

Task-Augmented Knowledge Distillation for Multi-Task Personalized Recommendation

Presenter : Chenxiao Yang
May 31st, 2021

Background | Multi-Task Learning (MTL)

Heterogenous User Feedback as Prediction Tasks: Click, Like, Add to Cart, etc.

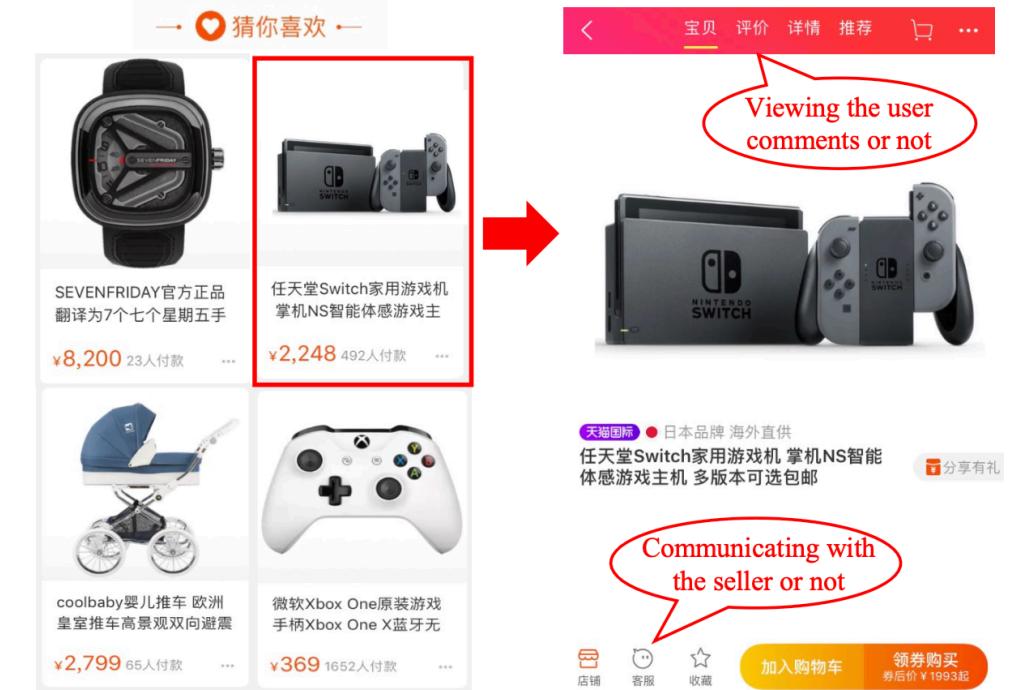
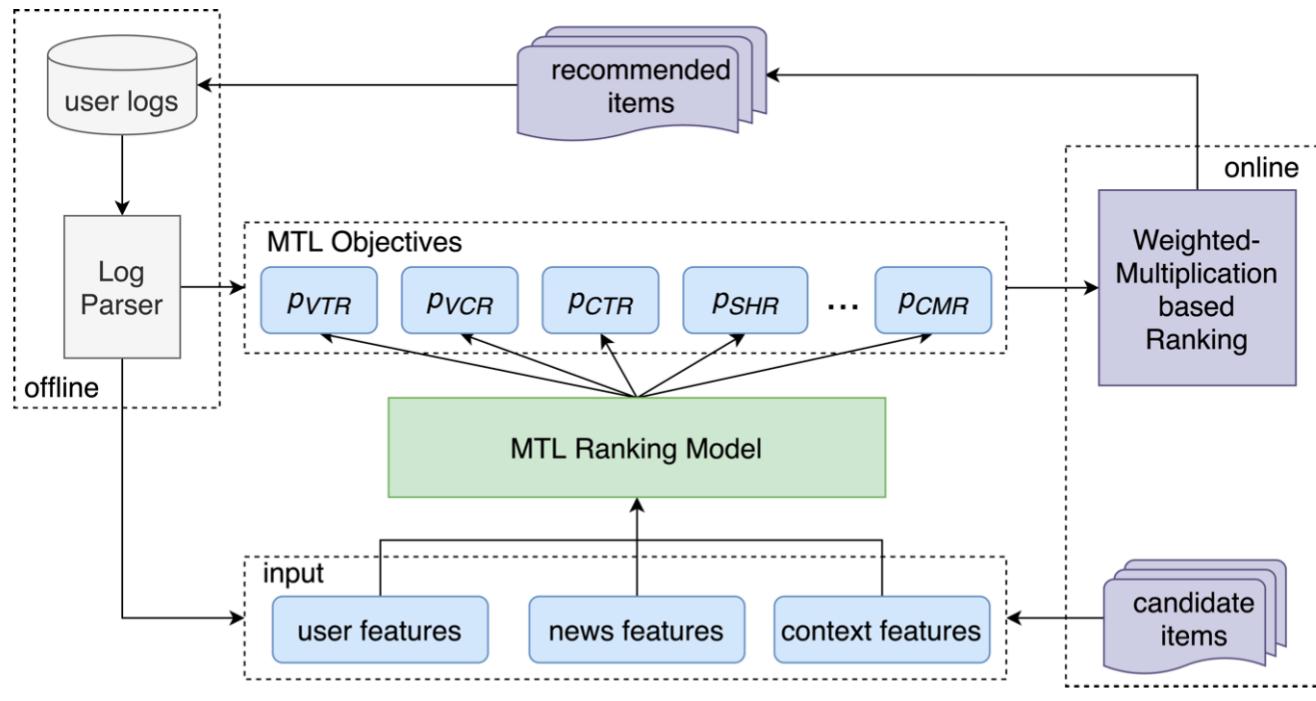


Fig. Illustration of Multi-Task Recommender Systems (from PLE [1] and PFD [3])

Background | Multi-Task Learning (MTL)

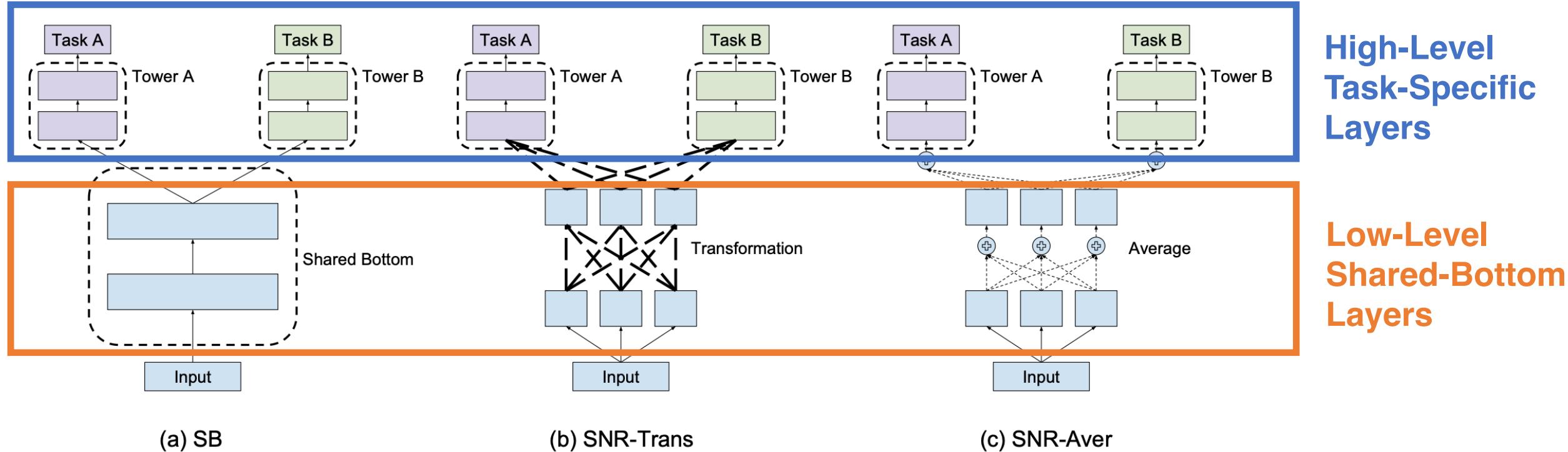


Fig. A Generic Paradigm for Multi-Task Learning Models (from SNR [2])

NO Knowledge Transfer in or above **Task-Specific Layers**

3

Implicit Knowledge Transfer in **Shared-Bottom Layers** by Branching (b), Gating (c), etc.

Background | Multi-Task Learning (MTL)

LIMITATIONS:

1. Bottleneck on the **ability of knowledge transfer**
2. A lack of **interpretability**
(will be discussed)
3. Unable to effectively solve **task conflicts**
(will be discussed)

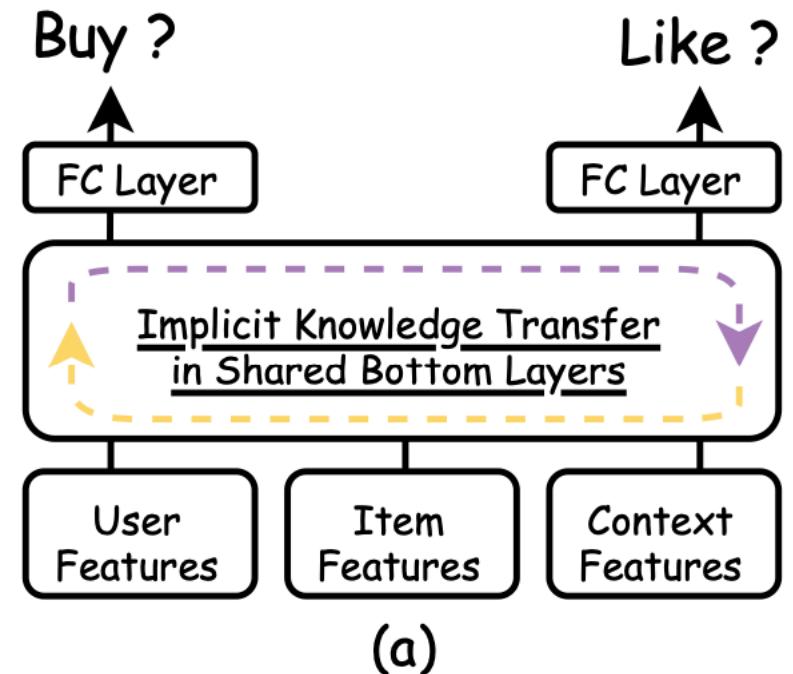


Fig. Common MTL framework

Background | Knowledge Distillation (KD)

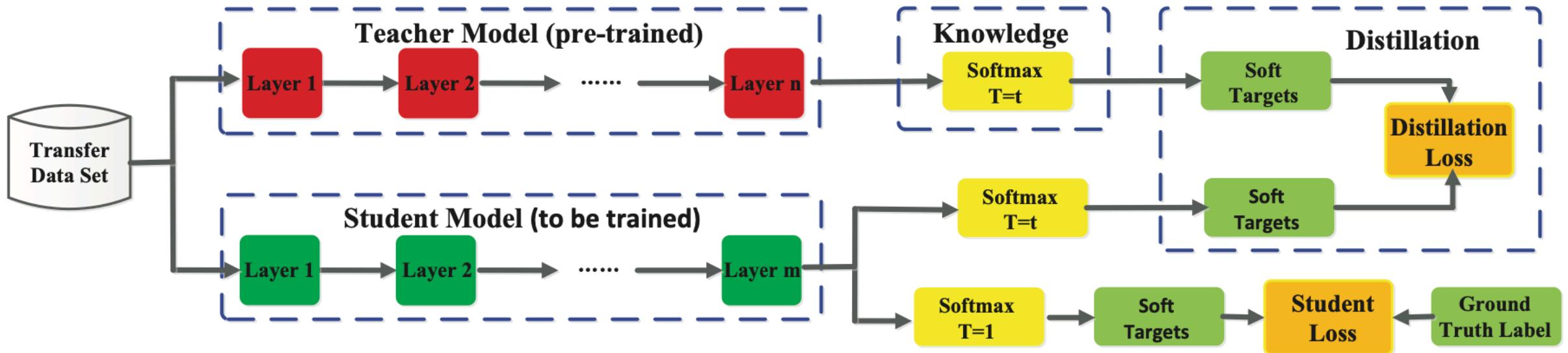


Fig. Generic Response-Based Knowledge Distillation (from a KD survey)

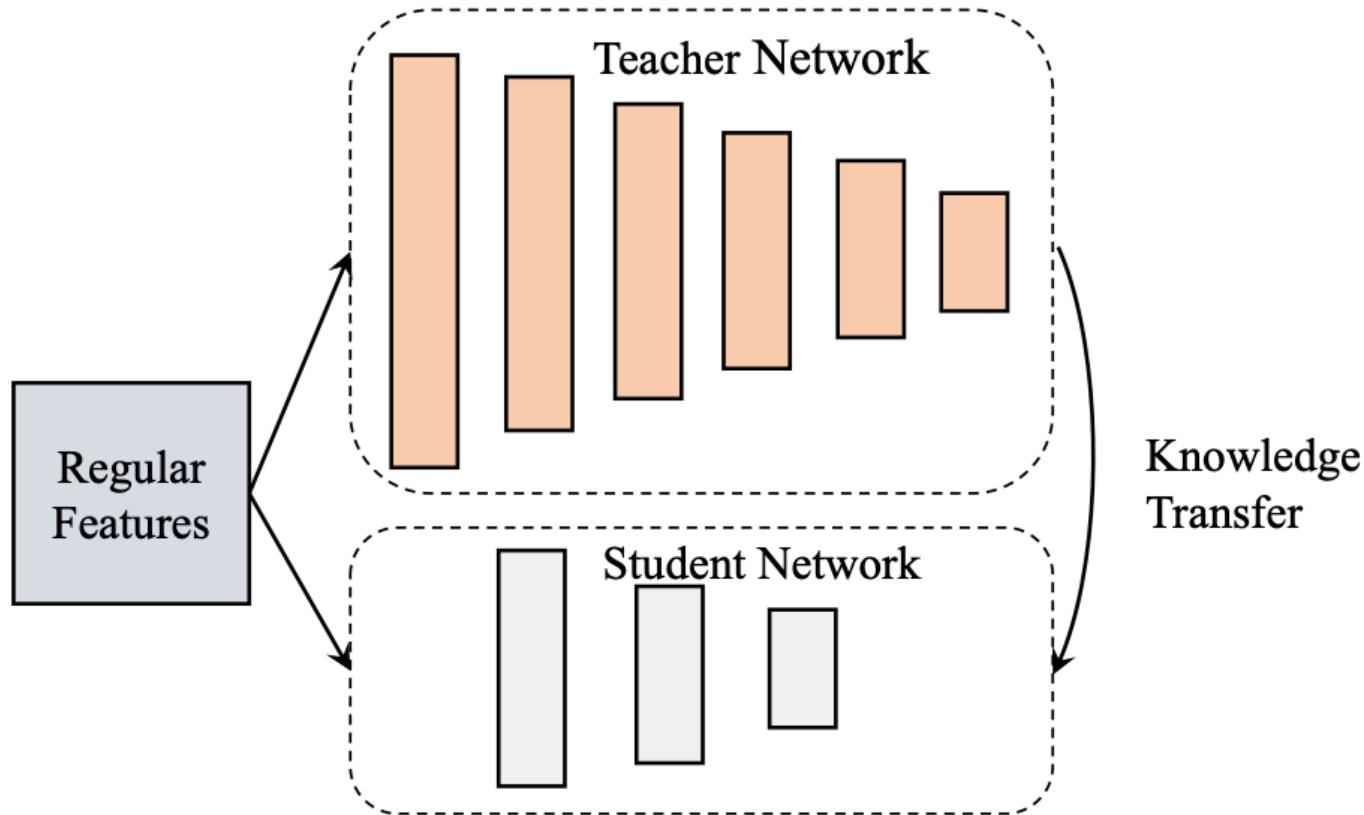
Hard targets	0	1	0	0
	cow	dog	cat	car
More Informative!				
Soft targets	10^{-6}	0.9	0.1	10^{-9}

$$L_D(p(z_t, T), p(z_s, T)) = \sum_i -p_i(z_{ti}, T) \log(p_i(z_{si}, T)) ,$$

$$L_S(y, p(z_s, T)) = \mathcal{L}_{CE}(y, p(z_s, T)) ,$$

$$\begin{aligned} L(x, W) = & \alpha * L_D(p(z_t, T), p(z_s, T)) \\ & + (1 - \alpha) * L_S(y, p(z_s, T)) , \end{aligned}$$

Background | Knowledge Distillation (KD)



Model Compression:
(Conventional Usage)

Cumbersome Teacher Model
->
Light-weight Student Model

Fig. Model Compression (from PFD [3])

Background | Knowledge Distillation (KD)

KD as Intelligent Label Smoothing Regularization for **Knowledge Transfer**
(Our usage)

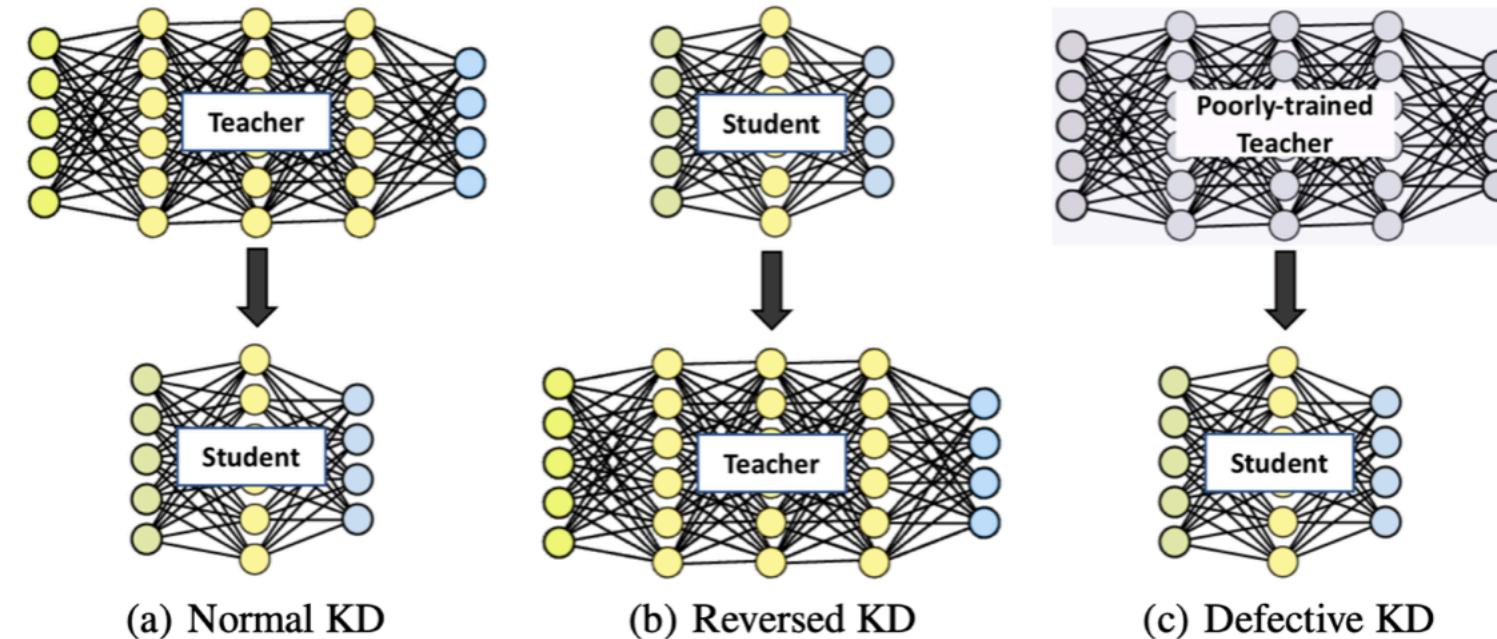


Fig. Exploratory Experiments (from [4])

Benefits of Knowledge Distillation:

1. Label smoothing
2. Example re-weighting
3. Prior knowledge of optimal output layer geometry

Method | Motivation

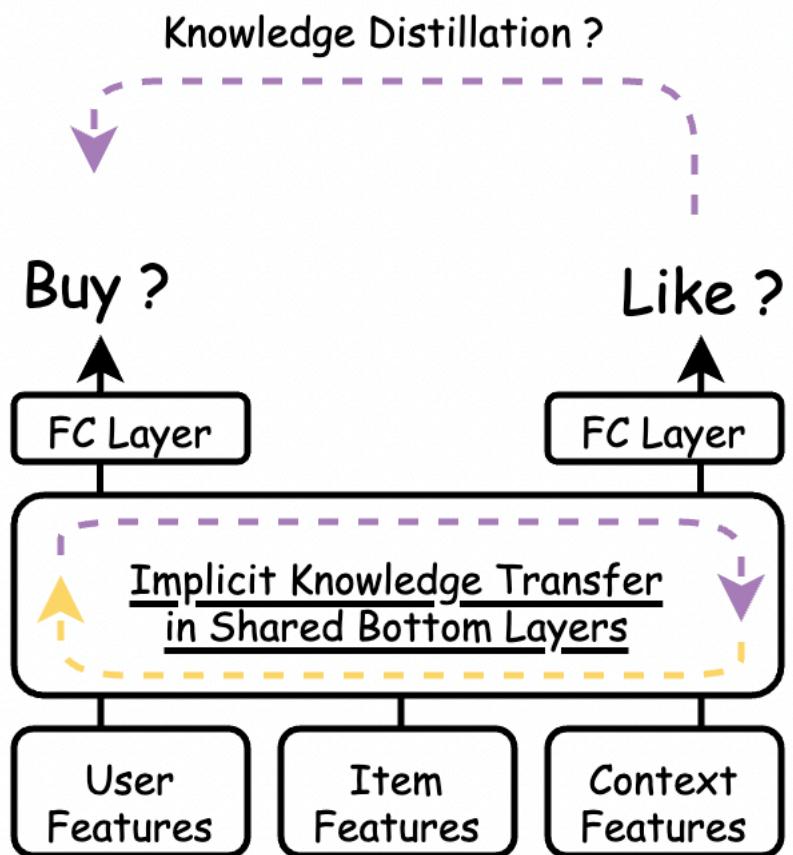


Fig. Intuitive Solution

Motivation:

Predictions of another task contains sample relation unseen to the main-task (which we want to transfer).

Intuitive Solution:

Directly Using Outputs of Another Task as Soft Labels

Problem:

TASK CONFLICTS !

Method | Motivation

Q: What kind of knowledge do we expect to **transfer** or elsewise to **avoid**?

Task A (Buy > Not Buy)

Buy & Like > Not Buy & Like	
Buy & Not Like > Not Buy & Like	
Buy & Like > Not Buy & Not Like	
Buy & Not Like > Not Buy & Not Like	

Task B (Like > Not Like)

	→ Buy & Like ? Not Buy & Like
*	Buy & Not Like < Not Buy & Like
*	Buy & Like < Not Buy & Not Like
	→ Buy & Not Like ? Not Buy & Not Like

Fig. Knowledge Disentanglement

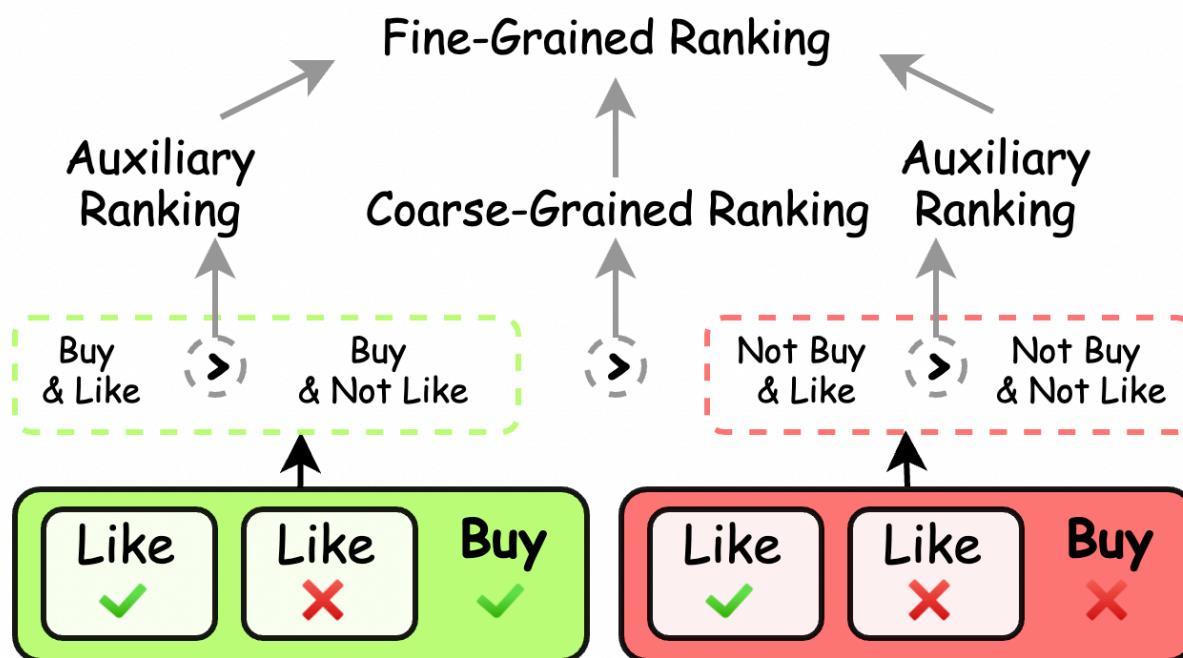
Solution: To Disentangle **Positive** and **Negative** Knowledge

ONLY transfer **Positive** Knowledge to Avoid Task Conflicts in Essence

Method | Motivation

UNCONFLICTED Fine-Grained Ranking

= Coarse-Grained Ranking + Auxiliary Ranking (**Positive Knowledge**)



Coarse-Grained Ranking

1. Only contain knowledge about the main task
2. Conflicted across tasks

V.S.

Fine-Grained Ranking

1. Contain across-task knowledge
2. Unconflicted with the main task

Fig. Unconflicted Fine-Grained Ranking Knowledge

Method | Motivation

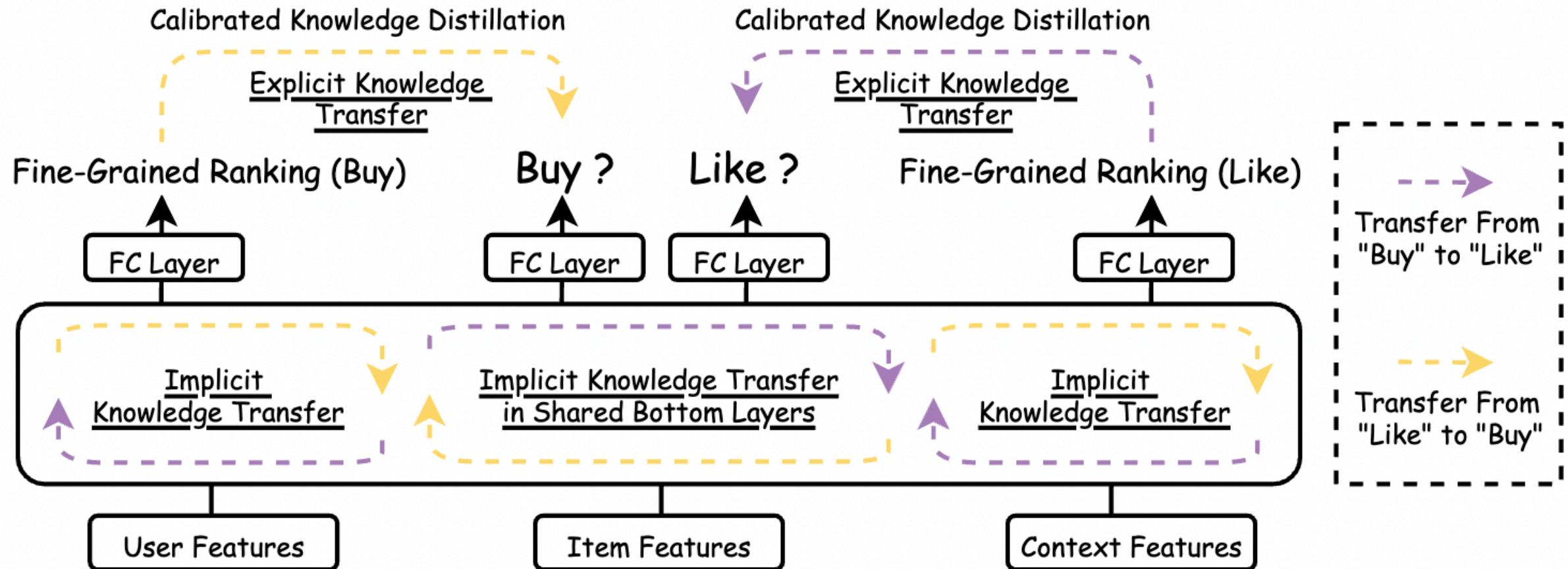


Fig. Simple Framework for TAKD based on Aforementioned Motivation

Method | Task Augmentation

Binary Cross Entropy Loss (Regression-Based): learning to fit the label

$$\mathcal{L}^A = CE(y^A, \hat{y}^A) = \sum_{\mathbf{x}_i \in \mathcal{D}} -y_i^A \ln \hat{y}_i^A - (1 - y_i^A) \ln(1 - \hat{y}_i^A)$$

$$\mathcal{L}^B = CE(y^B, \hat{y}^B) = \sum_{\mathbf{x}_i \in \mathcal{D}} -y_i^B \ln \hat{y}_i^B - (1 - y_i^B) \ln(1 - \hat{y}_i^B),$$

BPR Pair-wise Loss (Ranking-Based) [5]:

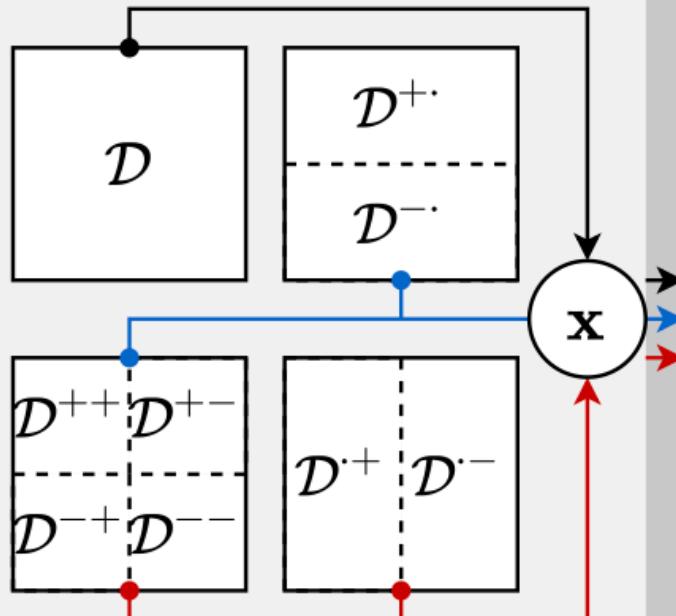
learning a bipartite ranking relation (positive sample > negative sample)

$$\mathcal{L}^{BPR} = \sum_{x_i \in \mathcal{D}^+} \sum_{x_j \in \mathcal{D}^-} -\log \sigma(r_i - r_j),$$

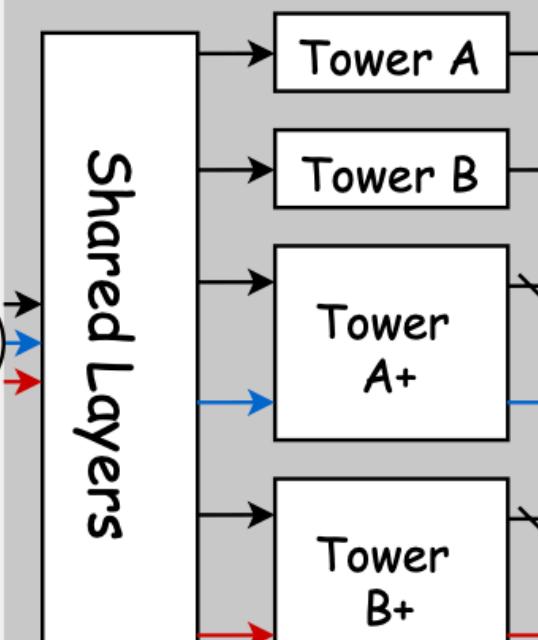
- Sigmoid(r) is strictly NOT a probability for BPR loss

Method | Task Augmentation

1. Dataset Sampling



2. Multi-Task Model



Split Training Dataset According to Permutations of Tasks' Labels

$$\mathcal{D}^{+-} = \{\mathbf{x}_i \in \mathcal{D} \mid y_i^A = 1, y_i^B = 0\},$$

$$\mathcal{D}^{-+} = \{\mathbf{x}_i \in \mathcal{D} \mid y_i^A = 0, y_i^B = 1\},$$

$$\mathcal{D}^{--} = \mathcal{D}^{--} \cup \mathcal{D}^{-+}, \quad \mathcal{D}^{++} = \mathcal{D}^{+-} \cup \mathcal{D}^{++},$$

$$\mathcal{D}^{--} = \mathcal{D}^{--} \cup \mathcal{D}^{+-}, \quad \mathcal{D}^{+-} = \mathcal{D}^{-+} \cup \mathcal{D}^{++}.$$

Fig. Training Set Split

Method | Task Augmentation

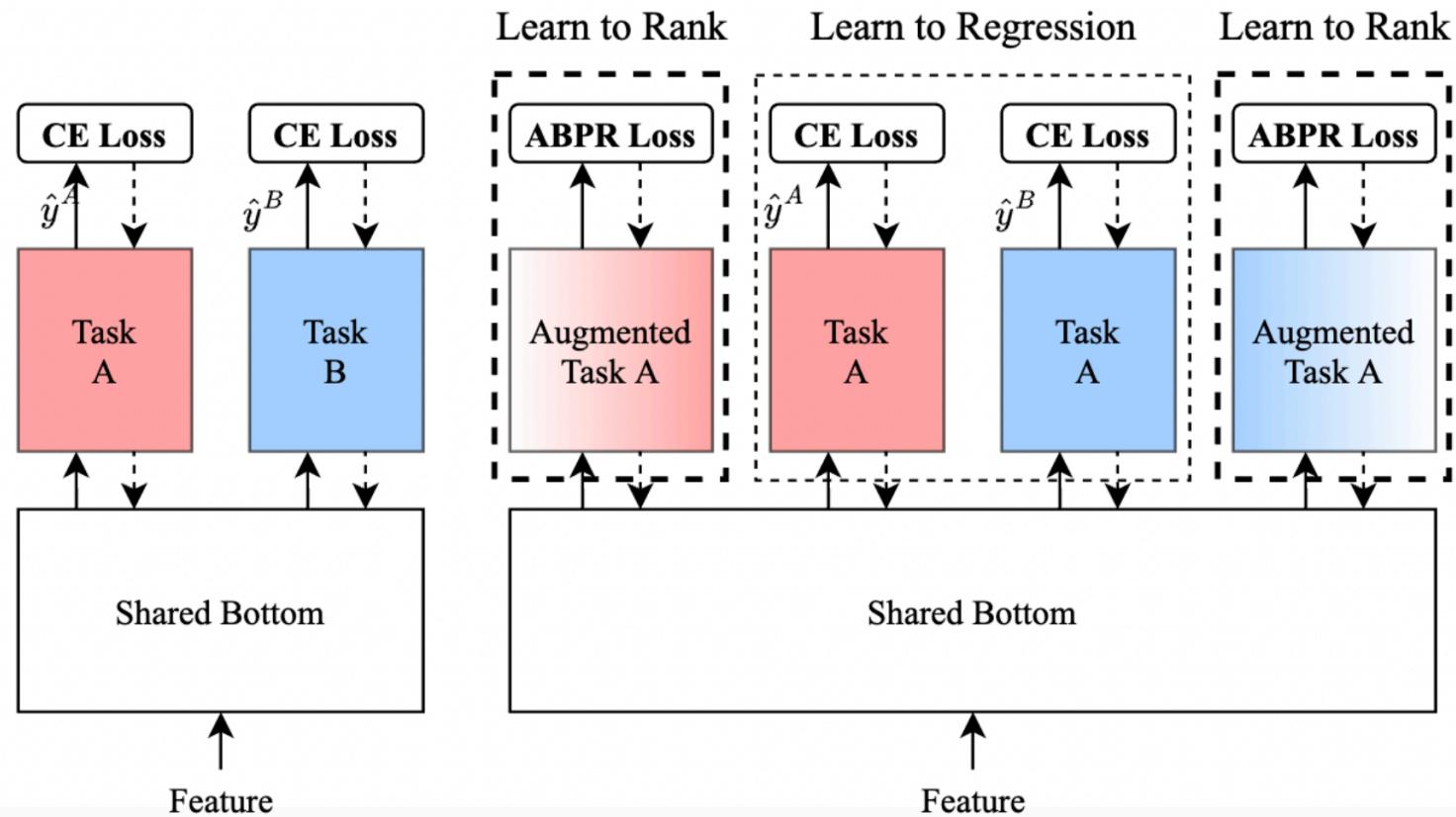


Fig. Task Augmentation

Augmented Ranking-Based Tasks

Aim to learn fine-grained ranking $++ > +- > -+ > --$ with the **optimization target** :

$$\begin{aligned} & \ln p(\Theta | \succ) \\ &= \ln p(\mathbf{x}_{++} \succ \mathbf{x}_{+-} | \Theta) \cdot p(\mathbf{x}_{-+} \succ \mathbf{x}_{--} | \Theta) \cdot p(\Theta) \\ &= \sum_{(\mathbf{x}_{++}, \mathbf{x}_{+-}, \mathbf{x}_{-+}, \mathbf{x}_{--})} \ln \sigma(\hat{r}_{++\succ+-}) + \ln \sigma(\hat{r}_{-+\succ--}) - \\ & \quad - Reg(\Theta), \end{aligned}$$

Method | Task Augmentation

Combined with Original BPR Loss:

$$\mathcal{L}^{A+} = \sum_{(\mathbf{x}_{++}, \mathbf{x}_{+-}, \mathbf{x}_{-+}, \mathbf{x}_{--})} -\beta_1^A \ln \sigma(\hat{r}_{++\succ+-}) - \beta_2^A \ln \sigma(\hat{r}_{-+\succ--})$$

Auxiliary Rankings

$$- \sum_{(\mathbf{x}_{+..}, \mathbf{x}_{-..})} -\ln \sigma(\hat{r}_{+.\succ-}).$$

Coarse-Grained Ranking

$$\mathcal{L}^{B+} = \sum_{(\mathbf{x}_{++}, \mathbf{x}_{+-}, \mathbf{x}_{-+}, \mathbf{x}_{--})} -\beta_1^B \ln \sigma(\hat{r}_{++\succ-+}) - \beta_2^B \ln \sigma(\hat{r}_{+-\succ--})$$
$$+ \sum_{(\mathbf{x}_{.+}, \mathbf{x}_{.-})} -\ln \sigma(\hat{r}_{.+}\succ-.).$$

Method | Calibrated Knowledge Distillation

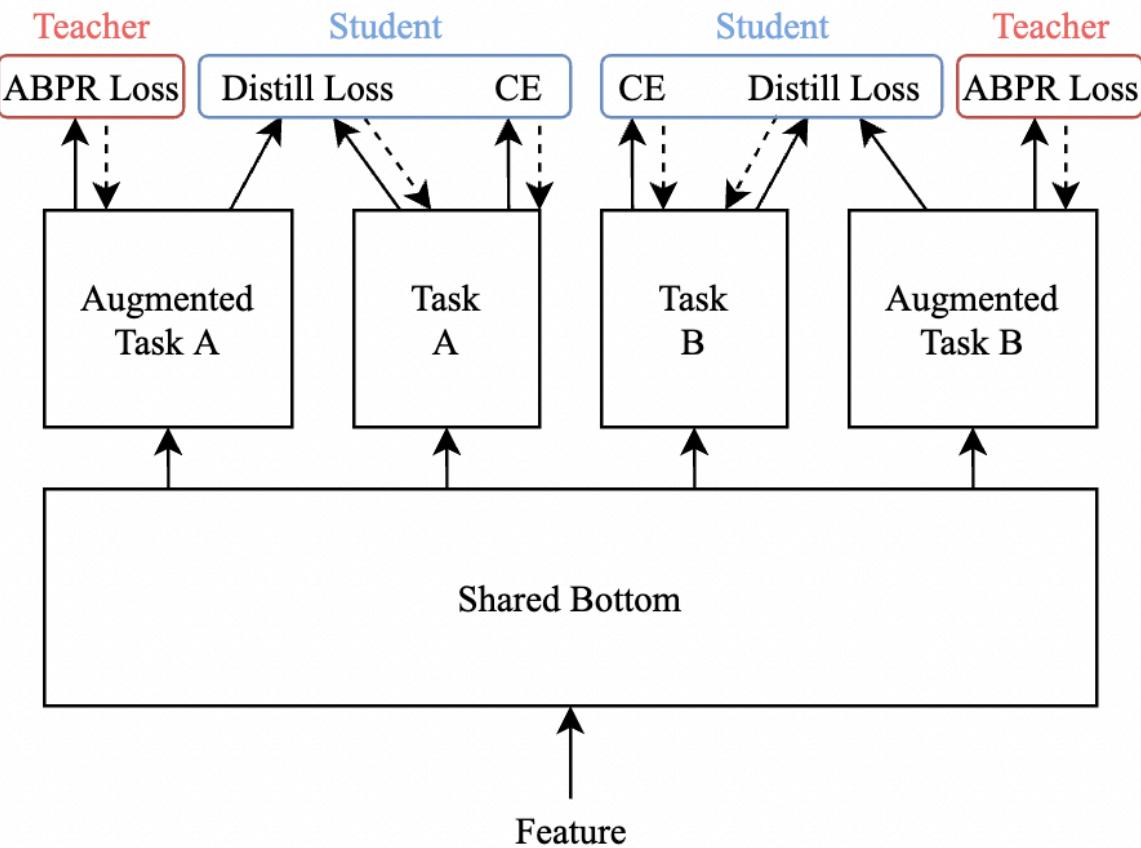


Fig. Directly Using Teacher Outputs

Directly Using Outputs of Teacher for Distillation

Distillation Loss:

$$\mathcal{L}^{A-KD} = CE(\sigma(\hat{r}^{A+}/\tau), \sigma(\hat{r}^A/\tau)),$$

$$\mathcal{L}^{B-KD} = CE(\sigma(\hat{r}^{B+}/\tau), \sigma(\hat{r}^B/\tau)),$$

Student Loss:

$$\mathcal{L}^{A-Stu} = (1 - \alpha^A)\mathcal{L}^A + \alpha^A \mathcal{L}^{A-KD},$$

$$\mathcal{L}^{B-Stu} = (1 - \alpha^B)\mathcal{L}^B + \alpha^B \mathcal{L}^{B-KD},$$

Method | Calibrated Knowledge Distillation

Challenge : Sigmoid(r) is strictly **NOT** a probability for BPR loss. **Calibration** is required !

$$\mathcal{L}^{BPR} = \sum_{x_i \in \mathcal{D}^+} \sum_{x_j \in \mathcal{D}^-} -\log \sigma(r_i - r_j),$$

A Simple Solution : Platt Scaling

$$\tilde{r}^{A+} = P^A \cdot \hat{r}^{A+} + Q^A, \quad \tilde{y}^{A+} = \frac{1}{1 + \exp \tilde{r}^{A+}}$$

Calibration Loss $\mathcal{L}^{Cal} = \mathcal{L}^{A-Cal} + \mathcal{L}^{B-Cal} = CE(y^A, \tilde{y}^{A+}) + CE(y^B, \tilde{y}^{B+}).$

Method | Training Procedure

Algorithm 1: Training Algorithm for TAKD

```
Input: Training dataset  $\mathcal{D}$ , learning rate  $\gamma_1$  and  $\gamma_2$ , initial  
parameters  $\Theta$  and  $\Omega$ , task weights  $w_1, w_2, w_3, w_4$ .  
1 Construct set  $\mathcal{D}^{++}, \mathcal{D}^{+-}, \mathcal{D}^{-+}, \mathcal{D}^{--}, \mathcal{D}^{+\cdot}, \mathcal{D}^{-\cdot}, \mathcal{D}^{\cdot+}, \mathcal{D}^{\cdot-}$ ;  
2 while Not converged do  
3     Sample  $\mathbf{x}$  uniformly at random from  $\mathcal{D}$ ;  
4     Sample  $\mathbf{x}_{++}, \mathbf{x}_{+-}, \mathbf{x}_{-+}, \mathbf{x}_{--}$  uniformly at random  
       from  $\mathcal{D}^{++}, \mathcal{D}^{+-}, \mathcal{D}^{-+}, \mathcal{D}^{--}$  respectively;  
5     Sample  $\mathbf{x}_{+..}, \mathbf{x}_{-..}, \mathbf{x}_{..+}, \mathbf{x}_{..-}$  uniformly at random from  
        $\mathcal{D}^{+\cdot}, \mathcal{D}^{-\cdot}, \mathcal{D}^{\cdot+}, \mathcal{D}^{\cdot-}$  respectively;  
6     Model parameter  $\Theta$  optimization:  
7         Calculate  $\mathcal{L}^{A+}(\mathbf{x}_{+..}, \mathbf{x}_{-..}, \mathbf{x}_{++}, \mathbf{x}_{+-}, \mathbf{x}_{-+}, \mathbf{x}_{--}; \Theta)$ ;  
8         Calculate  $\mathcal{L}^{B+}(\mathbf{x}_{+..}, \mathbf{x}_{-..}, \mathbf{x}_{++}, \mathbf{x}_{+-}, \mathbf{x}_{-+}, \mathbf{x}_{--}; \Theta)$ ;  
9         Calculate  $\mathcal{L}^{A-Stu}(\mathbf{x}; \Theta), \mathcal{L}^{B-Stu}(\mathbf{x}; \Theta)$ ;  
10          $\mathcal{L}^{Model} \leftarrow w_1 \mathcal{L}^{A+} + w_2 \mathcal{L}^{B+} + w_3 \mathcal{L}^{A-Stu} + w_4 \mathcal{L}^{B-Stu}$ ;  
11          $\Theta \leftarrow \Theta - \gamma_1 \nabla_{\Theta} \mathcal{L}^{Model}$ ;  
12     Calibration parameter  $\Omega$  optimization:  
13         Calculate  $\mathcal{L}^{Cal}(\mathbf{x}; \Omega)$ ;  
14          $\Omega \leftarrow \Omega - \gamma_2 \nabla_{\Omega} \mathcal{L}^{Cal}$ ;  
15 end
```

Two Sets of Parameters:

Theta -> Model (Backbone) Parameters

Omega -> Calibration Parameters



End-to-End Training

Bi-Level Training Procedure:

1. Sampling
2. Updating Model Parameter
3. Updating Calibration Parameter

Method | Correcting Errors

Early Stage of Training – Teacher is NOT well-trained

frequent errors in soft labels may distract the training process of student, causing **slow convergence**

Later Stage of Training – Teacher is well-trained

it is still likely teacher model would occasionally provide mistaken predictions which is called **genetic errors** in another literature

Error Correction Mechanism (Align with the hard label)

$$r^{Teacher}(\mathbf{x}) \leftarrow \mathbb{1}[y] \cdot \text{Max} \{ \mathbb{1}[y] \cdot r^{Teacher}(\mathbf{x}), m \}$$

m: Error correction margin

$\mathbb{1}[y] = 1$ if ($y == 1$) else -1

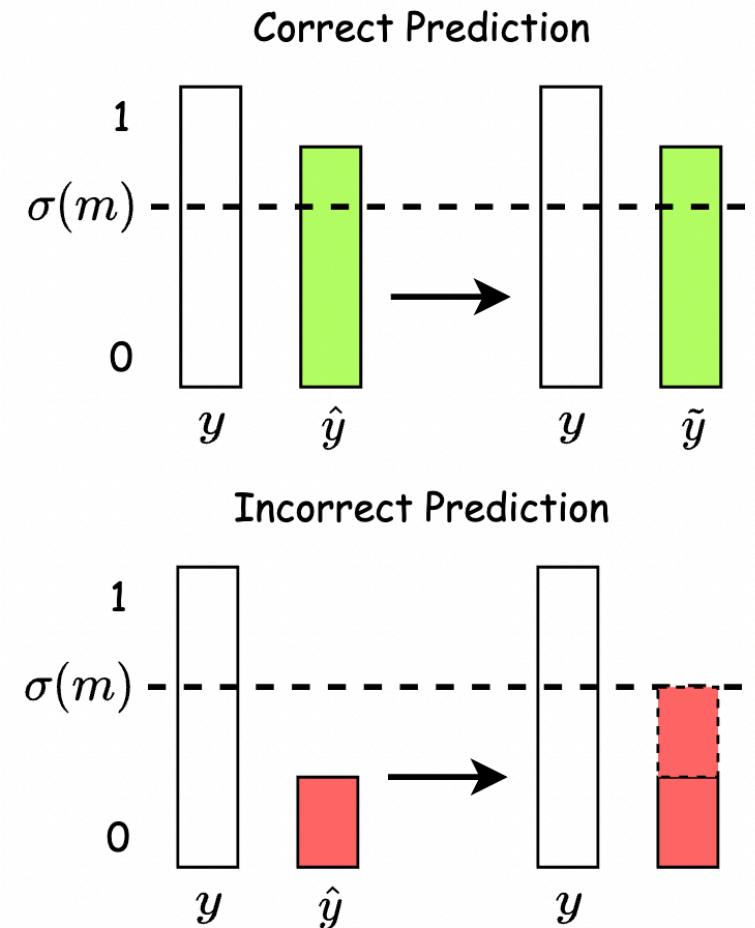


Fig. Error Correction Mechanism

Method | Framework Overview

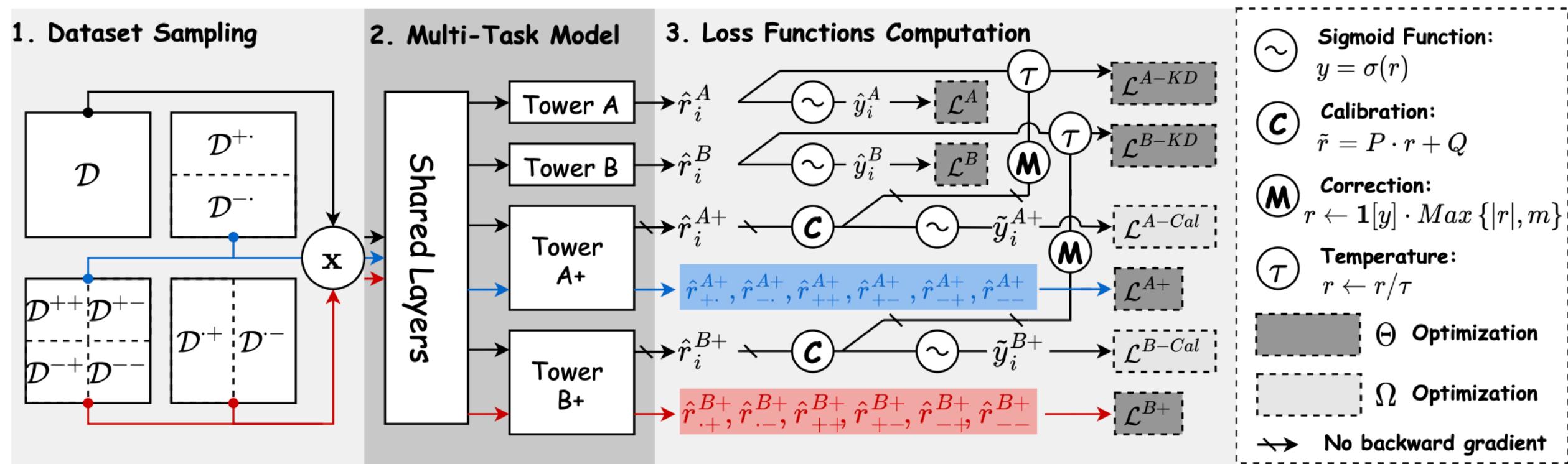


Fig. TAKD framework overview

Method | Adoptions

Example 1: A - Finish Watching, B - Like

Finish & Like > Finish & not Like > not Finish & Like > not Finish & not Like

Example 2: A - CVR, B - CTR

Click & Conversion > Click & not Conversion > not Click

Example 3: A - Click, B - not Interested

Click > not Click > not Interested

Example 4: A, B, C (priority A > B > C)

A & B & C > A & B & notC > A & notB & C > A & notB & notC >
notA & B & C > notA & B & notC > notA & notB & C > notA & notB & notC

Experiments | Settings

Dataset

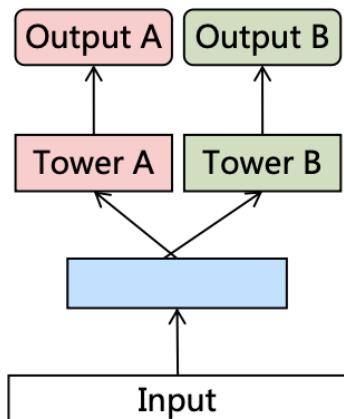
Table 1: Statistics of three datasets.

Datasets	#Samples	#Fields	#Features	Positive Ratio of Task A	Positive Ratio of Task B
TikTok Production	19622340	9	4691483	37.994%	1.101%
	9381820	10	447002	9.975%	1.510%

TikTok: Task A — Finish Watching Task B — Like

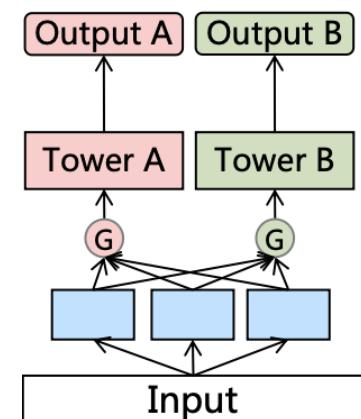
Production: Task A — pCTR Task B — Not Interested

Baselines

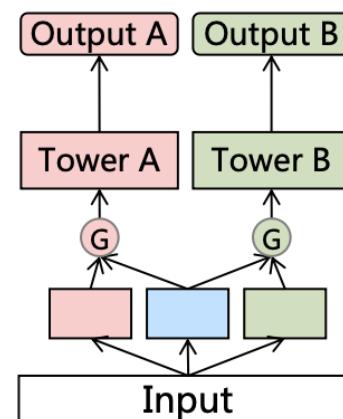


Shared-Bottom

$$\begin{bmatrix} \tilde{x}_1^i \\ \tilde{x}_2^i \end{bmatrix} = \begin{bmatrix} \alpha_{11} & \alpha_{12} \\ \alpha_{21} & \alpha_{22} \end{bmatrix} \begin{bmatrix} x_1^i \\ x_2^i \end{bmatrix}$$



Cross-Stitch



MMOE

PLE

Experiments | Settings

Metrics

AUC : Evaluating Bipartite Ranking Performance

$$\text{AUC} = \frac{1}{N^+ N^-} \sum_{\mathbf{x}_i \in D^+} \sum_{\mathbf{x}_j \in D^-} (\mathcal{I}(p(\mathbf{x}_i) > p(\mathbf{x}_j)))$$

Multi-AUC : Evaluating Multipartite Ranking Performance

$$\text{Multi-AUC} = \frac{2}{c(c-1)} \sum_{j=1}^c \sum_{k>j}^c p(j \cup k) \cdot AUC(k, j),$$

Experiments | Ablation Study

TABLE IV
ABLATION ANALYSIS FOR TASK A ON TIKTOK DATASET.

Variants	AUC	Multi-AUC
w/o AuxiliaryRank	0.7488 _(-0.0006)	0.6510 _(-0.0901)
w/o Calibration	0.7478 _(-0.0016)	0.7396 _(-0.0015)
w/o Correction	0.7486 _(-0.0008)	0.7399 _(-0.0012)
KD	0.7489 _(-0.0005)	0.6901 _(-0.0510)
Baseline	0.7494	0.7411

TABLE V
ABLATION ANALYSIS FOR TASK B ON TIKTOK DATASET.

Variants	AUC	Multi-AUC
w/o AuxiliaryRank	0.9501 _(-0.0012)	0.8005 _(-0.0336)
w/o Calibration	0.9504 _(-0.0009)	0.8312 _(-0.0029)
w/o Correction	0.9508 _(-0.0005)	0.8310 _(-0.0031)
KD	0.9505 _(-0.0008)	0.8014 _(-0.0327)
Baseline	0.9513	0.8341

w/o AuxiliaryRank : removing BPR losses for leaning auxiliary ranking relations

w/o Calibration : directly employing the teacher model outputs

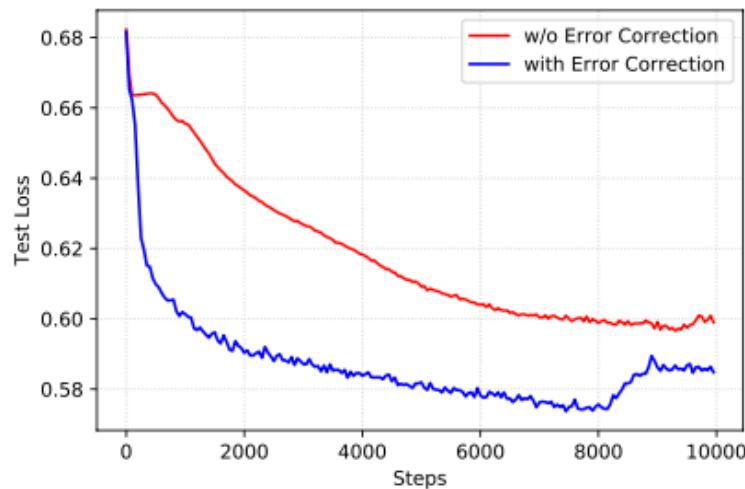
w/o Correction : not applying the error correction mechanism

KD : replacing teachers with regression-based models and using the vanilla knowledge distillation

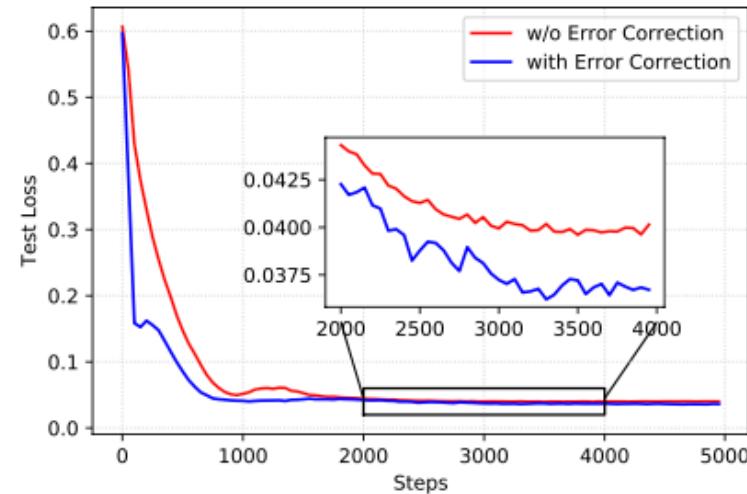
Baseline : our method

Experiments | Ablation Study

Q: Does **Error Correction Mechanism** Help to Accelerate Convergence and Enhance Knowledge Quality?



(a) Task A (Finish Watching).

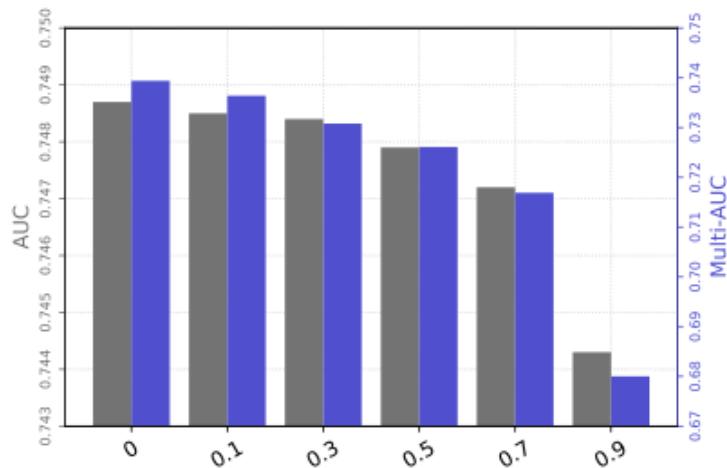


(b) Task B (Like).

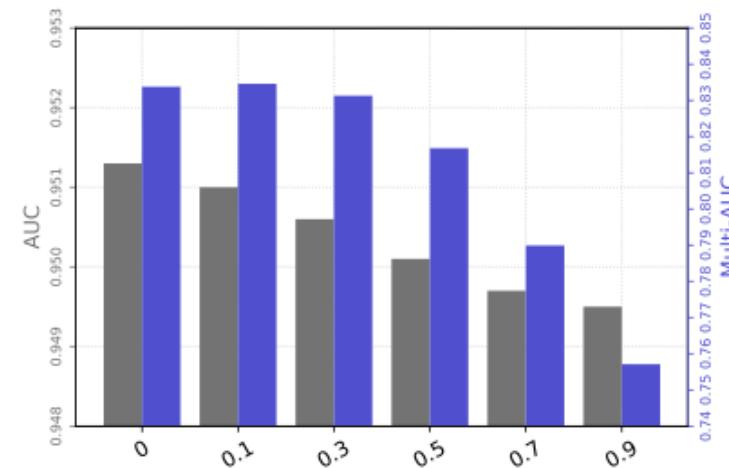
Fig. 4. Learning curves of TAKD with and without error correction mechanism on TikTok dataset.

Experiments | Ablation Study

Q: Is Student Really Benefiting from **Auxiliary Ranking Knowledge** from Other Tasks?



(a) Task A (Finish Watching).



(b) Task B (Like).

Fig. 5. Impact of corrupted auxiliary ranking information on student model performance for TikTok dataset.

Experiments | Hyper-Parameter Study

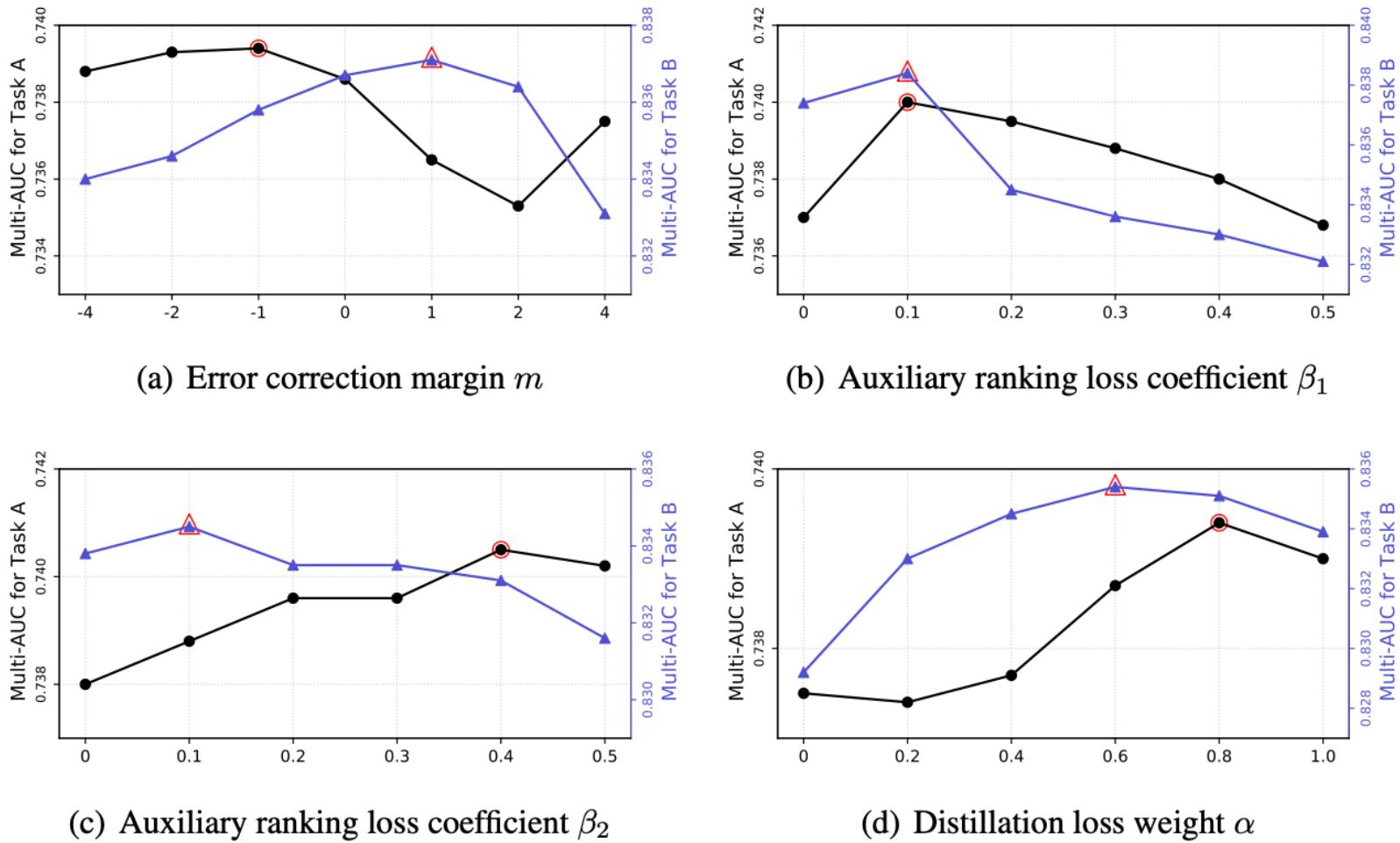


Figure 5.3 Multi-AUC performance on TikTok dataset for TaskA and Task B w.r.t. different hyper-parameters.

References

- [1] H. Tang, J. Liu, M. Zhao, and X. Gong, “Progressive layered extraction (ple): A novel multi-task learning (mtl) model for personalized recommendations,” in Fourteenth ACM Conference on Recommender Systems, 2020, pp. 269–278.
- [2] J. Ma, Z. Zhao, J. Chen, A. Li, L. Hong, and E. H. Chi, “Snr: Sub- network routing for flexible parameter sharing in multi-task learning,” in Proceedings of the AAAI Conference on Artificial Intelligence, vol. 33, no. 01, 2019, pp. 216–223.
- [3] C. Xu, Q. Li, J. Ge, J. Gao, X. Yang, C. Pei, F. Sun, J. Wu, H. Sun, and W. Ou, “Privileged features distillation at taobao recommendations,” in Proceedings of the 26th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining, 2020, pp. 2590–2598.
- [4] L. Yuan, F. E. Tay, G. Li, T. Wang, and J. Feng, “Revisit knowledge distillation: a teacher-free framework,” arXiv preprint arXiv:1909.11723, 2019.
- [5] S. Rendle, C. Freudenthaler, Z. Gantner, and L. Schmidt-Thieme, “Bpr: Bayesian personalized ranking from implicit feedback,” arXiv preprint arXiv:1205.2618, 2012.