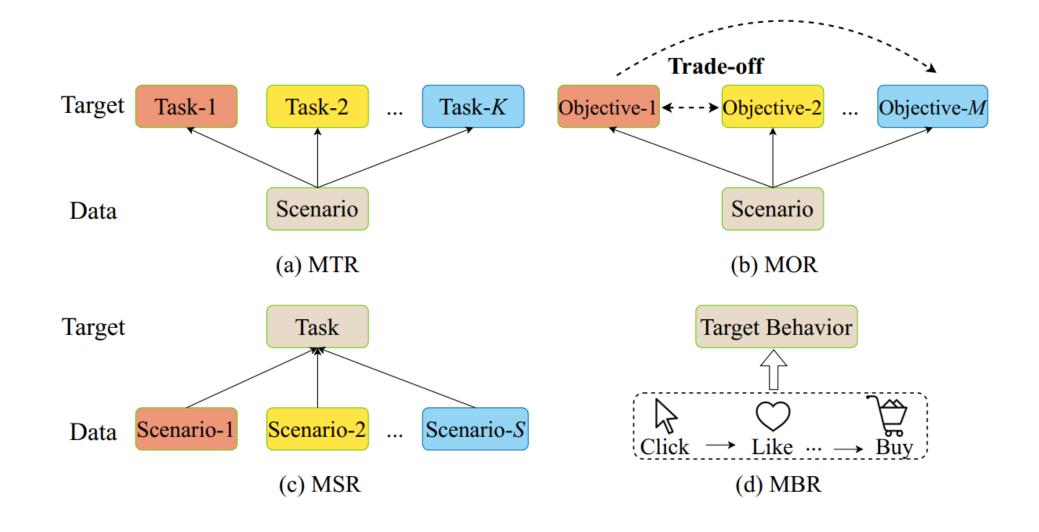
MTR, MOR, MSR, MBR











Comparison with CV & NLP





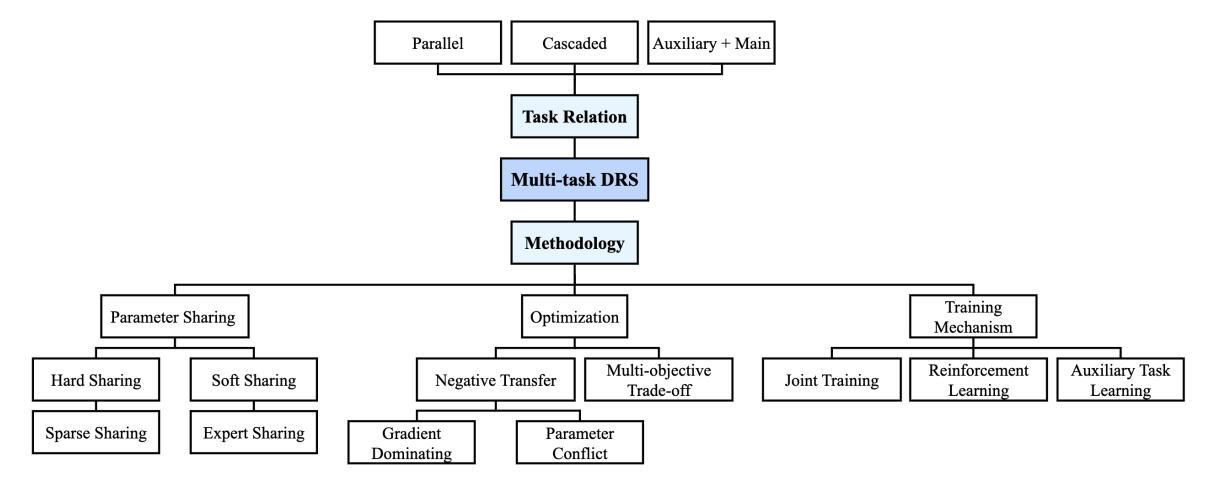




Task	Description	Explanation
CV	Multi-target segmentation and further classification for each object	Utilizing feature transformation to represent common features based on a multi-layer feed-forward network
NLP	Mostly focus on the design of MTL architectures	Based on RNN because of the sequence pattern Can be divided into word-, sentence-, and document-level by granularity





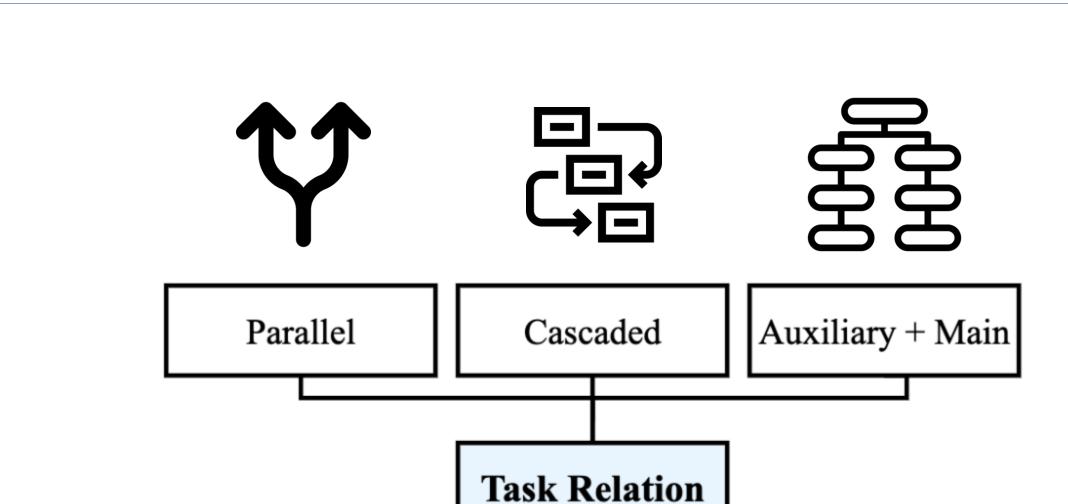


Task Relation









Parallel









Tasks independently calculated without sequential dependency

➤ Objective function: Weighted sum with constant loss weights

Cascaded









- > Cascaded task relationship: sequential dependency
- > Computation of current task depends on **previous** ones
 - E.g. CTCVR = CTR \times CVR
- > General formulation:

$$\hat{y}_n^k \left(\theta^s, \theta^k\right) - \hat{y}_n^{k-1} \left(\theta^s, \theta^k\right) = P\left(\epsilon_k = 0, \epsilon_{k-1} = 1\right)$$

- ϵ_k : Indicator variable for task k
- Difference is the probability of the task k not happening while the task k-1 is observed

Cascaded







Model	Problem	Behavior Sequence
ESMM [Ma et al., 2018b]	SSB & DS	$impression \rightarrow click \rightarrow conversion$
ESM ² [Wen et al., 2020]	SSB & DS	$impression \rightarrow click \rightarrow D(O)Action \rightarrow purchase$
Multi-IPW & DR [Zhang et al., 2020]	SSB & DS	$exposure \rightarrow click \rightarrow conversion$
ESDF [Wang et al., 2020b]	SSB & DS & time delay	$impression \rightarrow click \rightarrow pay$
HM ³ [Wen <i>et al.</i> , 2021]	SSB & DS & micro and macro behavior modeling	$impression \rightarrow click \rightarrow micro \rightarrow macro \rightarrow purchase$
AITM [Xi et al., 2021]	sequential dependence in multi-step conversions	$impression \rightarrow click \rightarrow application \rightarrow approval \rightarrow activation$
MLPR [Wu et al., 2022]	sequential engagement & vocabulary mismatch in product ranking	$impression \rightarrow click \rightarrow add\text{-to-cart} \rightarrow purchase$
ESCM ² [Wang et al., 2022a]	inherent estimation bias & potential independence priority	$impression \rightarrow click \rightarrow conversion$
HEROES [Jin et al., 2022]	multi-scale behavior & unbiased learning-to-rank	$observation \rightarrow click \rightarrow conversion$
APEM [Tao et al., 2023]	sample-wise representation learning in SDMTL	$impression \rightarrow click \rightarrow authorize \rightarrow conversion$
DCMT [Zhu et al., 2023]	SSB & DS & potential independence priority (PIP)	$exposure \rightarrow click \rightarrow conversion$
	<i>'</i>	

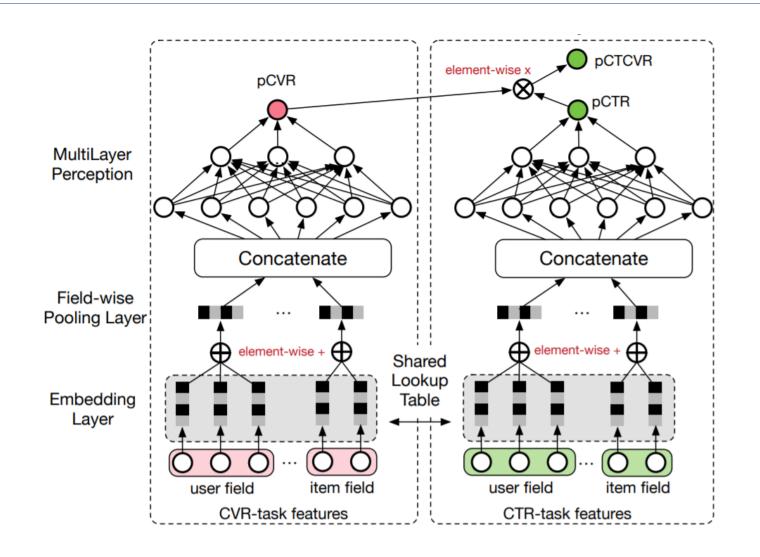
SSB: Sample Selection Bias DS: Data Sparsity











Auxiliary with Main Task







- ➤ A task specified as the main task
 while associated auxiliary tasks help to improve performance
- > Probability estimation for main task the probability of auxiliary tasks
- ➤ Provide richer information across entire space

Auxiliary with Main Task







Model	References	Method
ESDF Multi-IPW and Multi-DR DMTL Metabalance	[Wang et al., 2020b] [Zhang et al., 2020] [Zhao et al., 2021] [He et al., 2022]	Adopt the original recommendation tasks as auxiliaries
MTRec PICO MTAE Cross-Distill	[Li et al., 2020a] [Lin et al., 2022] [Yang et al., 2021] [Yang et al., 2022a]	Manually design various auxiliary tasks
CSRec	[Bai et al., 2022]	Contrastive learning as the auxiliary
Self-auxiliary*	[Wang et al., 2022b]	Under-parameterized self-auxiliaries

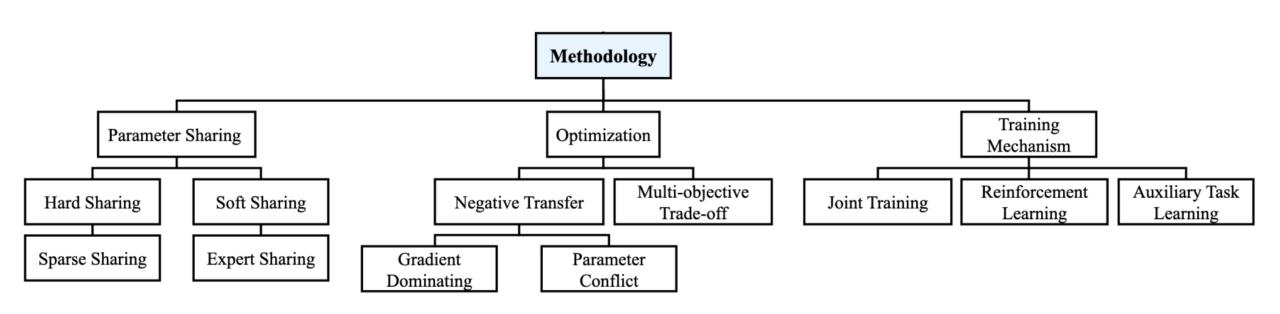
Methodology











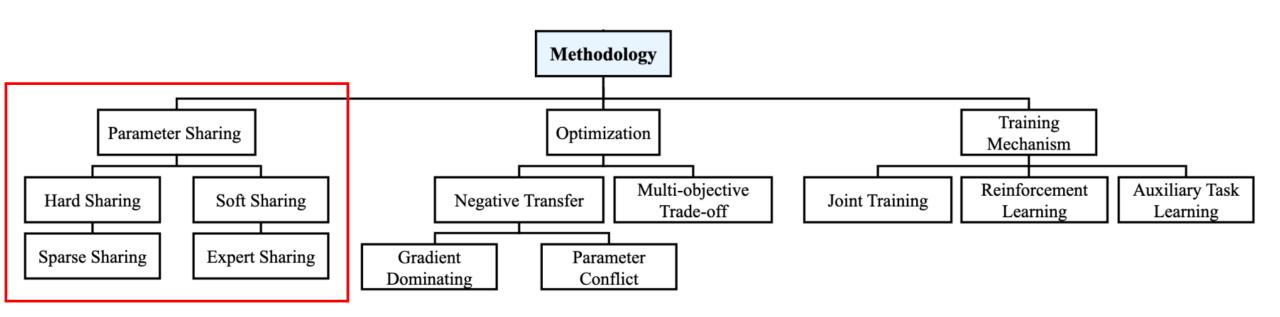
Parameter Sharing











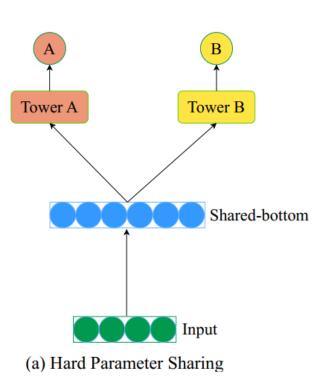
Parameter Sharing

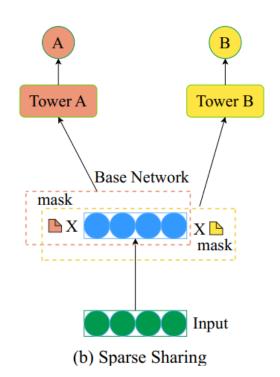


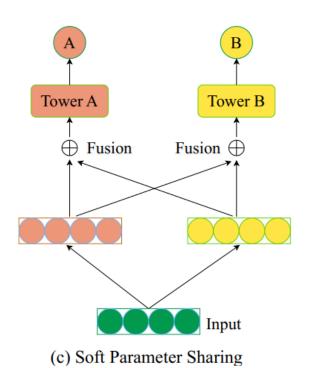


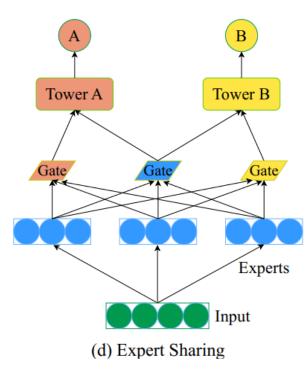










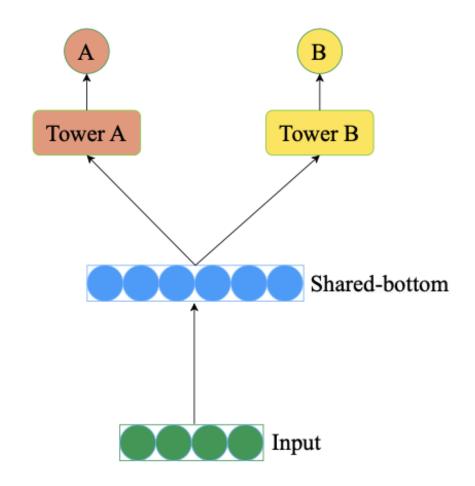


Hard Sharing









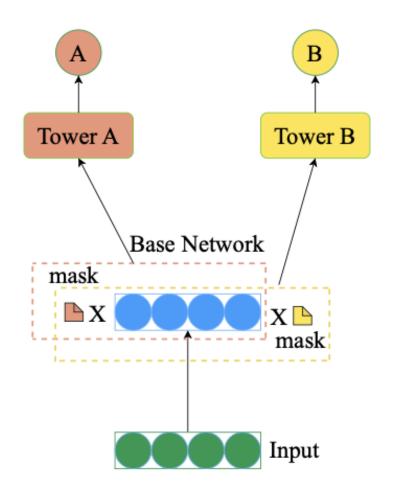
- Shared bottom layers extract the same information for different tasks,
- Task-specific top layers are trained individually
- ✓ Improving computation efficiency and alleviating over-fitting
- X Limited capacity of the shared parameter space → Weakly related tasks and noise

Sparse Sharing









- Extracting **sub-networks** for each task by parameter masks from a base network
 - Special case of Hard Sharing
- ✓ Coping with the weakly related tasks flexibly
- X Negative transfer when updating shared parameters

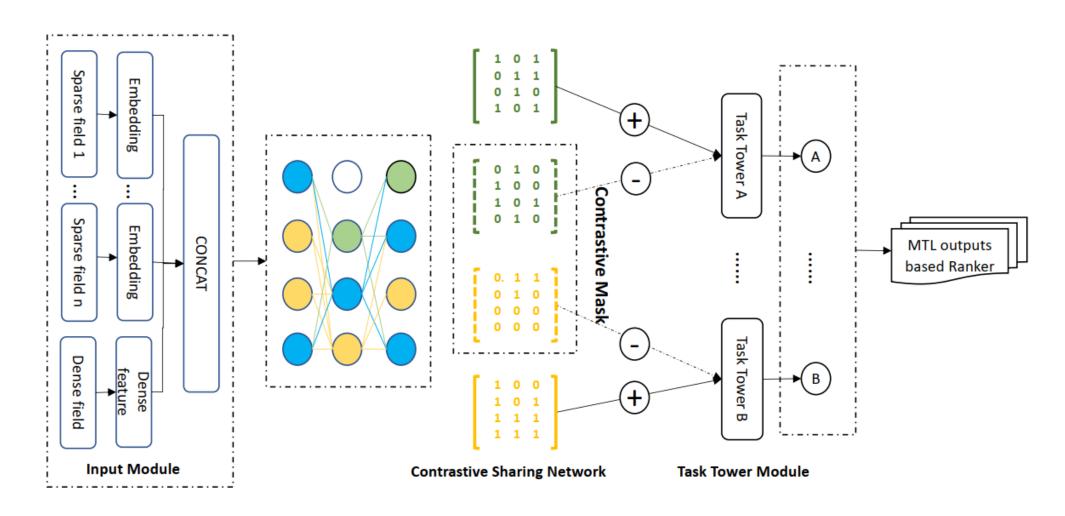
CSRec









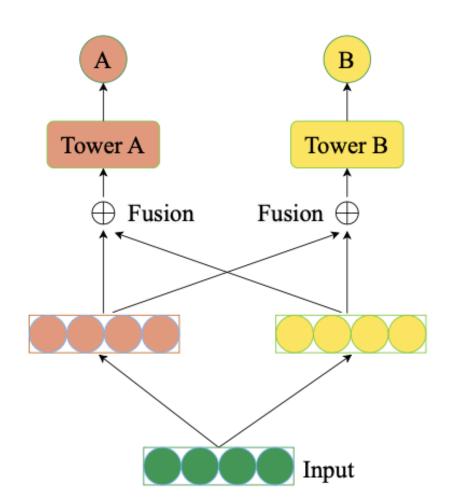


Soft Sharing









- Building separate models for tasks but the information among tasks is fused by weights of task relevance
- ✓ Relatively high **flexibility** in parameter sharing v.s. hard sharing
- X Can not reconcile the flexibility
- X Computation cost of the model

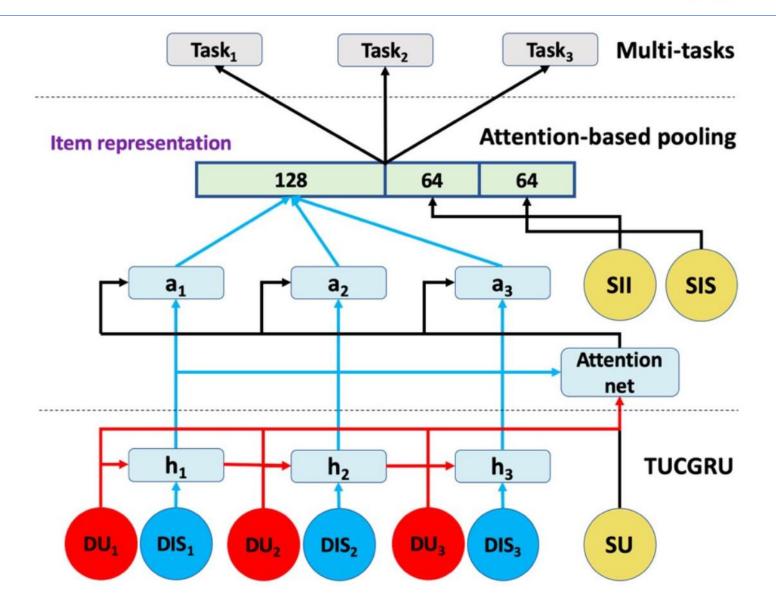
DINOP











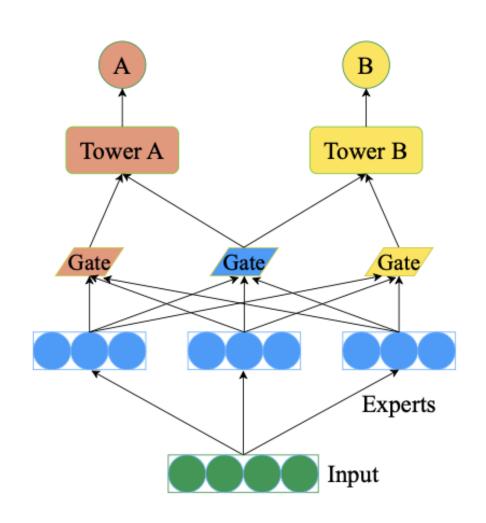
Expert Sharing











- Employing multiple expert networks to extract knowledge from shared bottom
 - → Fed into **task-specific** modules like gates
 - → Passed into the task-specific tower
- Mainly non-sequential input features
- Special case of Soft Sharing

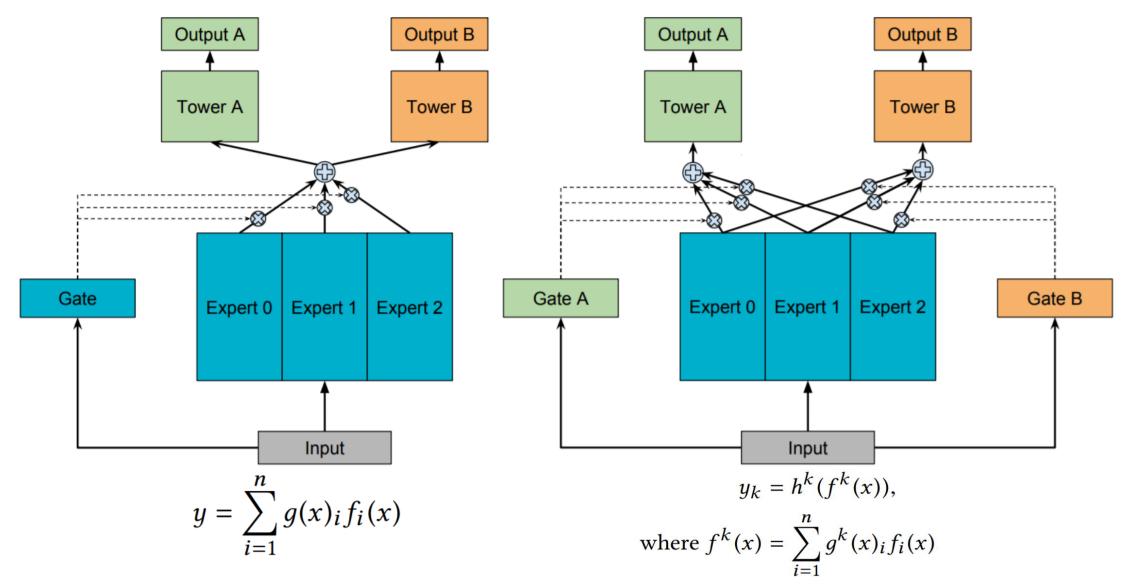
MMoE











Expert Sharing









Model	Reference
MMoE	[Ma et al., 2018a]
SNR	[Ma et al., 2019]
PLE	[Tang et al., 2020]
DMTL	[Zhao et al., 2021]
DSelect-k	[Hazimeh et al., 2021]
MetaHeac	[Zhu et al., 2021]
PFE	[Xin et al., 2022]
MVKE	[Xu et al., 2022]
FDN	[Zhou et al., 2023]
MoME	[Xu et al., 2024]
MoSE	[Qin et al., 2020]

Processing **non-sequential** input features, while the remaining models is ameliorated based on MMoE

Processing **sequential** input features utilizing LSTM & sequential experts

Special Case: Multi-Embedding Paradigm

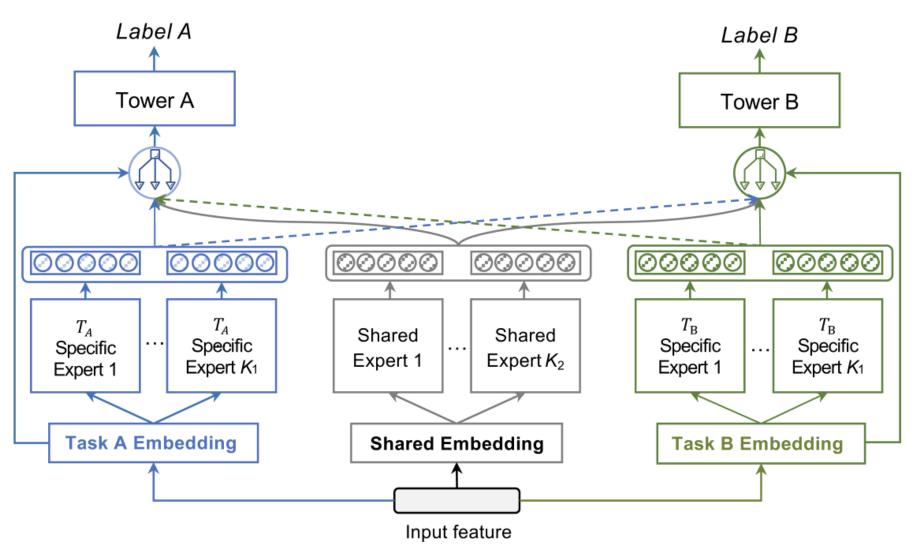


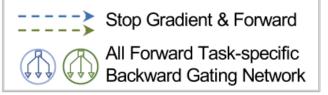


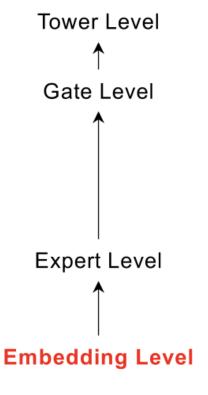




STEM (Shared and Task-specific EMbeddings)







Task-Specific Levels

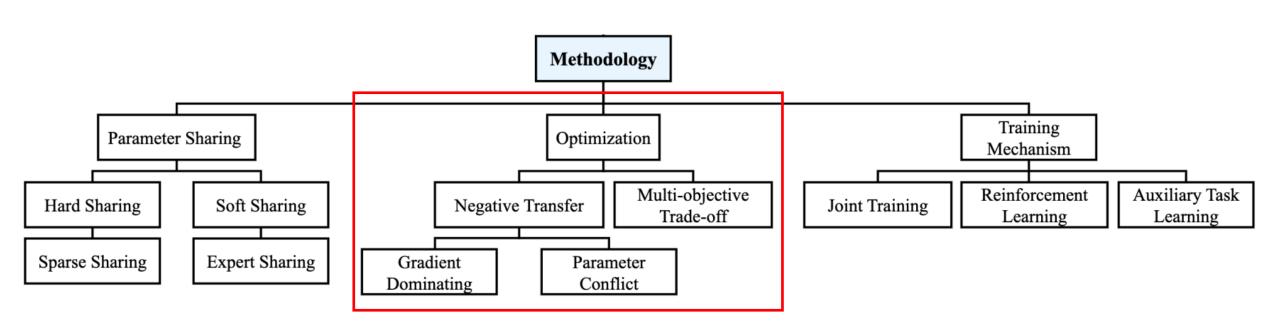
Optimization











Negative Transfer









Works	Approach
AdaTask [Yang et al., 2022b]	Quantifying task dominance of shared parameters, calculate task-specific accumulative gradients
MetaBalance [He et al., 2022]	Flexibly balancing the gradient magnitude proximity between auxiliary and target tasks by a relax factor

Opposite directions of gradient +- $\nabla_{\theta} L^k(\theta)$

Works	Approach
PLE [Tang et al., 2020]	Proposing customized gate control (CGC) separating shared and task-specific experts
CSRec [Bai et al., 2022]	Alternating training procedure and contrastive learning on parameter masks to reduce the conflict probability
GradCraft [Bai et al., 2024]	Adjusting gradient norm and deconflicting global direction through projection and combination 51

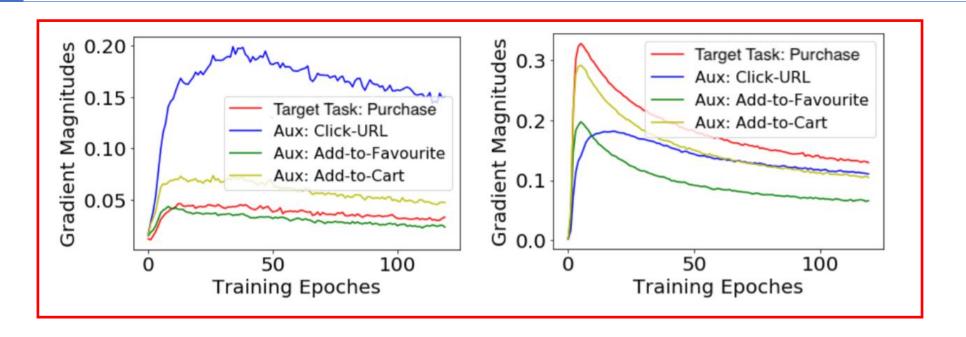
MetaBalance











$$\theta^{t+1} = \theta^t - \alpha * \mathbf{G}_{total}^t$$

$$\mathbf{G}_{total}^{t} = \nabla_{\theta} \mathcal{L}_{total}^{t} = \nabla_{\theta} \mathcal{L}_{tar}^{t} + \sum_{i=1}^{K} \nabla_{\theta} \mathcal{L}_{aux,i}^{t}$$

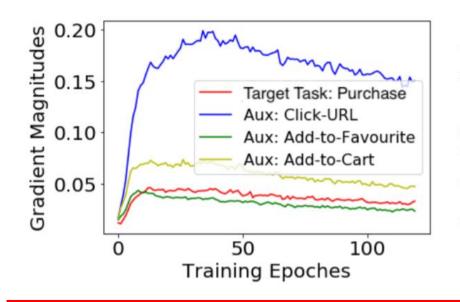
MetaBalance

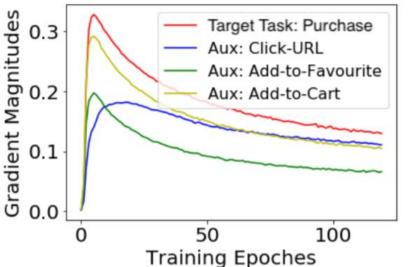












$$\begin{split} \boldsymbol{\theta}^{t+1} &= \boldsymbol{\theta}^{t} - \boldsymbol{\alpha} * \mathbf{G}_{total}^{t} \\ \mathbf{G}_{total}^{t} &= \nabla_{\boldsymbol{\theta}} \mathcal{L}_{total}^{t} = \nabla_{\boldsymbol{\theta}} \mathcal{L}_{tar}^{t} + \sum_{i=1}^{K} \nabla_{\boldsymbol{\theta}} \mathcal{L}_{aux,i}^{t} \end{split}$$

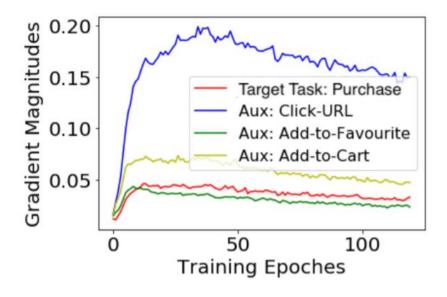
MetaBalance











$$\theta^{t+1} = \theta^t - \alpha * \mathbf{G}_{total}^t$$

$$G_{total}^{t} = \nabla_{\theta} \mathcal{L}_{total}^{t} = \nabla_{\theta} \mathcal{L}_{tar}^{t} + \sum_{i=1}^{K} \nabla_{\theta} \mathcal{L}_{aux,i}^{t}$$

$$\mathbf{G}_{aux,i}^{t} \leftarrow (\mathbf{G}_{aux,i}^{t} * \frac{\|\mathbf{G}_{tar}^{t}\|}{\|\mathbf{G}_{aux,i}^{t}\|}) * r + \mathbf{G}_{aux,i}^{t} * (1 - r)$$

Multi-objective Trade-off









Objectives optimized regardless of the potential conflict

Works	Trade-off
[Wang et al., 2021]	Group fairness and accuracy
[Wang et al., 2022b]	Minimizing task conflicts and improving multi-task generalization

Training Mechanism

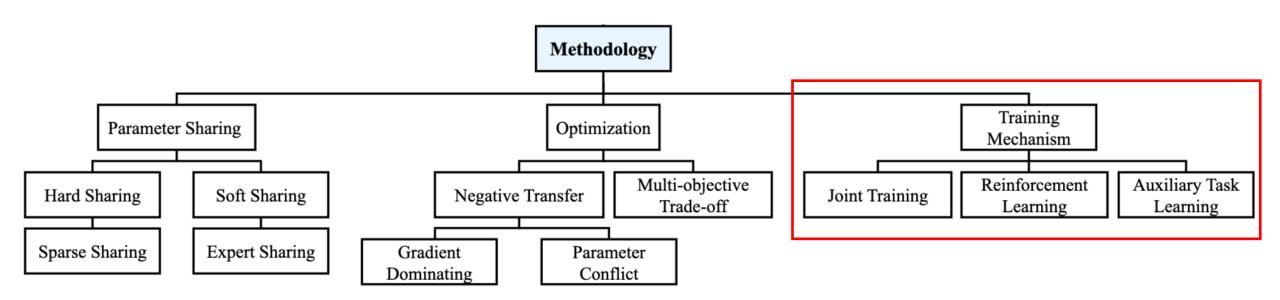








Training process & Learning strategy



Joint Training









Parallel manner

Category	Reference
Session-based RS	[Shalaby et al., 2022] [Qiu et al., 2021] [Meng et al., 2020]
Route RS	[Das, 2022]
Knowledge graph enhanced RS	[Wang et al., 2019]
Explainability	[Lu et al., 2018] [Wang et al., 2018]
Graph-based RS	[Wang et al., 2020a]

Reinforcement Learning









Sequential user behaviors as MDP

Summary	Reference
Formulating MTF as MDP and use batch RL to optimize long-term user satisfaction	[Zhang et al., 2022b]
Using an actor-critic model to learn the optimal fusion weight of tasks rather than greedy ranking strategies	[Han et al., 2019]
Using dynamic critic networks to adaptively adjust the fusion weight considering the session-wise property	[Liu et al., 2023]

Auxiliary Task Learning









Joint training & Others

Summary	Reference
Employing Expectation-Maximization (EM) algorithm for optimization	ESDF [Wang et al., 2020b]
Trained with task-specific sub- networks	Self-auxiliaries [Wang et al., 2022b]

Application Fields









- **E-commerce**: Main focus
- Advertising
 - **Utility & Cost**
 - i. MM-DFM [Hou et al., 2021]: Performing multiple conversion prediction tasks in different observation duration
 - ii. MetaHeac [Zhu et al., 2021]: Handling audience expansion tasks on contentbased mobile marketing
 - iii. MVKE [Xu et al., 2022]: Performing user tagging for online advertising

Social media

- i. MMoE [Zhao et al., 2019b]: YouTube engagement and satisfaction
- ii. LT4REC [Xiao et al., 2020]: Tencent Video
- iii. BatchRL-MTF [Zhang et al., 2022b]: Tencent short video platform

Datasets









Datasets	Stage	Tasks	Website
Ali-CCP [42]	Ranking	CTR, CVR	https://tianchi.aliyun.com/dataset/408/
Criteo [13]	Ranking	CTR, CVR	https://ailab.criteo.com/criteo-attribution-modeling-bidding-dataset/
AliExpress [32]	Ranking	CTR, CTCVR	https://tianchi.aliyun.com/dataset/74690/
MovieLens [23]	Recall & Ranking	Watch, Rating	https://grouplens.org/datasets/movielens/
Yelp	Recall & Ranking	Rating, Explanation	https://www.yelp.com/dataset/
Amazon [25]	Recall & Ranking	Rating, Explanation	http://jmcauley.ucsd.edu/data/amazon/
Kuairand [18]	Recall & Ranking	Click, Like, Follow, Comment,	https://kuairand.com/
Tenrec [77]	Recall & Ranking	Click, Like, Share, Follow,	https://github.com/yuangh-x/2022-NIPS-Tenrec/

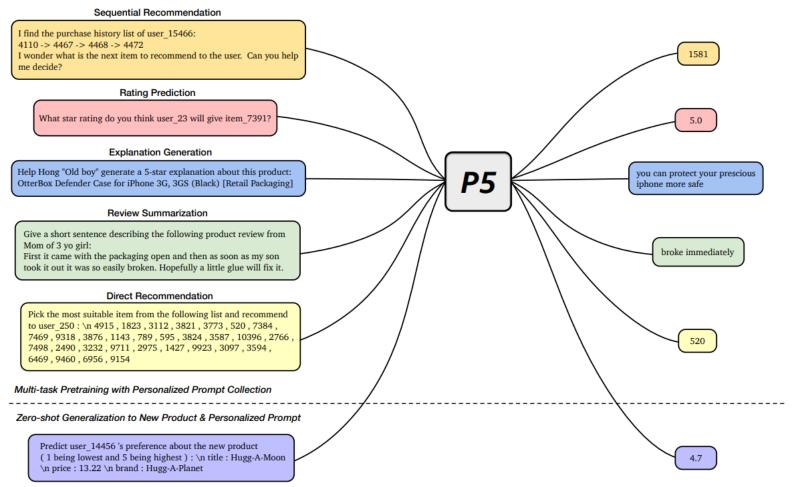








- ➤ P5: a unified recommendation model with pre-trained LLM model T5
- Fine-tuning with five commonly used tasks



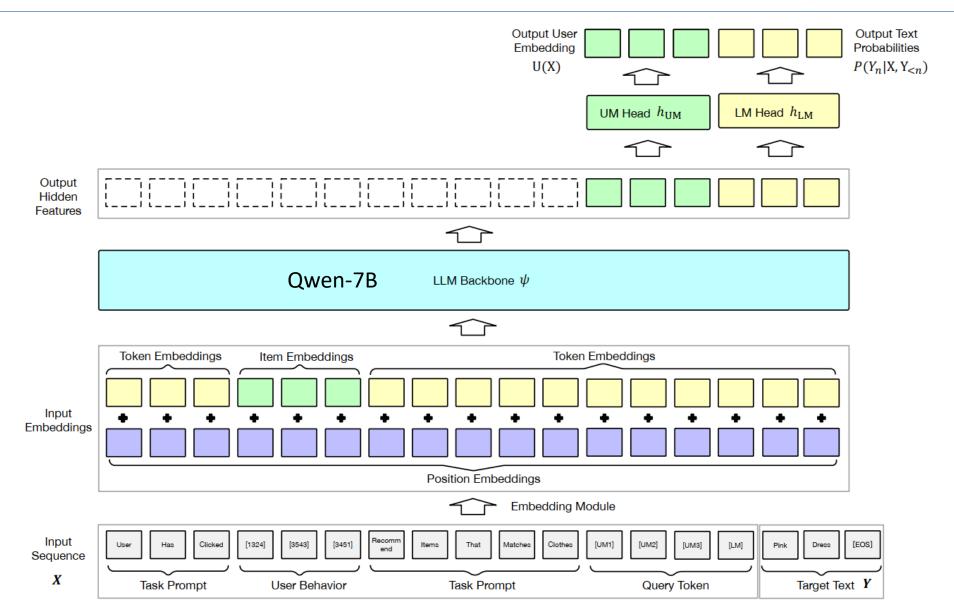
URM











Prompt Template









Multi-scenario Recommendation: The items the user has recently clicked on are as follows: {USER BEHAVIOR SE-QUENCE}. In scenario {SCENE}, please recommend items.

Multi-objective Recommendation: The items the user has recently clicked on are as follows: {USER BEHAVIOR SE-QUENCE}. Please find items that the user will {ACTION}.

Long-tail Item Recommendation: The items the user has recently clicked on are as follows: {USER BEHAVIOR SE-QUENCE}. Please recommend long-tail items.

Serendipity Recommendation: The items the user has recently clicked on are as follows: {USER BEHAVIOR SEQUENCE}.

Please recommend some new item categories.

Long-term Recommendation: The items the user has recently clicked on are as follows: {USER BEHAVIOR SEQUENCE}.

Please find items that match the user's long-term interests.

Search Problem: The items the user has recently clicked on are as follows: {USER BEHAVIOR SEQUENCE}.

Please recommend items that match {QUERY}.

Inputs: The items the user has recently clicked on are as follows: [7502][8308][8274][8380]. Please recommend items that match *Clothes*. [UM][LM]

Target Text: Swimwear & Beachwear for the Summer;

Casual Dresses for Every Occasion.

Target Items: [3632][1334]

Summary









➤ Multi-task Recommendation + Language Model

Model	Setting	Methods
P5	MTR+PLM	Prompt design;SFT
M6-Rec	MTR+PLM	Prompt design;SFT
UniMIND	MTR+PLM	Prompt design;SFT
URM	MTR+LLM	Prompt design;SFT
LUM	MTR+LLM	Next condition-item prediction







Topic	Challenge & future direction
Negative Transfer	 Extra complex inter-task correlation What, where, and when to transfer to alleviate negative transfer
AutoML	 Existing models only focus on the parameter sharing routing, while other components and hyper-parameters still under-explored
Explainability	Complex task relevance
Task-specific Biases	 Most existing models only focus on one specific bias Multiple bias should be tackled in future







Topic	Challenge & future direction
Negative Transfer	 Extra complex inter-task correlation What, where, and when to transfer to alleviate negative transfer
AutoML	 Existing models only focus on the parameter sharing routing, while other components and hyper-parameters still under-explored
Explainability	Complex task relevance
Task-specific Biases	 Most existing models only focus on one specific bias Multiple bias should be tackled in future







Topic	Challenge & future direction
Negative Transfer	 Extra complex inter-task correlation What, where, and when to transfer to alleviate negative transfer
AutoML	 Existing models only focus on the parameter sharing routing, while other components and hyper-parameters still under-explored
Explainability	Complex task relevance
Task-specific Biases	 Most existing models only focus on one specific bias Multiple bias should be tackled in future







Topic	Challenge & future direction
Negative Transfer	 Extra complex inter-task correlation What, where, and when to transfer to alleviate negative transfer
AutoML	 Existing models only focus on the parameter sharing routing, while other components and hyper-parameters still under-explored
Explainability	Complex task relevance
Task-specific Biases	 Most existing models only focus on one specific bias Multiple bias should be tackled in future

Conclusion









Parallel, Cascaded, Auxiliary with Main

➤ Methodology:

Parameter Sharing, Optimization, Training Mechanism

Survey









https://arxiv.org/abs/2302.03525

Multi-Task Deep Recommender Systems: A Survey

YUHAO WANG*, HA TSZ LAM*, and YI WONG*, City University of Hong Kong ZIRU LIU, City University of Hong Kong XIANGYU ZHAO[†], City University of Hong Kong YICHAO WANG, BO CHEN, HUIFENG GUO, and RUIMING TANG[†], Huawei Noah's Ark Lab

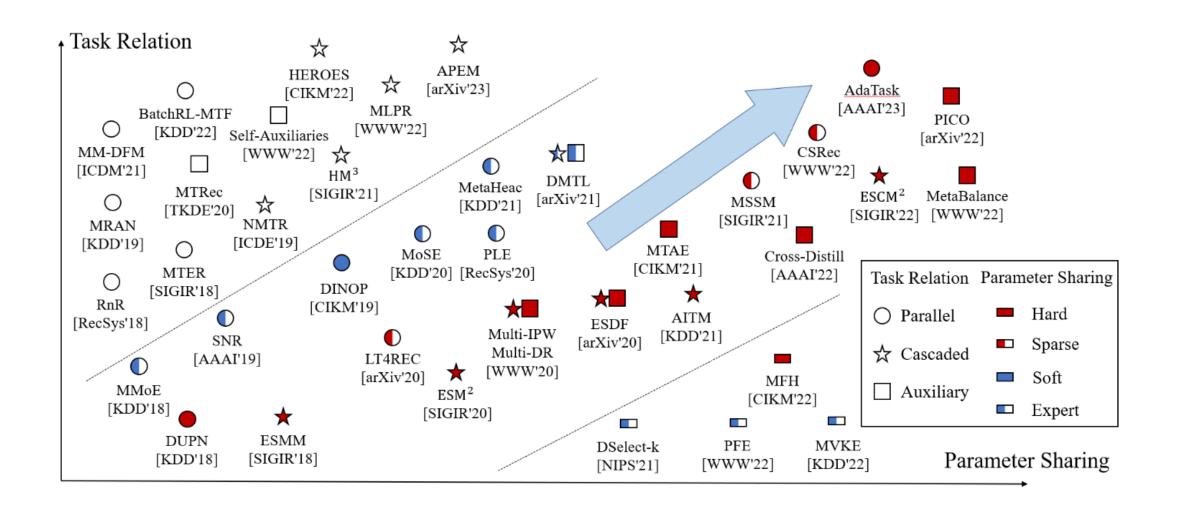
Trend of MTDRS











Agenda













Joint Modeling in RS





Multi-task Recommendation

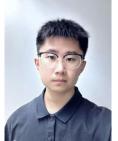


Yuhao Wang





Pengyue Jia Xiaopeng Li

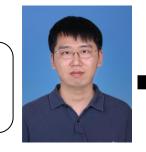


Multi-behavior Recommendation



Jingtong Gao

Multi-modal Recommendation



Qidong Liu



Future Work



Yichao Wang

Background









- Multi-Scenario Recommender Systems:
 - By using a unified model to simultaneously model multiple scenarios, the goal of improving the effects of different scenarios at the same time is achieved through information transfer between scenarios.
- > Importance:
 - Time/Memory efficiency; Maintenance cost
 - Accuracy

- Classification on Methods:
 - Shared-Specific network paradigm
 - Dynamic weight
 - Multi-scenario & Multi-task recommendation

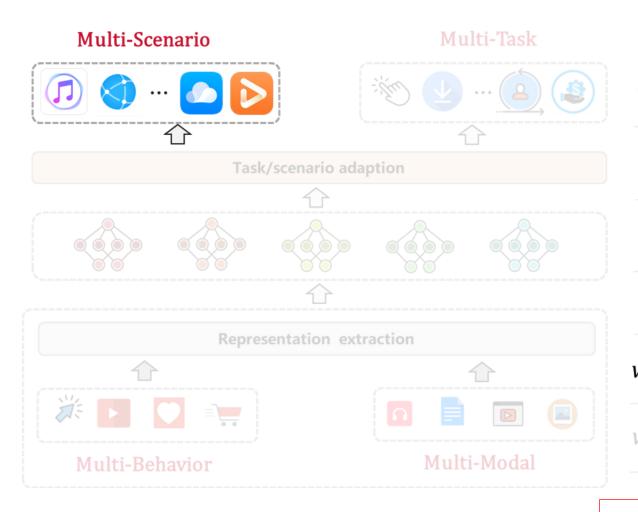
Multi-Scenario Modeling

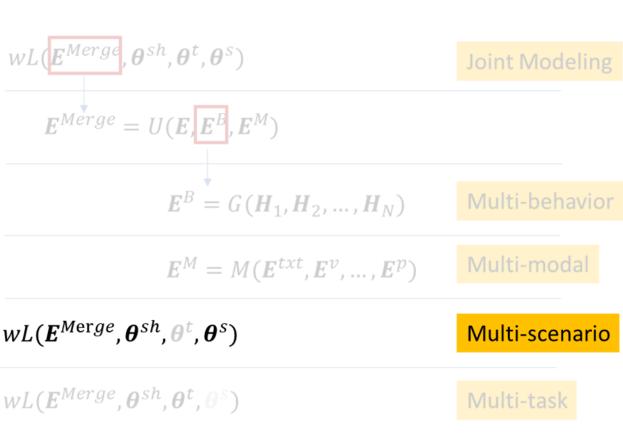












 θ^{sh} : shared parameters across scenarios θ^{s} : scenarios parameters of modeling

Recommendation Scenarios









➤ What is Scenario?

- Homepage, Searching page, Detailed page ...
- Food, Leisure and entertainment, ...
- Usually refers to different business scenarios

➤ Scenario and Domain?

- Generally do not make a distinction
- The same in this tutorial

Commonalities and Diversities

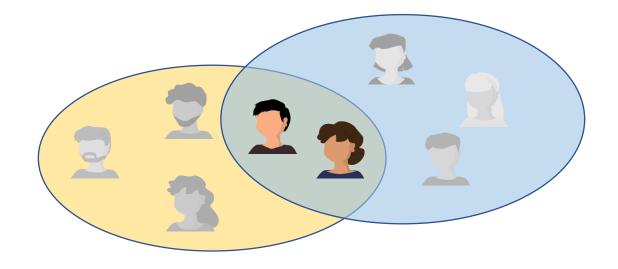




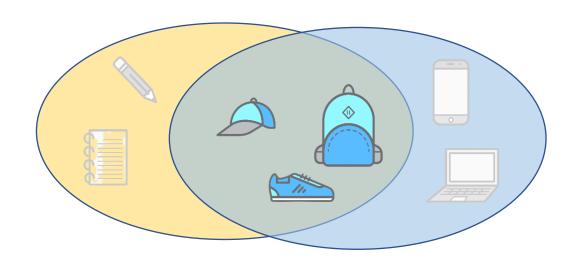




- **≻**Commonalities
 - User Overlap



- **≻**Commonalities
 - Item Overlap



Commonalities and Diversities









→ Diversities

- The specific user group may be different
- User's interest changes with the scenarios

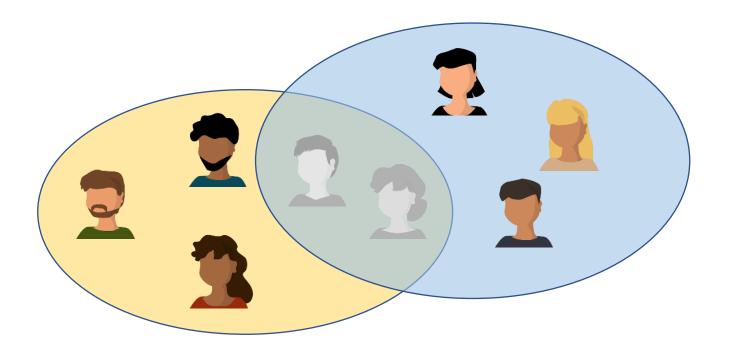


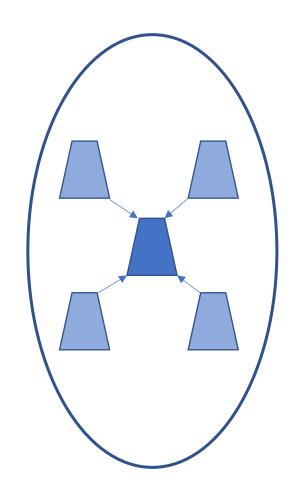
Table of Contents



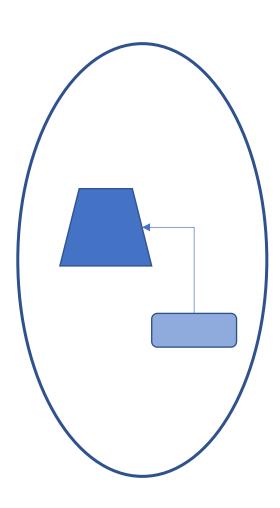




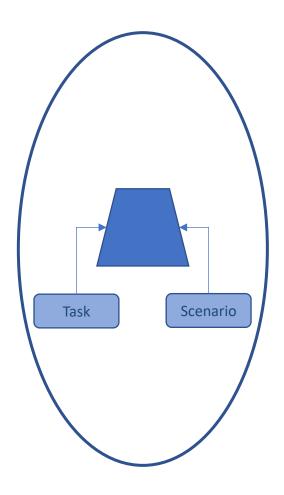




Shared-specific network paradigm $wL(E^{Merge}, \boldsymbol{\theta}^{sh}, \boldsymbol{\theta}^{t}, \boldsymbol{\theta}^{s})$



Dynamic weight $wL(\mathbf{E}^{Merge}, \boldsymbol{\theta}^{sh}, \boldsymbol{\theta}^{t}, \boldsymbol{\theta}^{s})$



Multi-Scenario & Multi-Task $wL(\mathbf{E}^{Merge}, \mathbf{\theta}^{sh}, \mathbf{\theta}^{t}, \mathbf{\theta}^{s})$

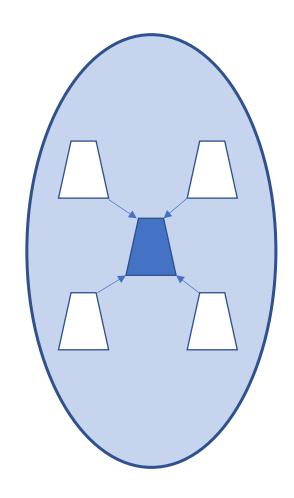
Table of Contents



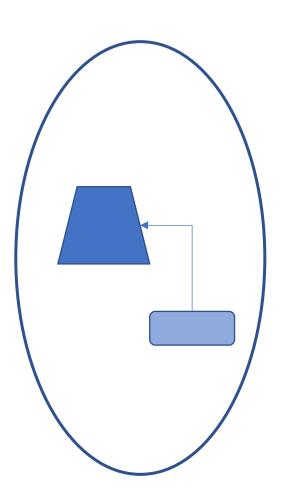




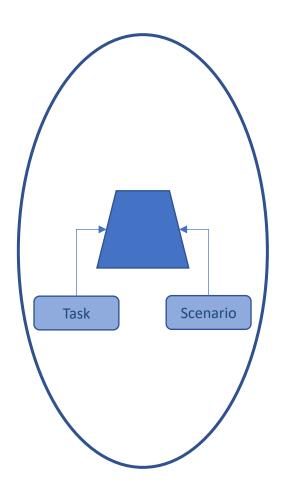




Shared-specific network paradigm $wL(E^{Merge}, \boldsymbol{\theta}^{sh}, \boldsymbol{\theta}^{t}, \boldsymbol{\theta}^{s})$



Dynamic weight $wL(\mathbf{E}^{Merge}, \boldsymbol{\theta}^{sh}, \boldsymbol{\theta}^{t}, \boldsymbol{\theta}^{s})$



Multi-Scenario & Multi-Task $wL(\mathbf{E}^{Merge}, \mathbf{\theta}^{sh}, \mathbf{\theta}^{t}, \mathbf{\theta}^{s})$

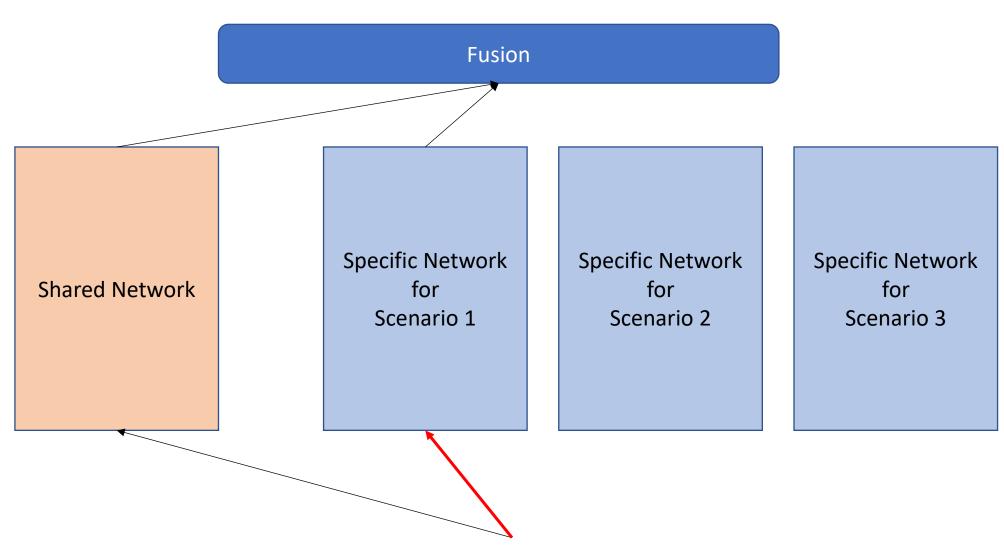
Shared-specific Network Paradigm









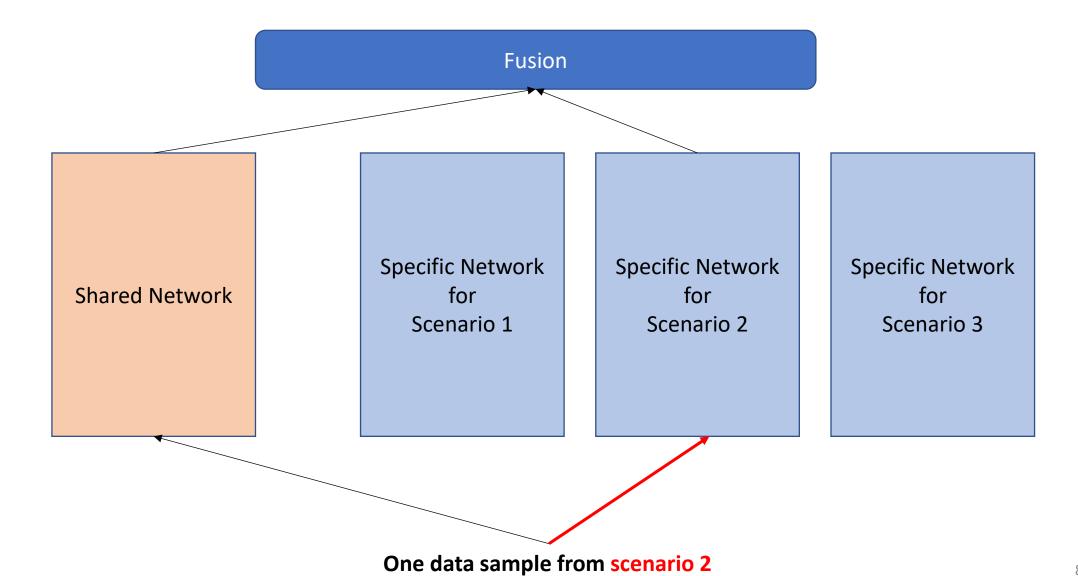


Shared-specific Network Paradigm









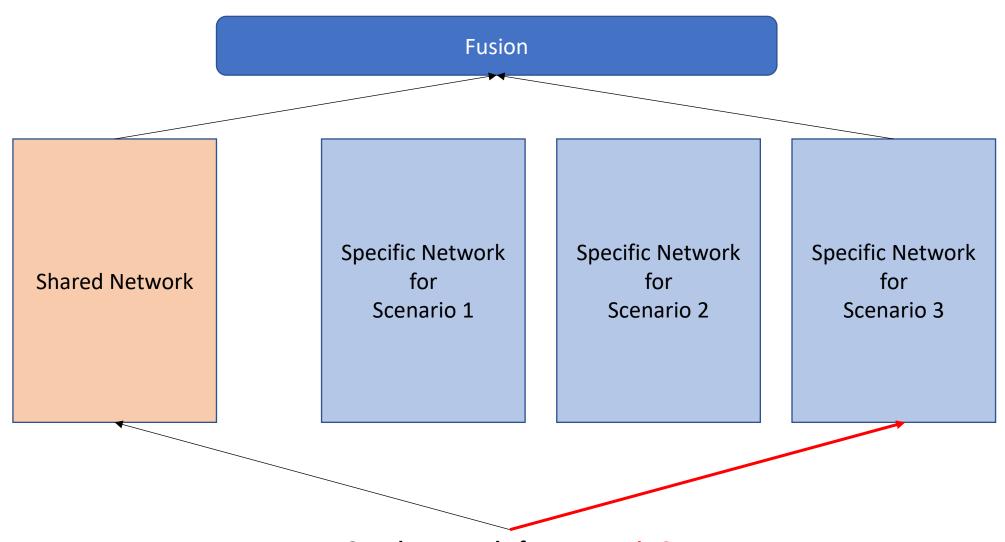
Shared-specific Network Paradigm











STAR









Motivation:

- Training individual models for each domain → does not fully use the data from all domains
- Data across domains owns commonalities and characteristics

> Target:

- Use a single model to serve multiple domains simultaneously
- Shared network → commonalities
- Specific network → characteristics

> Methods:

- Partitioned Normalization
- STAR Topology
- Auxiliary Network



Banner



Guess What You Like

STAR Details









- Partitioned Normalization (PN)
- > Training

$$z' = (\gamma * \gamma_p) rac{z - \mu}{\sqrt{\sigma^2 + \epsilon}} + (eta + eta_p)$$

Compared to BN

Testing

$$\mathbf{z}' = (\gamma * \gamma_p) rac{\mathbf{z} - E_p}{\sqrt{Var_p + \epsilon}} + (eta + eta_p)$$

- ➤ Batch Normalization (BN)
- > Training

$$\mathbf{z'} = \gamma \frac{\mathbf{z} - \mu}{\sqrt{\sigma^2 + \epsilon}} + \beta$$

Testing

$$\mathbf{z'} = \gamma \frac{\mathbf{z} - E}{\sqrt{Var + \epsilon}} + \beta$$

STAR Details









STAR Topology

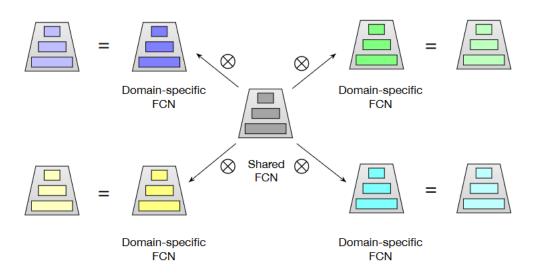
The final weight and bias for p-th domain is obtained by:

$$W_p^\star = W_p \otimes W, b_p^\star = b_p + b$$

The output for p-th domain is derived by:

$$out_p = \phi((W_p^\star)^ op in_p + b_p^\star)$$

⊗ element-wise product



SAR-Net









≻Motivation

• Traffic characteristics of different scenarios are significantly different (individual data scale or representative topic)

≻Target

Train a unified model to serve all scenarios

SAR-Net Details

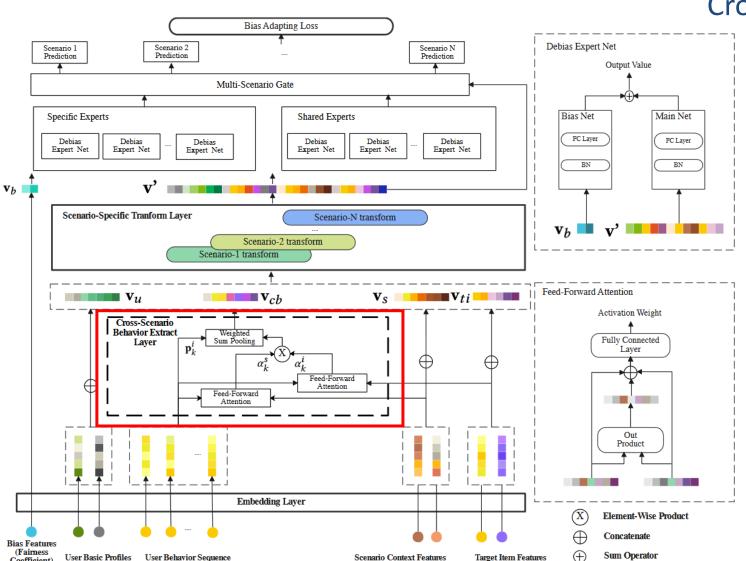








Cross-Scenario Behavior Extract Layer



How to aggregate the sequence?

 $p(B^i)$ is item behavior sequence

$$\mathbf{p}(B^i) = \{\mathbf{p}_1^i, \mathbf{p}_2^i, \cdots, \mathbf{p}_{|\mathbf{p}(B^i)|}^i\}$$

$$\mathbf{p}_{k}^{i} = [\mathbf{e}_{itemId} || \mathbf{e}_{destination} || \mathbf{e}_{category} || \cdots]$$

 $p(B^s)$ is scenario context sequence

 $\mathbf{p}_{k}^{s} = [\mathbf{e}_{scenarioId} || \mathbf{e}_{scenarioType} || \mathbf{e}_{behaviorTime} || \cdots]$

$$\mathbf{p}(B^s) = \{\mathbf{p}_1^s, \mathbf{p}_2^s, \cdots, \mathbf{p}_{|\mathbf{p}(B^s)|}^i\}$$

$$\alpha_k^i = \frac{\exp(\psi(\mathbf{p}_k^i, \mathbf{p}_t^i))}{\sum_{l=1}^{|\mathbf{p}(B^i)|} \exp(\psi(\mathbf{p}_l^i, \mathbf{p}_t^i))},$$

$$\alpha_k^s = \frac{\exp(\psi(\mathbf{p}_k^s, \mathbf{p}_t^s))}{\sum_{l=1}^{|\mathbf{p}(B^s)|} \exp(\psi(\mathbf{p}_l^s, \mathbf{p}_t^s))},$$

 α_k^i and α_k^s indicate the relevance between user's kth behavior item and the target item or target scenario

SAR-Net Details

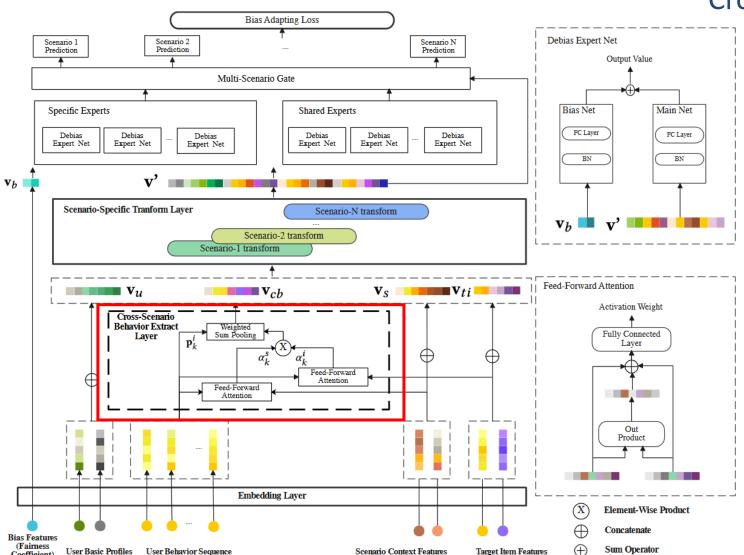








Cross-Scenario Behavior Extract Layer



How to aggregate the sequence?

$$\begin{split} \alpha_k^i &= \frac{\exp(\psi(\mathbf{p}_k^i, \mathbf{p}_t^i))}{\sum_{l=1}^{|\mathbf{p}(B^i)|} \exp(\psi(\mathbf{p}_l^i, \mathbf{p}_t^i))}.,\\ \alpha_k^s &= \frac{\exp(\psi(\mathbf{p}_k^s, \mathbf{p}_t^s))}{\sum_{l=1}^{|\mathbf{p}(B^s)|} \exp(\psi(\mathbf{p}_l^s, \mathbf{p}_t^s))}, \end{split}$$

$$\mathbf{p}_{k}^{i} = [\mathbf{e}_{itemId} || \mathbf{e}_{destination} || \mathbf{e}_{category} || \cdots]$$

$$\mathbf{v}_{cb} = \sum_{k=1}^t lpha_k^i st lpha_k^s st \mathbf{p}_k^i$$

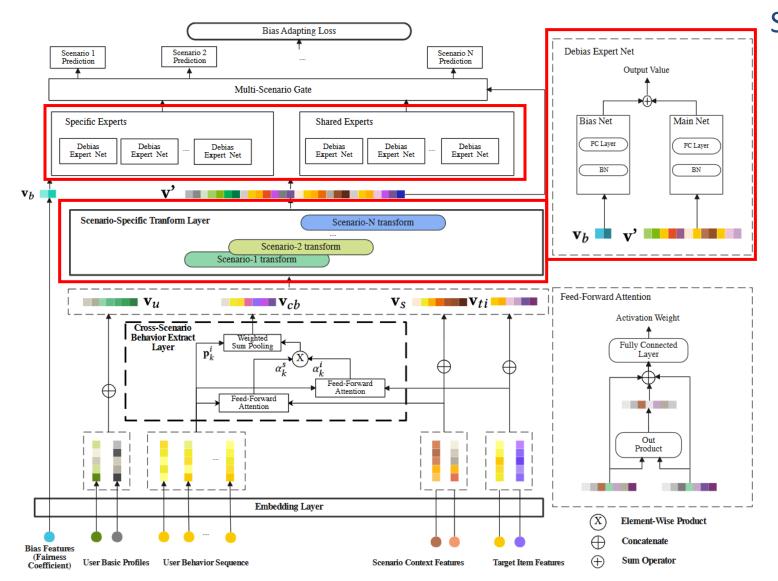
SAR-Net











Scenario-Specific Transform Layer

$$\mathbf{v'} = \mathbf{v} \otimes eta_i + \gamma_i$$

Mixture of Debias Experts

Multi-expert network. Each scenario has some scenario-specific experts and all the scenarios share several common experts.

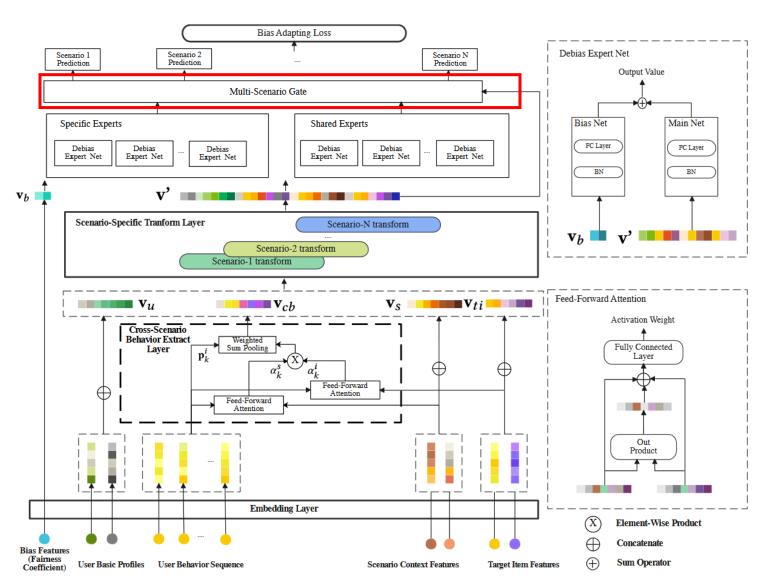
SAR-Net











Multi-Gate Network & Prediction

The output of experts:

$$S^{k}(x) = [o_{k,1}, o_{k,2}, \cdots, o_{k,m_{k}}, o_{s,1}, o_{s,2}, \cdots, o_{s,m_{s}}]^{T}$$

Final predicted score of scenario *k*

$$y^k(x) = w^k(x)S^k(x)$$

 $w^k(x)$ is derived by a single-layer feedforward network with a SoftMax activation function









► Motivation

- Separate model for each scenario, ignoring the cross-domain overlapping of user groups and items
- One shared model trained on mix data, model performance may decrease when different domains conflict

≻Target

- Modeling commonalities and diversities → common networks and domain-specific networks
- Tackle the feature-level domain adaptation → domain-specific batch normalization, domain interest adaptation layer



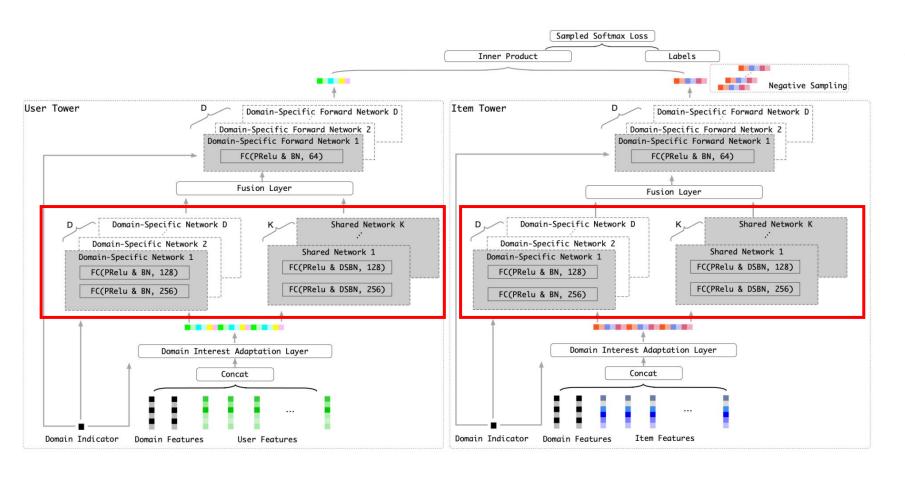






Backbone Network

Shared Network & Domain-Specific Network



$$egin{aligned} \mathbf{a}z_k &= rac{W_{shared}^k(f_{domain}) + b_{shared}^k}{\sum_{n=1}^K (W_{shared}^n(f_{domain}) + b_{shared}^n)} \ E_{shared} &= \sum_{k=1}^K lpha_k M L P_{shared}^k(\mathbf{F}) \ E_{spec}^{(d)} &= M L P_{spec}^{(d)}(\mathbf{F}^{(d)}) \end{aligned}$$

 f_{domain} Domain indicator embedding

 $\mathbf{F}^{(d)}$ Data from domain d

K hyperparameter, number of Shared Network

D domains, D Domain-Specific Network

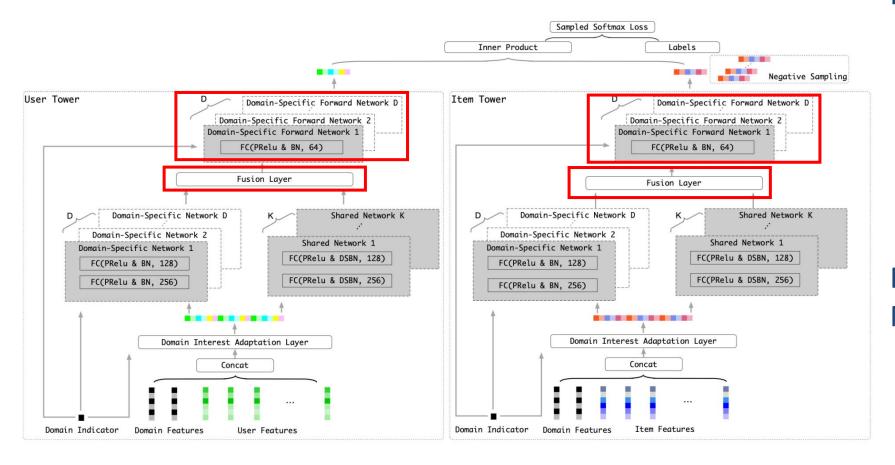








Backbone Network



Fusion Layer

$$eta_1^{(d)} = \sigma(W_{fusion_spec}^{(d)}(f_{domain})) \ eta_2^{(d)} = \sigma(W_{fusion_shared}^{(d)}(f_{domain}))$$

$$E_{fusion}^{(d)} = concat(eta_1^{(d)} E_{spec}^{(d)} \mid \ eta_1^{(d)} E_{spec}^{(d)} \odot eta_2^{(d)} E_{shared} \mid eta_2^{(d)} E_{shared})$$

Domain-Specific Forward Network

$$E = FC_{forward}^{(d)}(E_{fusion}^{(d)})$$



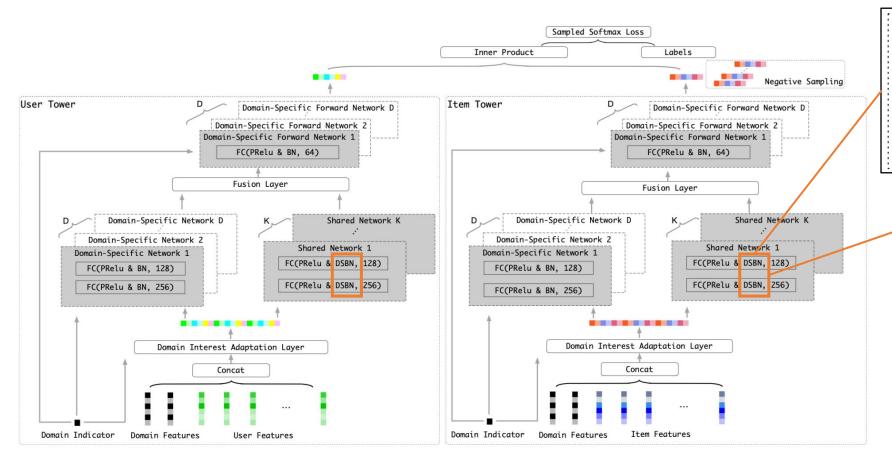


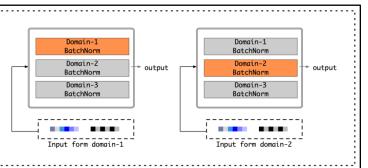




Domain Adaptation

Domain-Specific Batch Normalization (DSBN)





$$\hat{\mathbf{X}}^{(d)} = lpha^{(d)} rac{\mathbf{X}^{(d)} - \mu^{(d)}}{\sqrt{(\sigma^{(d)})^2 + \epsilon}} + eta^{(d)}$$



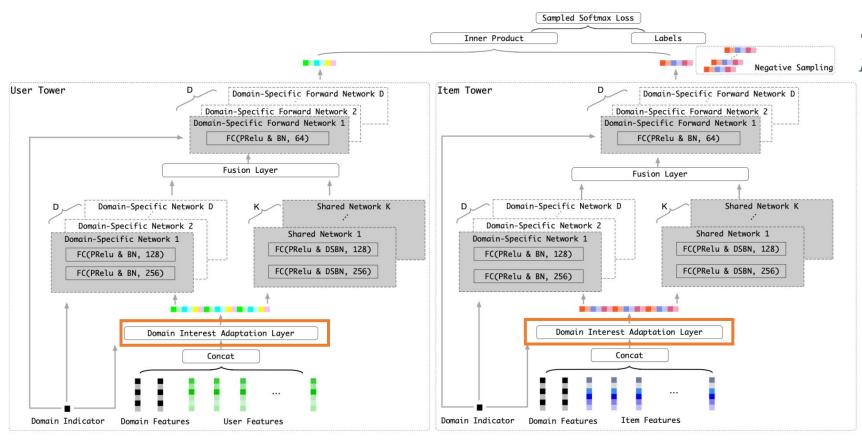






Domain Adaptation

Domain Interest Adaptation Layer



$$egin{aligned} lpha^{(d)} &= F_{se}(concat(F_{avg}(F_1^{(d)}) \mid \cdots \mid F_{avg}(F_N^{(d)}))) \ \hat{F}^{(d)} &= lpha^{(d)} \otimes concat(F_1^{(d)} \mid \cdots \mid F_N^{(d)}) \end{aligned}$$

 $F_i^{(d)}$ denotes ith feature of embedded input collected from domain d

 F_{se} denotes a (FC, Relu, FC) block and F_{avg} denotes average pooling operator.

Uni-CTR









- Due to varying data sparsity in different domains, models can easily be dominated by specific domains, leading to "seesaw phenomenon"
- Existing methods are difficult to handle newly added domain

≻Method

- Leveraging LLM to extract layer-wise representations to capture domain commonalities in order to migrate "seesaw phenomenon"
- Incorporating a pluggable domain-specific network to capture domain characteristics, ensuring scalability to new domains

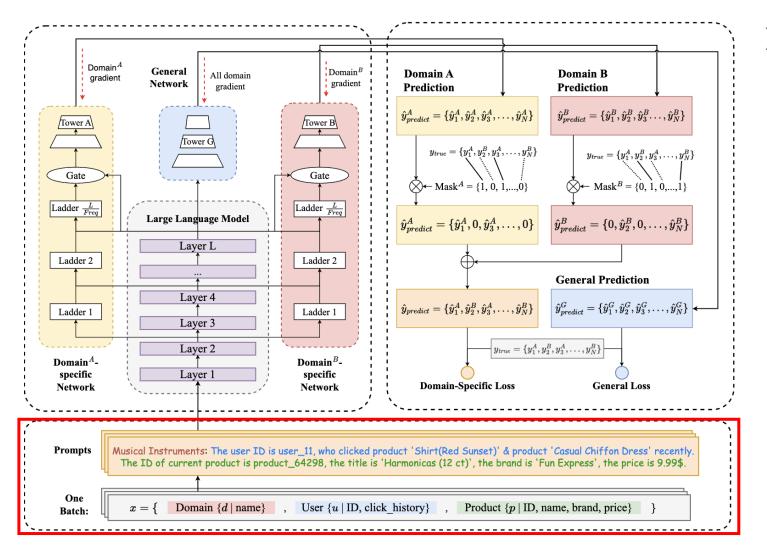
Uni-CTR Details











- Prompt-based Semantic Modeling
 - Capture rich semantic information via text-based features
 - Input
 - Domain Context
 - User Information
 - Product Informatio

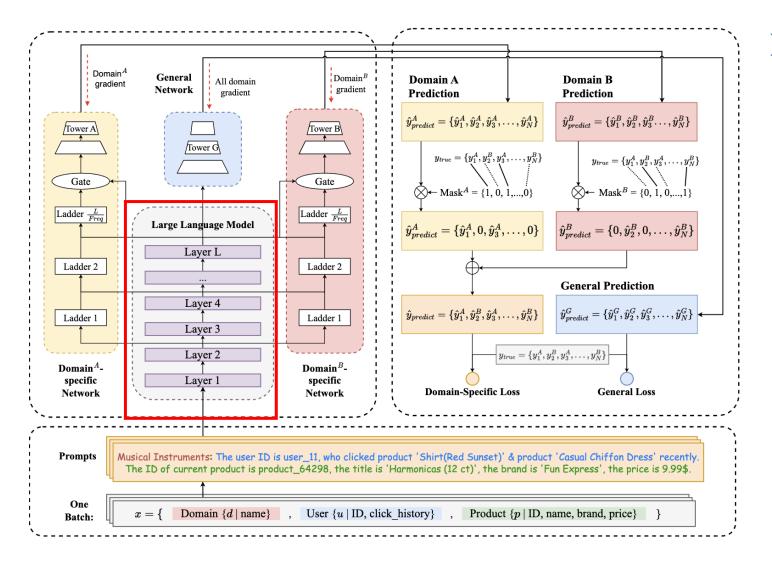
Uni-CTR Details











➤ Uni-CTR architecture

LLM Backbone

$$m{x}_{ ext{tokens}} = ext{Tokenizer}(x_{ ext{text}}) = \{t_0, t_1, \dots, t_J\},$$
 $m{h}_0 = m{E}_{embed}(m{x}_{ ext{tokens}}) = \{m{e}_0, m{e}_1, \dots, m{e}_J\}.$ $m{h}_l = ext{Transformer}_l(m{h}_{l-1}), l \in \{1, 2, \dots, L\},$

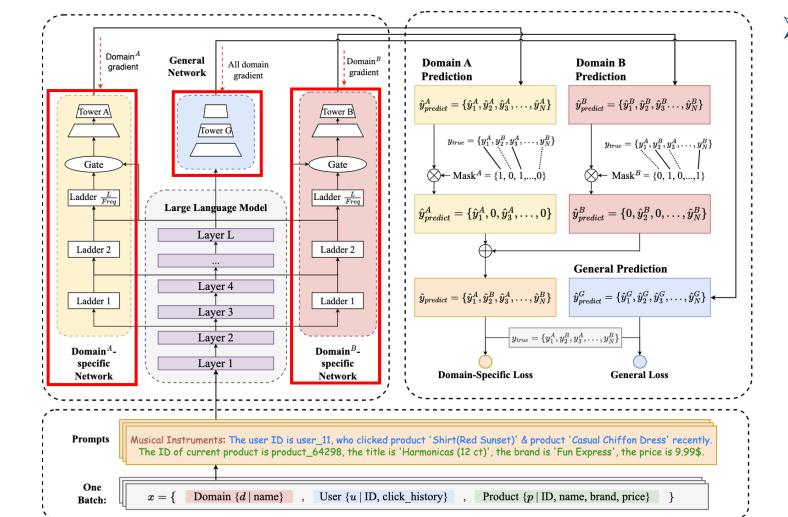
Uni-CTR Details











➤ Uni-CTR architecture

LLM Backbone

$$m{x}_{ ext{tokens}} = ext{Tokenizer}(x_{ ext{text}}) = \{t_0, t_1, \dots, t_J\},$$
 $m{h}_0 = m{E}_{embed}(m{x}_{ ext{tokens}}) = \{m{e}_0, m{e}_1, \dots, m{e}_J\}.$ $m{h}_l = ext{Transformer}_l(m{h}_{l-1}), l \in \{1, 2, \dots, L\},$

- Domain-Specific Network
 - Ladder Netowork
 - Gate Net
 - Tower Net
- General Network

$$\hat{m{y}}^G = ext{MLP}(m{h}_L; m{W}_\sigma^G, m{b}_\sigma^G)$$

M-LoRA







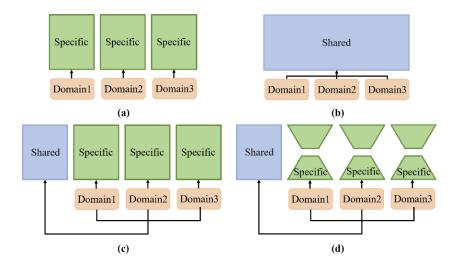


≻Motivation

- Suffering from data sparsity and ignoring domain relations
- Failing to capture domain diversity
- Suffering from a sharp increase in model parameters

≻Method

Incorporating Low-Rank Adaptor (LoRA) for multi-domain fine-tuning



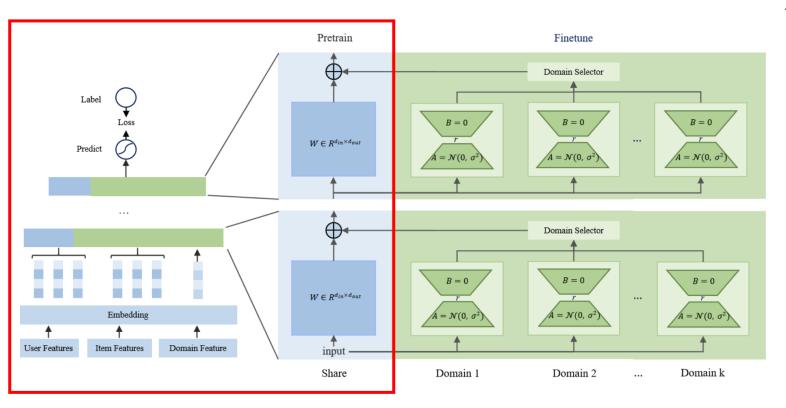
M-LoRA Details











- Pre-training
 - Lage-scale pre-training dataset
 - Shared Network and Embedding layer

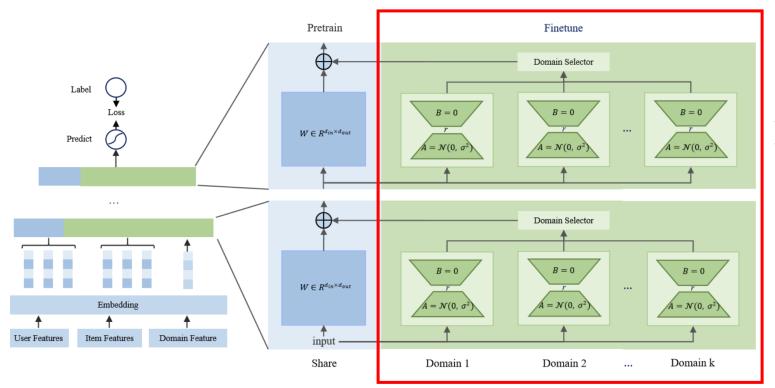
M-LoRA Details











- Pre-training
 - Lage-scale pre-training dataset
 - Shared Network and Embedding layer
- > Fine-tuning
 - LoRA module is integrated in each layer for each domain, including A and B

$$\mathbf{h}_t = \mathbf{W}\mathbf{x} + \Delta \mathbf{W}_t = \mathbf{W}\mathbf{x} + \mathbf{B}_t \mathbf{A}_t \mathbf{x},$$

where $\mathbf{B} \in R^{d_{out} \times r}$ and $\mathbf{A} \in R^{r \times d_{in}}$

r is not fixed for a more flexible representation

$$r = \max(\frac{d_{out}}{\alpha}, 1).$$

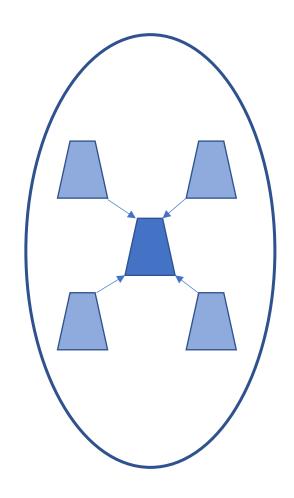
Table of Contents

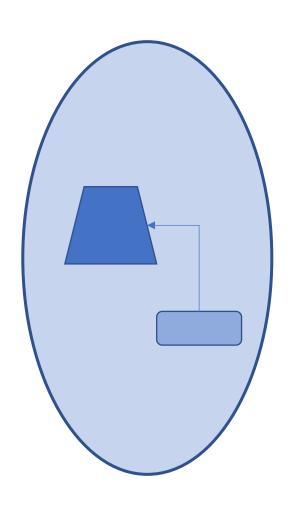


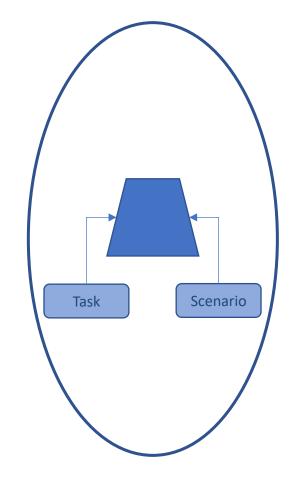












Shared-specific network paradigm

 $wL(E^{Merge}, \Theta, \Theta^t, (\Theta^{shared}, \Theta^{specific}))$

Dynamic weight

 $wL(E^{Merge}, \Theta, \Theta^t, \Theta^s)$

Multi-Scenario & Multi-Task

 $wL(E^{Merge}, \Theta, \Theta^t, \Theta^s, \Theta^T)$

Dynamic Weight

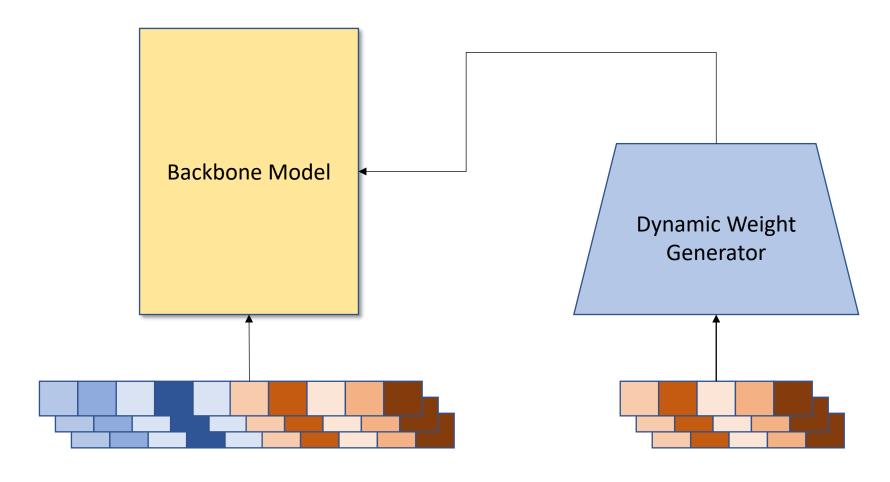








➤ Why Dynamic?



Input Features

Scenario Sensitive Features

HAMUR





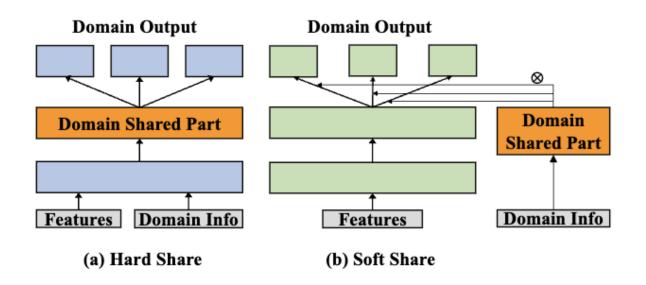




- Previous research relies on explicit sharing across different domains
- Static parameters constrain the representation of different domains

Methods

- Adapter for multi-domain dynamic adaptation
- Hyper-net for dynamic generation parameters for adapters



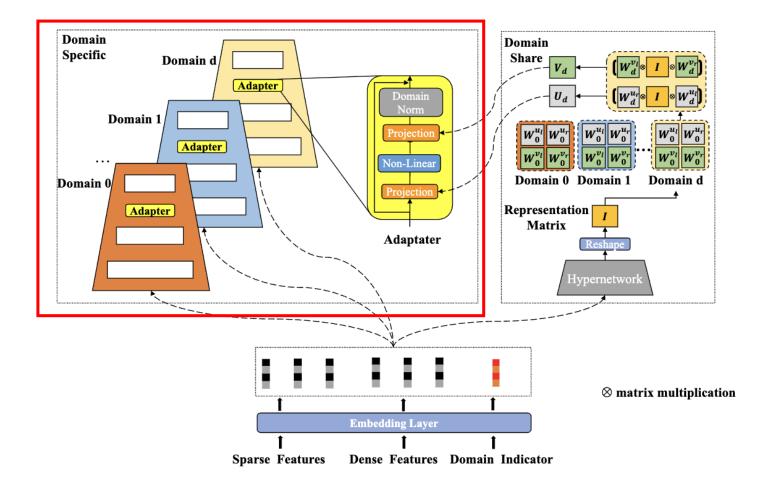
HAMUR











Domain-specific adapter

$$A_d(\mathbf{x}) = DN_d \Big(V_d(\sigma(U_d(\mathbf{x}))) \Big) + \mathbf{x}$$
$$DN_d = \gamma_d \odot \frac{\mathbf{x} - \mu}{\sqrt{\sigma^2 + \epsilon}} + \beta_d$$

- Domain Shared Hyper-Network
 - Parameters Generation

$$\mathbf{h}^i = \mathcal{H}(\mathbf{z}^i)$$
 $\mathbf{I}^i = reshape(\mathbf{h}^i)$

Low-Rank Decomposition

$$U_d^i = W_d^{u_l} \cdot I^i \cdot W_d^{u_r}$$
$$V_d^i = W_d^{v_l} \cdot I^i \cdot W_d^{v_r}$$

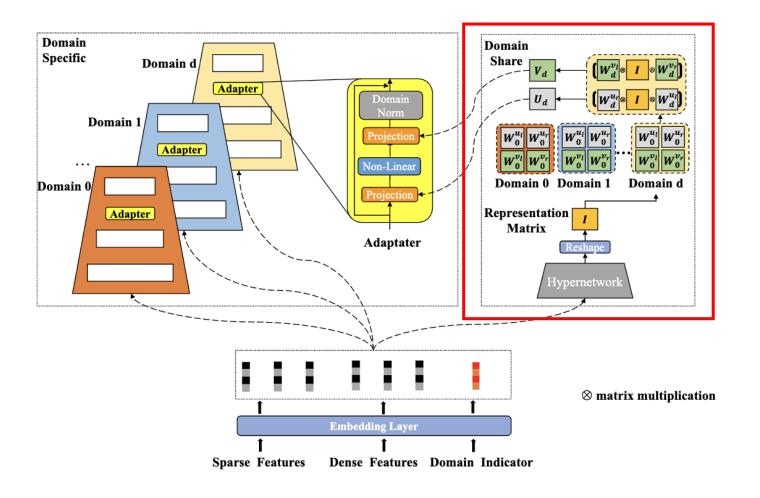
HAMUR











Domain-specific adapter

$$A_d(\mathbf{x}) = DN_d \Big(V_d(\sigma(U_d(\mathbf{x}))) \Big) + \mathbf{x}$$
$$DN_d = \gamma_d \odot \frac{\mathbf{x} - \mu}{\sqrt{\sigma^2 + \epsilon}} + \beta_d$$

- Domain Shared Hyper-Network
 - Parameters Generation

$$\mathbf{h}^i = \mathcal{H}(\mathbf{z}^i)$$
 $\mathbf{I}^i = reshape(\mathbf{h}^i)$

Low-Rank Decomposition

$$U_d^i = W_d^{u_l} \cdot I^i \cdot W_d^{u_r}$$
$$V_d^i = W_d^{v_l} \cdot I^i \cdot W_d^{v_r}$$

HierRec







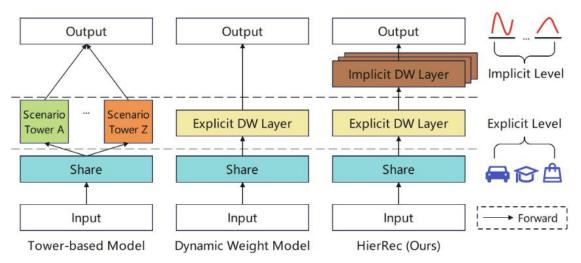


≻Motivation

- Current multi-scenario models mainly rely on explicit scenario modeling based on manually defined scenario IDs (like ad slots or channels).
- These manual rules are coarse-grained, rigid, and potentially biased, and they ignore internal variations within each scenario

≻Method

 Propose HierRec that models both explicit and implicit scenarios in a hierarchical and adaptive way



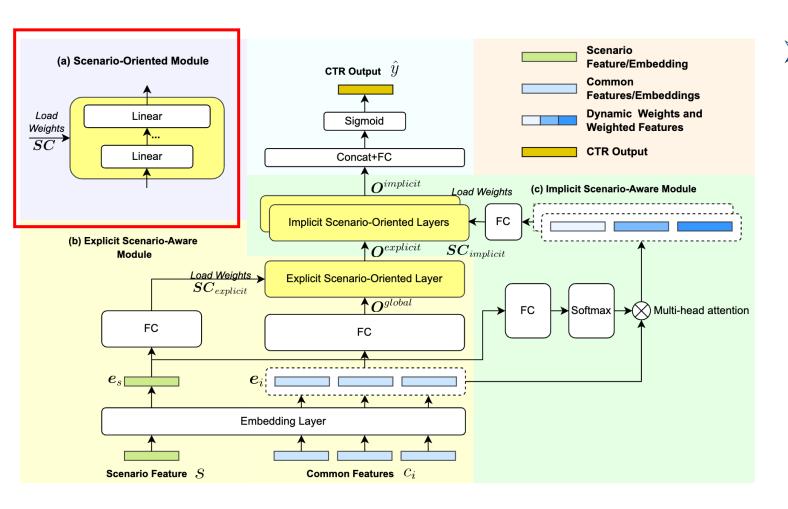
HierRec











- Scenario-Oriented Module
 - Adaptively generate parameters depending on scenario condition (SC)

$$\boldsymbol{W}_l, \boldsymbol{b}_l = Reshape(\boldsymbol{SC})[l] \quad l \in [1, L],$$

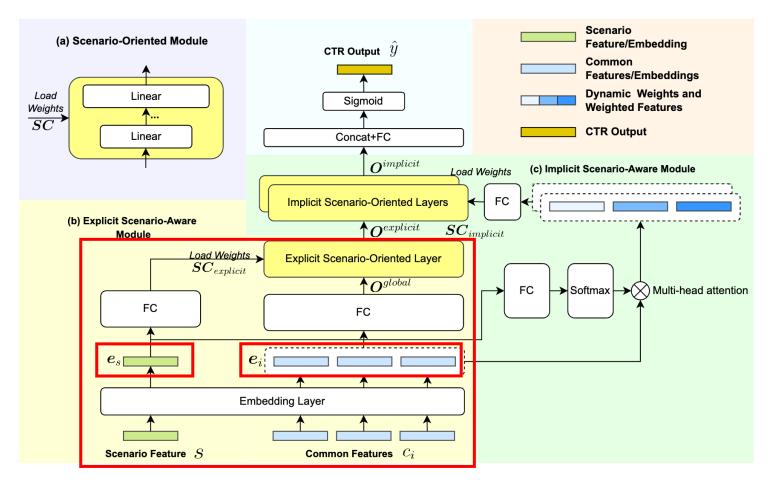
HierRec











- Scenario-Oriented Module
 - Adaptively generate parameters depending on scenario condition (SC)

$$\boldsymbol{W}_{l}, \boldsymbol{b}_{l} = Reshape(\boldsymbol{SC})[l] \quad l \in [1, L],$$

- Explicit Scenario-Aware Module
 - Model coarse-grained explicit scenario information

$$egin{aligned} egin{aligned} oldsymbol{e}_i &= oldsymbol{E}oldsymbol{M}_i \cdot Onehot(c_i), & i \in [1,I] \ oldsymbol{e}_s &= oldsymbol{E}oldsymbol{M}_s \cdot Onehot(s), \end{aligned}$$
 $oldsymbol{S}oldsymbol{C}_{explicit} &= FC(oldsymbol{e}_s).$

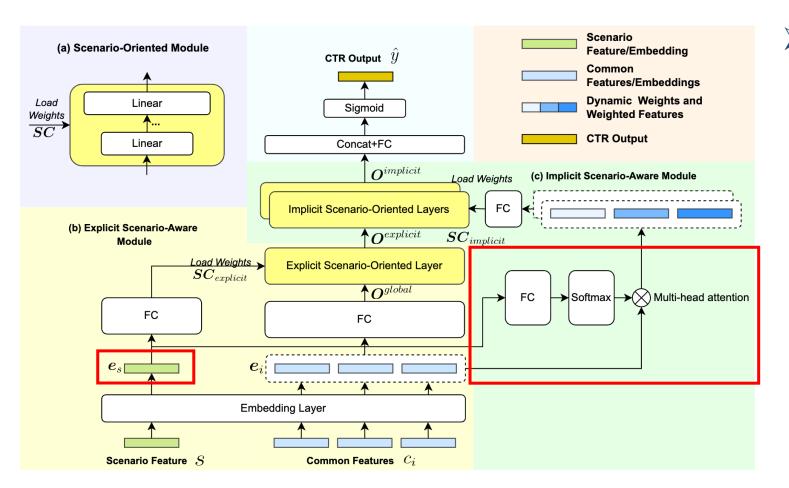
HierRec











- ➤ Implicit Scenario-Aware Module
 - Model fine-grained implicit scenario information

$$egin{cases} oldsymbol{weight}_{ori} = Reshape(FC(oldsymbol{e}_s)) \ oldsymbol{weight}_{norm}[g] = Softmax(oldsymbol{weight}_{ori}[g]), \ g \in [1, G] \end{cases}$$

$$oldsymbol{IE} = oldsymbol{weight}_{norm} \otimes oldsymbol{E}_c,$$

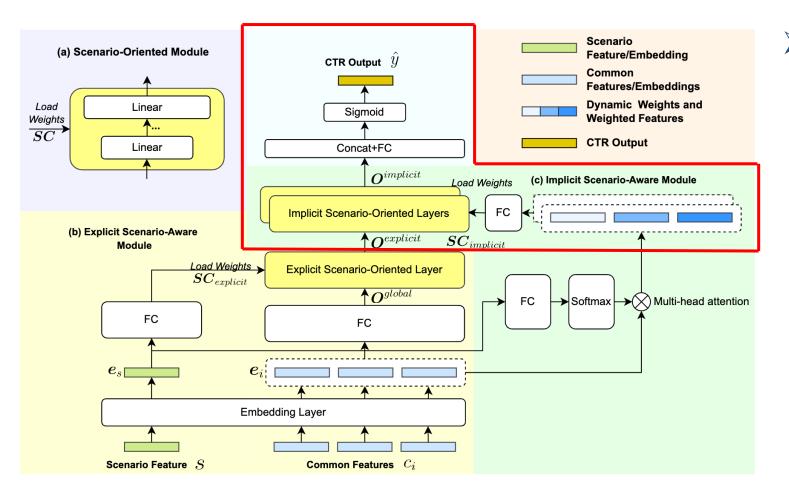
HierRec











- ➤ Implicit Scenario-Aware Module
 - Model fine-grained implicit scenario information

$$egin{cases} m{weight}_{ori} = Reshape(FC(m{e}_s)) \ m{weight}_{norm}[g] = Softmax(m{weight}_{ori}[g]), \ g \in [1,G] \end{cases}$$

$$oldsymbol{IE} = oldsymbol{weight}_{norm} \otimes oldsymbol{E}_c,$$

$$SC_{implicit}[g] = FC(IE[g]), \quad g \in [1, G].$$

$$\hat{y} = Sigmoid(FC(Concat(\boldsymbol{O}_{1}^{implicit}, ..., \boldsymbol{O}_{G}^{implicit}))).$$

LLM4MSR









≻ Motivation

- Insufficient scenario knowledge is incorporated
- Users' personalized preferences across scenarios tend to be ignored

≻Method

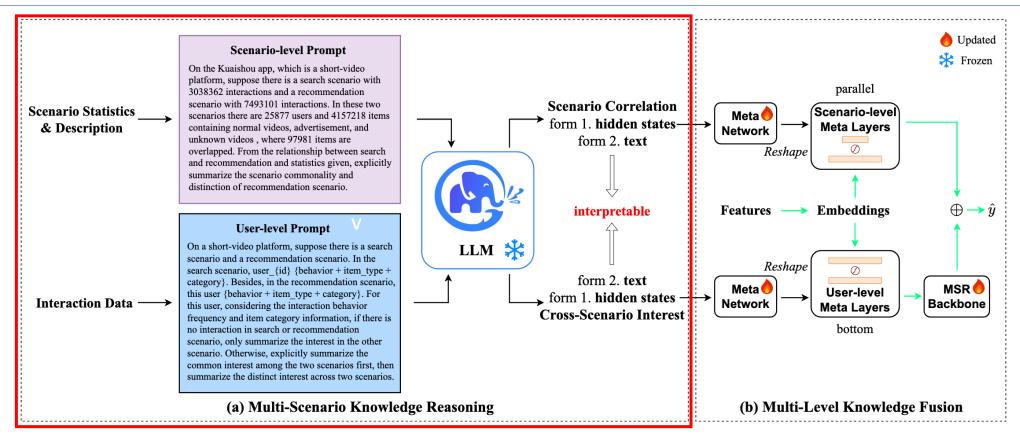
- Resorting to use LLM and hierarchical meta-networks
- Using LLM to grasp cross-scenario correlation and personalized preferences
- Using meta-network as a bridge connecting the semantic space in LLM and recommendation space in the multi-scenario backbone model

LLM4MSR Details









- Multi-scenario Knowledge Reasoning
 - Scenario-level prompt construction
 - User-level prompt construction

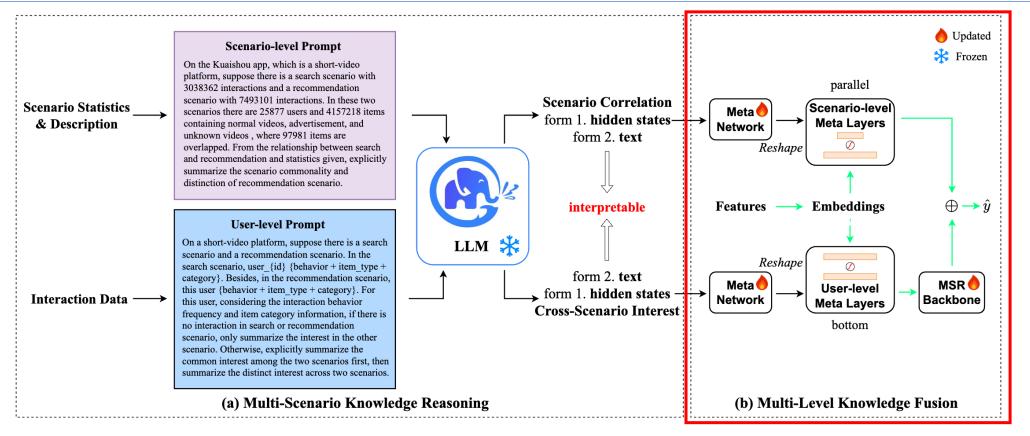
LLM4MSR Details











- Multi-Level Knowledge Fusion
 - Meta-net generates meta layers to fuse the scenario- and user-level knowledge

$$m{h}_{mw}, m{h}_{mb} = ext{Meta Network}(m{h}_{LLM}),$$
 $m{W}_l^{(i)} = Reshape(m{h}_{mw})$
 $m{b}_l^{(i)} = Reshape(m{h}_{mb}), i \in \{1, 2, \dots, K\}$

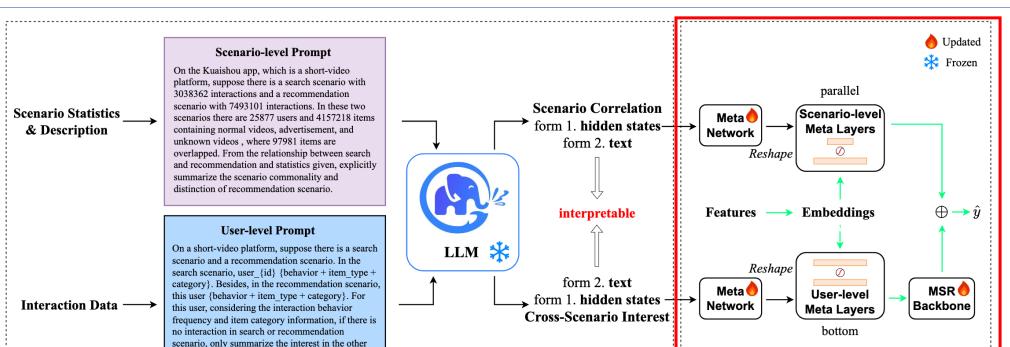
LLM4MSR Details



(b) Multi-Level Knowledge Fusion







Multi-Level Knowledge Fusion

scenario. Otherwise, explicitly summarize the common interest among the two scenarios first, then summarize the distinct interest across two scenarios.

• Prediction
$$\begin{aligned} \boldsymbol{h}^{(i)} &= \sigma(\boldsymbol{W}_l^{(i)} \boldsymbol{h}^{(i-1)} + \boldsymbol{b}_l^{(i)}), i \in \{1, 2, \dots, K\} \\ \boldsymbol{h} &= \mathrm{MSR}(\boldsymbol{h}_u^{(K)}), \\ \hat{y} &= \sigma'(\alpha \cdot \boldsymbol{h}_s^{(K)} + (1 - \alpha) \cdot \boldsymbol{h}), \end{aligned}$$

(a) Multi-Scenario Knowledge Reasoning

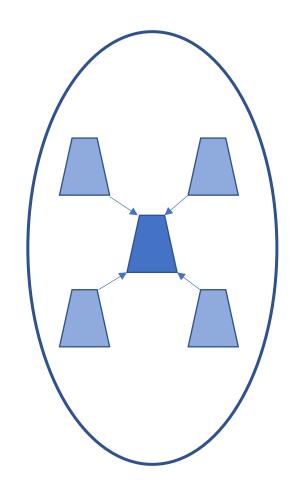
Table of Contents

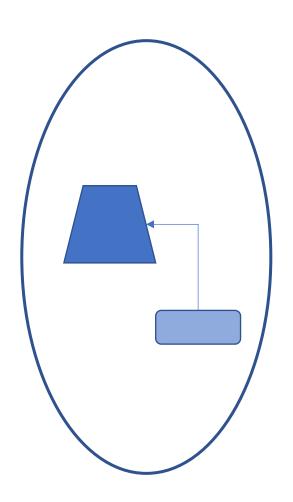


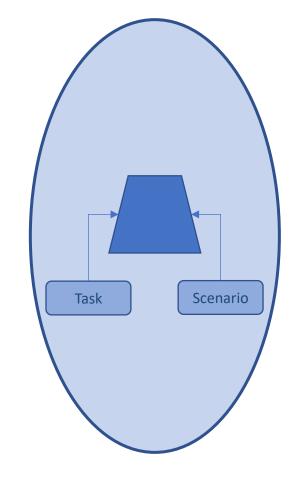












Shared-specific network paradigm

 $wL(E^{Merge}, \Theta, \Theta^t, (\Theta^{shared}, \Theta^{specific}))$

Dynamic weight

 $wL(E^{Merge}, \Theta, \Theta^t, \Theta^s)$

Multi-Scenario & Multi-Task

 $wL(E^{Merge}, \Theta, \Theta^t, \Theta^s, \Theta^T)$

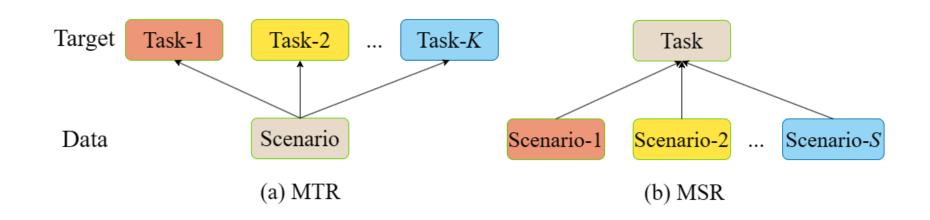
Multi-Scenario & Multi-Task Studies











PEPNet









Motivation

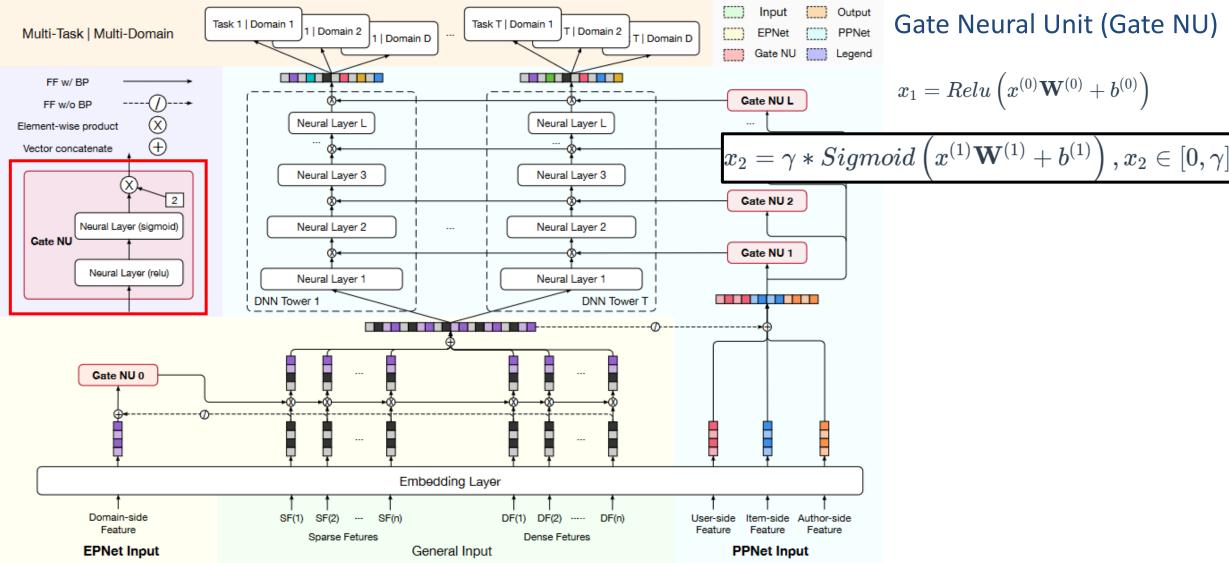
- The imperfectly double seesaw phenomenon
- More accurate personalization estimates can alleviate the imperfectly double seesaw problem

> Target

- Jointly model multi-domain and multi-task
- an efficient, low-cost deployment and plug-andplay method that can be injected in any network.







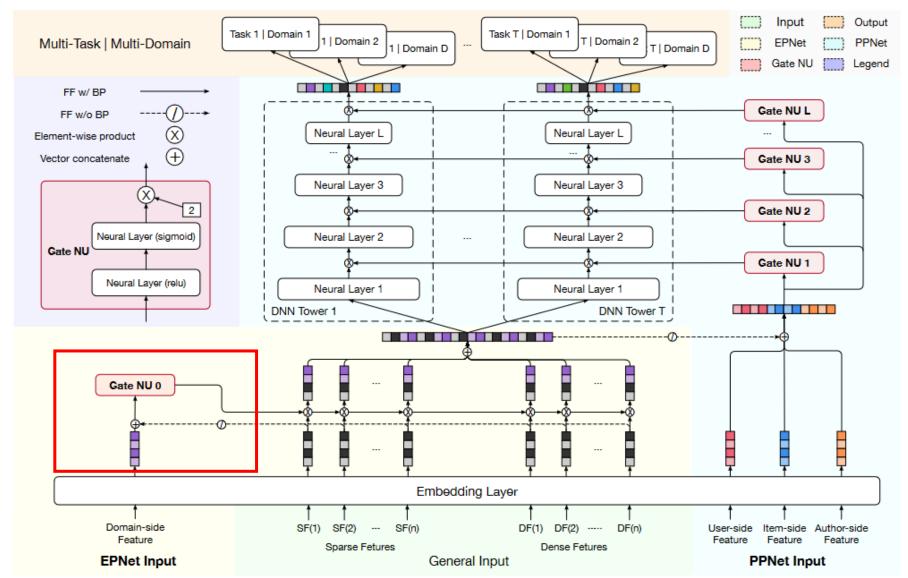
PEPNet Details











EPNet

$$\mathbf{E} = E(\mathcal{F}_S) \oplus E(\mathcal{F}_D)$$

Embeddings of sparse features and dense features

$$\delta_{domain} = \mathrm{U}_{ep}(E(\mathcal{F}_d) \oplus (\oslash(\mathbf{E})))$$

$$\mathbf{O}_{ep} = \delta_{domain} \otimes \mathbf{E}$$

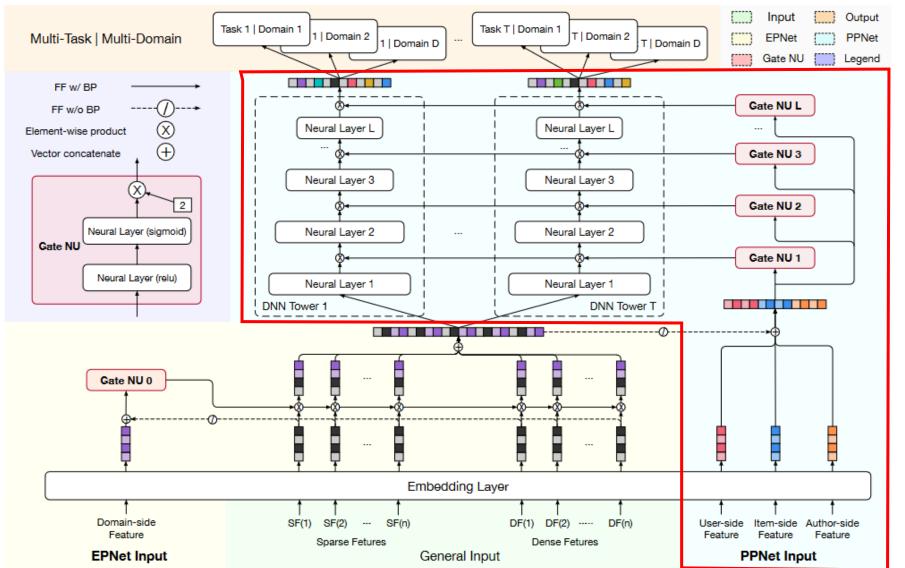
PEPNet Details











PPNet

$$egin{aligned} &0_{prior} = E(uf) \oplus E(if) \oplus E(af) \ &\delta_{task} = \mathbf{U}_{pp}(\mathbf{O}_{prior} \oplus (\oslash(\mathbf{O}_{ep}))) \ &\mathbf{O}_{pp}^{(l)} = oldsymbol{\delta}_{task}^{(l)} \otimes \mathbf{H}^{(l)}, \ &\mathbf{H}^{(l+1)} = f(\mathbf{O}_{pp}^{(l)} \mathbf{W}^{(l)} + b^{(l)}), l \in \{1,...,L\} \end{aligned}$$

M₂M









Motivation

- Less attention has been drawn to advertisers
- Major e-commerce platforms provide multiple marketing scenarios.

Methods

- Meta unit
- Meta attention module
- Meta tower module



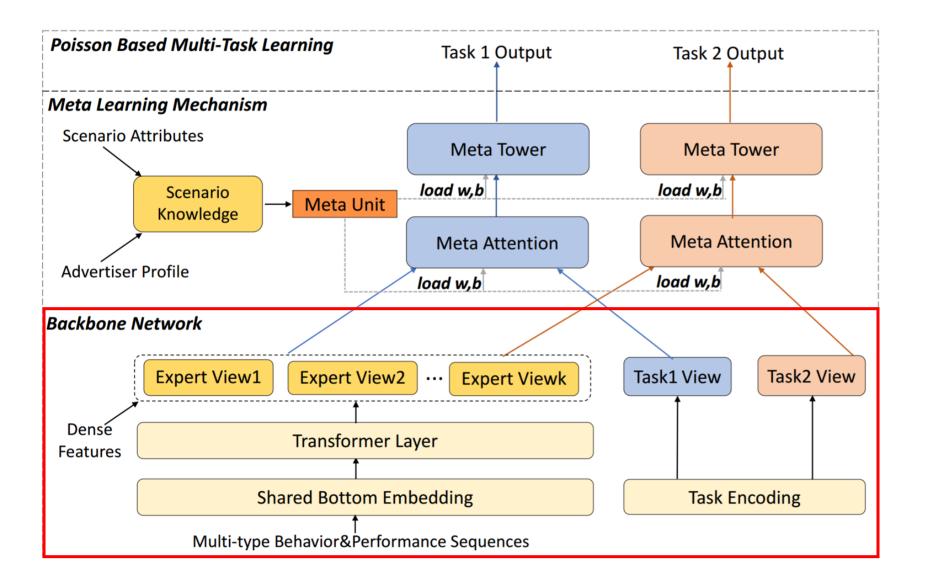
M2M Overview









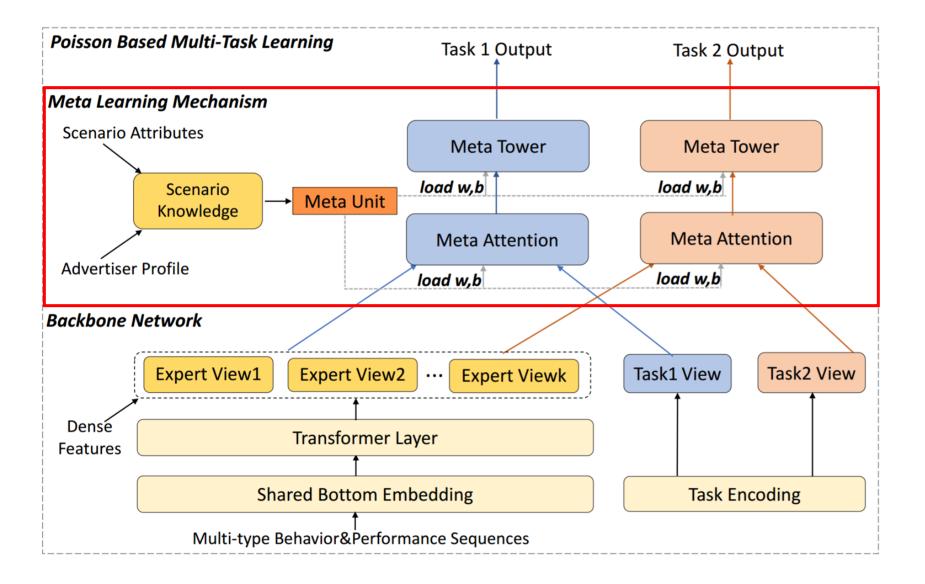


M2M Overview









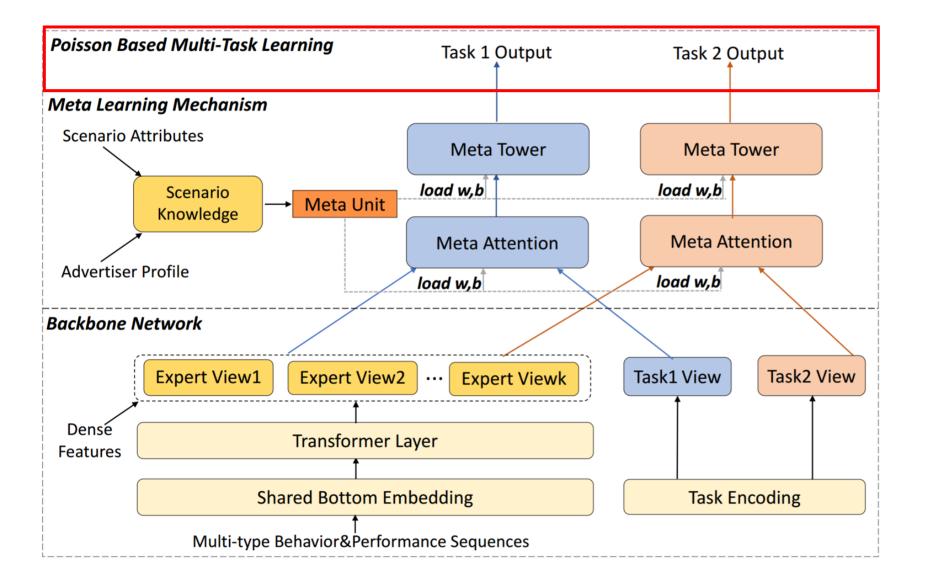
M2M Overview









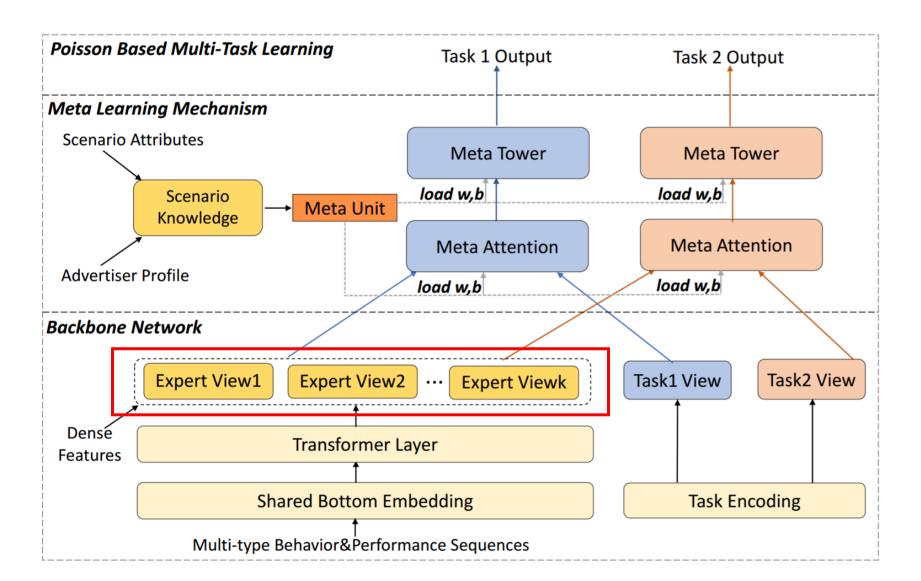












Backbone Network Expert View Representation

$$ext{E}_{ ext{i}} = f_{m{M}LP}(\mathbf{F}), orall i \in 1, 2, .., k$$

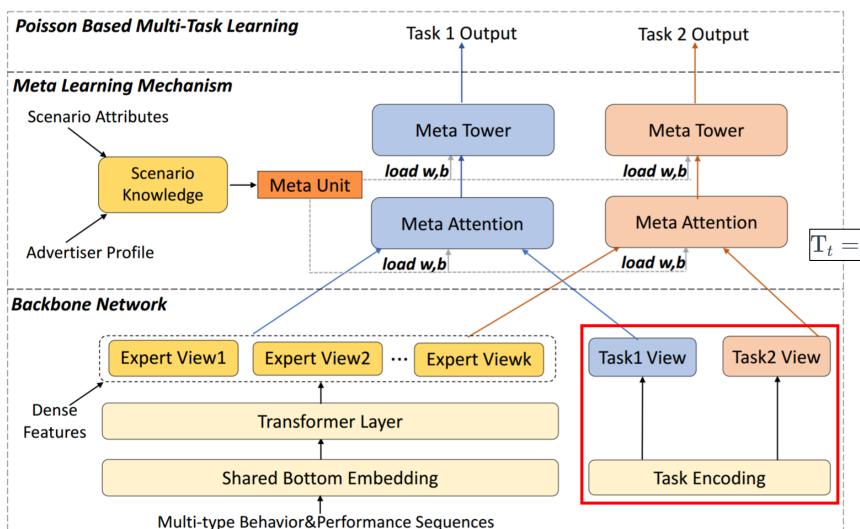
F is the output of transformer layer











Backbone Network

Expert View Representation

$$ext{E}_{ ext{i}} = f_{m{M}LP}(\mathbf{F}), orall i \in {1,2,..,k}$$

F is the output of transformer layer

Task View Representation

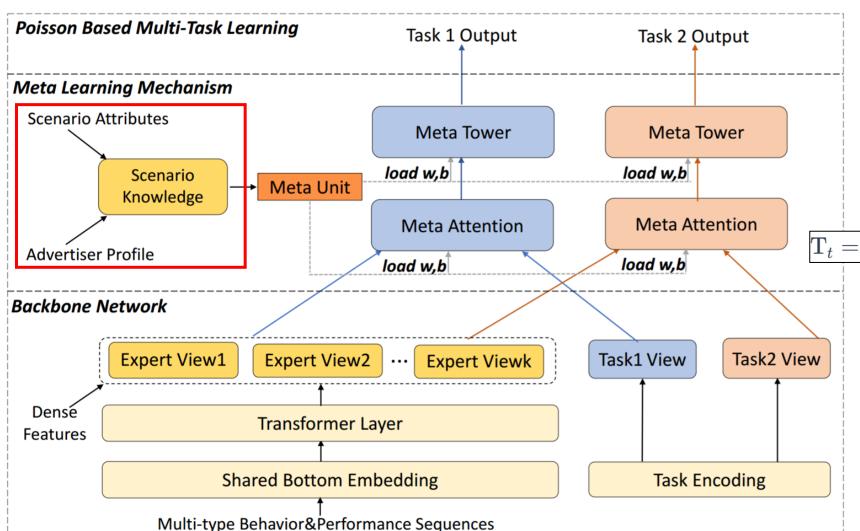
$$ext{T}_t = f_{MLP}(Embedding(t)), orall t \in 1, 2, .., m$$











Backbone Network

Expert View Representation

$$ext{E}_{ ext{i}} = f_{m{M}LP}(\mathbf{F}), orall i \in {1,2,..,k}$$

F is the output of transformer layer

Task View Representation

$$ext{T}_t = f_{m{M}LP}(Embedding(t)), orall t \in 1, 2, .., m$$

Scenario Knowledge Representation

$$ilde{ ext{S}} = f_{MLP}(ext{S}, \Lambda)$$



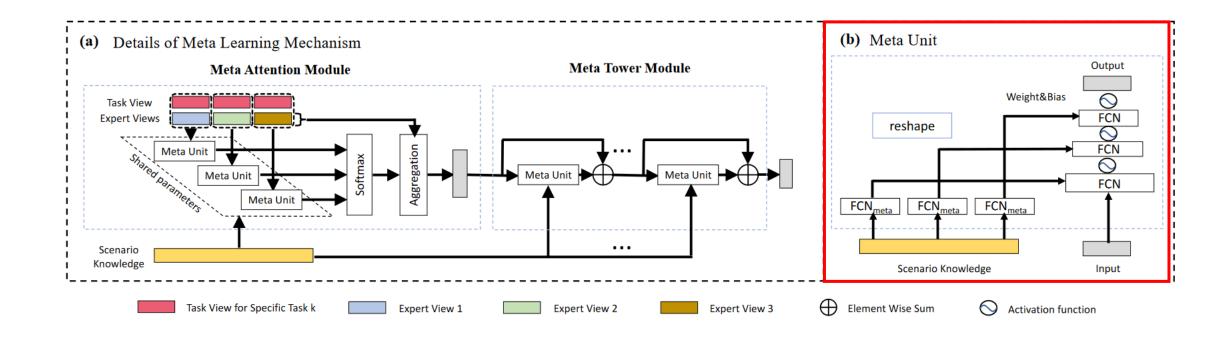






Meta Unit

$$\mathbf{h}_{output} = \mathbf{h}^K = Meta(\mathbf{h}_{input})$$





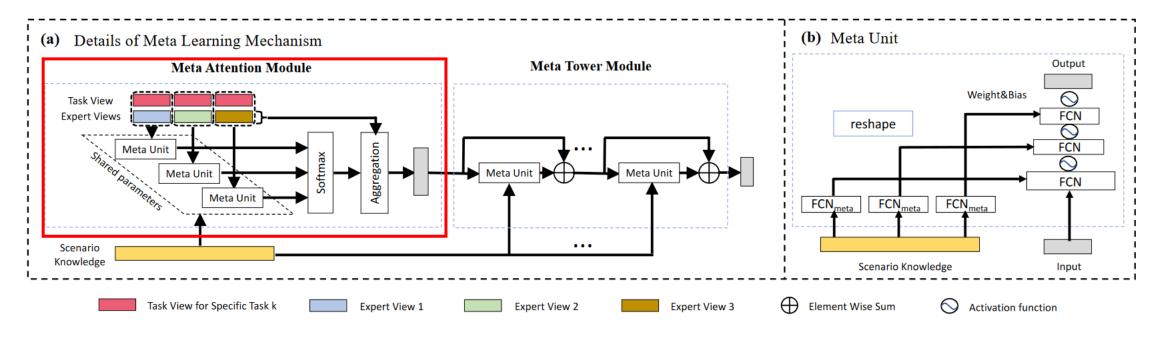




Meta Attention Module

$$a_{t_i} = \operatorname{v}^T Meta_t([\operatorname{E_i} \parallel \operatorname{T}_t])$$

$$lpha_{t_i} = rac{exp(a_{t_i})}{\sum_{j=1}^{M} exp(a_{t_j})}, \qquad \mathrm{R}_t = \sum_{i=1}^{k} lpha_{t_i} \mathrm{E}_i.$$











Meta Attention Module

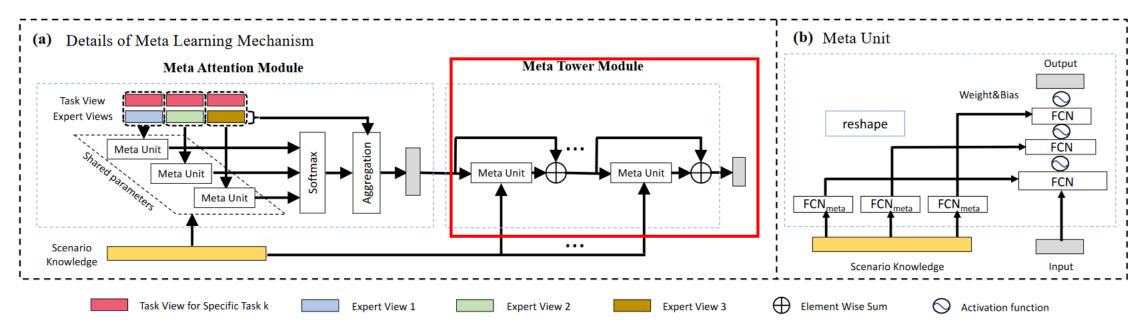
$$a_{t_i} = \operatorname{v}^T Meta_t([\operatorname{E_i} \parallel \operatorname{T}_t])$$

$$lpha_{t_i} = rac{exp(a_{t_i})}{\sum_{j=1}^{M} exp(a_{t_j})}, \qquad \mathrm{R}_t = \sum_{i=1}^{k} lpha_{t_i} \mathrm{E}_i.$$

Meta Tower Module

$$\mathrm{L}_t^{(0)}=\mathrm{R}_t$$

$$\mathbf{L}_{t}^{(j)} = \sigma(Meta^{(j-1)}(\mathbf{L}_{t}^{(j-1)}) + \mathbf{L}_{t}^{(j-1)}), orall j \in 1, 2, .., L$$



HiNet

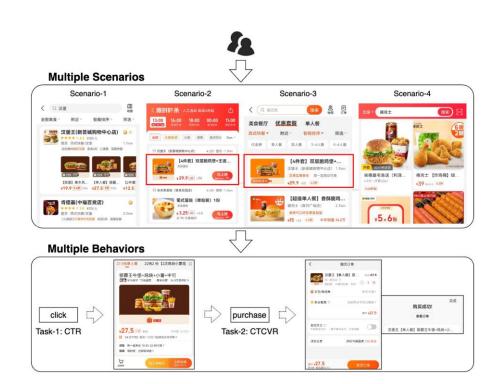








- Motivation
 - Multi-scenario and multi-task (CTR, CTCVR) optimization
- Methods
 - Proposing Hierarchical information extraction
 Network (HiNet) for multi-scenario & multi-task
 - Scenario Extraction Layer: Sharing information among scenarios and extracting scenario-specific characteristic
 - Task Extraction Layer: Resolving negative transfer problem in multi-task learning



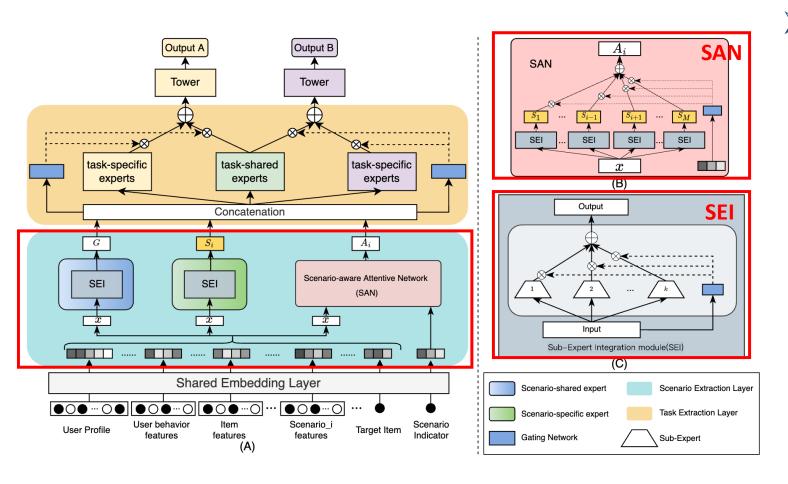
HiNet











- Scenario Extraction Layer
 - Scenario-shared expert network (SEI)
 - Scenario-specific expert network (SEI)
 - Scenario-aware attentive network (SAN)

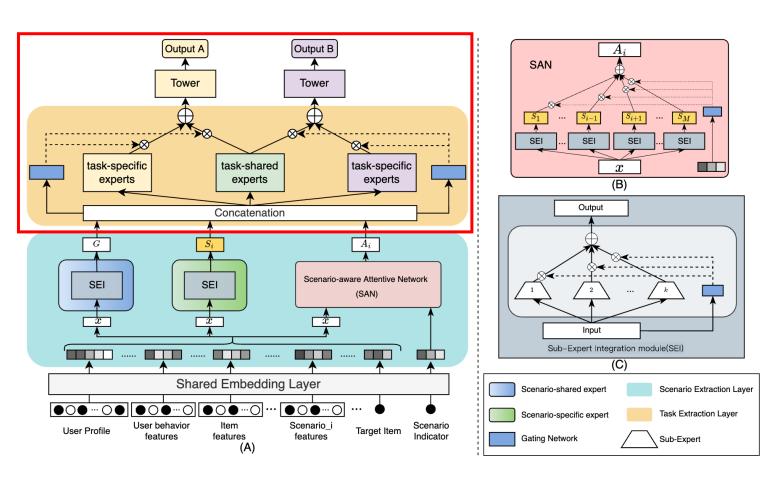
HiNet











- Scenario Extraction Layer
 - Scenario-shared expert network (SEI)
 - Scenario-specific expert network (SEI)
 - Scenario-aware attentive network (SAN)
- > Task Extraction Layer
 - Task-shared expert networks
 - Task-specific expert networks







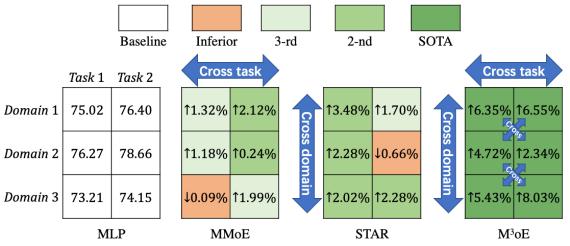


Motivation

- Multi-Domain Mult-Task (MDMT) seesaw problem
- The same multi-domain information transfer method may not generalize to different tasks
- The same multi-task optimization balancing strategy may not generalize to different domains.

Methods

- Domain representation extraction layer
- Multi-view expert learning layer
- MDMT objective prediction layer

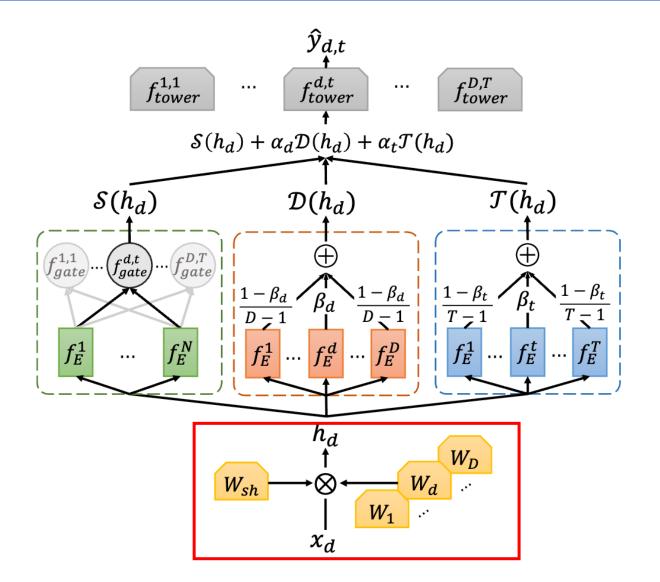












Domain Representation Extraction Layer

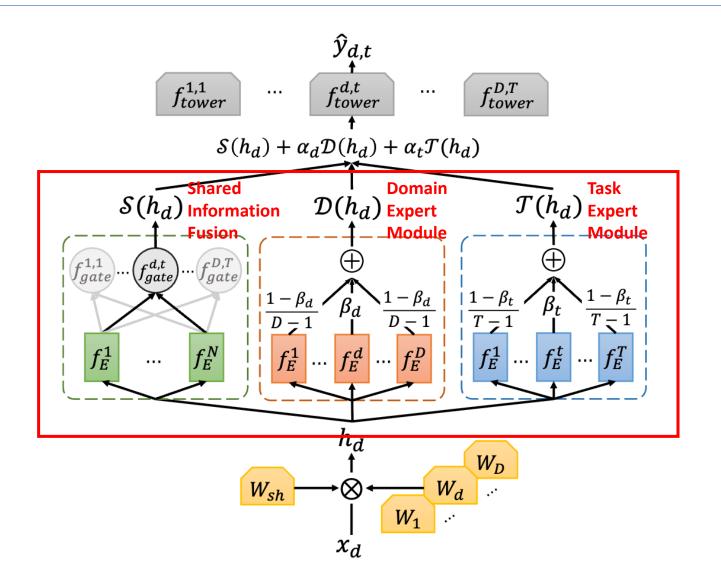
$$\begin{aligned} \widehat{\boldsymbol{W}}_d &= \boldsymbol{W}_d \otimes \boldsymbol{W}_{\mathrm{sh}} \\ f_{\mathrm{DR}}(\boldsymbol{x}_d) &= \widehat{\boldsymbol{W}}_d \boldsymbol{x}_d + \boldsymbol{b}_d + \boldsymbol{b}_{\mathrm{sh}} \\ \boldsymbol{h}_d &= \boldsymbol{W}_c f_{\mathrm{DR}}(\boldsymbol{x}_d) + \boldsymbol{b}_c + f_{\mathrm{DA}}(\boldsymbol{x}_d) \end{aligned}$$











Domain Representation Extraction Layer

$$\begin{split} \widehat{\boldsymbol{W}}_d &= \boldsymbol{W}_d \otimes \boldsymbol{W}_{\mathrm{sh}} \\ f_{\mathrm{DR}}(\boldsymbol{x}_d) &= \widehat{\boldsymbol{W}}_d \boldsymbol{x}_d + \boldsymbol{b}_d + \boldsymbol{b}_{\mathrm{sh}} \\ \boldsymbol{h}_d &= \boldsymbol{W}_c f_{\mathrm{DR}}(\boldsymbol{x}_d) + \boldsymbol{b}_c + f_{\mathrm{DA}}(\boldsymbol{x}_d) \end{split}$$

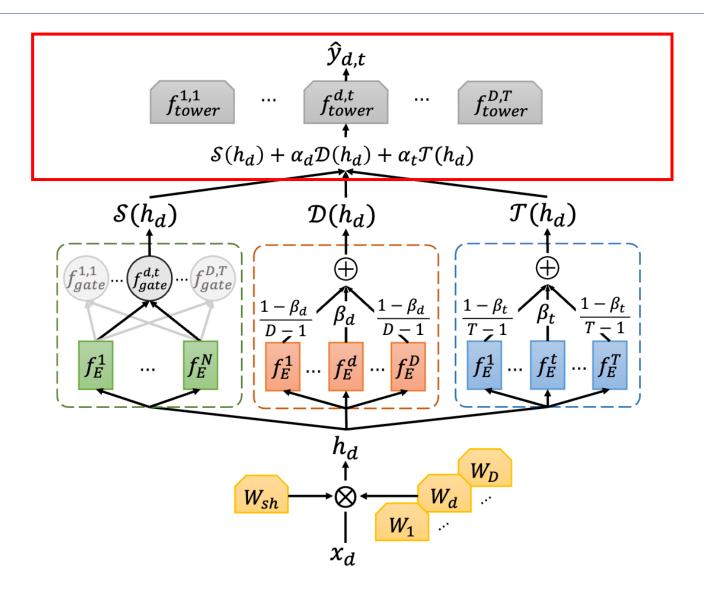
- ➤ Multi-View Expert Learning Layer
 - Shared Information Fusion
 - Domain Expert Module
 - Task Expert Module











Domain Representation Extraction Layer

$$\begin{split} \widehat{\boldsymbol{W}}_d &= \boldsymbol{W}_d \otimes \boldsymbol{W}_{\mathrm{sh}} \\ f_{\mathrm{DR}}(\boldsymbol{x}_d) &= \widehat{\boldsymbol{W}}_d \boldsymbol{x}_d + \boldsymbol{b}_d + \boldsymbol{b}_{\mathrm{sh}} \\ \boldsymbol{h}_d &= \boldsymbol{W}_c f_{\mathrm{DR}}(\boldsymbol{x}_d) + \boldsymbol{b}_c + f_{\mathrm{DA}}(\boldsymbol{x}_d) \end{split}$$

- Multi-View Expert Learning Layer
 - Shared Information Fusion
 - Domain Expert Module
 - Task Expert Module
- Multi-View Representation Balancing

$$\begin{split} \overline{\boldsymbol{h}}_{d} &= \mathcal{S}(\boldsymbol{h}_{d}) + \alpha_{d} \cdot \mathcal{T}(\boldsymbol{h}_{d}) + \alpha_{t} \cdot \mathcal{D}(\boldsymbol{h}_{d}) \\ f_{tower}^{d,t}(\overline{\boldsymbol{h}}_{d}) &= \boldsymbol{W}_{d,t}^{2} \text{ReLU}(\boldsymbol{W}_{d,t}^{1} \overline{\boldsymbol{h}}_{d} + \boldsymbol{b}_{d,t}^{1}) + \boldsymbol{b}_{d,t}^{2} \\ \widehat{\boldsymbol{y}}_{d,t} &= \text{Sigmoid}(f_{tower}^{d,t}(\overline{\boldsymbol{h}}_{d})) \end{split}$$

Benchmark: Scenario-Wise Rec





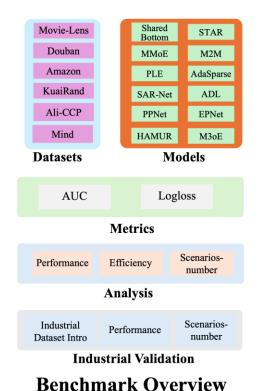






Scenario - Wise Rec

Benchmark for Multi-Scenario Recommendation



Dataset **Processing** Model Param **Training** Setting Data loading Epoch training **Evaluation** Early Stop Model Logs Results Saving Overall Scenario Result Results

Overall Pipeline





Paper Repo

GitHub Repo

- Highlight
 - A comprehensive benchmark exclusively for MSR
 - Complete data loading, training, and evaluation process
 - Providing 6 datasets and 12 MSR models
 - Compresentative analysis and step-by-step tutorial

Summary









Model	Setting	Methods
STAR	Multi-Scenario	Shared-Specific
SAR-Net	Multi-Scenario	Shared-Specific; Experts
ADI	Multi-Scenario	Shared-Specific
Uni-CTR	Multi-Scenario	Shared-Specific; LLMs
M-LoRA	Multi-Scenario	Shared-Specific; LoRAs
HAMUR	Multi-Scenario	Dynamic Weight
HierRec	Multi-Scenario	Dynamic Weight

Model	Setting	Methods
LLM4MSR	Multi-Scenario	Dynamic Weight; LLMs
PEPNet	Multi-Scenario & Multi-Task	Dynamic Weight
M2M	Multi-Scenario & Multi-Task	Dynamic Weight; Experts
HiNet	Multi-Scenario & Multi-Task	Experts
МЗоЕ	Multi-Scenario & Multi-Task	Experts
Scenario-Wise Rec	Multi-Scenario	Benchmark

Future Directions









Topic	Challenge & future direction
LLM-based multi-scenario & multi-task modeling	 Explore quantification or compression techniques for handling large-scale scenarios. More fine-grained modeling to bridge semantic gaps between LLM and MSR models.
Robustness	• Scenarios with different available information (multimodal)
Privacy	 Data need to be shared between different scenarios to build a unified model. Methods to protect user privacy should be proposed.
Fairness and Bias	The issue of fairness in recommendation scenarios.

Joint Modeling in Recommendations









Coffee Break



Huawei Noah's Ark Lab



WWW25 Huawei Noah's Ark Lab Chat Group



AML Lab CityU

Agenda













Joint Modeling in RS





Multi-task Recommendation



Yuhao Wang

Multi-scenario Recommendation





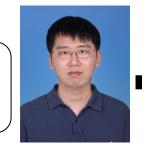
Pengyue Jia Xiaopeng Li

Multi-behavior Recommendation



Jingtong Gao

Multi-modal Recommendation



Qidong Liu



Future Work

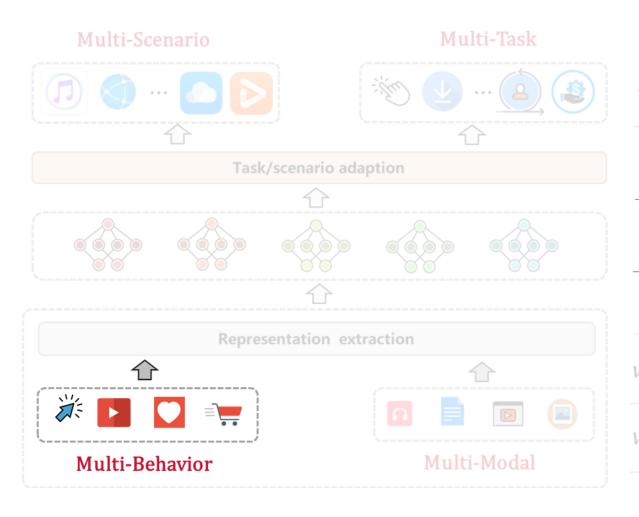


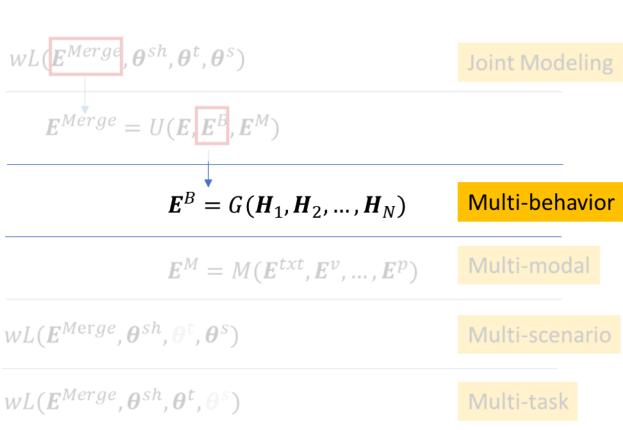
Multi-Behavior Modeling











Multi-Behavior Modeling

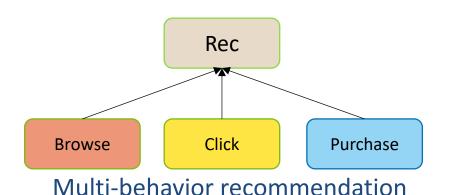








- >Understanding behavior patterns and behavior correlations at a fine-grained granularity
- Explicitly considering the different behavior types as they convey subtle differences in user interest modeling



Behavior Type Definition







- ➤ An open question
- ➤ Roughly three categories:
 - Macro behaviors: interaction with different items
 E.g. user 1 interact with item 1, then item 22, then item 81.
 - Micro behaviors: actions taken on this item E.g. click, add to cart,...
 - Behaviors from different domains or scenarios
 - E.g. Same behavior in two domains => different behaviors (highlight the distinctions)



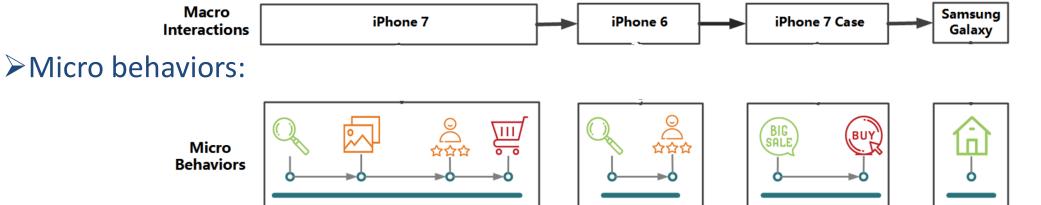
Behavior Type Definition





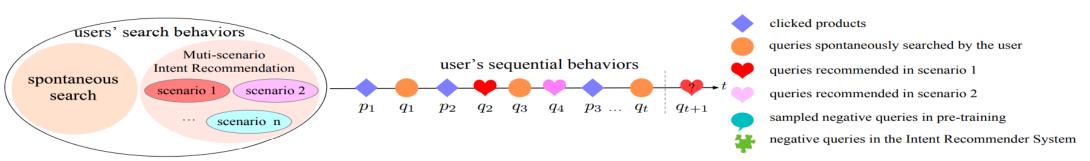






> Behaviors from different domains or scenarios

E.g. Same behavior in two domains => different behaviors (highlight the distinctions)



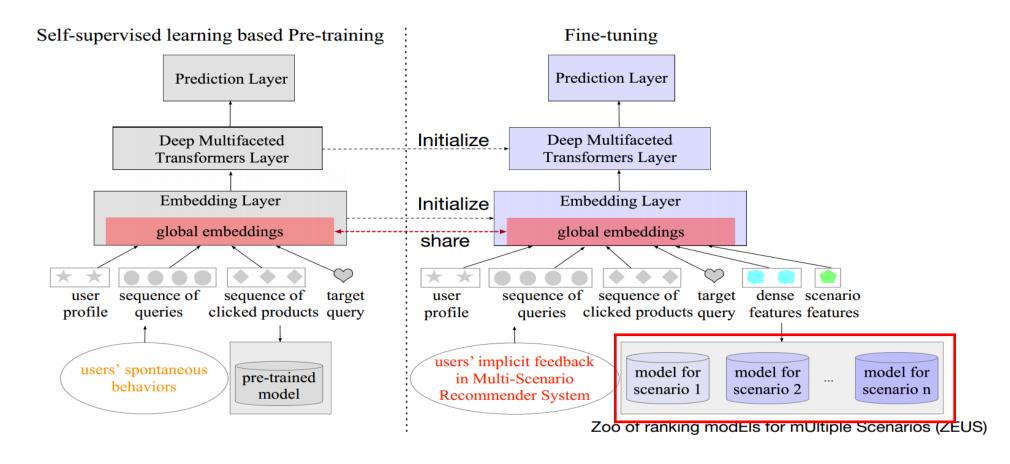
Multi-Behavior Fusion







➤ Modeling the complicated cross-scenario behavior dependencies



Example: pre-training and fine-tuning of ZEUS

Multi-Behavior Fusion

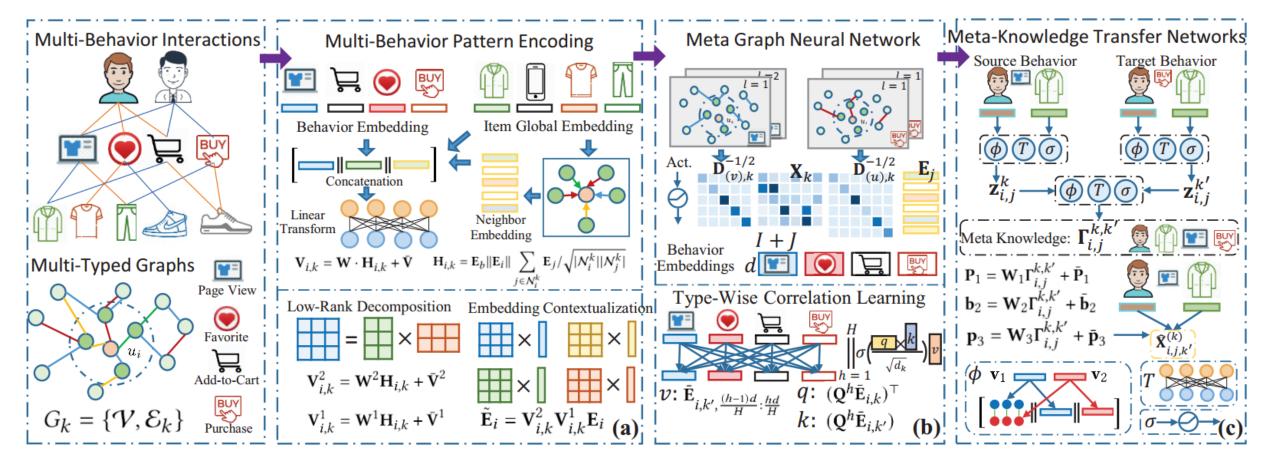








➤ Modeling the complicated cross-type behavior dependencies



Example: MB-GMN

Challenges









- ➤ Sequence modeling of heterogeneous behavioral feedback
 - How to model different behaviors and their feedbacks
- ➤ Modeling behavior relations
 - How to capture complicated behavior relations
- >Joint long-term and short-term preference modeling with heterogeneous behaviors of users
 - How to combining long-term preference and short-term preferences for better user modeling
- ➤ Avoiding noise and bias
 - How to solve the problem brought by noises and bias coming with different behaviors

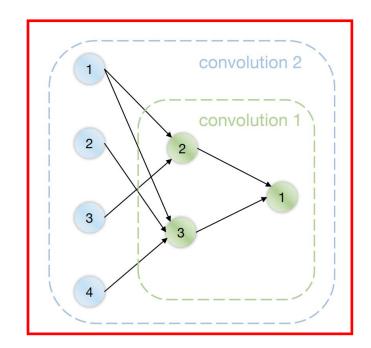


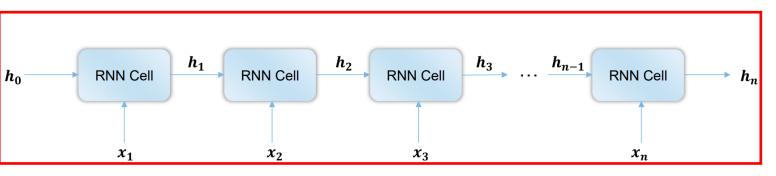


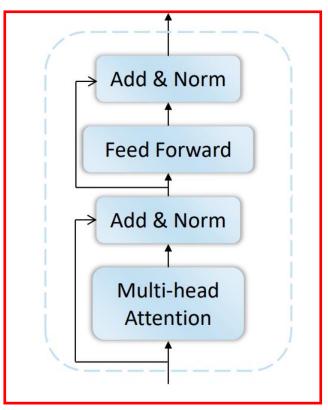




- ➤ RNN-based methods
- ➤ Graph-based methods
- >Transformer-based methods
- **≻**Other methods















- ➤ RNN-based methods
- ➤ Graph-based methods
- >Transformer-based methods
- **≻**Other methods

Works	Data Perspective	Model Perspective	Features
RLBL [14]	A sequence of (item,	Local	Capture the influence of heterogeneous behaviors by utilizing a be-
	behavior) pairs		havior transition matrix.
RIB [26]	A sequence of (item,	Local	Leverage GRU and attention mechanism simultaneously.
	behavior) pairs		
BINN [22]	A sequence of (item,	Local	Design the CLSTM and the Bi-CLSTM, where the behavior vector is
	behavior) pairs		as context in LSTM.
CBS [63]	Some behavior-specific	Local	Design of models with and without shared parameters for behaviors
	subsequences of items		simultaneously; towards the next-basket recommendation.
DIPN [64]	Some behavior-specific	Local	Leverage GRU and attention mechanism simultaneously; behaviors are
	subsequences of items		specific, including swipe, touch and browse interactive behavior.
HUP [27]	A sequence of (item,	Local	Design the Behavior-LSTM where adds behavior gate and time gate
	behavior) pairs		to the LSTM; leverage attention mechanism; take into account the
			category of the items.
IARS [28]	A sequence of (item,	Local	Propose Soft-MGRU (a multi-behavior gated recurrent unit) with
	behavior) pairs		sharing parameters between behaviors; leverage attention mechanism;
			take into account the category of the items.
DeepRec [62]		Local + Global	Utilizing multi-behavior sequence data to make privacy-preserving
	subsequences of items		recommendation.
MBN [65]	Some behavior-specific	Local	The overall Meta-RNN and the separate Behavior-RNN share the
	subsequences of items		learned potential representations by gathering and then scattering;
			towards the next-basket recommendation.

RNN-based Methods—RLBL





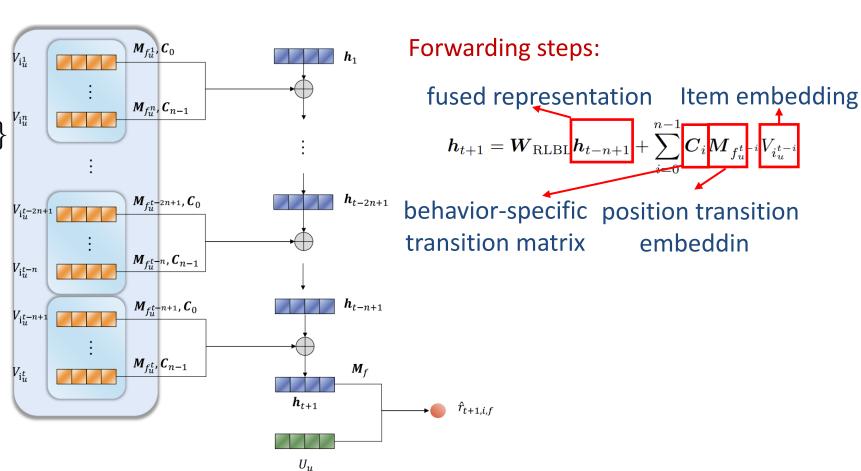




➤ Conducting item side behavior modeling via recurrent log-bilinear model

windowed representations
Item-behavior pairs

$$\{(i_u^{t-n+1}, \bar{f}_u^{t-n+1}), ..., (i_u^t, f_u^t)\}^{V_{i_u^n}}$$



RLBL

RNN-based Methods—MBN

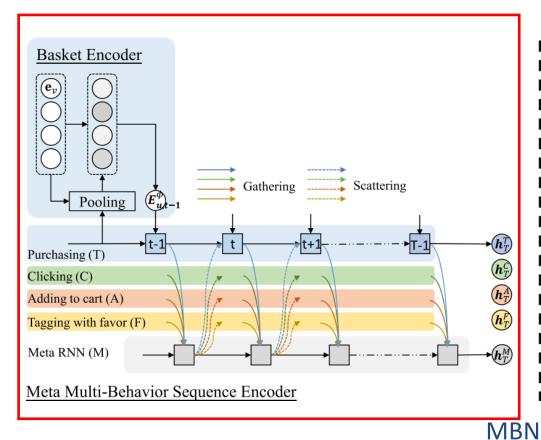


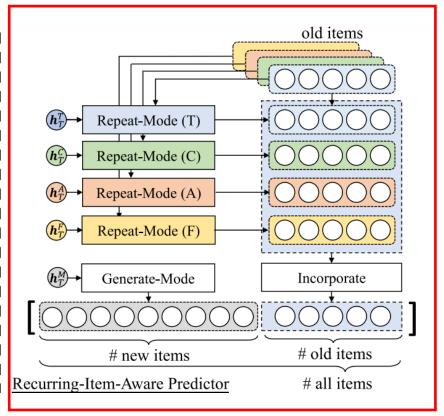






- ➤ Conducting Next basket recommendation with multi-behavior modeling
 - Encoder: Multiple Behavior-RNN and one Meta-RNN for behavior modeling and fusion
 - Predictor: Generating next items with both new items and old items













➤ RNN-based methods

- ➤ Graph-based methods
- >Transformer-based methods
- **≻**Other methods

			HOAVVEI
	Methods	Prons	Cons
S	RNN-based	 Suitable for sequence problems and can store short-term memories 	 Gradient disappearance & explosion problems Inefficient in predicting future sequences Rarely used currently
	Graph-based		
	Transformer- based		
	Others		









- >RNN-based methods
- ➤ Graph-based methods
- >Transformer-based methods
- **≻**Other methods

Works	Data Perspective	Model Perspective	Features
MGNN- SPred [24]	Some behavior-specific subsequences of items	Global	Modeling behavior from behavior transition relations, containing homogeneous behavior transitions intra each kind of behavior-specific subsequences.
DMBGN [71]	Some behavior-specific subsequences of items	Global	Focus on the task of voucher redemption rate prediction and model the relationship between multiple behaviors and vouchers effectively.
GPG4HSR [72]	A sequence of (item, behavior) pairs	Local + Global	Learn various behavior transition relations from the global graph and the personalized graph, respectively.
BGNN [73]	Some behavior-specific subsequences of items	Global	Construct directed graphs for different behavior transition (homogeneous and heterogeneous) information.
BA- GNN [<mark>74</mark>]	Some behavior-specific subsequence of items	Global	Construct directed graphs for different behavior-specific sequences respectively.



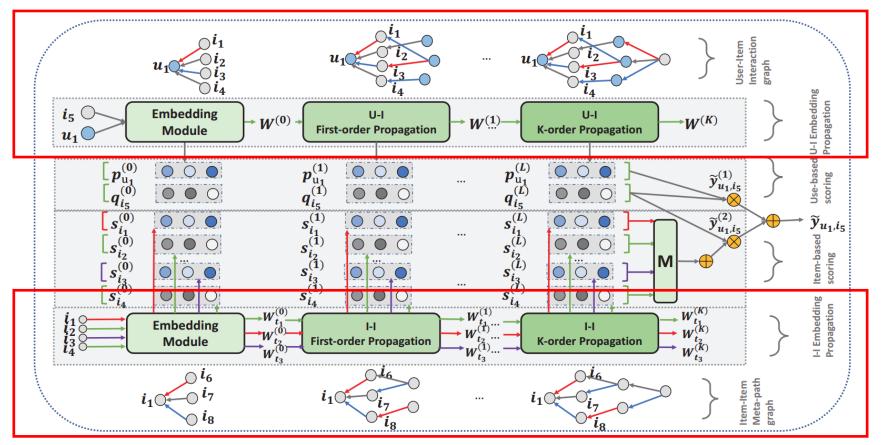






➤ Conducting behavior-aware user-item propagation and item-relevance aware item-item propagation in the user-item graph

U-I Embedding Propogation



I-I Embedding Propogation

MB-GCN









- Existing researches approach this task from two aspects
 - Utilizing multi-behavior data into the sampling process and builds multi-sampling pairs to reinforce the model learning process
 - Tryng to design model to capture multi-behavior information

➤ Their limitations

- The strength of multiple types of behaviors is not sufficiently utilized
- The semantics of multiple types of behaviors are not considered

➤Why?

• The limitations of existing methods lie in the fact that they cannot thoroughly address the above two challenges: modeling user-to-item based strength and item-to-item based semantics of multiple types of behaviors.



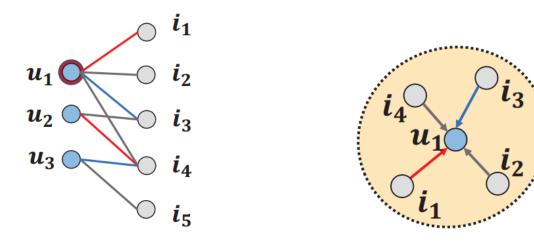






≻Solution

- Constructing a unified heterogeneous graph based on multiple types of behavioral data
- User/item represented as nodes and different types of behaviors represented as multiple types of edges of the graph



(a) U-I Interaction Graph

(b) Local Graph of u_1

An illustration of the user-item multi-behavior graph



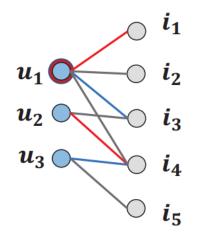


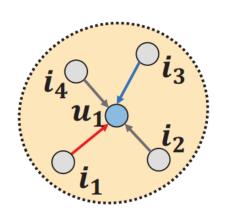




> Framework design

- Input: The user-item interaction data of T types of behaviors, $\{Y^1, Y^2, \dots, Y^T\}$
- Output: A recommendation model that estimates the probability that a user u will interact with an item i under the T-th behavior, i.e. target behavior.





(a) U-I Interaction Graph

(b) Local Graph of u_1

An illustration of the user-item multi-behavior graph

User Embedding Propagation

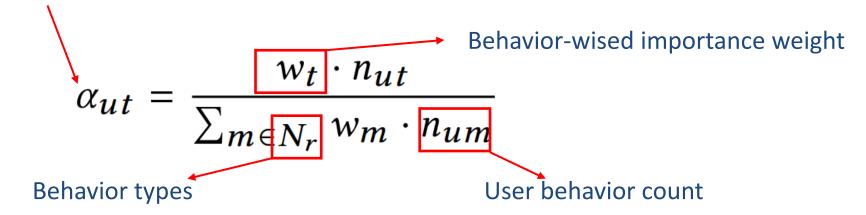








- ➤ User behavior propagation weight calculation
 - Different behavior contributes differently to the target behavior
 - The importance of each behavior to the target behavior cannot be measured artificially and should be learned by the model itself
 - A frequency-based propagation weight



User Embedding Propagation









- ➤ Neighbour item aggregation based on behavior
 - Items that are interacted under the same behavior reflect user's similar preference strength
 - Items that have the same behavior interaction with user are aggregated together so as to obtain one embedding for each behavior
 - Aggregation function: mean function, mean function with sampling, max pooling, etc.

$$p_{u,t}^{(l)} = \operatorname{aggregate}(q_i^{(l)}|i \in N_t^I(u))$$
 Item embedding Items the user interacted under behavior t

User Embedding Propagation

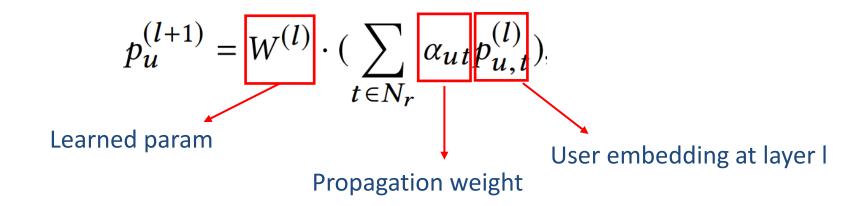








- ➤ Behavior-level Item Propagation for User
 - Summing neighbor item aggregation embedding together according to weight
 - Going through an encoder matrix to obtain the final neighbor item aggregation for users
 - A graph neural network to refine information based on multi-behavior



Item Embedding Propagation





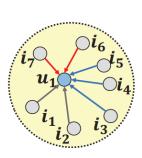


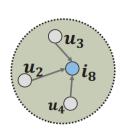


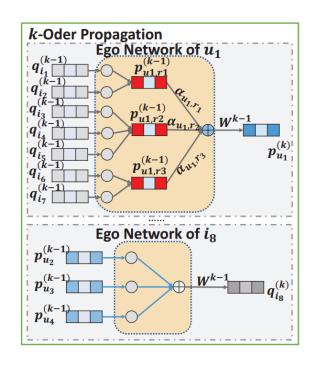
- The feature of the item is static
 - No importance needed. Assuming that different user has the same contribution to item

$$q_i^{(l+1)} = W^{(l)} \cdot \operatorname{aggregate}(p_j^{(l)}|j \in N^U(i)),$$

- ➤ Summary: Behavior-aware User-Item Propagation
 - User Embedding Propagation
 - Item Embedding Propagation







(a) Local Graph

(b) k-order Propagation Process

Item-Relevance Aware I-I Propagation



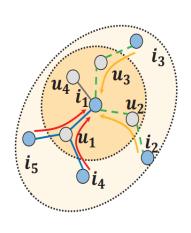


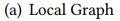


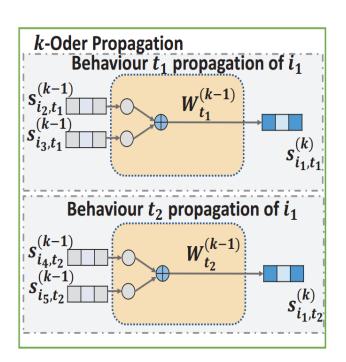


➤ Item information extracting

$$s_{it}^{(l+1)} = W_t^{(l)} \cdot \text{aggregate}(s_{jt}^{(l)}|j \in N_t^I(i))$$
 $s_{it}^{(0)} = q_i^{(0)}$







(b) k-order Propagation Process

Joint Prediction









≻Embedding aggregation

$$\begin{aligned} p_u^* &= p_u^{(0)}||...||p_u^{(L)}, \\ q_i^* &= q_i^{(0)}||...||q_i^{(L)}, \\ s_{it}^* &= s_{it}^{(0)}||...||s_{it}^{(L)}, t \in N_r \end{aligned}$$

➤ User-based CF scoring

$$y_1(u,i) = p_u^{*T} \cdot q_i^*$$

➤ Item-based CF scoring

$$y_2(u, i) = \sum_{t \in N_r} \sum_{j \in N_t^I(u)} \frac{s_{jt}^{*T} \cdot M_t \cdot s_{it}^*}{|N_t^I(u)|}$$

≻Combined Scoring

$$y(u, i) = \lambda \cdot y_1(u, i) + (1 - \lambda) \cdot y_2(u, i).$$

Overall performance







	Method	Recall@10	NDCG@10	Recall@20	NDCG@20	Recall@40	NDCG@40	Recall@80	NDCG@80
	MF-BPR	0.02331	0.01306	0.03161	0.01521	0.04239	0.01744	0.05977	0.02049
One-behavior	NCF	0.02507	0.01472	0.03319	0.01683	0.04502	0.01931	0.06352	0.02252
One-benavior	GraphSAGE-OB	0.01993	0.01157	0.02521	0.01296	0.03368	0.01474	0.04617	0.01693
	NGCF-OB	0.02608	0.01549	0.03409	0.01757	0.04612	0.02010	0.06415	0.02324
	MCBPR	0.02299	0.01344	0.03178	0.01558	0.04360	0.01813	0.06190	0.02132
	NMTR	0.02732	0.01445	0.04130	0.01831	0.06391	0.02279	0.09920	0.02891
Multi-behavior	GraphSAGE-MB	0.02094	0.01223	0.02805	0.01406	0.03804	0.01616	0.05351	0.01887
Muiti-benavior	NGCF-MB	0.03076	0.01754	0.04196	0.02042	0.05857	0.02389	0.08408	0.02833
	RGCN	0.01814	0.00955	0.02627	0.01165	0.03877	0.01426	0.05749	0.01750
	MBGCN	0.04006	0.02088	0.05797	0.02548	0.08348	0.03079	0.12091	0.03730
	Improvement	30.23%	19.04%	37.04%	24.78%	24.91%	28.88%	8.90%	26.40%

Comparison on Tmall

	Method	Recall@10	NDCG@10	Recall@20	NDCG@20	Recall@40	NDCG@40	Recall@80	NDCG@80
	MF-BPR	0.03873	0.02286	0.05517	0.02676	0.08984	0.03388	0.14137	0.04258
One-behavior	NCF	0.04209	0.02394	0.05609	0.02579	0.09118	0.03410	0.15426	0.04022
One-benavior	GraphSAGE-OB	0.034536	0.01728	0.06907	0.02594	0.11567	0.03547	0.18626	0.04747
	NGCF-OB	0.04112	0.02199	0.06336	0.02755	0.11051	0.03712	0.19524	0.05153
	MCBPR	0.03914	0.02264	0.04950	0.02525	0.09592	0.03467	0.15422	0.04462
	NMTR	0.03628	0.01901	0.06239	0.02559	0.10683	0.03461	0.18907	0.04855
Multi-behavior	GraphSAGE-MB	0.04204	0.02267	0.05862	0.02679	0.09707	0.03451	0.18272	0.04911
Multi-benavior	NGCF-MB	0.04241	0.02415	0.06152	0.02893	0.10370	0.03741	0.01771	0.04987
	RGCN	0.04204	0.02051	0.06354	0.02591	0.09859	0.03309	0.16121	0.04363
	MBGCN	0.04825	0.02446	0.07354	0.03077	0.11926	0.04005	0.20201	0.05409
	Improvement	13.77%	1.28%	11.76%	3.85%	7.68%	3.30%	6.58%	3.84%

Comparison on Beibei

Graph-based Methods—CML

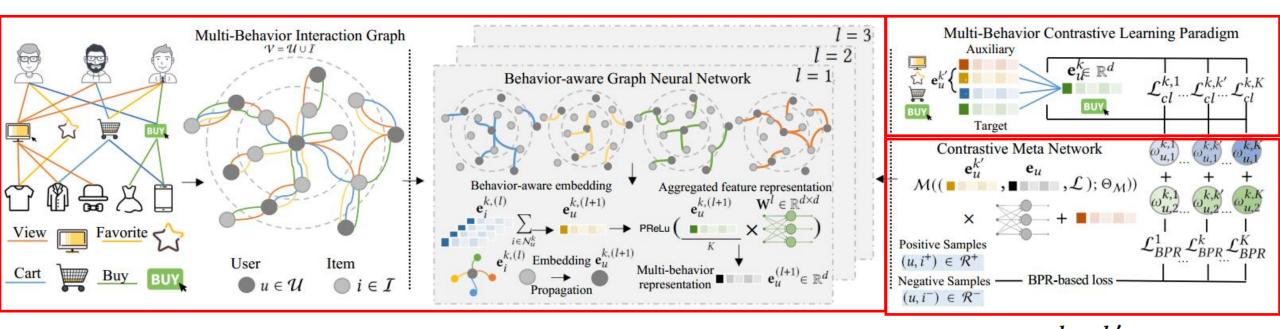








➤ Conducting constrastive learning among behaviors



$$\mathbf{e}_{u}^{k,(l+1)} = \sum_{i \in \mathcal{N}_{u}^{k}} \mathbf{e}_{i}^{k,(l)}; \ \mathbf{e}_{i}^{k,(l+1)} = \sum_{u \in \mathcal{N}_{i}^{k}} \mathbf{e}_{u}^{k,(l)}$$

$$\mathcal{L}_{cl}^{k,k'} = \sum_{u \in \mathcal{U}} -log \frac{\exp(\varphi(\mathbf{e}_{u}^{k}, \mathbf{e}_{u}^{k'})/\tau)}{\sum_{u' \in \mathcal{U}} \exp(\varphi(\mathbf{e}_{u}^{k}, \mathbf{e}_{u'}^{k'})/\tau)}$$









- ➤ RNN-based methods
- ➤ Graph-based methods
- >Transformer-based methods
- **≻**Other methods

			HOAWEI
	Methods	Prons	Cons
S	RNN-based	 Suitable for sequence problems and can store short-term memories 	 Gradient disappearance & explosion problems Inefficient in predicting future sequences Rarely used currently
	Graph-based	Detailed modeling for behavior relationsImproved performance	 Suffering from low efficiency
	Transformer- based		
	Others		









- ➤ RNN-based methods
- ➤ Graph-based methods
- >Transformer-based methods
- **≻**Other methods

-				
_	Works	Data Perspective	Model Perspective	Features
_	DMT [23]	Some behavior-specific subse-	Local + Global	Use target item as query; Consider implicit feedback bias by
-	_	quences of items		a bias deep neural network.
	DFN [84]	Some behavior-specific subse-	Local + Global	Use target item as query; Consider implicit negative feedback
		quences of items		noise by an attention network.
	DUMN [88]	Some behavior-specific subse-	Local	Consider implicit feedback noise; Use memory network to
		quences of items		obtain the long-term user preference.
	FeedRec [25]	Some behavior-specific subse-	Local + Global	Consider implicit feedback noise by an attention network;
		quences of items and a sequence		Consider multiple patterns of the multi-behavior sequences.
		of (item, behavior) pairs		
	NextIP [86]	Some behavior-specific subse-	Local + Global	Treat the problem as the item prediction task and the pur-
		quences of items and a sequence		chase prediction task; Consider multiple patterns of the multi-
_		of (item, behavior) pairs		behavior sequences.
	MB-	A sequence of (item, behavior)	Local	A novel positional encoding function to model multi-behavior
_	STR [87]	pairs		sequence relationships.
	FLAG [85]	A behavior-agnostic sequence of	Local + Global	Model user's local preference, local intention and global pref-
		items and a sequence of behaviors		erence simultaneously.

Transformer-based Methods—MB-STR

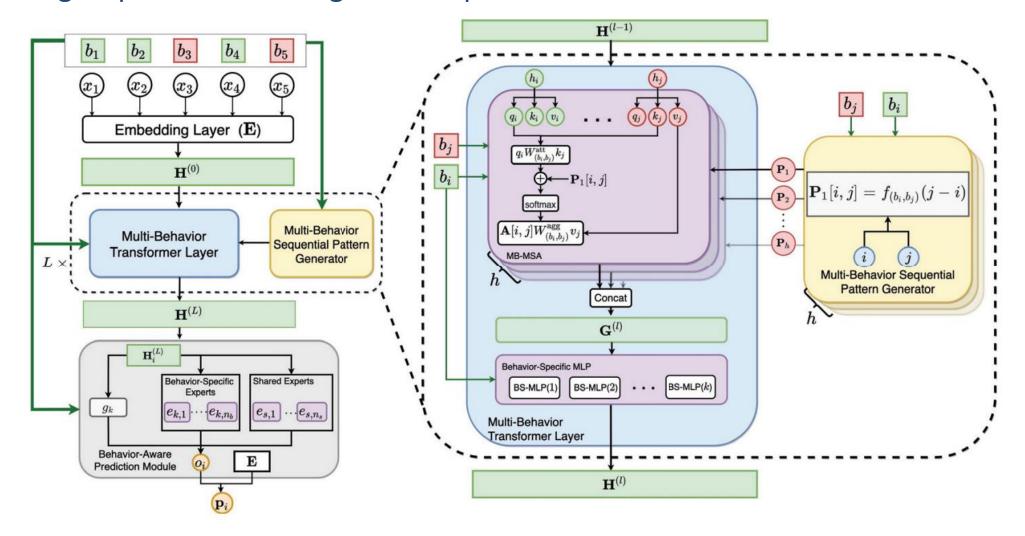








➤ Conducting sequential modeling for multiple behaviors



Transformer-based Methods—MB-STR









- ➤ Challenges in modeling of multi-behavior sequential recommendations
 - How to model heterogeneous multi-behavior dependencies at the fine-grained item-level
 - How to model diverse multi-behavior sequential patterns effectively
 - How to effectively mine users' multi-behavior sequence with multi-behavior supervision signals

	Multi-Behavior Modeling	Sequential Information	Behavior-Specific Prediction
MATN [39]	behavior-level	X	×
NMTR [9]	behavior-level	X	\checkmark
MBGCN [19]	behavior-level	X	×
MB-GMN [40]	behavior-level	×	\checkmark
DIPN [13]	behavior-level	fixed single behavior	\checkmark
DMT [11]	behavior-level	fixed single behavior	✓
MB-STR(our)	heterogeneous item-level	diverse multi-behavior	✓

Multi-Behavior Transformer Layer





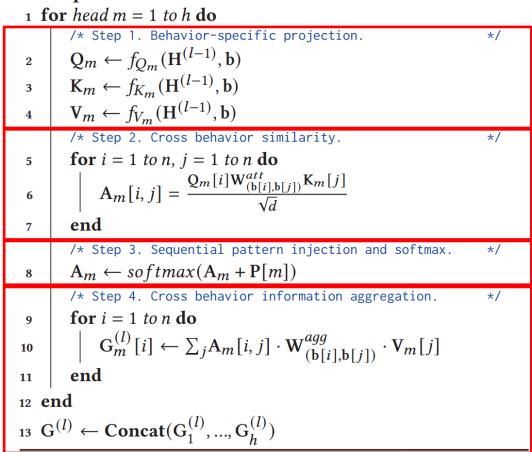




- ➤ Behavior-specific projection
 - n: number of behaviors in a sequence
- ➤ Cross behavior similarity
- ➤ Sequential pattern injection and softmax
 - P: multi-behavior sequential pattern matrix from MB-SPG
- ➤ Cross behavior information aggregation

Algorithm 1: Multi-Behavior Multi-head Self-Attention

Input: $\mathbf{H}^{(l-1)} \in \mathbb{R}^{n \times d}$, $\mathbf{b} \in \mathcal{B}^n$, $\mathbf{P}^{(l)} \in \mathbb{R}^{h \times n \times n}$ Output: $\mathbf{G}^{(l)} \in \mathbb{R}^{n \times d}$



Multi-Behavior Sequential Pattern Modeling





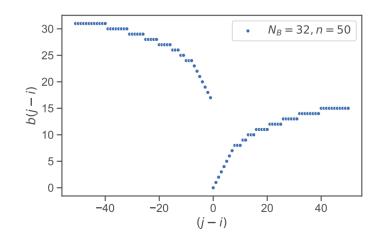




- >A designed position encoding function for balanced position pairs
 - k: heads
 - i,j: behavior position in a sequence

$$\mathbf{P}[k,i,j] = f_{(\mathbf{b}[i],\mathbf{b}[j])}(j-i)$$

$$b(j-i) = \begin{cases} B(j-i) & \text{if } (j-i) \ge 0 \\ B(-(j-i)) + \frac{N_B}{2} & \text{if } (j-i) < 0 \end{cases}$$



Behavior-Aware Masked Item Prediction









➤ Gating & expert

$$o_{i} = g_{k}(\mathbf{H}^{(L)}[i])^{\top} E_{k}(\mathbf{H}^{(L)}[i])$$

$$g_{k}(x) = softmax(\mathbf{W}_{g}^{k}x) \quad k = \mathbf{b}[i]$$

$$E_{k}(x) = \left[e_{k,1}(x), e_{k,2}(x), ..., e_{k,n_{b}}(x), e_{s,1}(x), e_{s,2}(x), ..., e_{s,n_{s}}(x)\right]$$

$$\mathbf{p}_{i}(v) = softmax(o_{i} \cdot \mathbf{E}^{\top})$$

Overall performance









O/S: One/multiple behavior

NS/S: Non-sequential/sequential

		lataget .	l V	ılı,	Taobao		IJCAI	
Dataset			16	elp	1a0	Dao	l ijC.	A1
Metrics		HR	NDCG	HR	NDCG	HR	NDCG	
		MF	0.755	0.481	0.262	0.153	0.285	0.185
	NS	DMF	0.756	0.485	0.305	0.189	0.392	0.250
O	1100	NGCF	0.789	0.500	0.302	0.185	0.461	0.292
O		LightGCN	0.810	0.513	0.373	0.235	0.443	0.283
	${s}$	SASRec	0.796	0.504	0.372	0.221	0.597	0.406
	8	BERT4Rec	0.816	0.531	0.385	0.234	0.605	0.431
	NS	NGCF _M	0.793	0.492	0.374	0.221	0.481	0.307
		$LightGCN_M$	0.872	0.585	0.391	0.243	0.486	0.317
		NMTR	0.790	0.478	0.332	0.179	0.481	0.304
		MATN	0.826	0.530	0.354	0.209	0.489	0.309
M		MBGCN	0.796	0.502	0.369	0.222	0.463	0.277
1V1		MB-GMN	0.87	0.582	0.491	0.300	0.532	0.345
		DIPN	0.791	0.500	0.317	0.178	0.475	0.296
	S	$SASRec_{M}$	0.819	0.531	0.637	0.442	0.795	0.611
	٥	$BERT4Rec_{M}$	0.838	0.558	0.675	0.476	0.816	0.632
		DMT	0.652	0.515	0.666	0.415	0.682	0.513
Our MB-STR		0.882*	0.624*	0.768*	0.608*	0.879*	0.713*	
	Rela	, Improv.	1.15%	6.67%	13.78%	27.73%	7.72%	12.82%

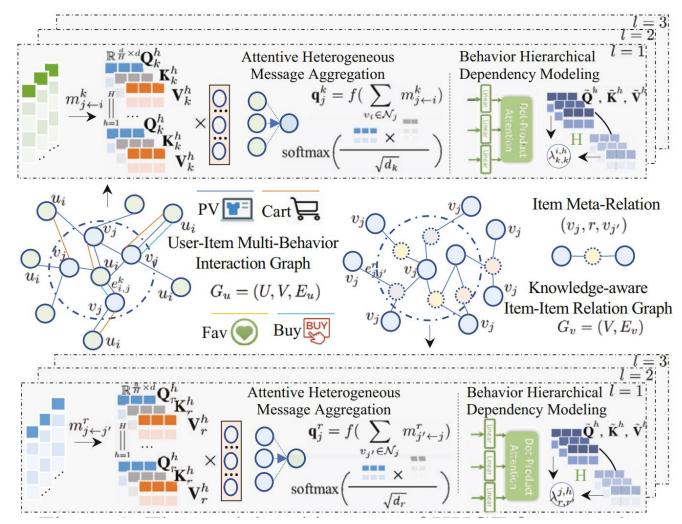
Transformer-based Methods—KHGT



















- >RNN-based methods
- ➤ Graph-based methods
- >Transformer-based methods
- **≻**Other methods

			HUAWEI
	Methods	Prons	Cons
S	RNN-based	 Suitable for sequence problems and can store short-term memories 	 Gradient disappearance & explosion problems Inefficient in predicting future sequences Rarely used currently
	Graph-based	Detailed modeling for behavior relationsImproved performance	 Suffering from low efficiency
	Transformer- based	 Exceptional performance from attention mechanism Superior parallel computing capabilities Enhanced ability to capture long-term dependencies Stronger explanability 	
	Others		









- >RNN-based methods
- ➤ Graph-based methods
- >Transformer-based methods
- **≻**Other methods

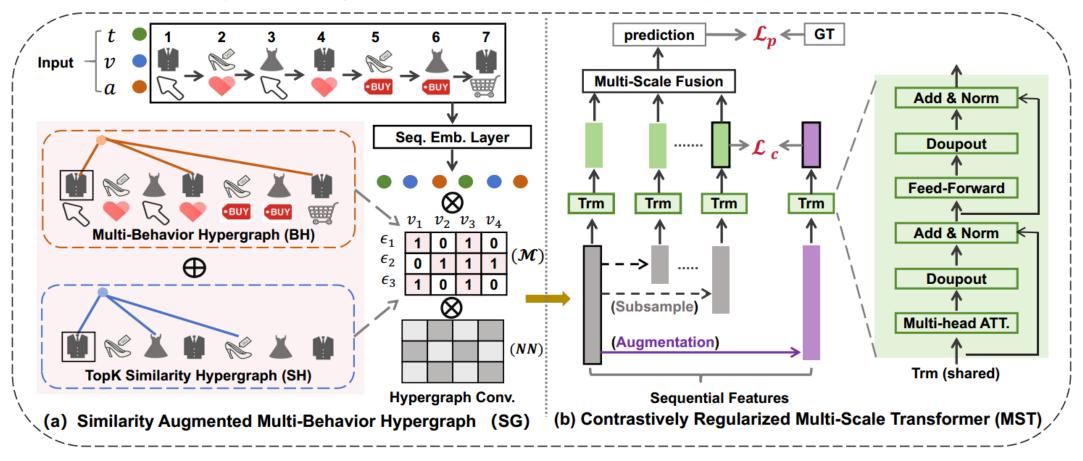
Other Methods—SG-MST







➤ Integrating a Similarity Augmented Multi-Behavior Hypergraph that captures complex behavior-aware dependencies among items and strengthens connections through item context similarities, producing more informative latent representations



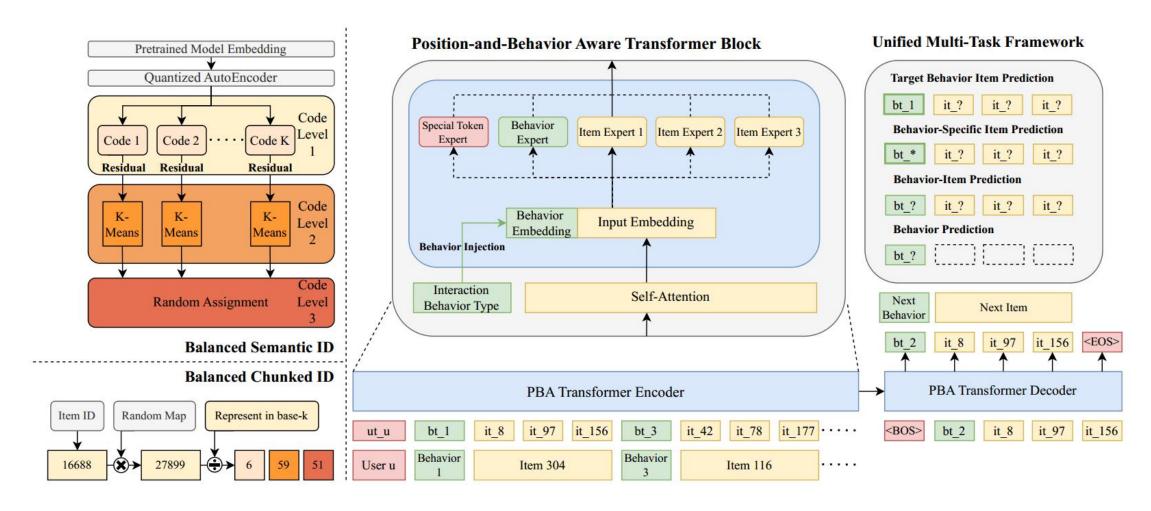
Other Methods—MBGen







> Generative modeling for multi-behavior sequential recommendations



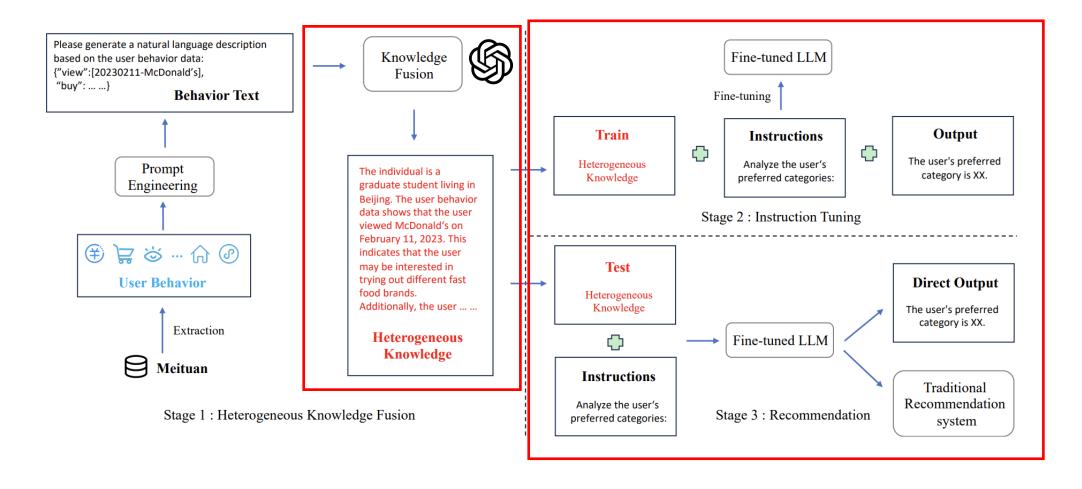
Other Methods—HKFR







Representation modeling for behaviors via LLM



Other Methods—HKFR







Experiments

- HKFRno-IT: Without instruct tuning
- HKFRno-HKF: Without heterogeneous knowledge fusion

Methods	Category				POIs			
	HR@5	NDCG@5	HR@10	NDCG@10	HR@5	NDCG@5	HR@10	NDCG@10
Caser	0.1152	0.1063	0.2147	0.1320	0.0897	0.0770	0.1842	0.1012
BERT4Rec	0.1217	0.1140	0.2196	0.1440	0.0875	0.0744	0.1811	0.0995
P5	0.1416	0.1384	0.2477	0.1589	0.1218	0.1159	0.2187	0.1260
ChatGLM-6B	0.1074	0.1019	0.2038	0.1254	0.0785	0.0720	0.1702	0.0872
$\overline{HKFR_{no-IT}}$	0.1241	0.1175	0.2267	0.1415	0.1014	0.0952	0.2050	0.1165
$HKFR_{no-HKF}$	0.1813	0.1308	0.2825	0.1580	0.1421	0.0975	0.2432	0.1270
HKFR	0.2160	0.1586	0.3007	0.1840	0.1726	0.1243	0.2610	0.1525

Taxonomy









- >RNN-based methods
- ➤ Graph-based methods
- >Transformer-based methods
- **≻**Other methods

			HUAWEI
	Methods	Prons	Cons
S	RNN-based	 Suitable for sequence problems and can store short-term memories 	 Gradient disappearance & explosion problems Inefficient in predicting future sequences Rarely used currently
	Graph-based	Detailed modeling for behavior relationsImproved performance	 Suffering from low efficiency
	Transformer- based	 Exceptional performance from attention mechanism Superior parallel computing capabilities Enhanced ability to capture long-term dependencies Stronger explanability 	
	Others	Better modelGenerative mModeling with	

Conclusion









Methods	Prons	Cons			
RNN-based	 Suitable for sequence problems and can store short-term memories 	 Gradient disappearance & explosion problems Inefficient in predicting future sequences Rarely used currently 			
Graph-based	Detailed modeling for behavior relationsImproved performance	 Suffering from low efficiency 			
Transformer- based	 Exceptional performance from attention mechanism Superior parallel computing capabilities Enhanced ability to capture long-term dependencies Stronger explanability 				
Others	Better modelGenerative mModeling wit	nodeling			

Model	Methods
RLBL	RNN-based
MBN	RNN-based
MB-GCN	Graph-based
CML	Graph-based
MB-STR	Transformer-based
KHGT	Transformer-based
SG-MST	Other
MBGen	Other
HKFR	Other

Future Directions









- **▶** Deeper information fusion
 - Better representation modeling
- **➤** More efficient learning method
 - Better modeling for better behavior modeling and lighter computation burden
- **➤ More explainable user representations**
 - Improving explanability
- > Fine-grained modeling with LLM
 - Conducting fine-grained modeling with LLM

Agenda









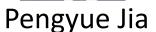




Yejing Wang

Joint Modeling in RS









Yuhao Wang







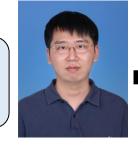
Pengyue Jia Xiaopeng Li

Multi-behavior Recommendation



Jingtong Gao

Multi-modal Recommendation



Qidong Liu



Future Work



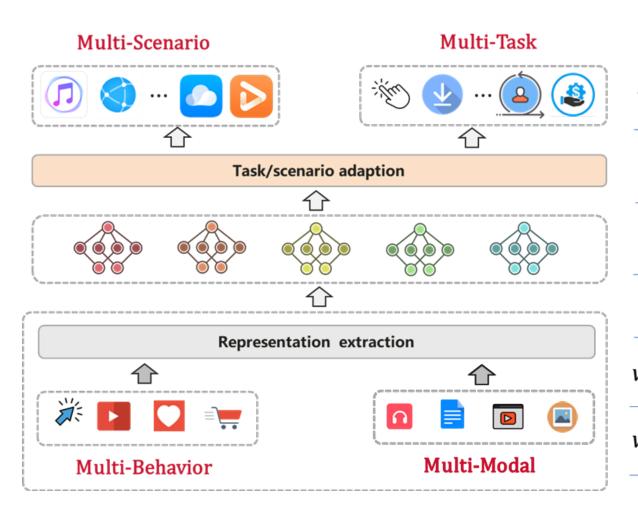
Yichao Wang

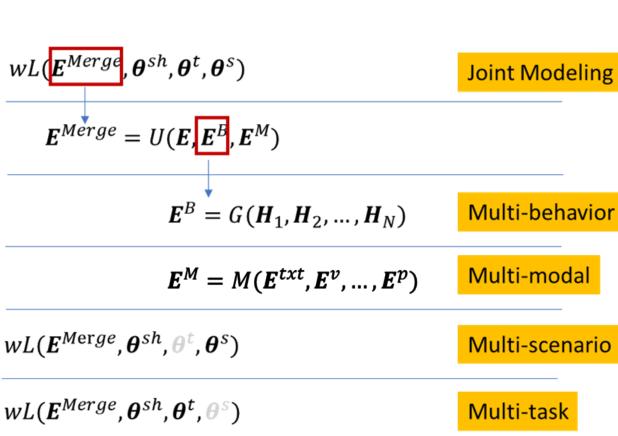
Multi-Modal Modeling











Formulation









- Learning a unified multimodal representations for users and items by various raw multimodal features $(x^{txt}, x^v, ..., x^p)$
- > Optimization problem:

$$E^{M} = M\left(E^{txt}, E^{v}, ..., E^{p}\right) = M\left(\mathcal{E}_{txt}(x^{txt}), \mathcal{E}_{v}(x^{v}), ..., \mathcal{E}_{p}(x^{p})\right)$$

- $\mathcal{E}_*(\cdot)$: the corresponding modality encoder
- $M(\cdot)$: feature interaction function, enabling combination of various modalities

Multimodal Recommender Systems (MRS)

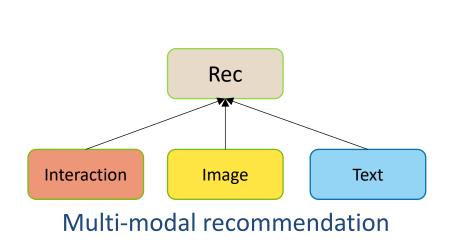


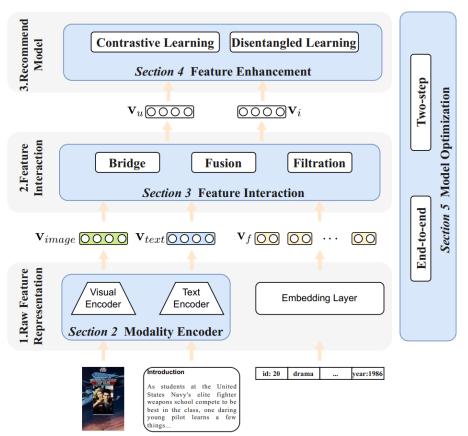






- >Using various types of information generated by multimedia applications and services to enhance recommender systems' performance
- Making use of multimodal features simultaneously, such as image, audio, and text





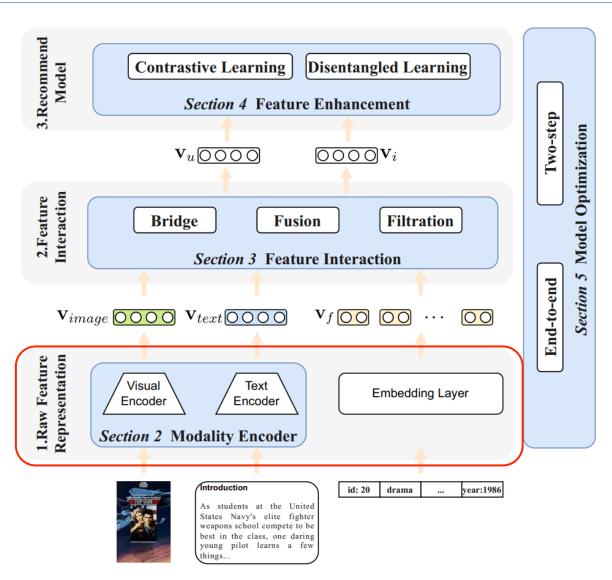








- ➤ Modality Encoder
 - Challenge: how to extract representations from complex raw features
 - Specialty: various encoders



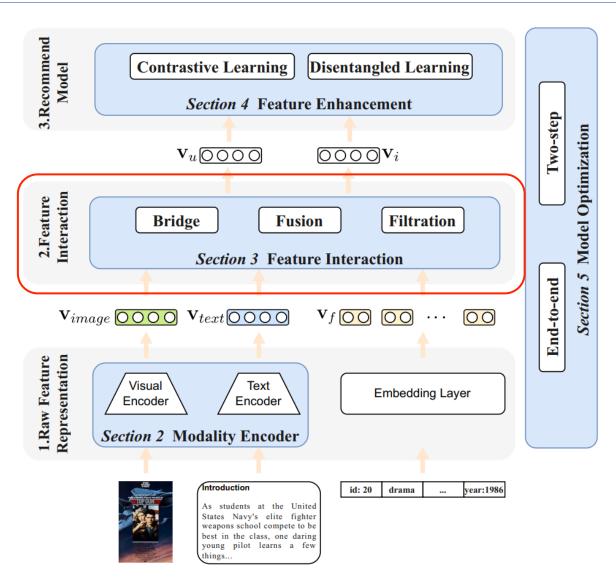








- ➤ Modality Encoder
- > Feature Interaction
 - Challenge: how to fuse the modality features in different semantic spaces
 - Specialty: modality alignment and fusion



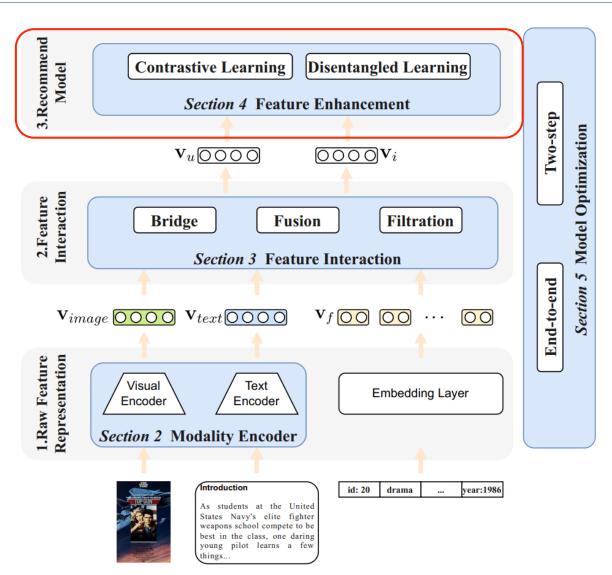








- ➤ Modality Encoder
- > Feature Interaction
- > Feature Enhancement
 - Challenge: how to get comprehensive representations for recommendation models under the data-sparse condition
 - Specialty: multimodal enhancement



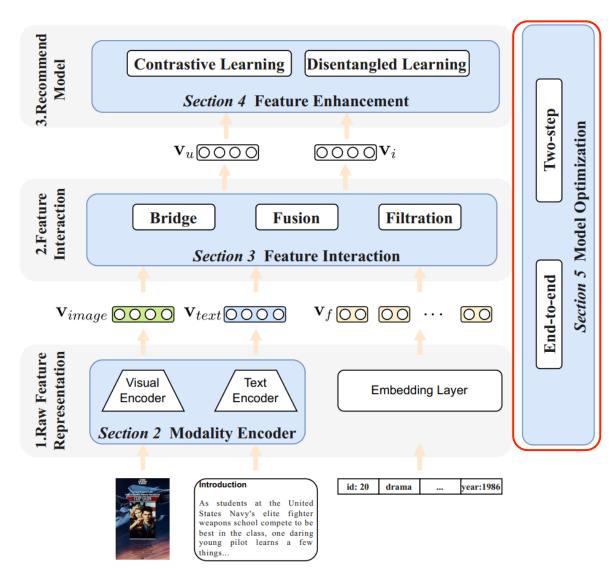








- ➤ Modality Encoder
- > Feature Interaction
- > Feature Enhancement
- ➤ Model Optimization
 - Challenge: how to optimize the lightweight recommendation models and parameterized modality encoder
 - Specialty: parameterized modality encoder



Modality Encoder



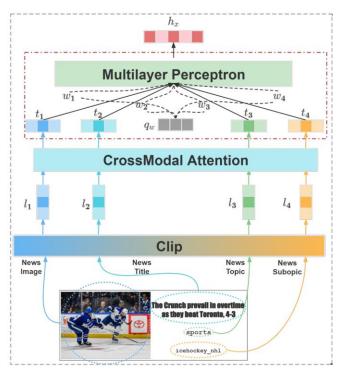






- ➤ Target: encoding different multimodal features
- ➤ Taxonomy:
 - Visual: CNN-based, ViT / Transformer-based
 - Textual: Word2Vec, CNN-based, RNN-based, Transformer-based
 - Others: E.g., converting acoustic and video data into text or visual information

Modality	Category
Visual Encoder	CNN ResNet Transformer
Textual Encoder	Word2vec RNN CNN Sentence-transformer Bert
Other Modality Encoder	Published Feature



Example:
Multimodal encoder
in VLSNR: Clip+ViT

Feature Interaction

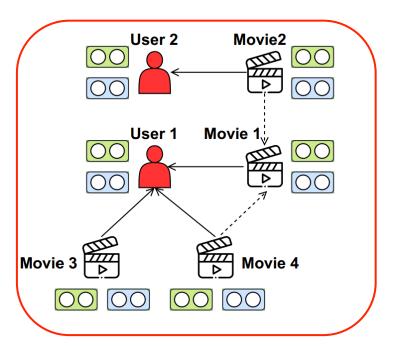


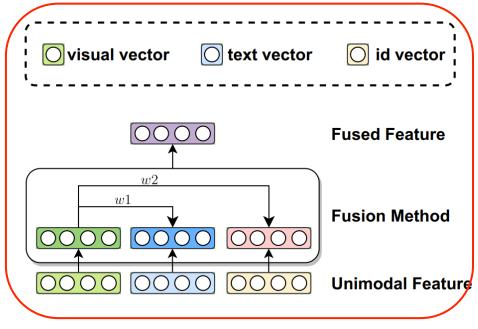


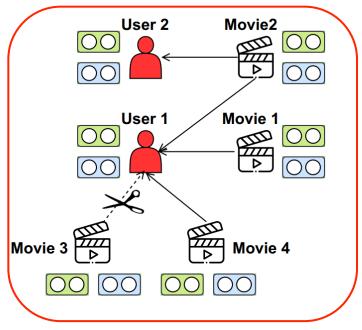




- Target: connecting various modalities
- **►**Taxonomy:
 - Bridge: capturing inter-relationship between users and items considering modalities
 - Fusion: capturing multimodal intra-relationships of items
 - Filtration: filtering out noisy data in interaction graph or multimodal features







Feature Interaction: Bridge









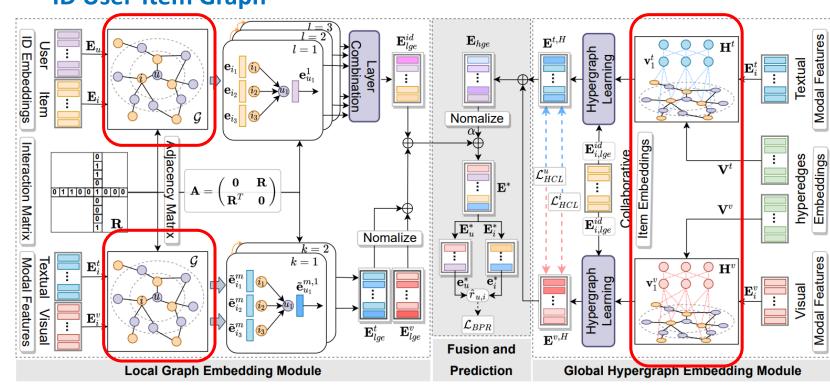
➤ User-Item Graph

Leveraging the information exchange between users and items to capture user's

multimodal preferences ID User-Item Graph

➤ LGMRec (AAAI'24)

- User-Item Graph: capture local preference
- Hypergraph: capture global preference
- Aggregating local and global preferences



Modal User-Item Graph

Hypergraph

Feature Interaction: Bridge









>Item-Item Graph

 Capturing latent semantic item-item structures to better learn item representations and improve model performance

➤MICRO (TKDE'22)

- Latent Structure Learning connect items by similarity in each modal
- Graph Convolutions: encodes two graphs
- Contrastive Fusion: learn fine-grained multimodal representations

Contrastive loss \mathcal{L}_C (a) Sports fleece hoodie... (b) Casual summer T-shirt... (c) Loose-fit summer shorts... Textual (d) Shirt for business, embedding wedding... Contrastive H^{t} (e) Summer casual tank top... Structure (f) Business pants... Textual structure ATextual features Fusion Visual Learning embedding Collaborative filtering $oldsymbol{H}^{ ext{v}}$ (MF, NGCF, ...)

Visual structure A

Textual Item-Item Graph

Visual features

Feature Interaction: Bridge









➤ Knowledge Graph

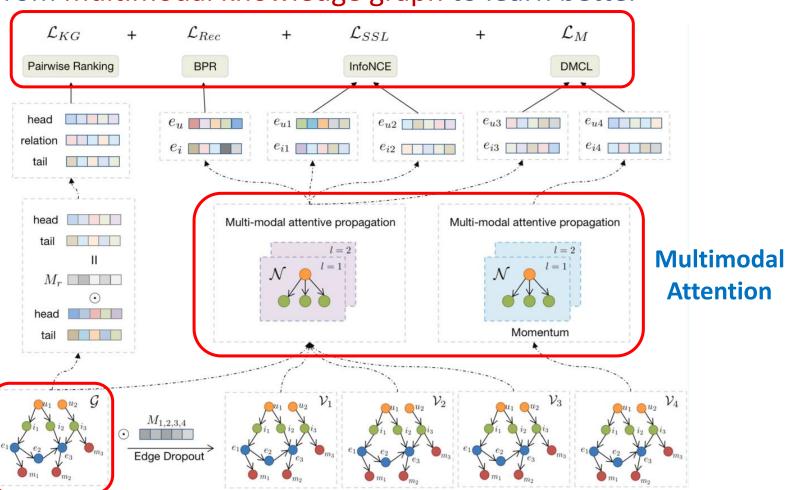
Utilizing auxiliary information from multimodal knowledge graph to learn better

multimodal representations

➤ M3KGR (IS'24)

- CLIP+Multimodal Attention: get aligned entity embedding
- Momentum Updating: increase quality of negatives
- Multi-task Learning: better item representations

Multimodal KG



Feature Interaction: Fusion





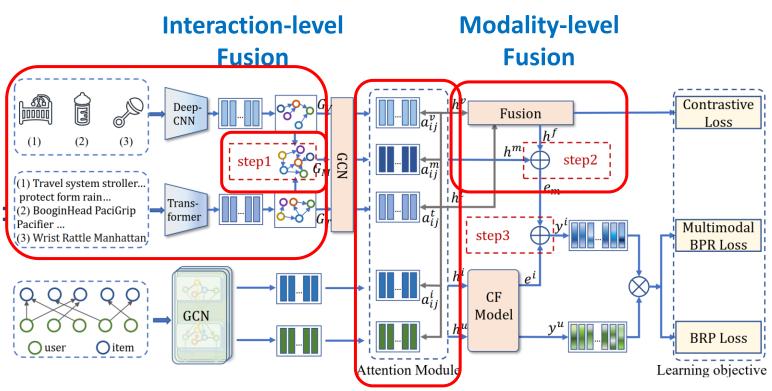




- ➤ Coarse-grained Attention
 - Capturing the multimodal relationships between interactions

>TMFUN (SIGIR'23)

- Graph Construction: itemitem graph for each modal
- Attention Relationship Mining: extract user's preferences
- Multi-step Fusion: interaction level and modality level



$$G_M = \lambda G_V + (1 - \lambda) G_T$$

Feature Interaction: Fusion









Fine-grained Attention

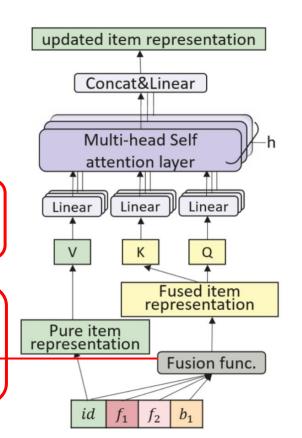
Capturing the multimodal relationships between modalities

➤ NOVA (AAAI'21)

- Non-invasive Attention: Q and K are multimodal feature embeddings, V is the ID embedding

$$\text{NOVA}(R, R^{(ID)}) = \sigma(\frac{QK^T}{\sqrt{d_k}})V$$

Fusion Operations: gating fusor
$$\mathcal{F}_{\mathrm{gating}}(f_1,\ldots,f_m) = \sum_{i=1}^m G^{(i)}f_i$$
 with trainable coefficients $G = \sigma(FW^F)$



Feature Interaction: Fusion







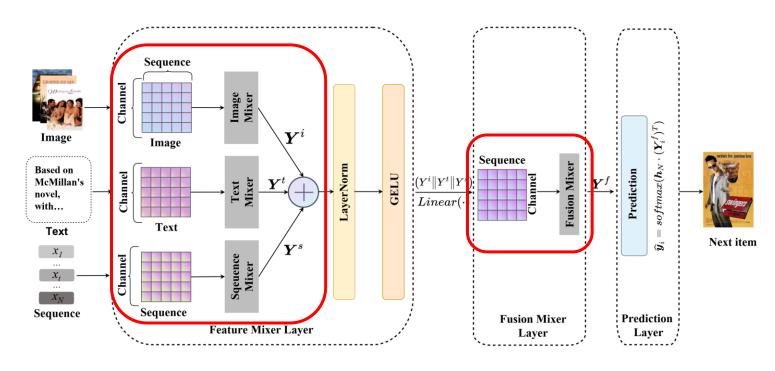


≻Other Methods

 Applying other simple methods, including concatenation operations, gating mechanism and etc.

>MMMLP (WWW'23)

- Feature Mixer Layer: learn sequential patterns
- Fusion Mixer Layer: fuse various modalities



Feature Mixer

Fusion Mixer

Feature Interaction: Filtration







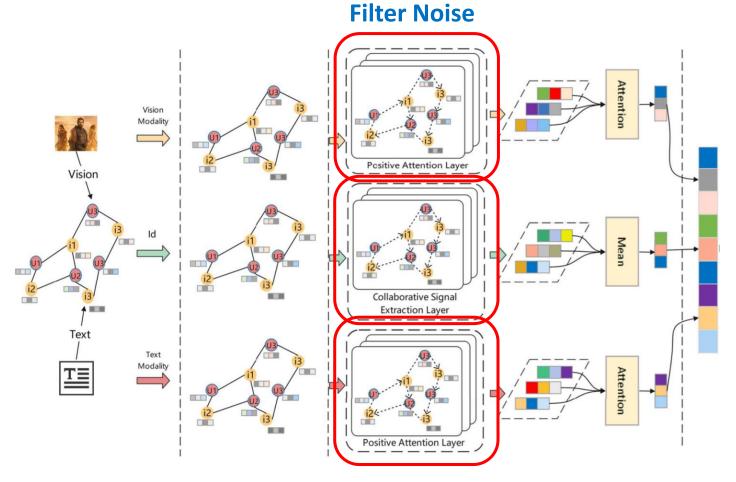


> Filtration

➤ Aiming at filtering out noisy data (data that is unrelated to user preferences)

➤ PMGCRN (APPL INTELL'23)

- Positive Attention Layer: remove implicit noisy edges by node similarity
- Collaborative Signal
 Extraction Layer: integrate
 original user-item
 interaction information



Feature Enhancement





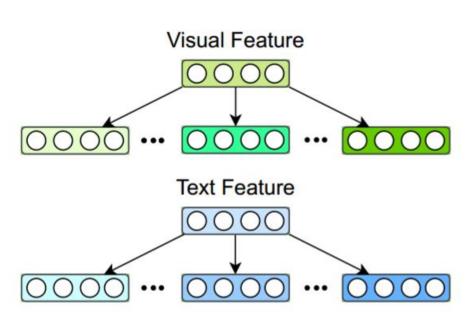




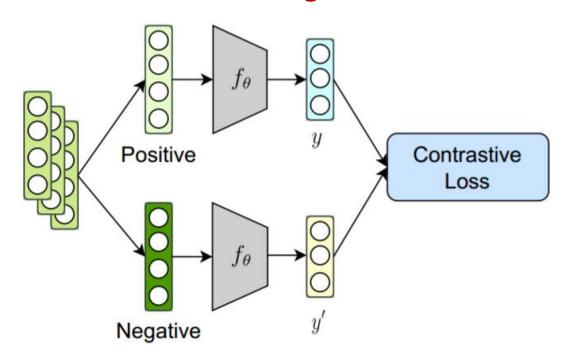
Target: distinguishing the unique and common characteristics of the multimodal features to improve the performance and generalization of MRS

≻Taxonomy:

Disentangled Representation Learning and Contrastive Learning



(a) Disentangled Representation Learning



(b) Contrastive Learning

Feature Enhancement









➤ Disentangled Representation Learning (DRL)

• introducing decomposition learning techniques to dig out the meticulous factors in user

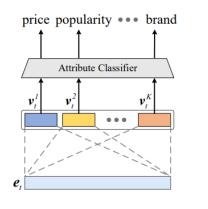
preference

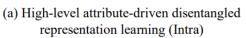
Low-level

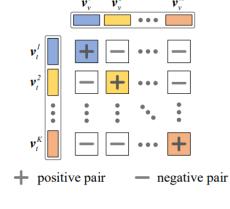
High-level

>AD-DRL (MM'24)

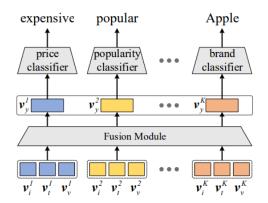
- Low-level Arrtribute-driven DRL:
 predicts attributes by raw modality
 features (intra- and inter disentanglement)
- High-level Arrtribute-driven DRL: predicts attributes by fused modality features







(b) High-level attribute-driven disentangled representation learning (Inter)



(c) Low-level attribute-driven disentangled representation learning

Feature Enhancement







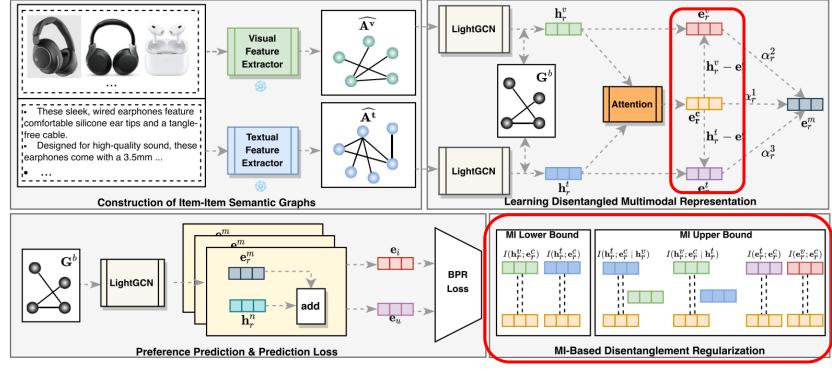


Contrastive Learning (CL)

 minimizing the distance between modalities of one item and maximizing the distance **Modality-shared** between the modalities of different ones **Modality-specific**

➤ CMDL (TOIS'25)

- **Disentangled Multimodal Represnetation:** modality-invariant part and modality-specific part
- **MI-based** Disentanglement **Regularization**: extended contrastive loss Extended CL Loss: $\mathcal{L}^{fub}(X,Y) = \mathbb{E}_{(x,y)\sim p(x,y)}[\hat{f}(x,y)] - \mathbb{E}_{x\sim p(x)}\mathbb{E}_{y\sim p(y)}[\hat{f}(x,y^-)],$



Model Optimization



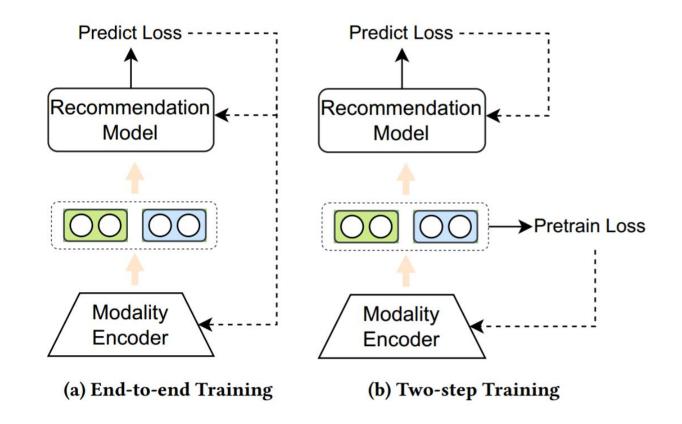






➤ Target: meeting the computational requirements for training MRS, including the multimodal encoders and RS model

- **≻**Taxonomy:
 - End-to-end Training and Two-step Training



Model Optimization







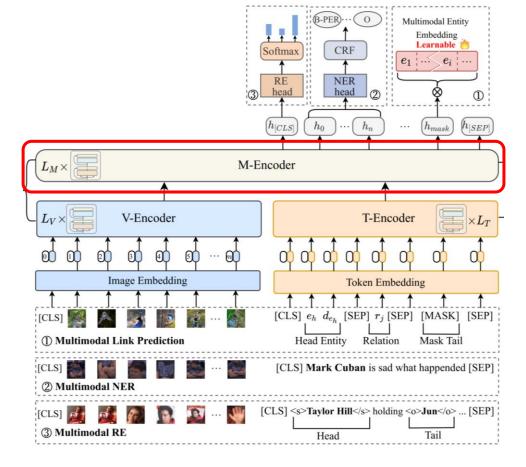


► End-to-end Training

Training the RS model and modality encoders together, making the encoders better adaptable to RS tasks

➤ MKGformer (SIGIR'22)

- Modality-specifc Parameters: each modality maintain their own layers, i.e., V-Encoder and T-Encoder
- **Modality-shared Parameters**: both modalities share serveal layers, i.e., M-Encoder



Model Optimization









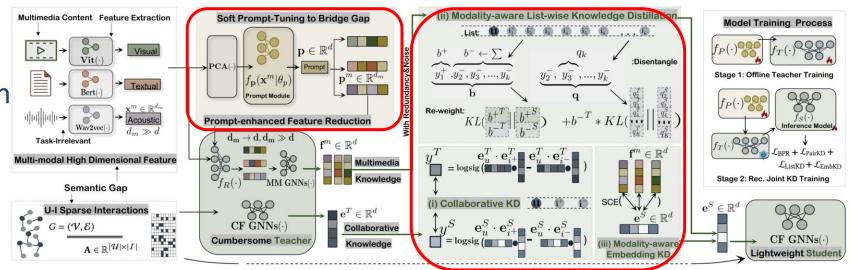
➤ Two-step Training

Training the RS model and modality encoders separately, making the training process more efficiently

➤ PromptMM (WWW'24)

- Soft Prompt Tuning: train the modality adapter for better alignment
- List-wise Knowledge
 Distillation: distill the multimodal knowledge to ID embedding

Soft Prompt Tuning Knowledge Distillation



Large Language Model for MRS









> LLM for Multimodal Recommendation

Model	LLM Type	Modality	RS Model
Rec-GPT4V (arXiv'24)	MLLM	Image→Text	LLM
MLLM-MSR (arXiv'24)	MLLM	Image→Text	LLM
Bundle-LLM (arXiv'24)	LLM	Image→Emb	LLM
QARM (arXiv'24)	MLLM	Image→Emb	Conventional RS
NoteLLM-2 (KDD'25)	LLM	Image→Emb	Conventional RS
UniMP (ICLR'25)	LLM	Image→Emb	LLM
UTGRec (arXiv'25)	MLLM	Image→Semantic ID	LLM

Large Language Model for MRS







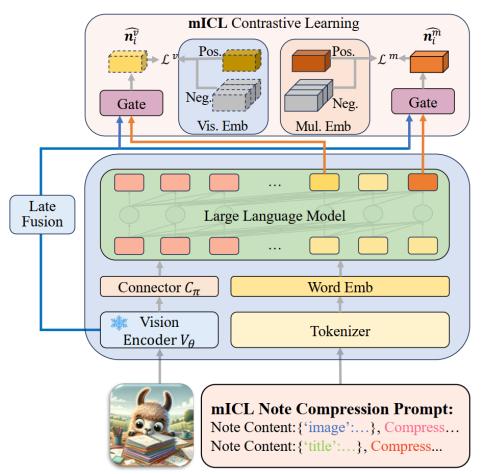


➤ Multimodal Large Language Model

Utilizing the powerful understanding abilities of MLLM

➤ NoteLLM-2 (KDD'25)

 Multimodal In-context Learning (mICL): separates multimodal content into visual and textual components, subsequently compressing the content into two modality-compressed words



Applications and Datasets









Data	Field	Modality	Scale	Link
Tiktok	Micro-video	V,T,M,A	726K+	https://paperswithcode.com/dataset/tiktok-dataset
Kwai	Micro-video	V,T,M	1 million+	https://zenodo.org/record/4023390#.Y9YZ6XZBw7c
Movielens + IMDB	Movie	V,T	100k~25m	https://grouplens.org/datasets/movielens/
Douban	Movie,Book,Music	V,T	1 million+	https://github.com/FengZhu-Joey/GA-DTCDR/tree/main/Data
Yelp	POI	V,T,POI	1 million+	https://www.yelp.com/dataset
Amazon	E-commerce	V,T	100 million+	https://cseweb.ucsd.edu/ jmcauley/datasets.html#amazon_reviews
Book-Crossings	Book	V,T	1 million+	http://www2.informatik.uni-freiburg.de/ cziegler/BX/
Amazon Books	Book	V,T	3 million	https://jmcauley.ucsd.edu/data/amazon/
Amazon Fashion	Fashion	V,T	1 million	https://jmcauley.ucsd.edu/data/amazon/
POG	Fashion	V,T	1 million+	https://drive.google.com/drive/folders/1xFdx5xuNXHGsUVG2VIohFTXf9S7G5veq
Tianmao	Fashion	V,T	8 million+	https://tianchi.aliyun.com/dataset/43
Taobao	Fashion	V,T	1 million+	https://tianchi.aliyun.com/dataset/52
Tianchi News	News	T	3 million+	https://tianchi.aliyun.com/competition/entrance/531842/introduction
MIND	News	V,T	15 million+	https://msnews.github.io/
Last.FM	Music	V,T,A	186 k+	https://www.heywhale.com/mw/dataset/5cfe0526e727f8002c36b9d9/content
MSD	Music	T,A	48 million+	http://millionsongdataset.com/challenge/

¹ 'V', 'T', 'M', 'A' indicate the visual data, textual data, video data and acoustic data, respectively.

Future Directions









>A Universal Solution

• Designing a universal solution with the combinations of the techniques mentioned before

≻Model Interpretability

Making the recommended items from MRS interpretable

> Computational Complexity

 Shrinking the computational cost and time required by MRS, due to parameterized modality encoder

≻Privacy

Protecting user's privacy under condition of affluent multimodal information

Conclusion









> Challenges:

Raw feature representation, feature interaction, recommendation

>Methodology:

Modality encoders, feature interaction, feature enhancement, model optimization

Survey









[CSUR'24] https://arxiv.org/abs/2302.03883

Multimodal Recommender Systems: A Survey

QIDONG LIU*, Xi'an Jiaotong University & City University of Hong Kong, China JIAXI HU*, City University of Hong Kong, China YUTIAN XIAO*, City University of Hong Kong, China XIANGYU ZHAO[†], City University of Hong Kong, China JINGTONG GAO, City University of Hong Kong, China WANYU WANG, City University of Hong Kong, China QING LI, The Hong Kong Polytechnic University, China JILIANG TANG, Michigan State University, USA



Agenda









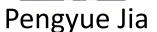




Yejing Wang

Joint Modeling in RS





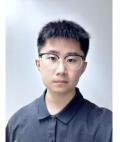




Yuhao Wang







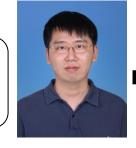
Pengyue Jia Xiaopeng Li

Multi-behavior Recommendation



Jingtong Gao

Multi-modal Recommendation



Qidong Liu



Future Work



Yichao Wang

Conclusion

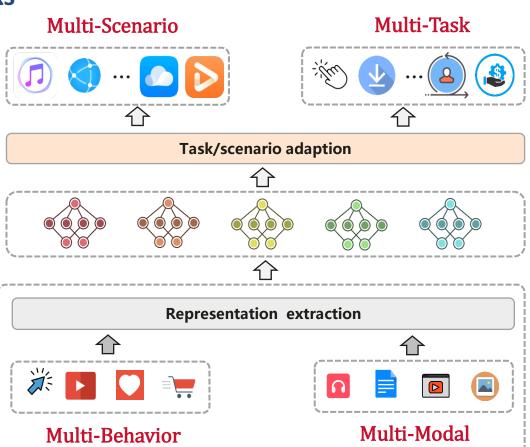








- ➤ Utilizing diverse user feedback signals from **different tasks**
- Extracting commonalities and diversities of user preferences from **different scenarios**
- Fusing heterogeneous information from different data modalities
- ➤ Acquiring multi-aspect user preferences from different type of **behaviors**
- Introducing open-world knowledge from large language models











➤ Multi-Task Recommendation

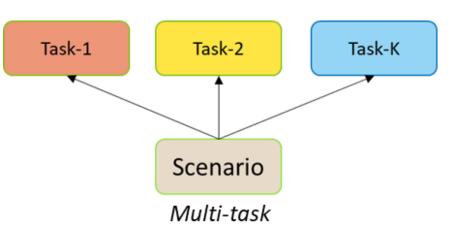
- Task relation:
 Parallel, Cascaded, Auxiliary with Main
- Methodology:
 Parameter Sharing, Optimization, Training Mechanism
- Trends:

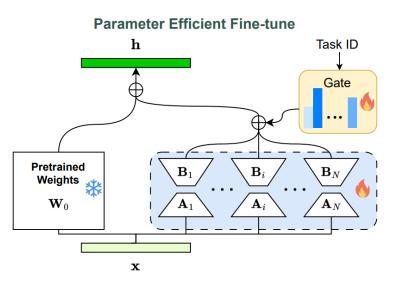
Using LLM Modeling the similarities and differences across tasks with MoE-LoRA to enable task generalization.

e.g., MOELoRA

Take Away:

- Future Direction:
 Mitigating negative transfer with LLM's world knowledge
- HOW to design the parameter sharing pattern with the purpose of capturing complex task relevance













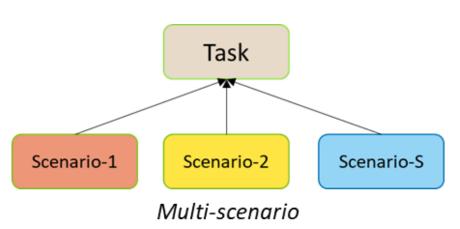
➤ Multi-Scenario Recommendation

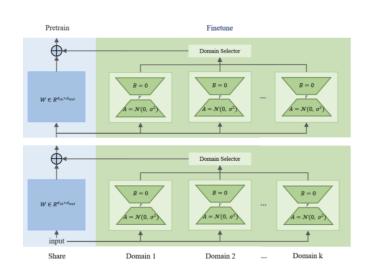
- Modeling methods: shared-specific network paradigm, and dynamic weight paradigm.
- Trends:

Incorporating Low-Rank Adaptor (LoRA) for domainspecific fine-tuning (e.g., M-LoRA); Incorporating LLM's world knowledge in domain understanding (e.g., LLM4MSR)

- Future Direction:
 Efficient M-LoRA architecture for large scale scenarios

 Scenario recommendation interpretability with LLMs
- Take Away:
 How to model the commonalities and diversities between different scenarios













Multi-Behavior Recommendation

• Behavior:

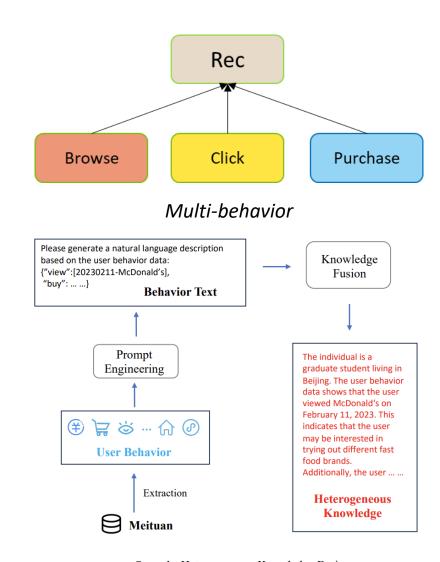
Macro behaviors, Micro behaviors, Behaviors from different domains or scenarios

- Methods:RNN, Graph, Transformer
- Trends:

Better modeling architecture; Generative Modeling; Modeling with LLM (e.g., knowledge fusion in HKFR)

- Future Direction:
 Fine-grained behavior understanding with LLMs
 Behavior debias with LLMs
- Take Away:

How to model the complex behavior patterns and behavior correlations from users



Stage 1 : Heterogeneous Knowledge Fusion









Multi-Modal Recommendation

Methods:

Modality Encoder, Feature Interaction, Feature Enhancement, Model Optimization

• Trends:

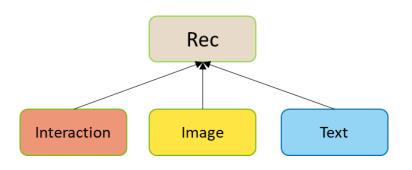
Utilizing the powerful understanding abilities of MLLM to enhance Multi-Modal Rec (e.g., NoteLLM)

Future Direction:

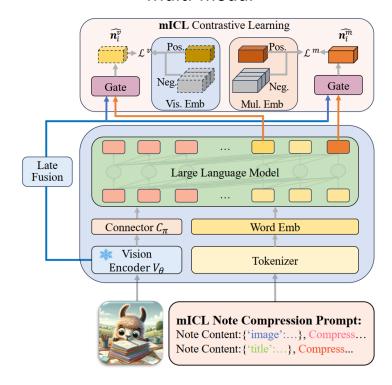
Recommendation representation align with MLLM representation; personalized multi-modal generation recommendation (e.g., image personalized generation)

Take Away:

How to encode different modalities to enhance recommender system



Multi-modal













We are hiring!



Huawei Noah's Ark Lab



WWW25 Huawei Noah's Ark Lab Chat Group



AML Lab
CityU

Tutorial Slides

https://zhaoxyai.github.io/paper/jointmodeling-www2025.pdf

- [1] Multi-Task Deep Recommendation Systems: A Survey. https://arxiv.org/abs/2302.03525
- [2] Scenario-Wise Rec: A Multi-Scenario Recommendation Benchmark. https://arxiv.org/abs/2412.17374
- [3] Multimodal Recommender Systems: A Survey. https://arxiv.org/abs/2302.03883