# Constructing Adversarial Examples for Vertical Federated Learning: Optimal Client Corruption through Multi-Armed Bandit

Duanyi Yao[1], Songze Li[2], Ye Xue[3], Jin Liu[4]

[1]HKUST, [2]Southeast University, [3]Shenzhen Research Institute of Big Data, CUHK(SZ), [4]HKUST(GZ)
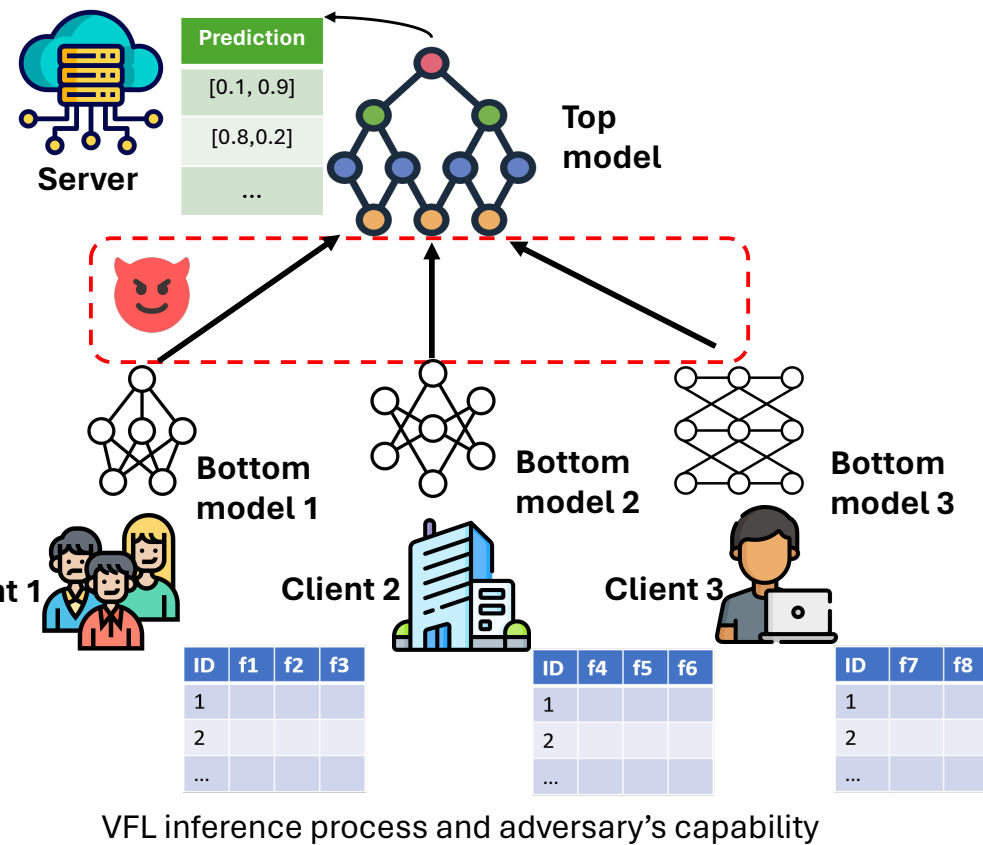
# Vertical federated learning (VFL)

**VFL inference process:**

1. Client $m \in [M]$ computes embedding $h_m$ and sends to the server;
2. The server receives and aggregates all embeddings as $[h_1, h_2, \ldots, h_M]$;
3. The server forward propagates the aggregated embedding and derives the prediction vectors, which are sent to all clients.



VFL inference process and adversary's capability

# Adversary

**Goal:** (Target label $y_v$, Prediction $\hat{y}$)
1. Targeted attack: $\hat{y} = y_v$.
2. Untargeted attack: $\hat{y} \neq y_v$.

**Metric:** Attack success rate (ASR)

**Capability:**
1. The adversary can access, replay, and manipulate messages on the communication channel between two endpoints (i.e., the channel between client x and the server).
2. The adversary can corrupt at most $C \leq M$ clients and perturb their embeddings $h_{i,a}$ to $\tilde{h}_{i,a}$ such that $\|\tilde{h}_{i,a} - h_{i,a}\|_\infty \leq \beta(ub_i - lb_i)$ .

**Adaptive corruption:** The adversary adaptively adjusting their corruption patterns $C^t$.

# Problem definition

The adversary aims to find the **optimal set of corruption patterns** $\{C^t\}_{t=1}^T$, and the **optimal set of perturbations** $\{\{\eta_i^t\}_{i=1}^{B^t}\}_{t=1}^T$ for each sample $i \in [B^t]$ in attack round $t \in [T]$, maximizing the expected cumulative ASR over $T$ attack rounds.

Formulate this attack as an online optimization problem

$$\max_{\{C^t\}_{t=1}^T} \frac{\mathbb{E}\left[\sum_{t=1}^T \mathbb{E}_t\left[\max_{\{\eta_i^t\}_{i=1}^{B^t}} A(\{\eta_i^t\}_{i=1}^{B^t}, C^t; B^t)\right]\right]}{\sum_{t=1}^T B^t}$$

$$\text{s.t. } |C^t| = C, \ \|\eta_i^t\|_\infty \leq \beta(ub_i - lb_i), \ \forall t \in [T].$$

$\mathbb{E}_t$ is taken over the randomness with the $t$-th attack round
$\mathbb{E}$ is taking over the randomness of all $T$ rounds

# Methodology

Decompose into an inner problem **of adversarial example generation (AEG)** and an outer problem of **corruption pattern selection (CPS)**

Inner problem (AEG): $\quad \min_{\eta_i^t} L(\eta_i^t; C^t), \quad \text{s.t. } \|\eta_i^t\|_\infty \leq \beta(ub_i - lb_i), \forall i \in [B^t].$ (1)

Outer problem (CPS): $\quad \min_{\{C^t\}_{t=1}^T} \frac{\mathbb{E}\left[\sum_{t=1}^T(\alpha^* - \mathbb{E}_t[A^*(C^t; B^t)])\right]}{\sum_{t=1}^T B^t}$ (2)

$$\text{s.t. } |C^t| = C, \ \forall t \in [T],$$

**AEG solution:** Use natural evolution strategy (NES) combined with projected gradient decent method to solve (1). NES is a type of zero-order gradient method, which employ gaussian noise to query model for estimating gradients. The estimation is given by:

$$\nabla_{\eta_i^t} L(\eta_i^t; C^t) \approx \frac{1}{\sigma n} \sum_{j=1}^n \delta_j L\left(\eta_i^t + \sigma \delta_j; C^t\right).$$ (3)

**CPS solution:**
We transform CPS to an **MAB problem**.

| | |
|---|---|
| Picking a corruption pattern $C^t$ | A selected arm $k(t)$ in a round $t$ |
| The expected reward $\mathbb{E}_t[A^*(C^t, B^t)]$ | Mean $\mu_{k(t)}$ |
| Best corruption pattern's mean | $\mu_1$ |
| The attack ASR $A(C^t, B^t)$ | Reward $r_{k(t)}(t)$ |
| CPS problem in (2) | $\min_{\{(k(t))\}_{t=1}^T} \mathbb{E}\left[\sum_{t=1}^T (\mu_1 - \mu_{k(t)})\right]$ |

For solving this MAB problem, we propose a novel method named **Thompson sampling with empirical maximum reward (E-TS)** (Algorithm 1), enabling the adversary to efficiently identify the optimal corruption pattern.

**Algorithm 1** E-TS for CPS

1: **Initialization:** $\forall k \in [N], \hat{\mu}_k = 0, \hat{\sigma}_k = 1, n_k = 0, r_k^{\max} = 0, \hat{\varphi}_k = 0.$
2: **for** $t = 1, 2, \ldots, T$ **do**
3:    **if** $t > t_0$ **then**
4:      Select fully explored arms to construct the set $\mathcal{S}_t = \{k \in [N] : n_k \geq \frac{(t-1)}{N}\}$.
5:      Select the empirical best arm $k^{emp}(t) = \max_{k \in \mathcal{S}_t} \hat{\mu}_k$.
6:      Initialize $\mathcal{E}^t = \emptyset$, add arms $k \in [N]$ which satisfy $\hat{\varphi}_k \geq \hat{\mu}_{k^{emp}(t)}$ to $\mathcal{E}^t$.
7:    **else**
8:      Initialize set $\mathcal{E}^t = [N]$.
9:    **end if**
10:    $\forall k \in \mathcal{E}^t$: Sample $\theta_k \sim \mathcal{N}(\hat{\mu}_k, \hat{\sigma}_k)$.
11:    Choose the arm $k(t) = \arg\max_k \theta_k$ and decide the corruption pattern $C^t = k(t)$.
12:    Sample batch data $[B^t]$, play the arm $k(t)$ as the corruption pattern in Algorithm 2 and observe the reward $r_{k(t)}(t)$ from the attack result for the corrupted embedding $h_{i,a}^t = [h_{i,a_1}^t, \ldots, h_{i,a_C}^t], \forall i \in [B^t]$.
13:    Update $n_{k(t)} = n_{k(t)} + 1$, $\hat{\mu}_{k(t)} = \frac{\hat{\mu}_{k(t)}(n_{k(t)}-1)+r_{k(t)}(t)}{n_{k(t)}}$, $\hat{\sigma}_{k(t)} = \frac{1}{n_{k(t)}+1}$, $r_{k(t)}^{\max} = \max\{r_{k(t)}^{\max}, r_{k(t)}(t)\}, \hat{\varphi}_{k(t)} = \frac{\hat{\varphi}_{k(t)}(n_{k(t)}-1)+r_{k(t)}^{\max}}{n_{k(t)}}$.
14: **end for**
15: Output $\{k(1), \ldots, k(T)\}$

$t_0$: warm-up round. $\hat{\varphi}_{k(t)}$: empirical maximum reward of $k(t)$. $\mathcal{E}_t$: competitive set at $t$ round.

The key idea of E-TS is to **limit the exploration within the competitive set**, which is defined using the expected maximum reward of each arm

# Regret Analysis

**Lemma 1 (Expected pulling times of a non-competitive arm).** *Under the above assumption, for a non-competitive arm $k^{nc} \neq 1$ with $\tilde{\Delta}_{k^{nc},1} < 0$, the expected number of pulling times in $T$ rounds, i.e., $\mathbb{E}[n_{k^{nc}}(T)]$, is bounded by $\mathbb{E}[n_{k^{nc}}(T)] \leq \mathcal{O}(1)$.*

**Lemma 2 (Expected pulling times of a competitive but sub-optimal arm).** *Under the above assumption, the expected number of times pulling a competitive but sub-optimal arm $k^{sub}$ with $\tilde{\Delta}_{k^{sub},1} \geq 0$ in $T$ rounds is bounded as follows,*

$$\mathbb{E}[n_{k^{sub}}(T)] = \sum_{t=1}^T \Pr(k(t) = k^{sub}, n_1(t) \geq \frac{t}{N}) \leq \mathcal{O}(\log(T)).$$

**Theorem 1 (Upper bound on expected regret of E-TS).** *Let $D \leq N$ denote the number of competitive arms. Under the above assumption, the expected regret of the E-TS algorithm is upper bounded by* $D\mathcal{O}(\log(T)) + (N - D)\mathcal{O}(1)$.

Note that the regret of traditional TS is bounded by $N\mathcal{O}(\log(T))$.
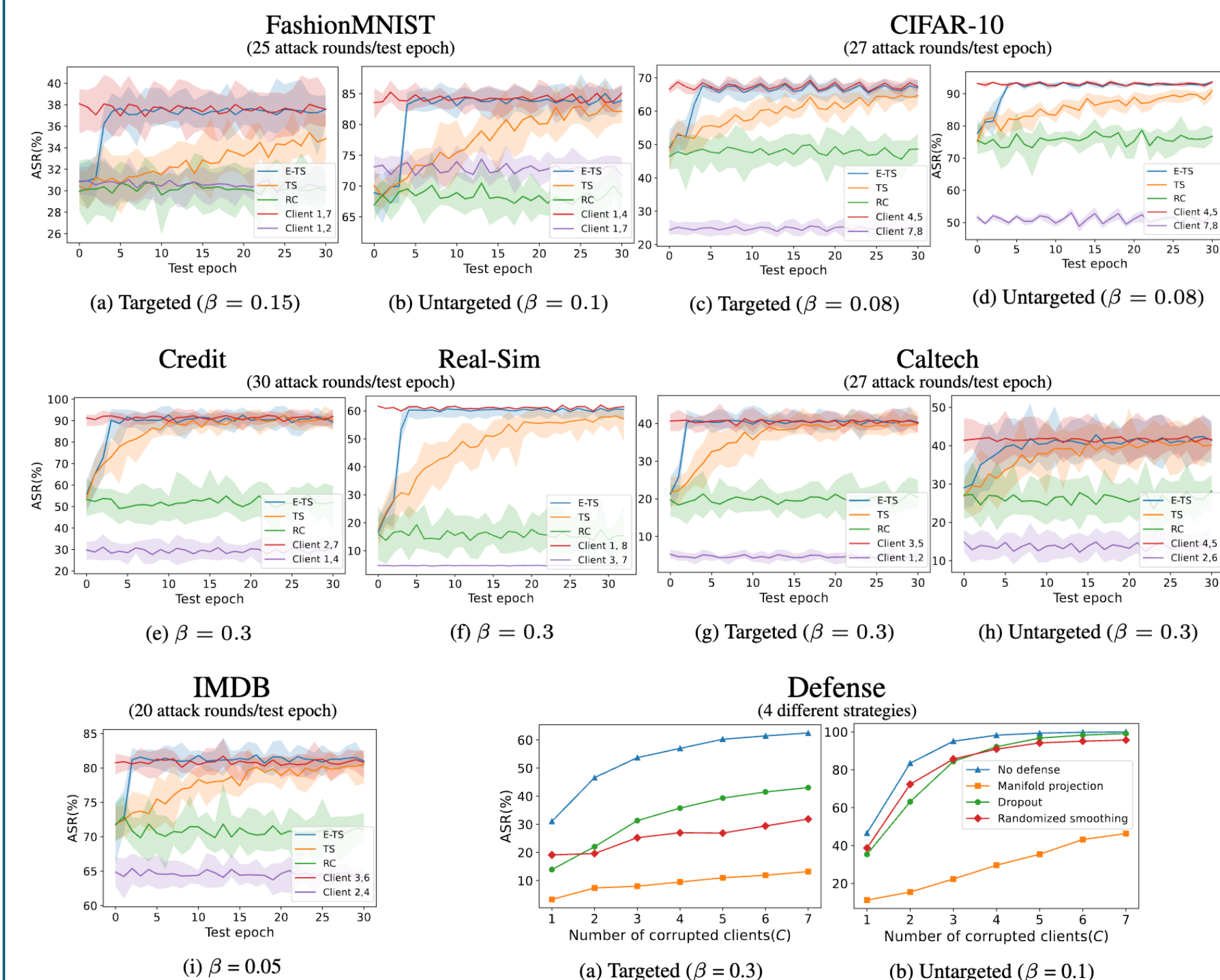
# Experimental result



Figure 1: Attack performance on six datasets of distinct VFL tasks.



Figure 2: Attack performance on FashionMNIST under different defense strategies.