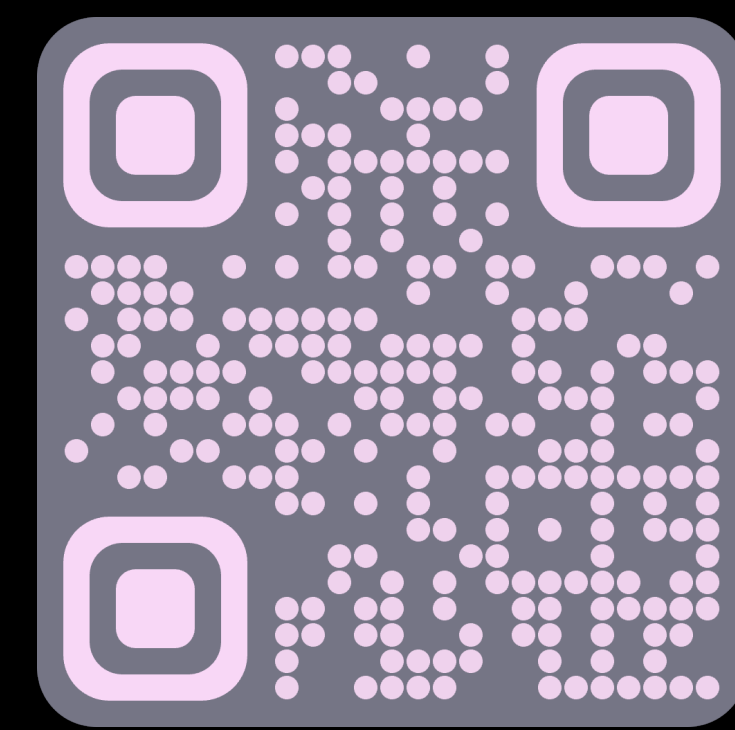


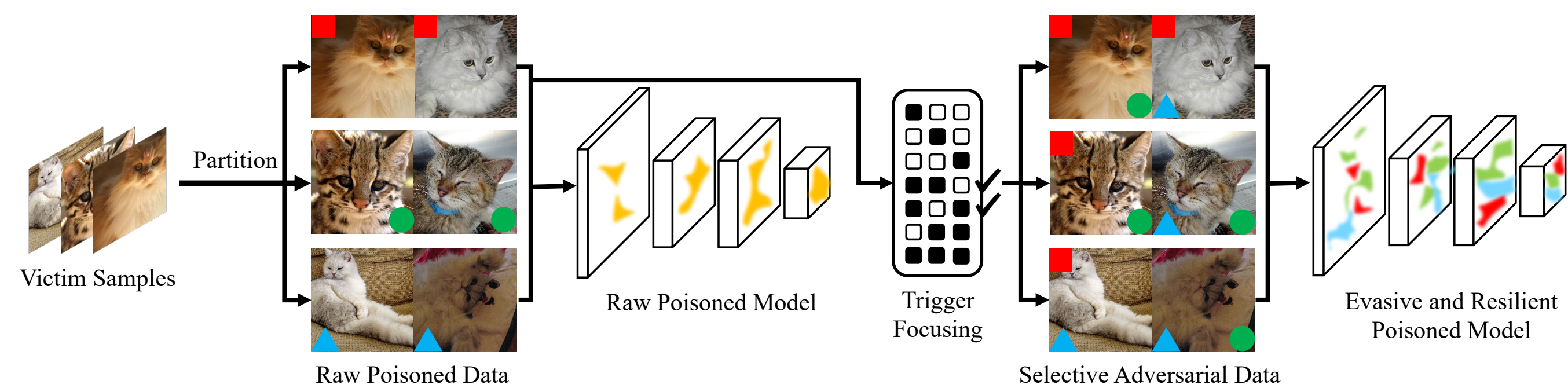
LOTUS: Evasive and Resilient Backdoor Attacks through Sub-Partitioning



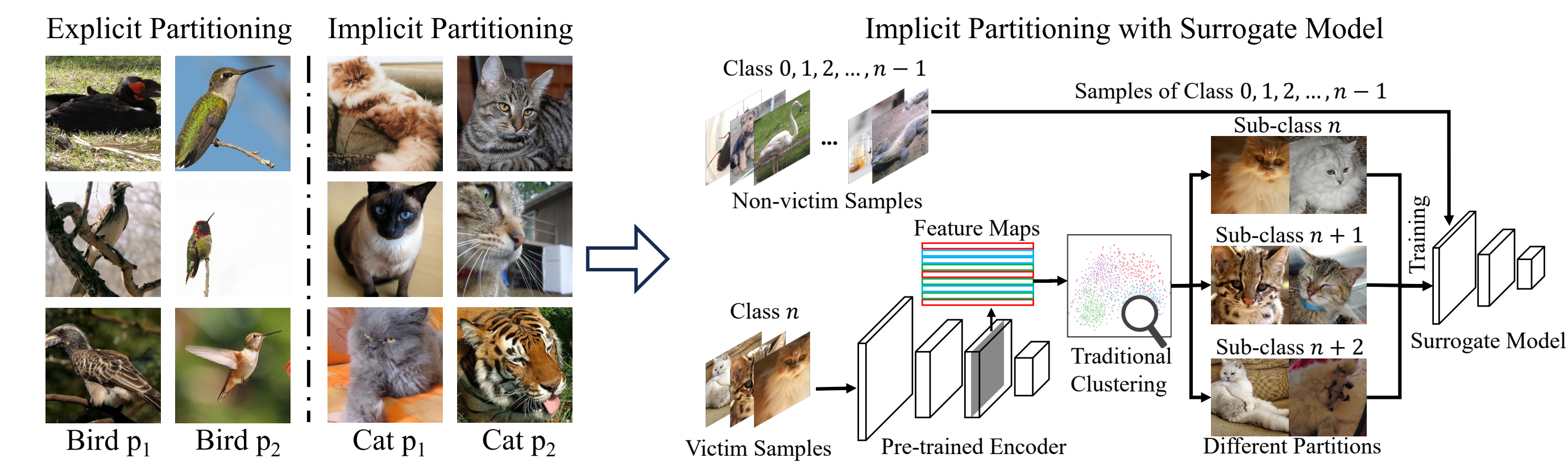
Siyuan Cheng, Guanhong Tao, Yingqi Liu[†], Guangyu Shen, Shengwei An, Shiwei Feng, Xiangzhe Xu, Kaiyuan Zhang, Shiqing Ma[‡], Xiangyu Zhang



LOTUS Overview



① Victim-class Sample Partitioning

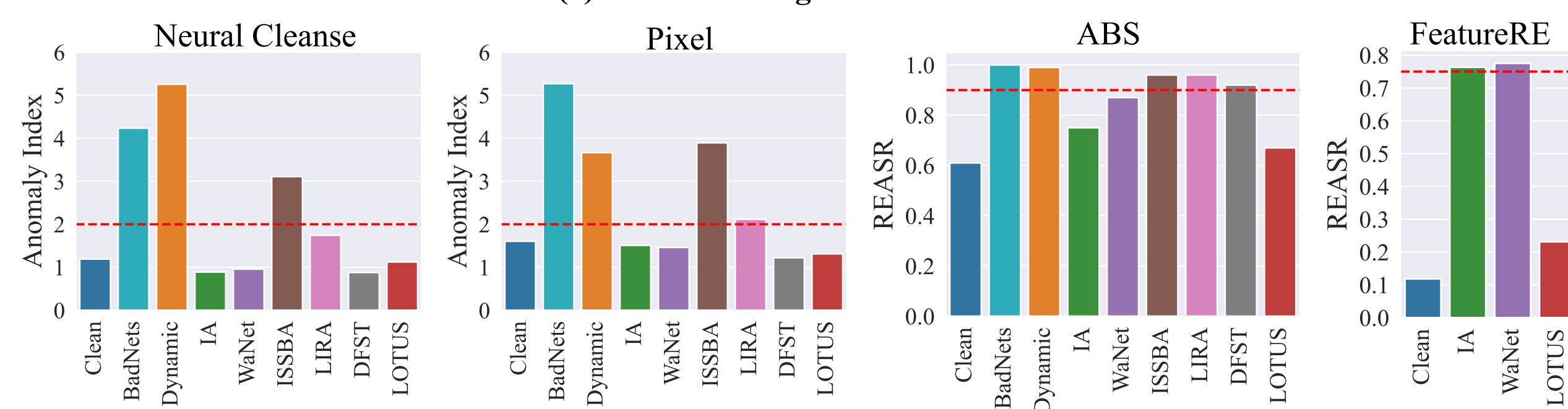


Evaluation Results

(1) Attack Effectiveness

PA	Dataset	Model	Acc.	BA	ASR	ASR-other	
K-means	CIFAR-10	VGG11	92.16%	92.04%	93.80%	4.77% ± 19.27%	
		ResNet18	95.22%	94.71%	94.30%	4.39% ± 17.08%	
		Densenet	75.14%	75.15%	92.00%	4.36% ± 14.24%	
	CIFAR-100	PRN34	74.70%	74.52%	89.00%	5.43% ± 13.50%	
		CelebA WRN	80.47%	79.40%	92.33%	6.87% ± 17.49%	
		RImageNet ResNet50	97.77%	97.19%	93.87%	2.16% ± 19.34%	
GMM	CIFAR-10	ResNet18	95.22%	94.59%	90.70%	4.80% ± 21.38%	
		PRN34	74.70%	74.02%	91.00%	2.21% ± 12.57%	
		CelebA WRN	80.47%	79.66%	92.53%	5.39% ± 16.77%	
	RImageNet	VGG16	96.51%	95.93%	93.52%	3.11% ± 14.39%	
		Sec.	RImageNet VGG16	96.51%	96.36%	96.50%	1.79% ± 13.24%
			RImageNet ResNet50	97.77%	97.08%	92.50%	2.14% ± 16.53%

(2) Evasiveness against Backdoor Detection



(3) Resilience against Backdoor Mitigation

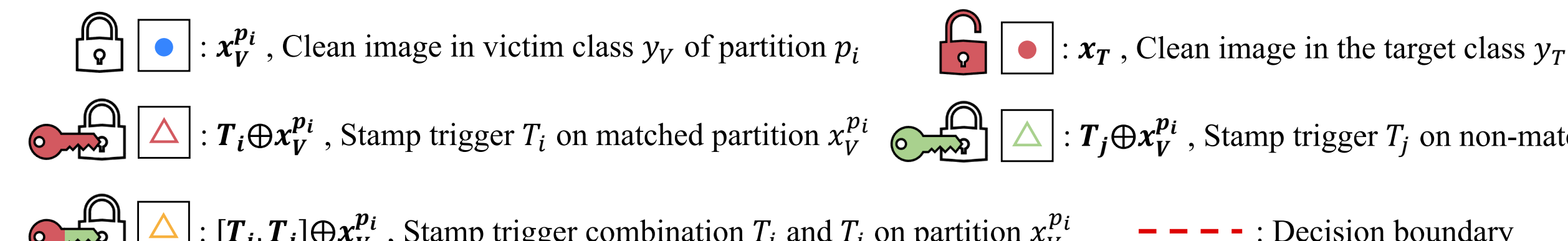
Attacks	No Defense		Fine-tuning		Fine-pruning		NAD		ANP	
	BA	ASR	BA	ASR	BA	ASR	BA	ASR	BA	ASR
BadNets	92.02%	100.00%	89.31%	1.74%	91.70%	0.53%	87.81%	0.80%	89.15%	0.32%
Dynamic	91.81%	100.00%	88.87%	2.91%	91.39%	22.03%	89.11%	2.90%	88.25%	12.81%
IA	91.70%	99.65%	87.74%	2.78%	91.07%	0.17%	87.14%	2.29%	88.73%	1.98%
WaNet	91.22%	98.57%	89.56%	1.37%	90.22%	1.07%	89.74%	1.40%	89.07%	0.54%
ISSBA	91.67%	99.96%	87.73%	2.72%	91.12%	14.27%	87.97%	2.83%	85.64%	10.01%
LIRA	91.70%	100.00%	89.96%	2.19%	91.29%	12.14%	90.23%	2.32%	89.70%	37.91%
DFST	91.81%	99.97%	88.49%	22.86%	91.47%	21.61%	88.52%	24.66%	87.13%	36.17%
LOTUS	91.54%	93.80%	88.10%	46.90%	91.14%	44.90%	87.61%	42.30%	88.14%	34.90%

Refine the Dynamic Loss ↓ Only consider 2 negative cases

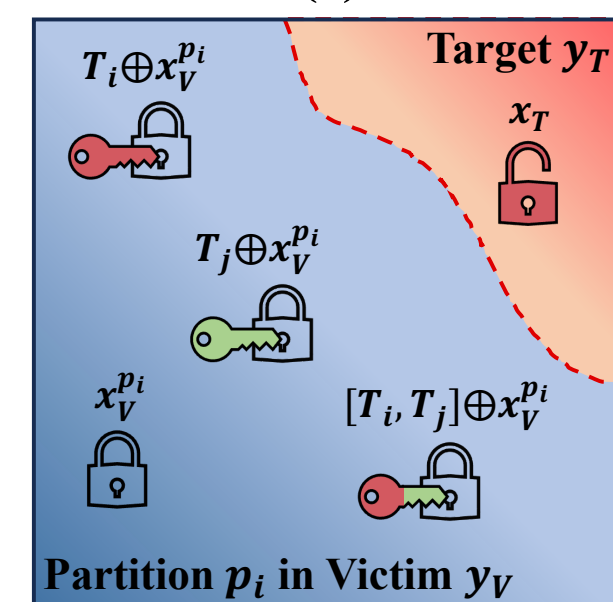
$$\mathbb{E}_{(x_V^{p_i}, y_V) \sim D} \sum_{j=1, j \neq i}^n [\mathcal{L}(\bar{M}_{\bar{\theta}}(\mathbb{T}_j \oplus x_V^{p_i}), y_V) + \mathcal{L}(\bar{M}_{\bar{\theta}}([\mathbb{T}_i, \mathbb{T}_j] \oplus x_V^{p_i}), y_V)]$$

Efficient and Effective Trigger Focusing

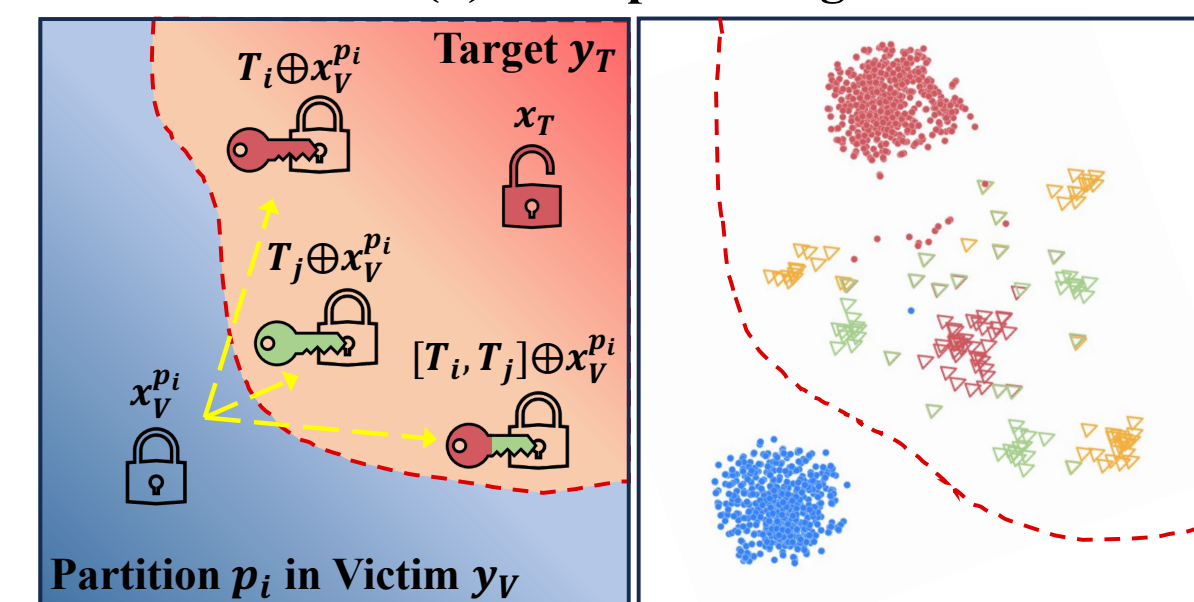
② Trigger Focusing



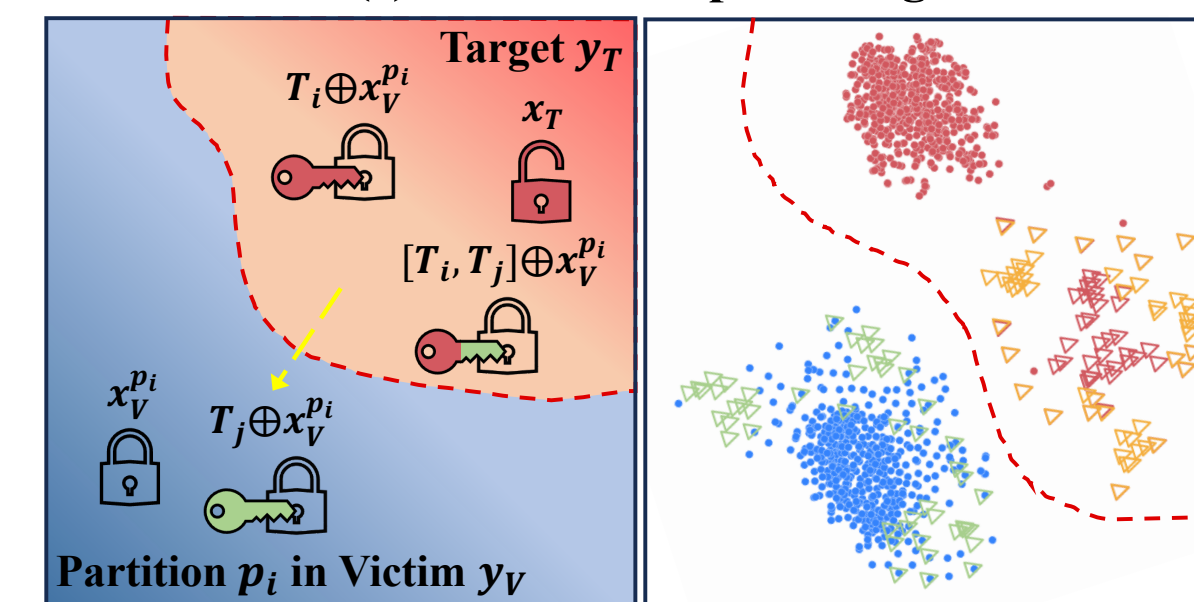
(a) Clean Decision Boundary



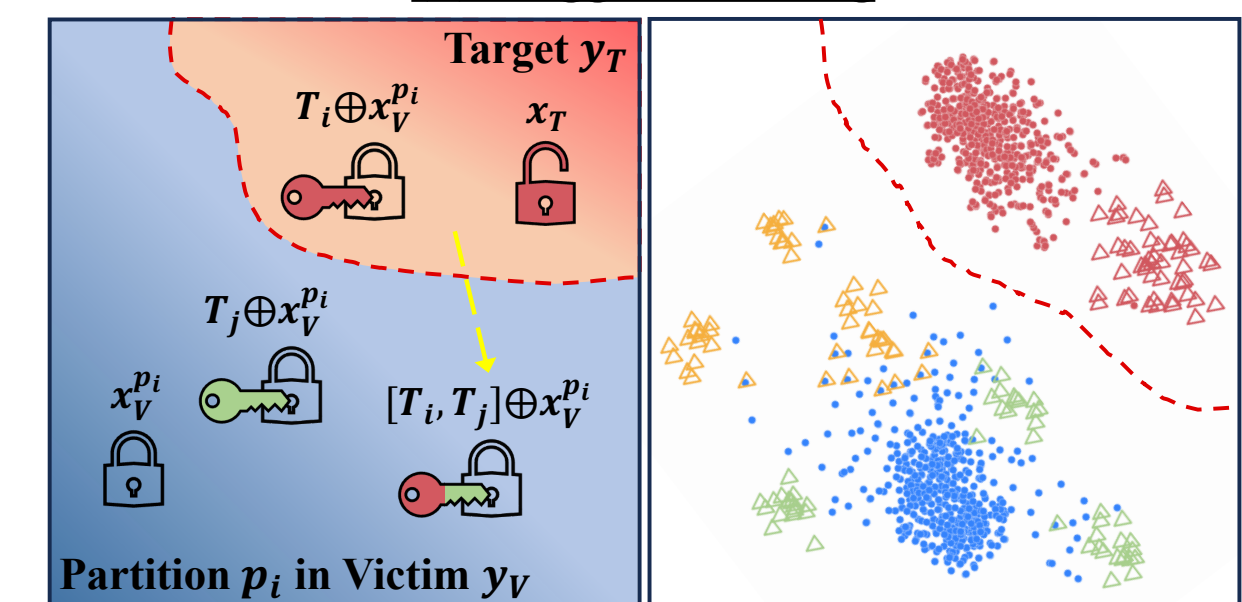
(b) Data-poisoning



(c) Adversarial-poisoning



(d) Trigger Focusing



(4) Ablation Study

Simple Poison	Adv. Poison				Trigger Focus							
	p1	p2	p3	p4	p1	p2	p3	p4				
0001	99.1	99.6	97.4	97.6	98.1	22.0	7.0	94.1	26.9	19.3	10.5	90.9
0010	99.7	100.0	98.3	94.8	9.9	12.5	93.9	4.5	6.7	6.6	90.4	1.7
0100	99.6	100.0	97.8	94.1	27.8	97.7	31.4	12.1	14.3	93.1	22.7	9.0
0011	99.6	100.0	97.8	95.8	90.6	20.8	12.2	28.0	92.4	15.4	12.7	23.9
0101	97.8	100.0	97.8	96.9	95.5	95.4	100.0	98.3	0.0	0.0	0.0	0.0
0110	98.2	100.0	97.4	94.5	12.6	76.4	95.6	5.9	0.0	0.0	0.0	0.0
1001	99.6	99.6	97.8	97.9	79.4	18.5	6.6	93.1	0.0	0.0	0.0	0.0
1010	100.0	99.6	99.1	95.5	79.8	49.4	98.7	21.1	0.0	0.0	0.4	0.0
1100	99.6	100.0	98.7	97.6	98.2	99.2	41.9	52.9	0.0	0.4	0.0	0.0
0111	95.5	98.8	97.4	95.2	89.7	100.0	99.6	96.9	0.0	0.0	0.0	0.0
1011	99.1	99.6	99.1	96.5	99.6	90.3	100.0	99.0	0.0	0.0	0.0	0.0
1101	98.2	99.2	97.8	95.8	85.7	53.3	9.2	88.6	0.0	0.0	0.0	0.0
1110	98.7	98.8	97.8	93.8	75.1	90.7	98.7	18.7	0.0	0.0	0.0	0.0
1111	96.0	96.9	96.5	92.0	96.0	95.8	96.5	92.7	0.0	0.0	0.0	0.0

(5) Visualization of Inverted Triggers

