



# Untargeted Attack Against Federated Recommendation Systems via Poisonous Item Embeddings and the Defense

**Yang Yu<sup>1,2</sup>, Qi Liu<sup>1,2\*</sup>, Likang Wu<sup>1,2</sup>, Runlong Yu<sup>1,2</sup>, Sanshi Lei Yu<sup>1,2</sup>, Zaixi Zhang<sup>1,2</sup>**

<sup>1</sup>Anhui Province Key Laboratory of Big Data Analysis and Application,  
School of Computer Science and Technology, University of Science and Technology of China

<sup>2</sup>State Key Laboratory of Cognitive Intelligence

- Recommender systems are widely used to alleviate the information overload problem.
- Most existing recommender systems are trained on **centralized user data**.
  - Risk of data leakage.
  - Privacy concerns.
- Privacy regulations (e.g., GDPR, CCPA) make it more difficult to collect user data for centralized model training.



- **Federated learning (FL)** enables multiple clients to collaboratively learn a global model without sharing their local data.
- Several studies have applied FL to train privacy-preserving **federated recommendation (FedRec)** systems.
- Unfortunately, FL is known to be vulnerable to **poisoning attacks**.
  - **Targeted Attack**
    - Increase the exposure rate of certain target items.
  - **Untargeted Attack**
    - Degrade the overall performance of the FedRec system.
    - Also known as the denial-of-service attack.
    - Continuously disrupt the user experience → Severe losses of customers and revenue.

## □ Challenges

- The attack method must be effective even with a small fraction of malicious clients.
- The attacker can only access a small set of data stored on the malicious clients.
- The attack needs to manipulate the model output on arbitrary inputs.
- Many recommenders are naturally robust to malicious perturbation to a certain degree.

## □ In this work

- **ClusterAttack**: an effective and covert untargeted model poisoning attack method.
- **UNION**: a general uniformity-based defense mechanism.

## □ Federated Recommendation Systems

□ The parameters of the recommendation model  $\Theta = [\Theta_{\text{item}}; \Theta_{\text{user}}; \Theta_{\text{pred}}]$ .

□ Standard FL procedure

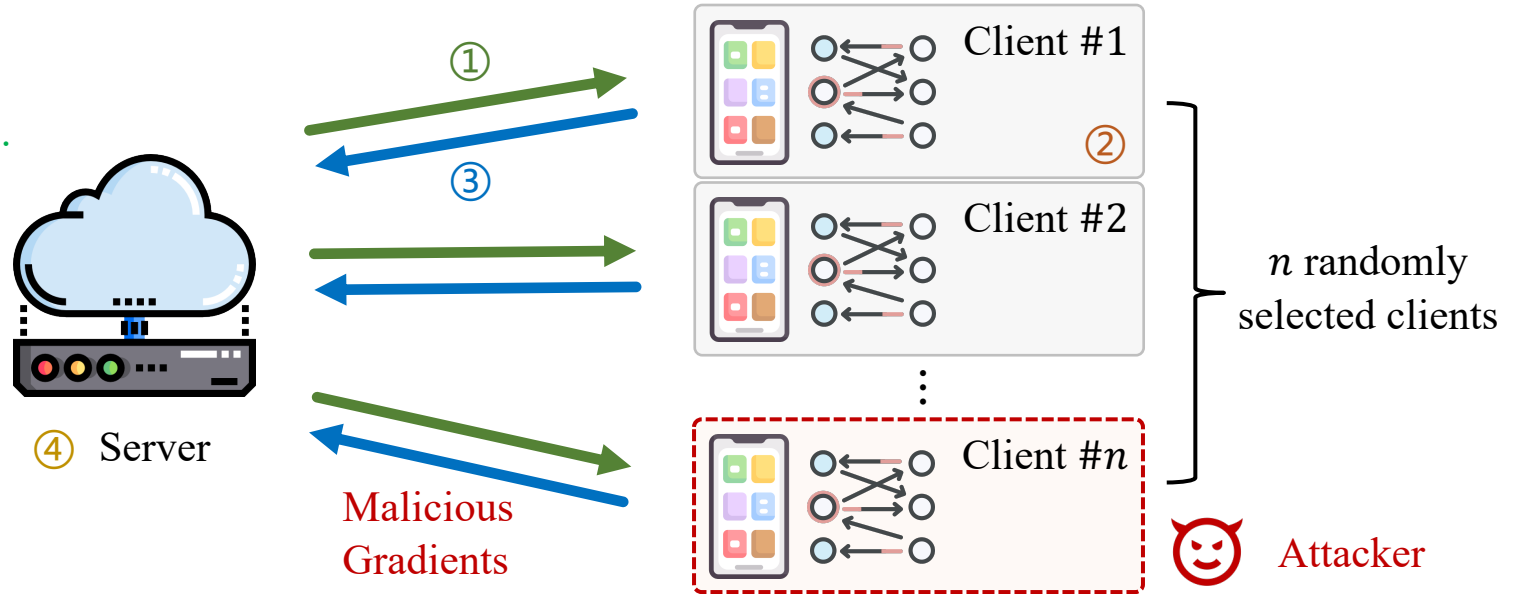
① Distribute global model  $[\Theta_{\text{item}}; \Theta_{\text{pred}}]$ .

② Compute local gradients

$$\mathbf{g} = [\mathbf{g}_{\text{item}}; \mathbf{g}_{\text{user}}; \mathbf{g}_{\text{pred}}]$$

③ Upload  $[\mathbf{g}_{\text{item}}; \mathbf{g}_{\text{pred}}]$  and update local  $\Theta_{\text{user}}$  with  $\mathbf{g}_{\text{user}}$ .

④ Aggregate and update global model.



## □ Threat Model

□  $m\%$  (typically small, e.g., 1%) of clients are controlled by the attacker.

□ The attacker knows the training code, local model, and user data of malicious clients.

□ The attacker cannot access the data or gradients of other benign clients.

## ClusterAttack

- The recommendation model generally predicts the ranking score based on the user embedding and the item embedding.
- Upload malicious gradients that converge item embeddings into several dense clusters.

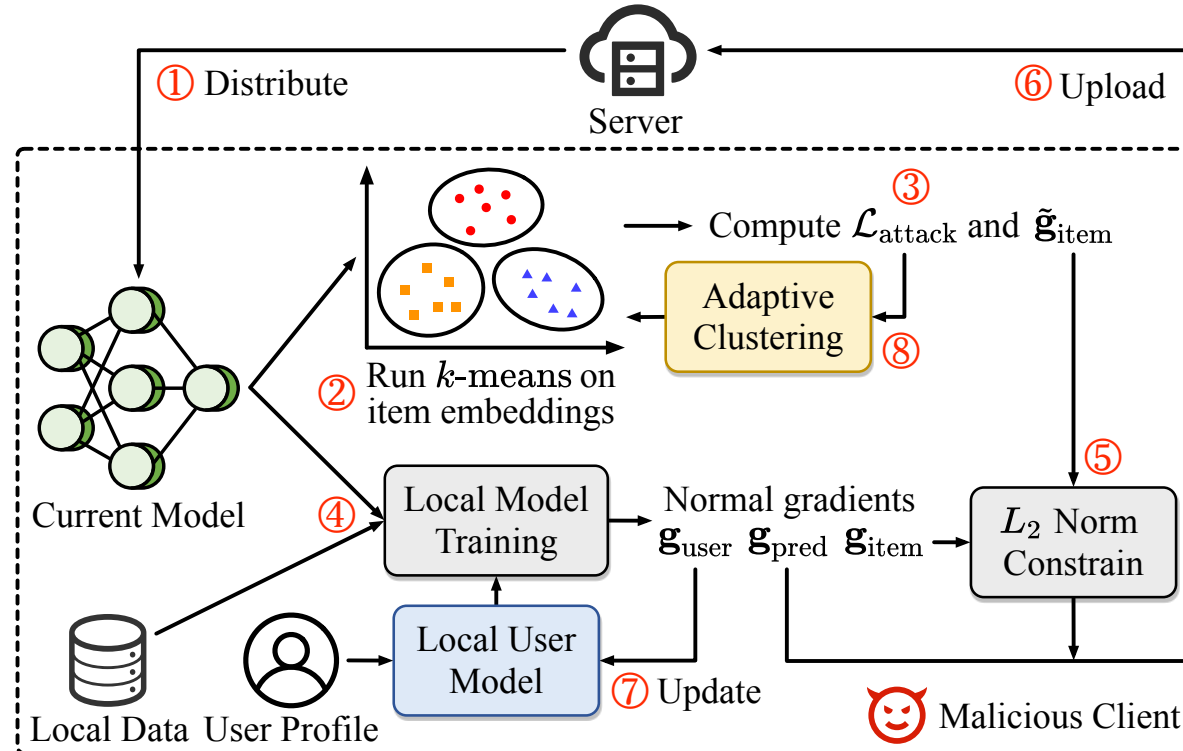


Figure 1: The procedure of ClusterAttack.

## ClusterAttack

- Apply  $k$ -means to split the item embeddings  $\{v_i\}_{i=1}^M$  into  $K$  clusters  $\{C_i\}_{i=1}^K$  with centroids  $\{c_i\}_{i=1}^K$ .
- Compute the within-cluster variance and the malicious gradient.

$$\mathcal{L}_{attack} = \sum_{i=1}^K \sum_{v_j \in C_i} \|v_j - c_i\|_2^2 \quad \tilde{\mathbf{g}}_{v_i} = \frac{\partial \mathcal{L}_{attack}}{\partial v_i}$$

## Gradient clipping

- Compute the normal gradient of each malicious client.
- Calculate the mean  $\mu$  and standard deviation  $\sigma$  of the  $L_2$  norms of all normal item embedding gradients.

$$b_i^{(j)} = \mu + \lambda_i^{(j)} \sigma, \lambda_i^{(j)} \in [0,3] \quad \hat{\mathbf{g}}_{v_i}^{(j)} = \frac{\tilde{\mathbf{g}}_{v_i}}{\max(1, \|\tilde{\mathbf{g}}_{v_i}\|_2 / b_i^{(j)})}$$

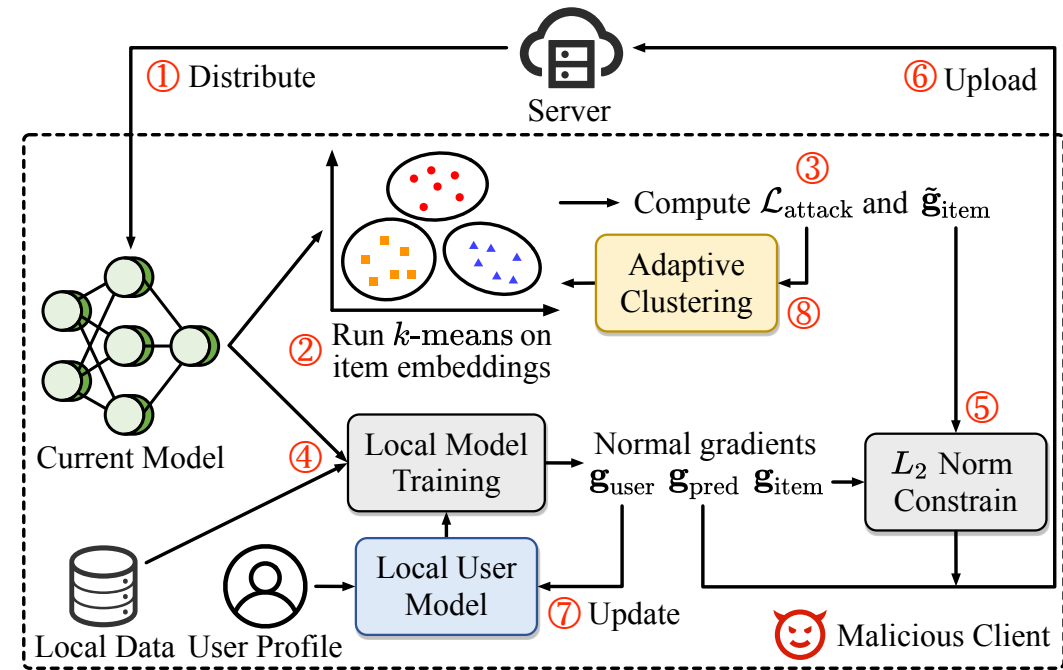


Figure 1: The procedure of ClusterAttack.

## ClusterAttack

### Adaptive clustering

- Adjust the number of clusters  $K$  after each round of attack.
- Use the change of  $\mathcal{L}_{attack}$  as the feedback.
- $\mathcal{L}_{attack}$  keeps increasing  $\rightarrow \mathcal{L}_{attack}$  cannot converge well.
- $\mathcal{L}_{attack}$  keeps decreasing  $\rightarrow$  decrease  $K$  for stronger attack.

---

### Algorithm 1: Adaptive Clustering

---

**Input:** Number of clusters  $K$ , range of number of clusters  $[K_{\min}, K_{\max}]$ , threshold  $R$ , and decay rate  $\beta$ .

**Init:** Set  $\tilde{\mathcal{L}}_{attack}^{(0)}$ ,  $n_{inc}$ ,  $n_{dec}$  and  $t$  as 0.  
// Repeat after each round of attack

- $t \leftarrow t + 1$ ;
- Calculate  $\mathcal{L}_{attack}^{(t)}$  with Equation (2);
- $\tilde{\mathcal{L}}_{attack}^{(t)} \leftarrow \beta \cdot \tilde{\mathcal{L}}_{attack}^{(t-1)} + (1 - \beta) \cdot \mathcal{L}_{attack}^{(t)}$ ;
- $\hat{\mathcal{L}}_{attack}^{(t)} \leftarrow \tilde{\mathcal{L}}_{attack}^{(t)} / (1 - \beta^t)$ ;
- if**  $\hat{\mathcal{L}}_{attack}^{(t)} > \hat{\mathcal{L}}_{attack}^{(t-1)}$  **then**  $n_{inc} \leftarrow n_{inc} + 1$ ;
- else**  $n_{dec} \leftarrow n_{dec} + 1$ ;
- if**  $n_{inc} - n_{dec} \geq R$  **then**
  - $K \leftarrow \min(\lfloor K + \sqrt{K_{\max} - K} \rfloor, K_{\max})$ ;
  - Reset  $n_{inc}$ ,  $n_{dec}$  and  $t$  as 0;
- end if**
- if**  $n_{dec} - n_{inc} \geq R$  **then**
  - $K \leftarrow \max(\lfloor K - \sqrt{K - K_{\min}} \rfloor, K_{\min})$ ;
  - Reset  $n_{inc}$ ,  $n_{dec}$  and  $t$  as 0;
- end if**

---



## □ UNION Mechanism

### □ Client Side

- Train the local recommendation model with an **additional contrastive learning task**.
- Denote the item set interacted by the user as  $\mathcal{V}_u = \{v_i\}_{i=1}^L$  and the entire item set as  $\mathcal{V}$ .
- For each  $v_i \in \mathcal{V}_u$ , randomly select another positive item  $v_i^+ \in \mathcal{V}_u$  and  $P$  negative items  $\{v_i^-\}_{i=1}^P \subseteq \mathcal{V} \setminus \mathcal{V}_u$ .

$$\mathcal{L}_{\text{cl}} = - \sum_{i=1}^L \log \frac{e^{f(v_i)^T f(v_i^+)}}{e^{f(v_i)^T f(v_i^+)} + \sum_{j=1}^P e^{f(v_i)^T f(v_i^-)}} \quad \mathcal{L} = \mathcal{L}_{\text{rec}} + \alpha \mathcal{L}_{\text{cl}}$$

- $\mathcal{L}_{\text{cl}}$  can regularize the item embeddings toward a uniform distribution in the space [1] while training with the recommendation task (opposite to the goal of ClusterAttack).

## □ UNION Mechanism

### □ Server Side

- Estimate the uniformity of updated item embeddings for each received gradient.

$$d_i = \mathbb{E}_{x, y \sim p_{\text{data}}^{\text{i.i.d}}} \|f(x) - f(y)\|_2^2$$

- Use the **Gap Statistics algorithm** [2] to estimate the number of clusters in  $\{d_i\}_{i=1}^n$ .
- If the algorithm estimates that there is more than one cluster, we apply  $k$ -means to split  $\{d_i\}_{i=1}^n$  into two clusters and remove all the gradients belonging to the minor one.

### □ Note

- UNION is a general mechanism that aims to preserve the distribution of item embeddings.
- It can be combined with existing Byzantine-robust FL methods (e.g., MultiKrum, NormBound) to provide more comprehensive protection for FedRec systems.

[2] Tibshirani et al. Estimating the Number of Clusters in a Data Set via the Gap Statistic. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*. 2001.

## □ Datasets

- MovieLens-1M
- Gowalla

## □ Base Recommendation Model

- MF
- SASRec

## □ Metrics

- Hit Ratio (HR)
- Normalized Discounted Cumulative Gain (NDCG)
- Only calculated on benign clients using the all-ranking protocol.

Dataset	#Users	#Items	#Actions	Avg. length	Density
ML-1M	6,040	3,706	1,000,209	165.6	4.47%
Gowalla	29,858	40,981	1,585,043	53.1	0.13%

Table 2: Detailed statistics of the two datasets.

## □ Attack Performance with No Defense

Model	Attack Method	ML-1M		Gowalla	
		HR@5	NDCG@5	HR@5	NDCG@5
MF	No Attack	0.03549 (-)	0.02226(-)	0.02523 (-)	0.01697 (-)
	LabelFlip	0.03561 (-0.34%)	0.02238 (-0.54%)	0.02541 (-0.71%)	0.01711 (-0.82%)
	FedAttack	0.03358 (5.38%)	0.02118 (4.85%)	0.02371 (6.02%)	0.01585 (6.60%)
	Gaussian	0.03555 (-0.17%)	0.02224 (0.09%)	0.02528 (-0.20%)	0.01701 (-0.24%)
	LIE	0.03259 (8.17%)	0.02062 (7.37%)	0.02316 (8.20%)	0.01571 (7.42%)
	Fang	0.03038 (14.40%)	0.01897 (14.78%)	0.02131 (15.54%)	0.01448 (14.67%)
	ClusterAttack	<b>0.02451 (30.94%)</b>	<b>0.01545 (30.59%)</b>	<b>0.01664 (34.05%)</b>	<b>0.01117 (34.18%)</b>
SASRec	No Attack	0.10810 (-)	0.07053 (-)	0.03251 (-)	0.02217 (-)
	LabelFlip	0.10857 (-0.43%)	0.07071 (-0.26%)	0.03270 (-0.58%)	0.02222 (-0.23%)
	FedAttack	0.10013 (7.37%)	0.06572 (6.82%)	0.03054 (6.06%)	0.02087 (5.86%)
	Gaussian	0.10769 (0.38%)	0.07055 (-0.03%)	0.03226 (0.77%)	0.02222 (-0.23%)
	LIE	0.09677 (10.48%)	0.06281 (10.95%)	0.03008 (7.47%)	0.02021 (8.84%)
	Fang	0.08964 (17.08%)	0.05909 (16.22%)	0.02797 (13.96%)	0.01883 (15.07%)
	ClusterAttack	<b>0.06547 (39.44%)</b>	<b>0.04130 (41.44%)</b>	<b>0.02223 (31.62%)</b>	<b>0.01544 (30.36%)</b>

Table 1: Model performance under different untargeted attack methods with no defense. The percentages in parentheses indicate the relative performance degradation compared with the no-attack scenario.

## □ Attack Performance with No Defense

Model	Attack Method	ML-1M		Gowalla	
		HR@5	NDCG@5	HR@5	NDCG@5
MF	No Attack	0.03549 (-)	0.02226(-)	0.02523 (-)	0.01697 (-)
	LabelFlip	0.03561 (-0.34%)	0.02238 (-0.54%)	0.02541 (-0.71%)	0.01711 (-0.82%)
	FedAttack	0.03358 (5.38%)	0.02118 (4.85%)	0.02371 (6.02%)	0.01585 (6.60%)
	Gaussian	0.03555 (-0.17%)	0.02224 (0.09%)	0.02528 (-0.20%)	0.01701 (-0.24%)
	LIE	0.03259 (8.17%)	0.02062 (7.37%)	0.02316 (8.20%)	0.01571 (7.42%)
	Fang	0.03038 (14.40%)	0.01897 (14.78%)	0.02131 (15.54%)	0.01448 (14.67%)
	ClusterAttack	<b>0.02451 (30.94%)</b>	<b>0.01545 (30.59%)</b>	<b>0.01664 (34.05%)</b>	<b>0.01117 (34.18%)</b>
SASRec	No Attack	0.10810 (-)	0.07053 (-)	0.03251 (-)	0.02217 (-)
	LabelFlip	0.10857 (-0.43%)	0.07071 (-0.26%)	0.03270 (-0.58%)	0.02222 (-0.23%)
	FedAttack	0.10013 (7.37%)	0.06572 (6.82%)	0.03054 (6.06%)	0.02087 (5.86%)
	Gaussian	0.10769 (0.38%)	0.07055 (-0.03%)	0.03226 (0.77%)	0.02222 (-0.23%)
	LIE	0.09677 (10.48%)	0.06281 (10.95%)	0.03008 (7.47%)	0.02021 (8.84%)
	Fang	0.08964 (17.08%)	0.05909 (16.22%)	0.02797 (13.96%)	0.01883 (15.07%)
	ClusterAttack	<b>0.06547 (39.44%)</b>	<b>0.04130 (41.44%)</b>	<b>0.02223 (31.62%)</b>	<b>0.01544 (30.36%)</b>

Table 1: Model performance under different untargeted attack methods with no defense. The percentages in parentheses indicate the relative performance degradation compared with the no-attack scenario.

## □ Attack Performance with No Defense

Model	Attack Method	ML-1M		Gowalla	
		HR@5	NDCG@5	HR@5	NDCG@5
MF	No Attack	0.03549 (-)	0.02226(-)	0.02523 (-)	0.01697 (-)
	LabelFlip	0.03561 (-0.34%)	0.02238 (-0.54%)	0.02541 (-0.71%)	0.01711 (-0.82%)
	FedAttack	0.03358 (5.38%)	0.02118 (4.85%)	0.02371 (6.02%)	0.01585 (6.60%)
	Gaussian	0.03555 (-0.17%)	0.02224 (0.09%)	0.02528 (-0.20%)	0.01701 (-0.24%)
	LIE	0.03259 (8.17%)	0.02062 (7.37%)	0.02316 (8.20%)	0.01571 (7.42%)
	Fang	0.03038 (14.40%)	0.01897 (14.78%)	0.02131 (15.54%)	0.01448 (14.67%)
	<b>ClusterAttack</b>	<b>0.02451 (30.94%)</b>	<b>0.01545 (30.59%)</b>	<b>0.01664 (34.05%)</b>	<b>0.01117 (34.18%)</b>
SASRec	No Attack	0.10810 (-)	0.07053 (-)	0.03251 (-)	0.02217 (-)
	LabelFlip	0.10857 (-0.43%)	0.07071 (-0.26%)	0.03270 (-0.58%)	0.02222 (-0.23%)
	FedAttack	0.10013 (7.37%)	0.06572 (6.82%)	0.03054 (6.06%)	0.02087 (5.86%)
	Gaussian	0.10769 (0.38%)	0.07055 (-0.03%)	0.03226 (0.77%)	0.02222 (-0.23%)
	LIE	0.09677 (10.48%)	0.06281 (10.95%)	0.03008 (7.47%)	0.02021 (8.84%)
	Fang	0.08964 (17.08%)	0.05909 (16.22%)	0.02797 (13.96%)	0.01883 (15.07%)
	<b>ClusterAttack</b>	<b>0.06547 (39.44%)</b>	<b>0.04130 (41.44%)</b>	<b>0.02223 (31.62%)</b>	<b>0.01544 (30.36%)</b>

Table 1: Model performance under different untargeted attack methods with no defense. The percentages in parentheses indicate the relative performance degradation compared with the no-attack scenario.

## Attack Performance under Defense (Left five groups)

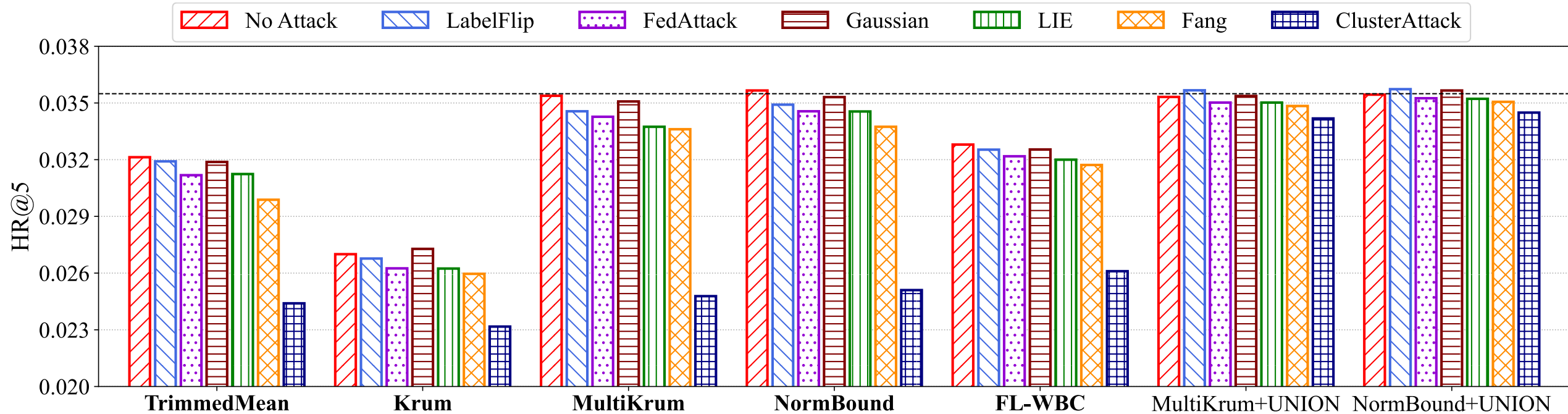


Figure 2: Model performance under different untargeted attack methods with different defense mechanisms. The black dashed line represents the model performance without any attack or defense.

## Attack Performance under Defense (Left five groups)

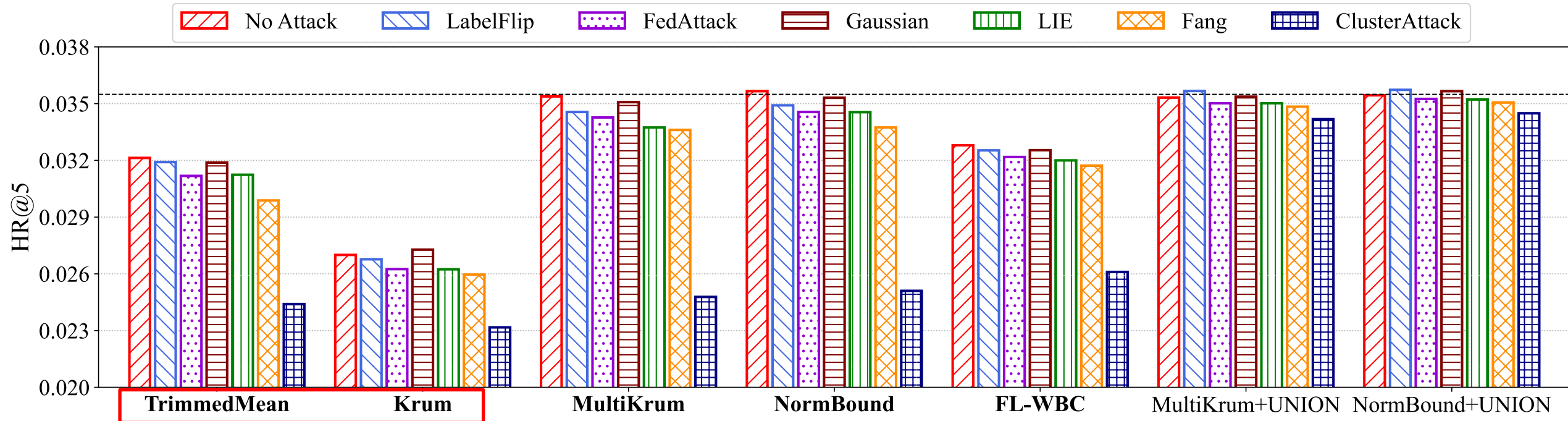


Figure 2: Model performance under different untargeted attack methods with different defense mechanisms. The black dashed line represents the model performance without any attack or defense.



## Attack Performance under Defense (Left five groups)

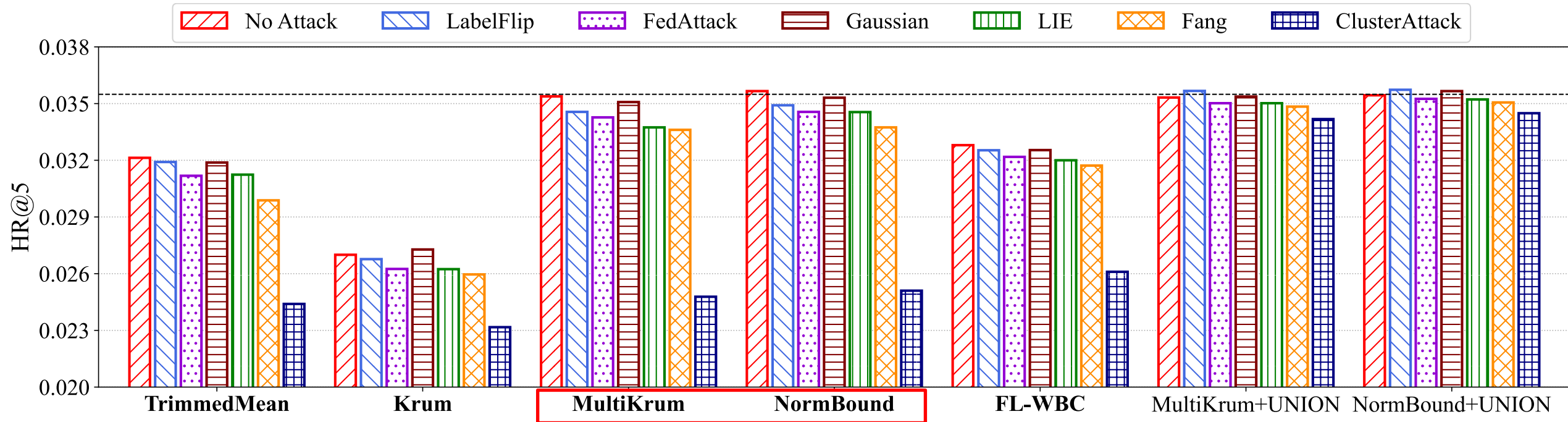


Figure 2: Model performance under different untargeted attack methods with different defense mechanisms. The black dashed line represents the model performance without any attack or defense.

## Defense Performance (Right two groups)

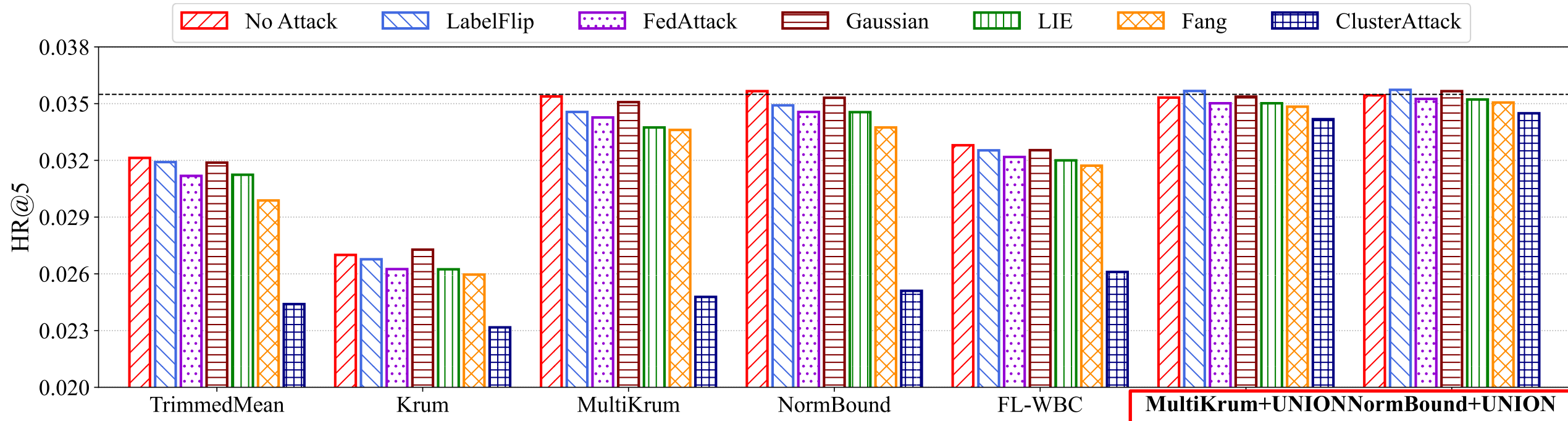


Figure 2: Model performance under different untargeted attack methods with different defense mechanisms. The black dashed line represents the model performance without any attack or defense.

## □ Can the attacker evade UNION?

□  $\mathcal{L}'_{attack} = \mathcal{L}_{attack} + \alpha \cdot \mathcal{L}_{cl}$

□ The extra contrastive learning task weakens the attack effect of ClusterAttack.

Defense Method	Attack Method	HR@5
MultiKrum+UNION	ClusterAttack	0.03378 (4.82%)
	ClusterAttack+CL	0.03525 (0.68%)
NormBound+UNION	ClusterAttack	0.03449 (2.82%)
	ClusterAttack+CL	0.03566 (-0.48%)

Table 3: Attack performance of ClusterAttack+CL.

## □ Impact of Adaptive Clustering

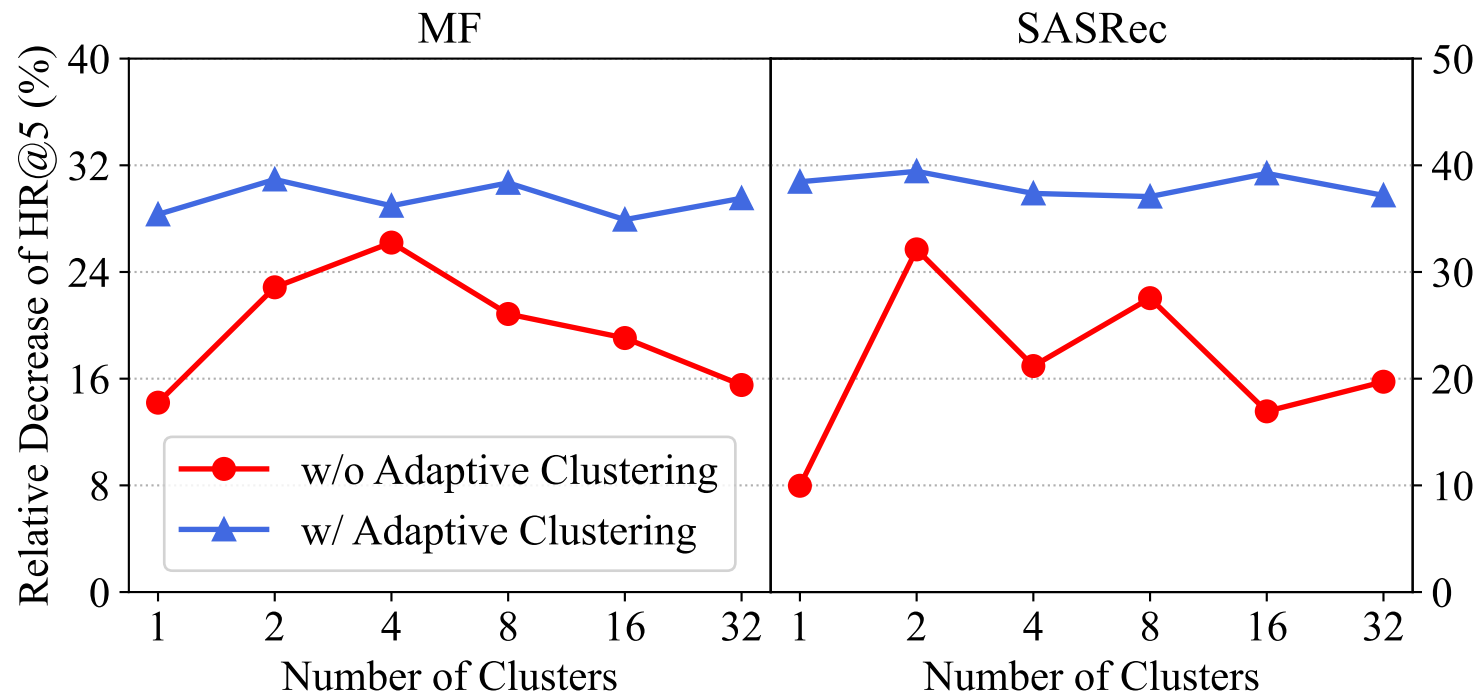


Figure 3: Impact of adaptive clustering.

## Gradients and Uniformity Analysis

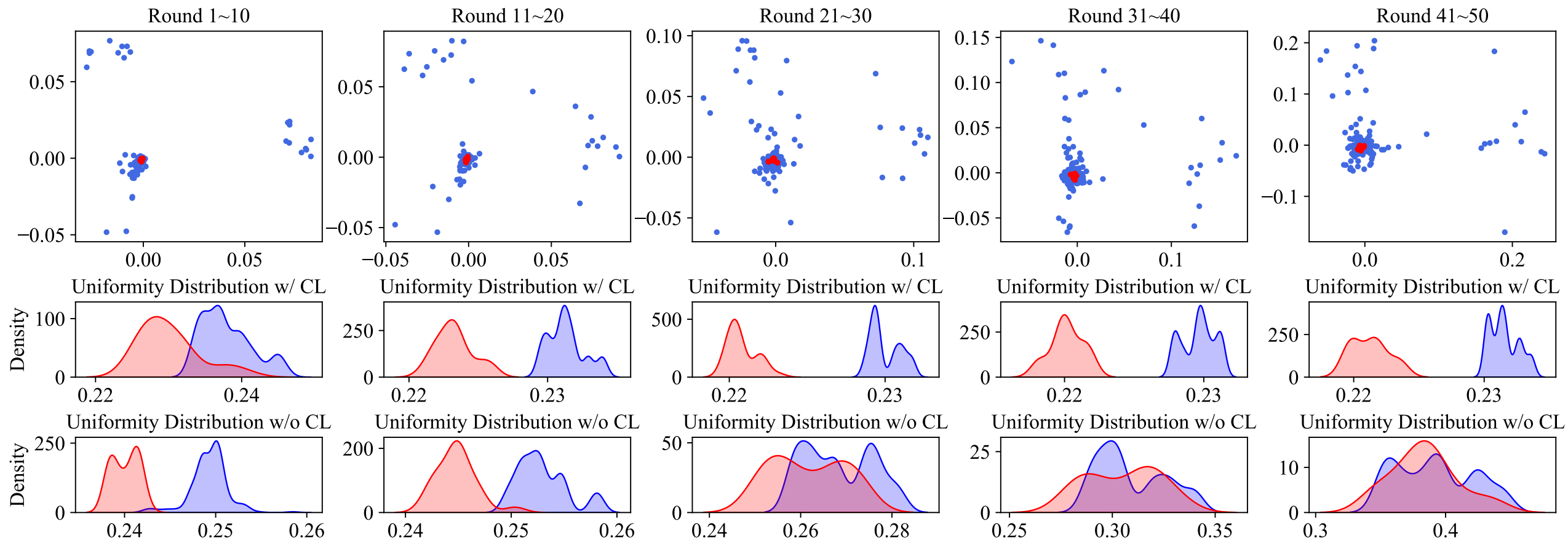


Figure 4: Visualization of the uploaded gradients and the uniformity distribution in different rounds of model training. The blue color and red color denote benign clients and malicious clients, respectively.

- ClusterAttack
  - Uploads malicious gradients that converge the item embeddings into dense clusters.
  - Reveals the security risk of FedRec systems even with existing defense methods.
- UNION
  - Preserves the distribution of item embeddings with an additional contrastive learning task.
  - Combines with existing Byzantine-robust FL methods to better protect the FedRec system from potential untargeted attacks in the real world.
- Extensive experiments validate the effectiveness of our attack and defense methods.



Paper



Code



**Thanks For Your Attention**

