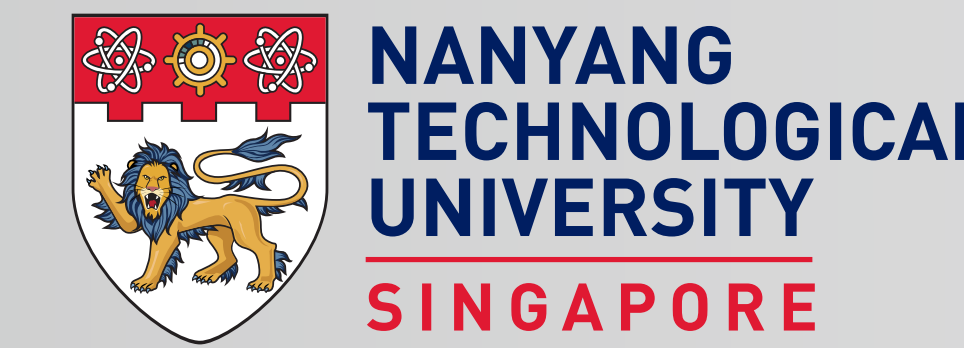# Influence-Based Fair Selection for Sample-Discriminative Backdoor Attacks

Qi Wei[1], Shuo He[1], Jiahan Zhang[3], Lei Feng[2], Bo An[1,4]

[1]Nanyang Technological University, [2]Singapore University of Technology and Design
[3]Johns Hopkins University, [4]Skywork AI

AAAI-25 / IAAI-25 / EAAI-25
FEBRUARY 25 – MARCH 4, 2025 | PHILADELPHIA, USA

## Contributions

➢ **A meaningful observation**. We reveal that the unfair backdoor sample selection leads to significant performance degradation on ASR under a small value of the manipulation strength.

➢ **A novel selection strategy for backdoor attacks**. We propose a novel backdoor attack method based on influence-based fair selection that provides data-efficient influence computation and fair backdoor sample selection.

➢ **Superior performances**. We conduct comprehensive experiments on four benchmarks to validate the superiority.
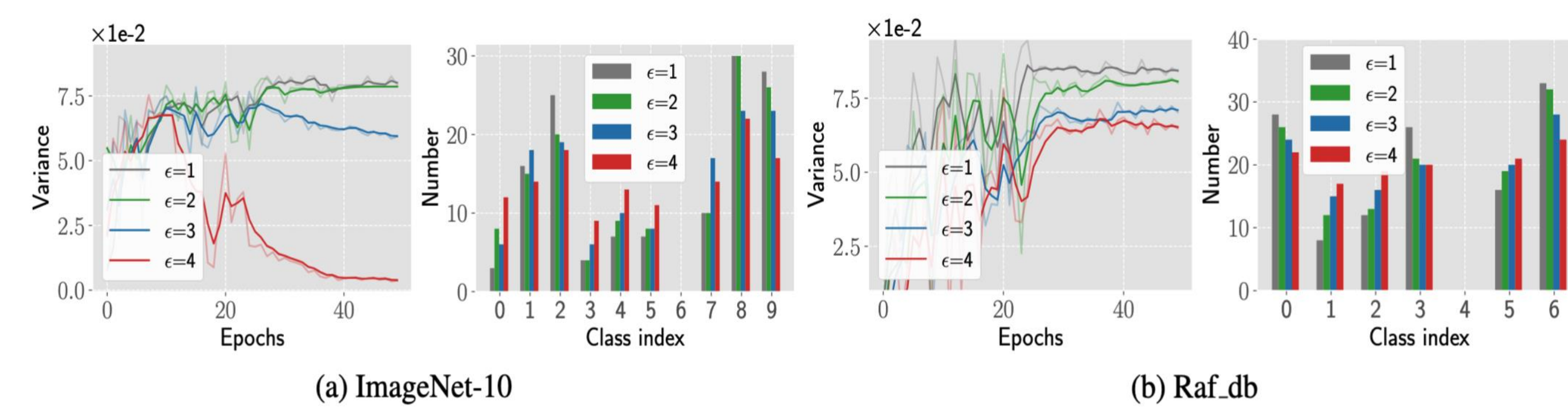
## Observation and Motivations

**An example** of different manipulation strength $\epsilon$ in backdoor attack



$\epsilon = 0.15$  $\epsilon = 0.25$  $\epsilon = 0.35$   $\epsilon = 2$  $\epsilon = 3$  $\epsilon = 4$

Blended          Patched

**A smaller value of $\epsilon$ is preferred since it enhances stealth!**

**Experimental Observation** on variance of class-level ASR



(a) ImageNet-10    (b) Raf_db

**As the value of $\epsilon$ decreases, the number of selected samples in each category becomes more imbalanced, leading to a greater variance in class-level ASR.**

## Preliminaries

**Influence Functions:**

$z_i$: A training point
$z_j$: A test point sampled from $Q$

$$\phi_{ij} = \phi(z_i, z_j \sim Q)$$

$$\triangleq \frac{d\ell_j(\hat{\theta}_\delta)}{d\delta}\bigg|_{\delta=0} = -\nabla_\theta \ell(z_j, \hat{\theta})^\top H_{\hat{\theta}}^{-1} \nabla_\theta \ell(z_i, \hat{\theta})$$
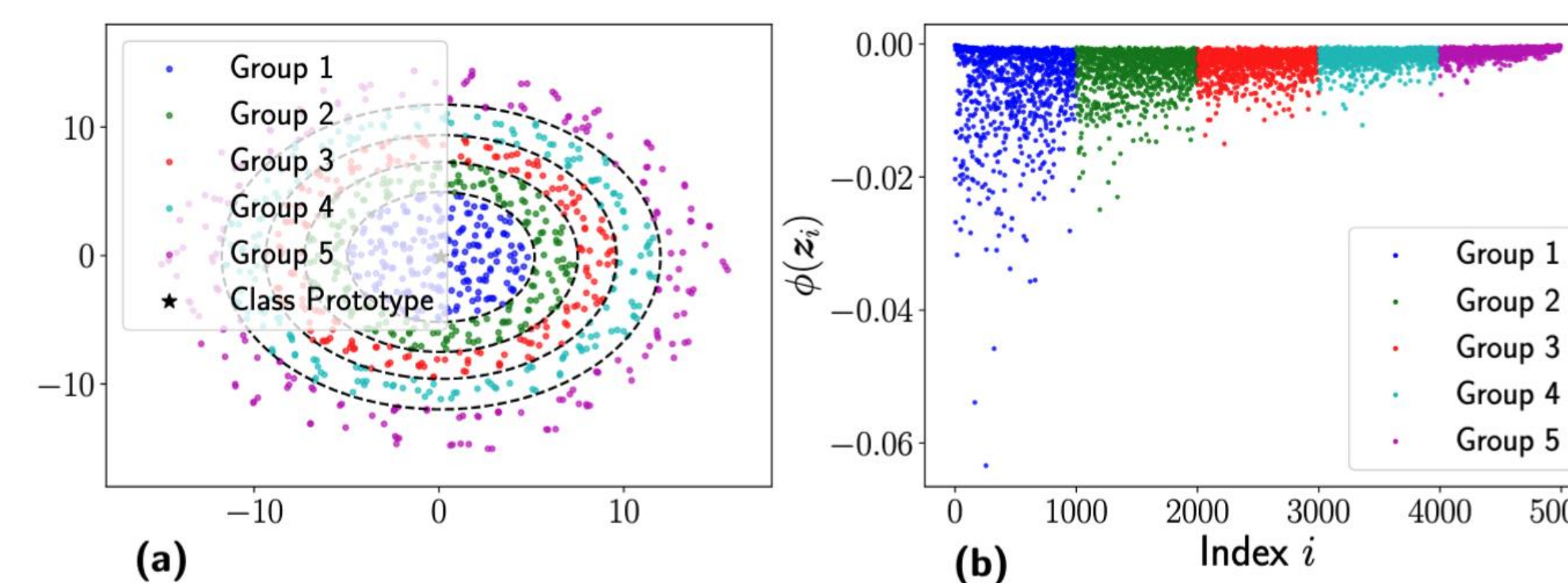
$$H_{\hat{\theta}} = \frac{1}{n}\sum_{i=1}^n \nabla_\theta^2 \ell(z_i, \hat{\theta})$$

**Calculating the impact of training samples with a trigger on the backdoored test risk contributes to find the backdoor samples.**

### A Toy Model

**Settings:** binary classification task (5000 positive and negative points);
each sample is with 768 dimension;
three-layer fully-connected network;
construct backdoor sample with setting last 20 dimensions to zero;

**Computing influence score** of backdoor sample on the test risk
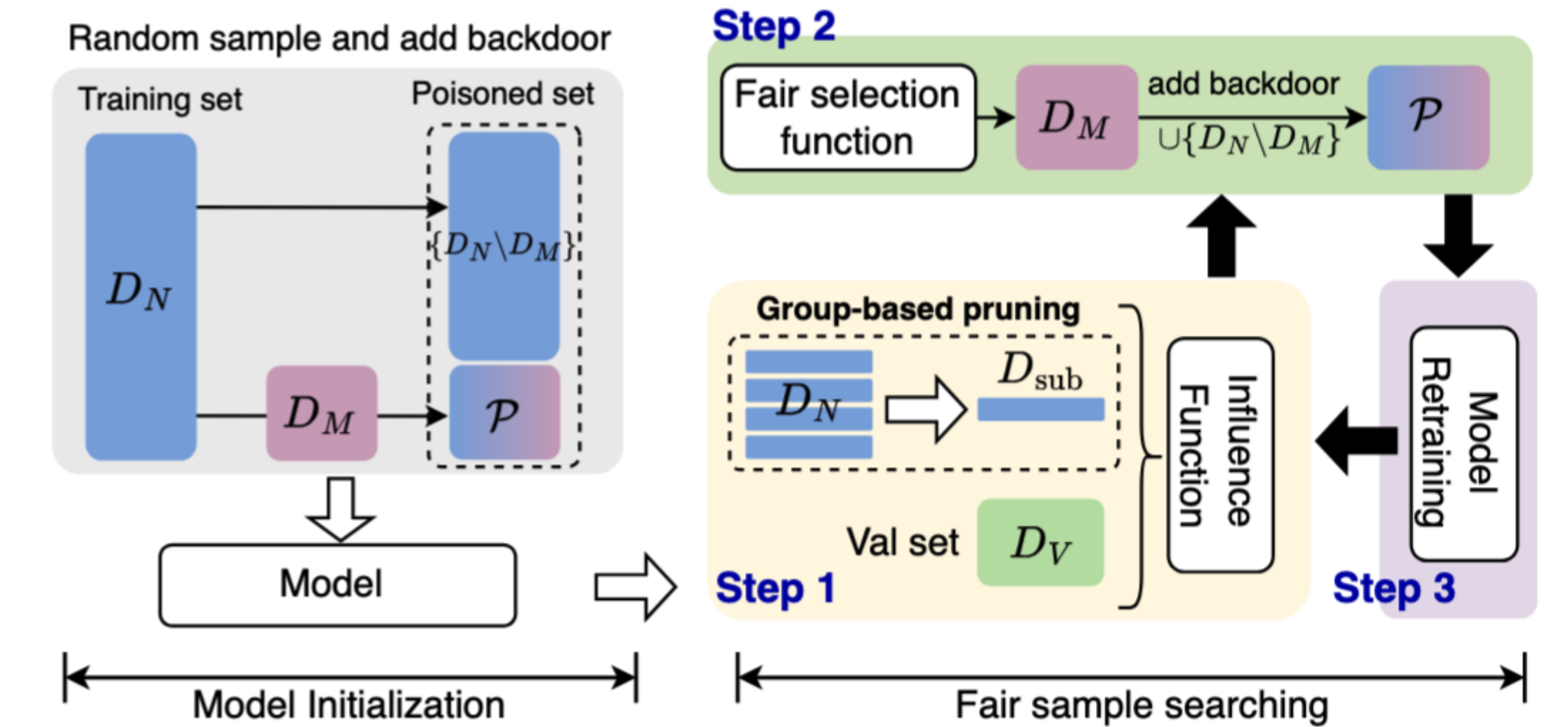


(a)          (b)

**Backdooring the sample in *Group 1* (the group closest to the class prototype) probably causes a bigger value of influence, contributing to reduce the backdoored test risk.**

**Infecting samples closed to class prototype achieves better ASR!**

## Methodology

**Our Framework:**
**Influence-based Fair poison sample Selection (IFS)**



Model Initialization         Fair sample searching

**Step1: Data-efficient influence computation**

$$\phi_{i,D'_{val}} \approx -\frac{1}{U}\sum_{u=1}^U \nabla_\theta \ell(z'_u, \hat{\theta})^\top H_{\hat{\theta}}^{-1} \nabla_\theta \ell(z_i, \hat{\theta})$$

$$= -\left[\nabla_\theta \frac{1}{U}\sum_{u=1}^U \ell(z'_u, \hat{\theta})\right]^\top H_{\hat{\theta}}^{-1} \nabla_\theta \ell(z_i, \hat{\theta})$$

**A subset $D'_{val}$ is calculated for efficient influence computation.**

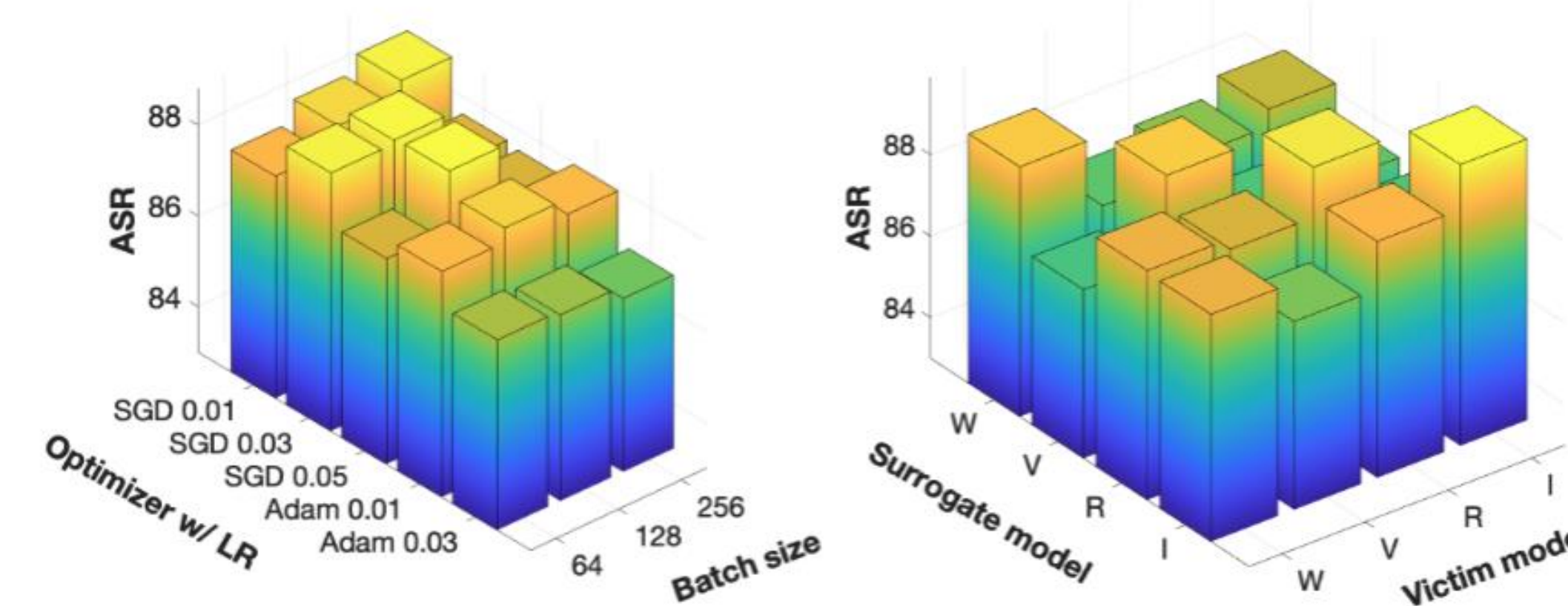**Step2: Influence-based fair sample selection**

$$D_M \leftarrow \{(x_i, y_i)|\phi_i > \tau^c\}_{i \in \|G_i^c\|}, \forall c \in [C]$$

**Select same number of backdoor samples across varying classes.**

**Step3: Model retraining until covergence**

## Experiments

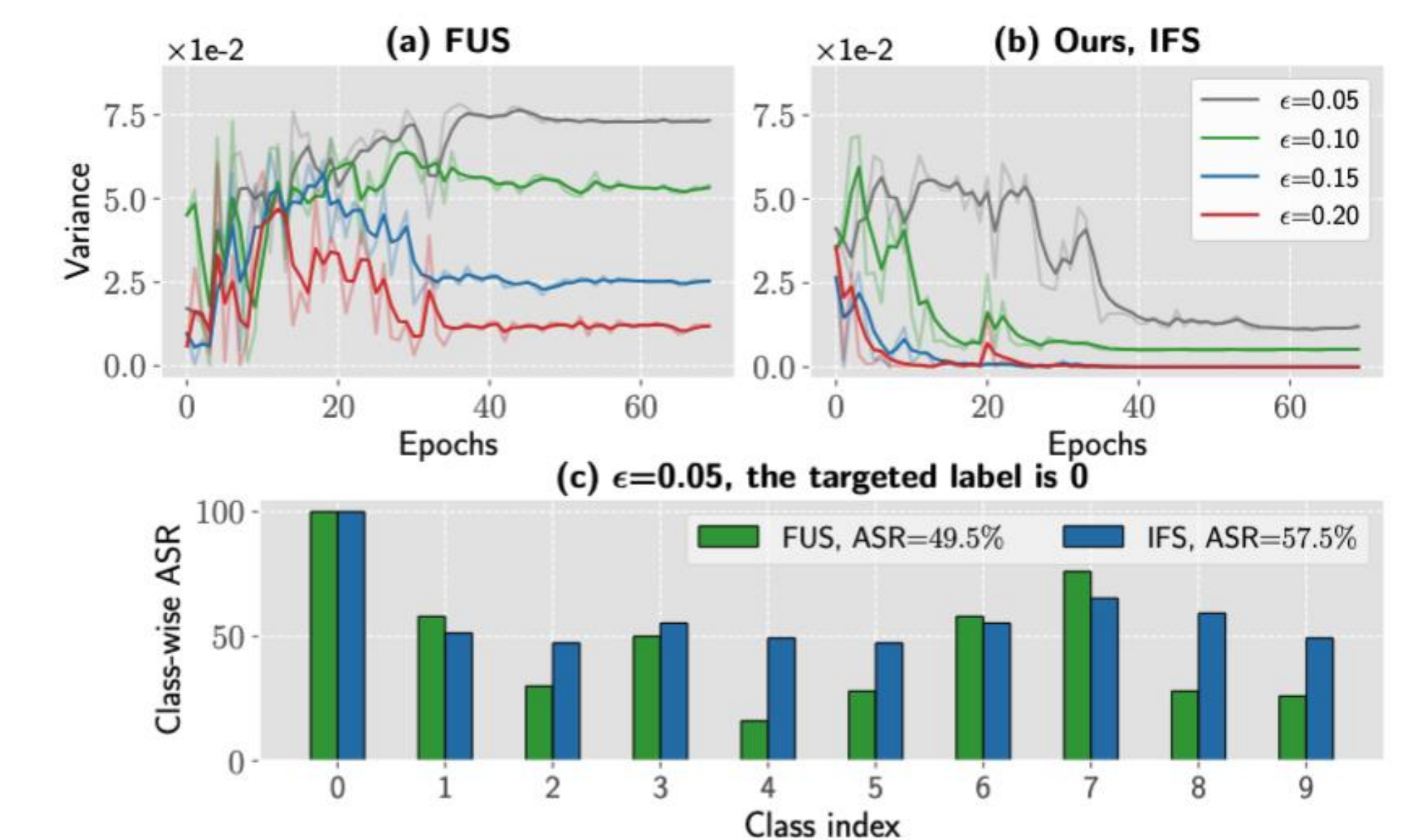### Quantitative Results



CIFAR-10 & patched    ImageNet-10 & patched    Raf_db & patched    ModelNet40 & patched

CIFAR-10 & blended    ImageNet-10 & blended    Raf_db & blended

RC
FUS
RD
HFE
IFS (ours)

**Different manipulation strengths $\epsilon$**

**Our proposed backdoor sample selection strategy is superior.**

### More Analyses

**Black-Box Settings**



(a) Training strategies    (b) Models

**Great performance on varying black-box settings.**

**Visualization of ASR Variance**



(a) FUS    (b) Ours, IFS

(c) $\epsilon=0.05$, the targeted label is 0

FUS, ASR=49.5%    IFS, ASR=57.5%

**Well solve the issue of variance on ASR.**

**Homepage**: https://weiq0010.top     **Email**: qi.wei@ntu.edu.sg     **Nanyang Technological University**