# AdaptSSR: Pre-training User Model with Augmentation-Adaptive Self-Supervised Ranking

Yang Yu, Qi Liu, Kai Zhang, Yuren Zhang, Chao Song, Min Hou, Yuqing Yuan, Zhihao Ye, Zaixi Zhang, Sanshi Lei Yu

## Introduction

### Background

- User modeling, which aims to capture the user's characteristics or interests, is critical for many user-oriented tasks, such as user profiling, personalized recommendation, and click-through rate prediction.
- Most existing methods heavily rely on task-specific labeled data and suffer from the data sparsity problem.
- Several recent studies tackled this issue by pre-training the user model on massive unlabeled user behavior sequences with a contrastive learning task.
- They assume different views of the same user behavior sequence constructed via data augmentation are **semantically consistent** (reflecting similar characteristics or interests of the user), thus maximizing their agreement in the feature space.

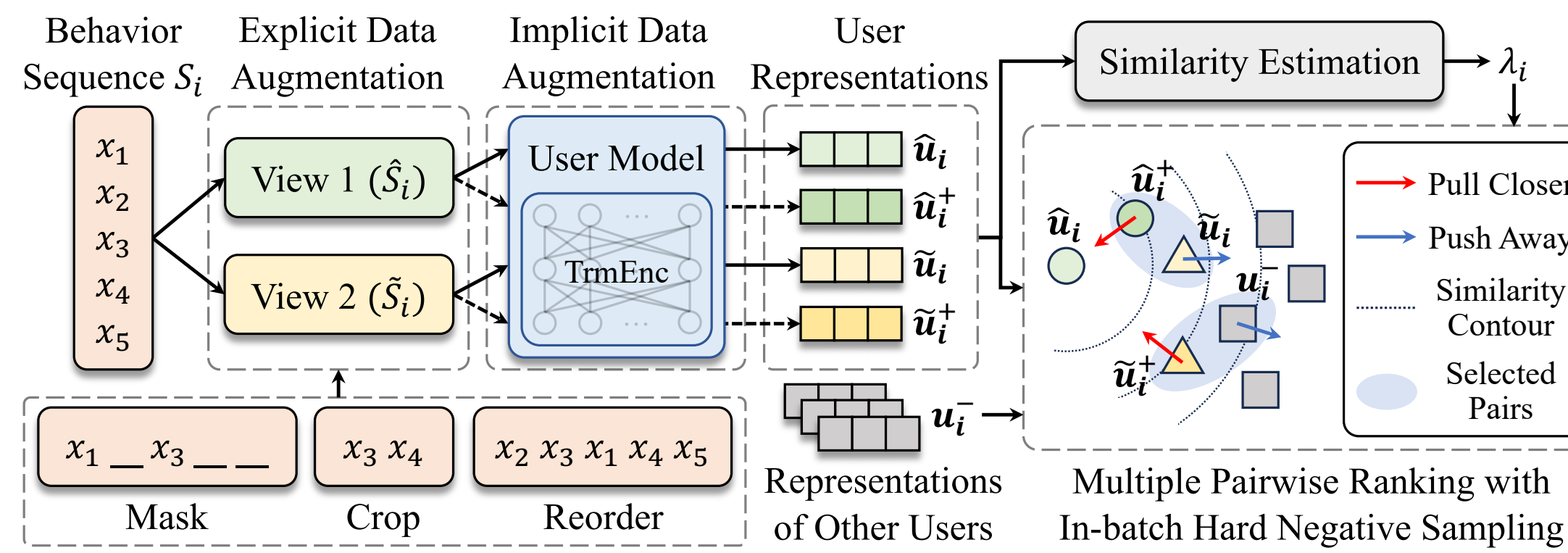### Challenge: The Semantic Inconsistency Problem



Each picture represents a news article clicked by the user. Pictures with the same color border reflect similar interests of the user. Dash borders are used to indicate behaviors replaced by the data augmentation method.

- Due to the diverse interests and heavy noise in user behaviors, existing data augmentation methods tend to lose certain characteristics of the user or introduce noisy interests that the user does not have.
- The impacts of data augmentation vary significantly across different user behavior sequences.
- Existing contrastive learning-based methods force the user model to maximize the agreement between the augmented views no matter whether they are similar or not, which may lead to a negative transfer for the downstream task.

## Methodology

### Main Idea: Self-Supervised Ranking

- Train the user model $\mathcal{M}$ to capture the similarity order between the implicitly augmented view, the explicitly augmented view, and views from other users.
- Given a user behavior sequence $S = \{x_1, x_2, ..., x_n\}$
  a) Input $S$ into $\mathcal{M}$ twice with different independently sampled dropout masks $\rightarrow \boldsymbol{u}, \boldsymbol{u}^+$ (implicit data augmentation)
  b) Input the augmented behavior sequence $\hat{S}$ into $\mathcal{M} \rightarrow \hat{\boldsymbol{u}}$ (explicit data augmentation)
  c) Input the behavior sequence of another user into $\mathcal{M} \rightarrow \boldsymbol{u}^-$
- Pre-training objective: $\text{sim}(\boldsymbol{u}, \boldsymbol{u}^+) \geq \text{sim}(\boldsymbol{u}, \hat{\boldsymbol{u}}) \geq \text{sim}(\boldsymbol{u}, \boldsymbol{u}^-)$



### Multiple Pairwise Ranking (MPR) with In-batch Hard Negative Sampling

- Given a batch of user behavior sequences $\{S_i\}_{i=1}^B$, apply two randomly selected explicit augmentation methods to each sequence $S_i \rightarrow \hat{S}_i$ and $\tilde{S}_i$
- Input $\hat{S}_i$ and $\tilde{S}_i$ into $\mathcal{M}$ twice $\rightarrow \hat{\boldsymbol{u}}_i, \hat{\boldsymbol{u}}_i^+$ and $\tilde{\boldsymbol{u}}_i, \tilde{\boldsymbol{u}}_i^+$
- MPR loss: extend the BPR loss to learn two pairwise ranking orders simultaneously.
- For the augmented sequence $\hat{S}_i$, the user representation $\hat{\boldsymbol{u}}_i, \hat{\boldsymbol{u}}_i^+$ and each $\boldsymbol{v} \in \{\tilde{\boldsymbol{u}}_i, \tilde{\boldsymbol{u}}_i^+\}$, $\boldsymbol{w} \in \mathbf{U}_i^- = \{\hat{\boldsymbol{u}}_j, \hat{\boldsymbol{u}}_j^+, \tilde{\boldsymbol{u}}_j, \tilde{\boldsymbol{u}}_j^+\}_{j=1, j\neq i}^B$ form a quadruple for model training.

$$\hat{\mathcal{L}}_i = -\frac{1}{2|\mathbf{U}_i^-|} \sum_{\boldsymbol{v} \in \{\tilde{\boldsymbol{u}}_i, \tilde{\boldsymbol{u}}_i^+\}} \sum_{\boldsymbol{w} \in \mathbf{U}_i^-} \log \sigma \left[ \lambda \left( \text{sim}(\hat{\boldsymbol{u}}_i, \hat{\boldsymbol{u}}_i^+) - \text{sim}(\hat{\boldsymbol{u}}_i, \boldsymbol{v}) \right) + (1 - \lambda) \left( \text{sim}(\hat{\boldsymbol{u}}_i, \boldsymbol{v}) - \text{sim}(\hat{\boldsymbol{u}}_i, \boldsymbol{w}) \right) \right]$$

- In-batch hard negative sampling: for each pairwise ranking order, select the pair with the smallest similarity difference to facilitate model training.

$$\hat{\mathcal{L}}_i = -\log \sigma \left[ \lambda \left( \text{sim}(\hat{\boldsymbol{u}}_i, \hat{\boldsymbol{u}}_i^+) - \max_{\boldsymbol{v} \in \{\tilde{\boldsymbol{u}}_i, \tilde{\boldsymbol{u}}_i^+\}} \text{sim}(\hat{\boldsymbol{u}}_i, \boldsymbol{v}) \right) + (1 - \lambda) \left( \min_{\boldsymbol{v} \in \{\tilde{\boldsymbol{u}}_i, \tilde{\boldsymbol{u}}_i^+\}} \text{sim}(\hat{\boldsymbol{u}}_i, \boldsymbol{v}) - \max_{\boldsymbol{w} \in \mathbf{U}_i^-} \text{sim}(\hat{\boldsymbol{u}}_i, \boldsymbol{w}) \right) \right]$$

- The loss function $\tilde{\mathcal{L}}_i$ for another augmented sequence $\tilde{S}_i$ is symmetrically defined and the overall loss is computed as $\mathcal{L} = \sum_{i=1}^B (\hat{\mathcal{L}}_i + \tilde{\mathcal{L}}_i)/2B$.

### Augmentation-Adaptive Fusion

- The effects of data augmentation vary significantly across different behavior sequences.
- The constant hyper-parameter $\lambda$ applies a fixed and unified constraint to all samples.
- Replace $\lambda$ with a dynamic coefficient $\lambda_i$, which is estimated based on the average similarity between the user representations generated from $\hat{S}_i$ and $\tilde{S}_i$.

$$\lambda_i = 1 - \frac{1}{4} \sum_{\hat{s} \in \{\hat{\boldsymbol{u}}_i, \hat{\boldsymbol{u}}_i^+\}} \sum_{\tilde{s} \in \{\tilde{\boldsymbol{u}}_i, \tilde{\boldsymbol{u}}_i^+\}} \max(\text{sim}(\hat{s}, \tilde{s}), 0)$$

- If $\hat{S}_i$ and $\tilde{S}_i$ are semantically similar, $\lambda_i$ will be small and $\hat{\mathcal{L}}_i$ will focus on maximizing the latter term $\min_{\boldsymbol{v} \in \{\tilde{\boldsymbol{u}}_i, \tilde{\boldsymbol{u}}_i^+\}} \text{sim}(\hat{\boldsymbol{u}}_i, \boldsymbol{v}) - \max_{\boldsymbol{w} \in \mathbf{U}_i^-} \text{sim}(\hat{\boldsymbol{u}}_i, \boldsymbol{w})$, which forces the user model to discriminate these similar explicitly augmented views from views of other users.
- Otherwise, $\lambda_i$ will be large and train the user model to pull the implicitly augmented view and these dissimilar explicitly augmented views apart.

## Experiments

### Datasets and Downstream Tasks

- TTL: users' recent 100 interactions on the QQ Browser platform.
- App: users' app installation behaviors on OPPO smartphones from 2022-12 to 2023-03.
- $\mathcal{T}_1$: age prediction
- $\mathcal{T}_2$: life status prediction
- $\mathcal{T}_3$: click recommendation
- $\mathcal{T}_4$: thumb-up recommendation
- $\mathcal{T}_5$: gender prediction
- $\mathcal{T}_6$: CVR prediction

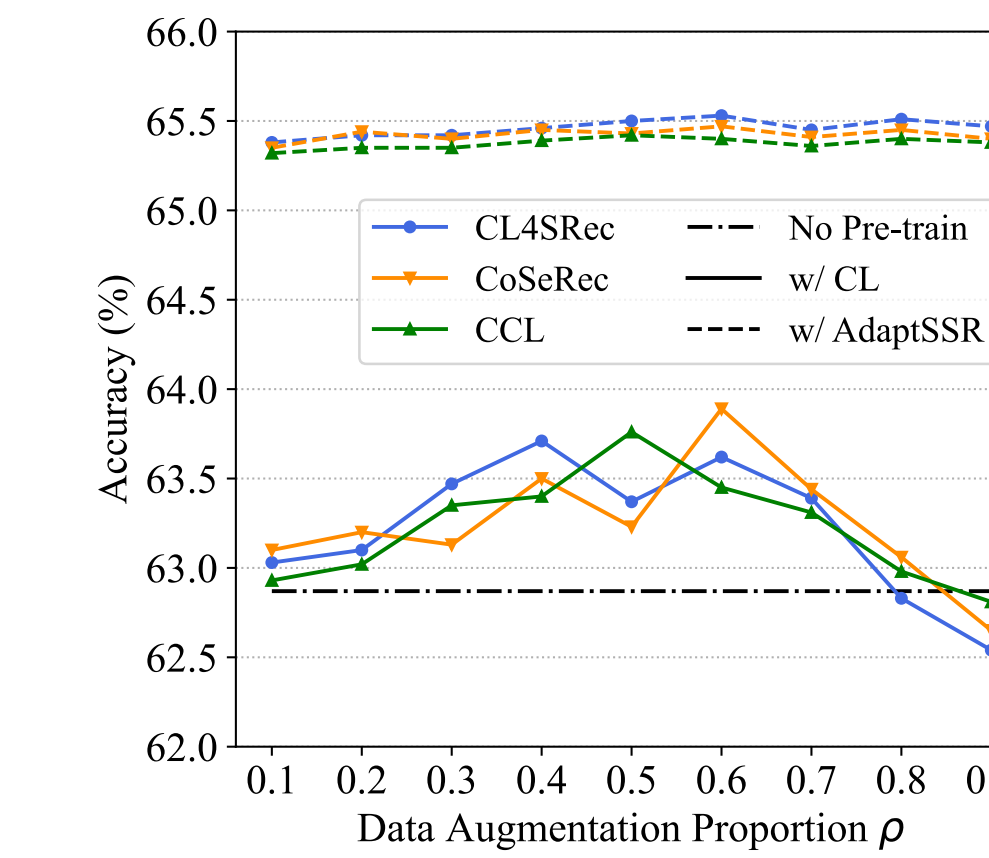| Dataset | TTL | | | | App | |
|---|---|---|---|---|---|---|
| # Behavior Sequences | 1,470,149 | | | | 1,575,837 | |
| # Different Behaviors | 645,972 | | | | 4,047 | |
| Avg. Sequence Length | 54.84 | | | | 44.13 | |
| Downstream Task | $\mathcal{T}_1$ | $\mathcal{T}_2$ | $\mathcal{T}_3$ | $\mathcal{T}_4$ | $\mathcal{T}_5$ | $\mathcal{T}_6$ |
| # Samples | 1,470,147 | 1,020,277 | 1,397,197 | 255,646 | 1,178,603 | 564,940 |
| # Labels/Items | 8 | 6 | 17,879 | 7,539 | 2 | 2 |

Detailed statistics of each dataset and downstream task.
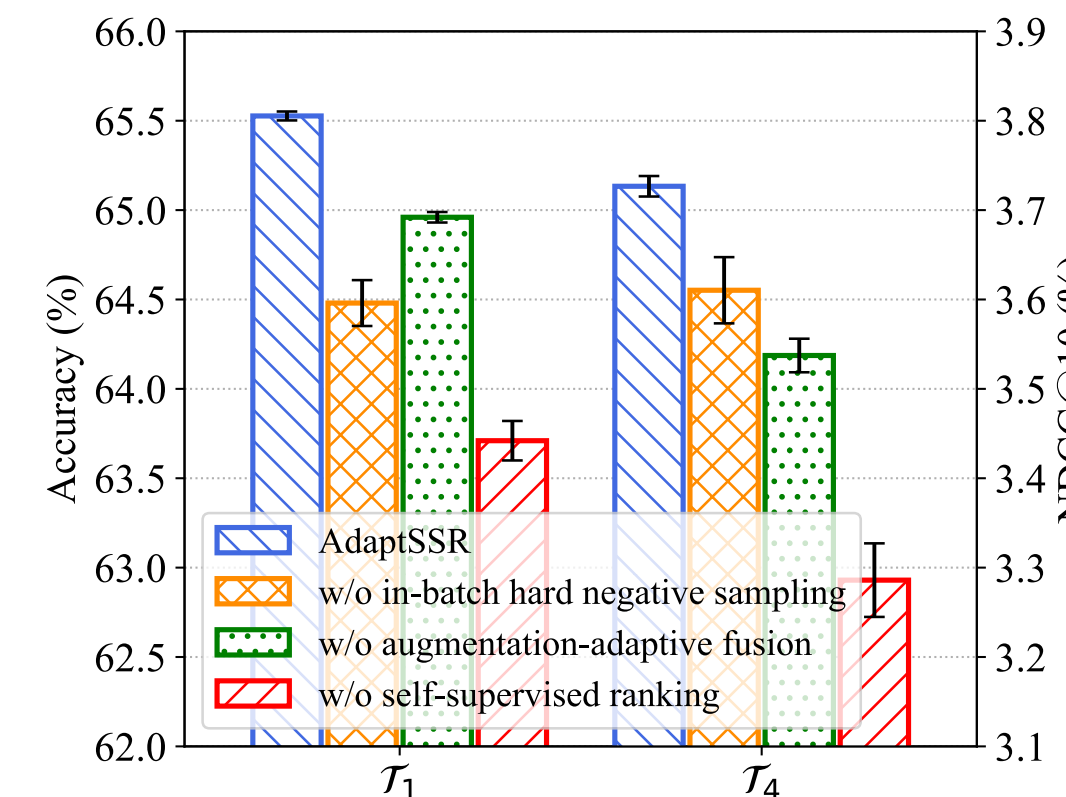
### Performance Comparison

Impr (%) indicates the relative improvement compared with the end-to-end training.

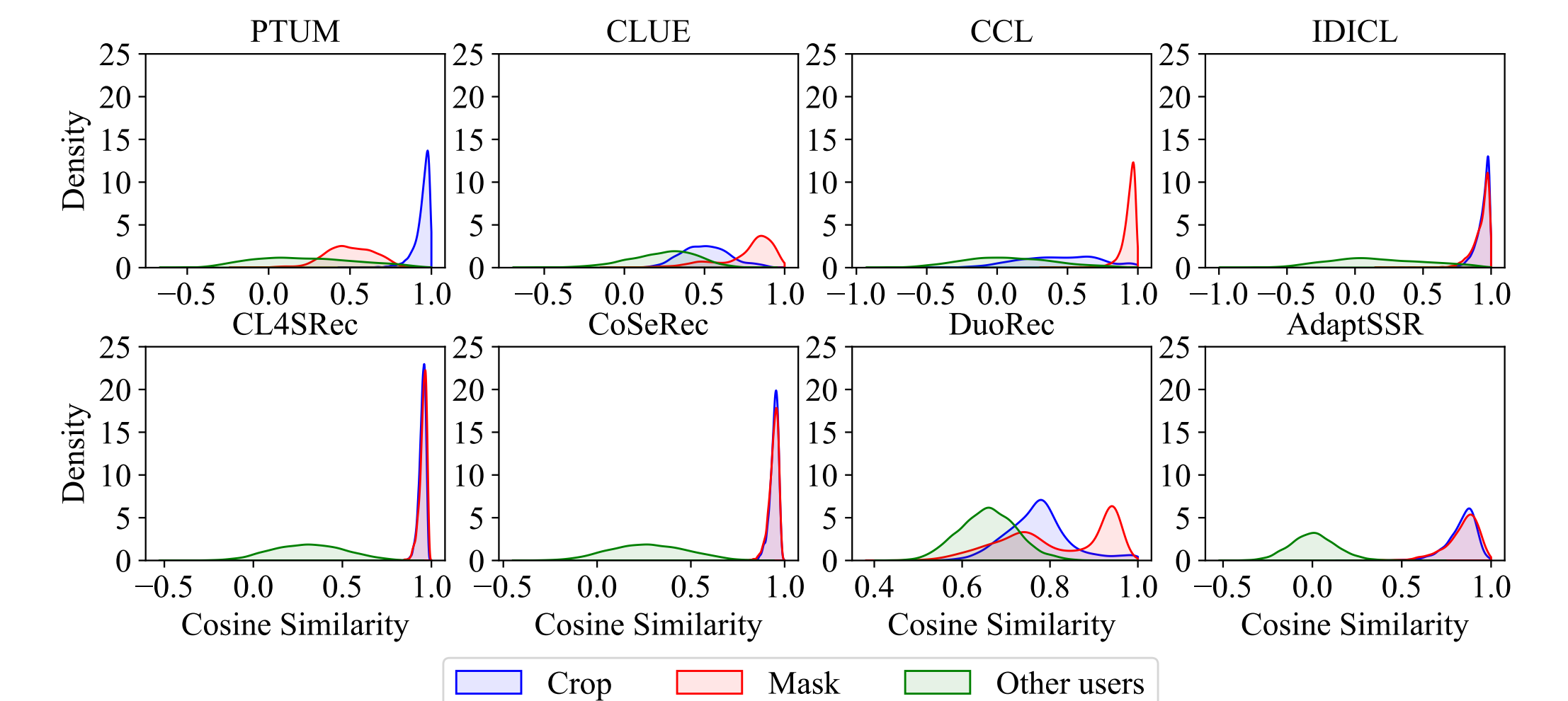| Pre-train Method | $\mathcal{T}_1$ | | $\mathcal{T}_2$ | | $\mathcal{T}_3$ | | $\mathcal{T}_4$ | | $\mathcal{T}_5$ | | $\mathcal{T}_6$ | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | Acc | Impr | Acc | Impr | NDCG@10 | Impr | NDCG@10 | Impr | AUC | Impr | AUC | Impr |
| None | 62.87±0.05 | - | 52.24±0.16 | - | 1.99±0.03 | - | 2.87±0.07 | - | 78.63±0.06 | - | 75.14±0.14 | - |
| PeterRec | 63.62±0.11 | 1.19 | 53.14±0.07 | 1.72 | 2.37±0.02 | 19.10 | 3.06±0.08 | 6.62 | 79.61±0.13 | 1.25 | 76.04±0.10 | 1.20 |
| PTUM | 63.21±0.14 | 0.54 | 53.05±0.04 | 1.55 | 2.29±0.03 | 15.08 | 2.96±0.03 | 3.14 | 79.48±0.11 | 1.08 | 75.82±0.13 | 0.90 |
| CLUE | 63.38±0.10 | 0.81 | 53.23±0.05 | 1.90 | 2.38±0.02 | 19.60 | 3.05±0.01 | 6.27 | 79.90±0.06 | 1.62 | 76.03±0.16 | 1.18 |
| CCL | 63.76±0.11 | 1.42 | 53.37±0.09 | 2.16 | 2.43±0.02 | 22.11 | 3.32±0.13 | 15.68 | 80.22±0.07 | 2.02 | 77.35±0.10 | 2.94 |
| IDICL | 63.88±0.04 | 1.61 | 53.45±0.05 | 2.32 | 2.46±0.02 | 23.32 | 3.42±0.04 | 19.16 | 80.34±0.05 | 2.17 | 77.92±0.08 | 3.70 |
| CL4SRec | 63.71±0.11 | 1.34 | 53.43±0.05 | 2.28 | 2.41±0.03 | 21.11 | 3.29±0.06 | 14.63 | 80.14±0.08 | 1.92 | 77.02±0.05 | 2.50 |
| CoSeRec | 63.89±0.03 | 1.62 | 53.53±0.09 | 2.47 | 2.44±0.02 | 22.61 | 3.33±0.16 | 16.03 | 80.48±0.06 | 2.35 | 77.71±0.09 | 3.42 |
| DuoRec | 63.50±0.09 | 1.00 | 53.26±0.06 | 1.95 | 2.39±0.01 | 20.10 | 3.11±0.16 | 8.36 | 80.03±0.09 | 1.78 | 76.85±0.09 | 2.28 |
| **AdaptSSR** | **65.53±0.04** | **4.23** | **54.41±0.02** | **4.15** | **2.61±0.03** | **31.16** | **3.73±0.03** | **29.97** | **82.30±0.03** | **4.67** | **79.92±0.05** | **6.36** |

### Performance with Different Augmentation Methods



### Ablation Study



### User Representation Similarity Distribution Analysis



Distributions of the cosine similarity between user representations generated from the original behavior sequence, different augmented behavior sequences, and behavior sequences of other users with various pre-training methods.