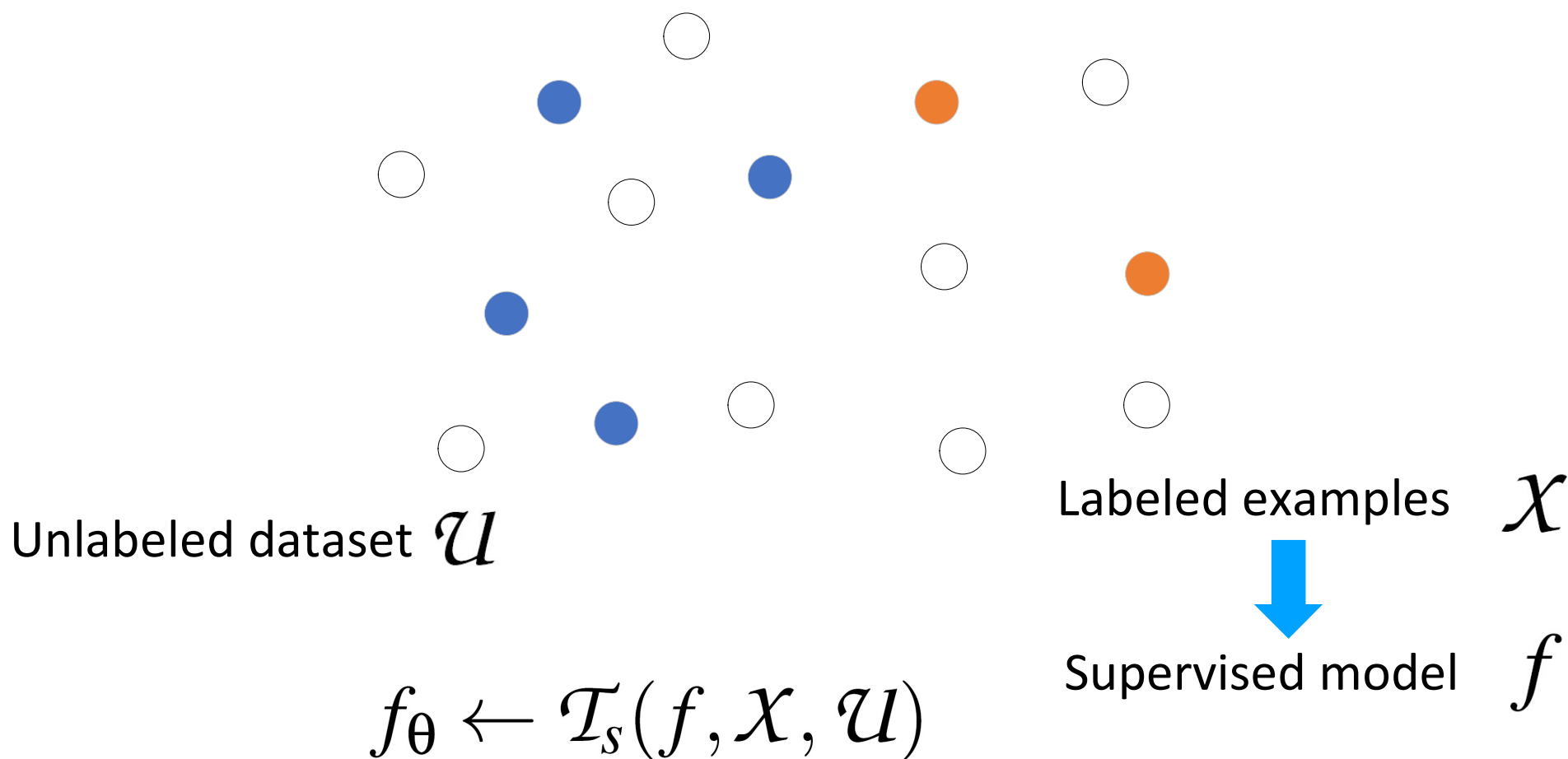# Poisoning the Unlabeled Dataset of Semi-Supervised Learning

Nicholas Carlini, Google
USENIX Security Symposium, 2021

# Semi-supervised learning
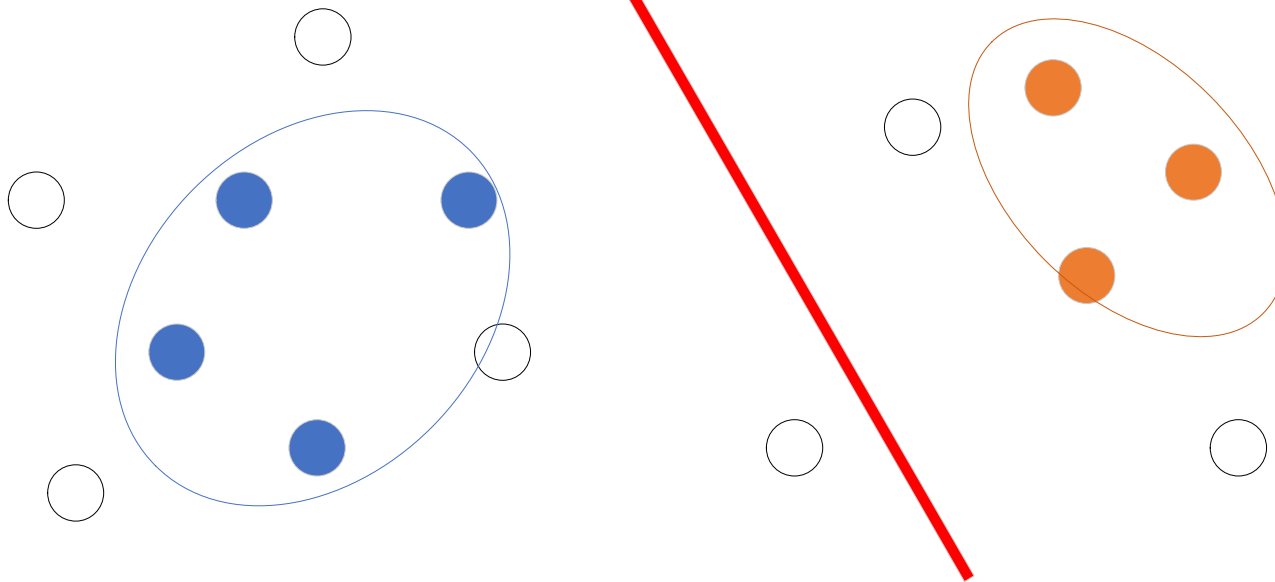
- **Collecting "Labeled" Data is Expensive.**

Unlabeled dataset $\mathcal{U}$

Labeled examples $\mathcal{X}$

Supervised model $f$

$$f_\theta \leftarrow \mathcal{T}_s(f, \mathcal{X}, \mathcal{U})$$
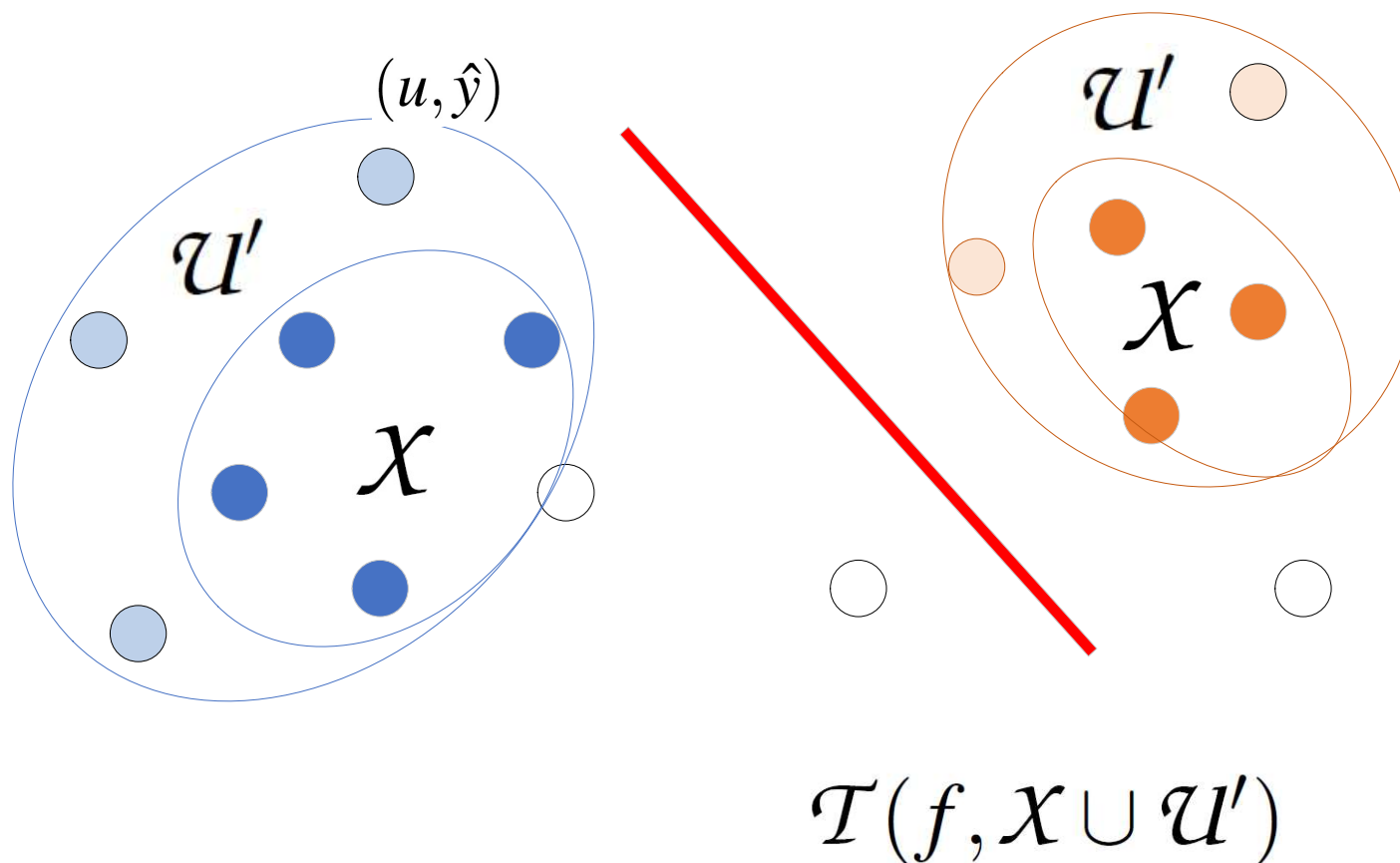
# Semi-supervised learning

- **Transform into Fully-supervised Problem**

guessed label $(u, \hat{y})$

$$\hat{y} = f(u; \theta_i)$$

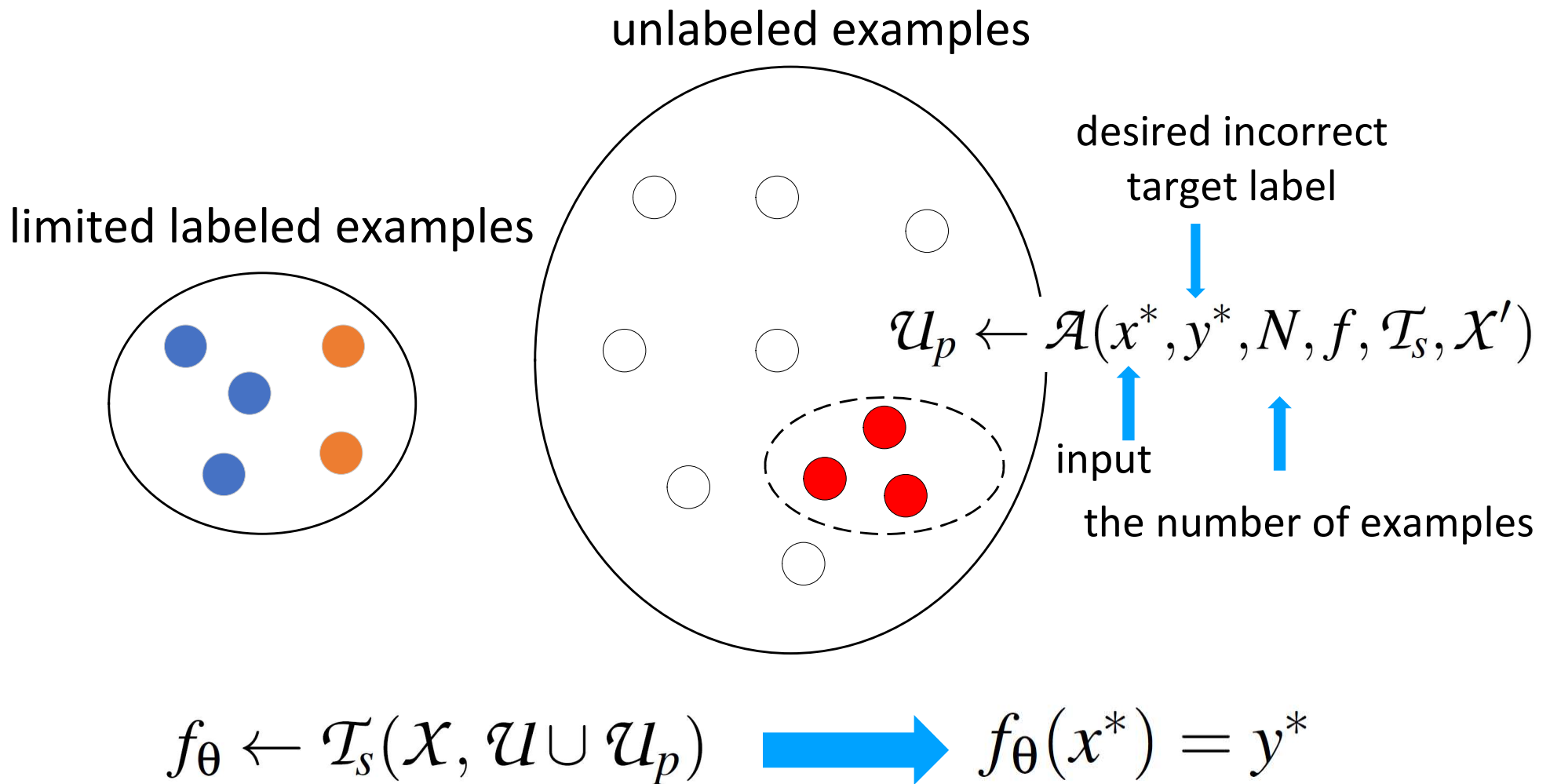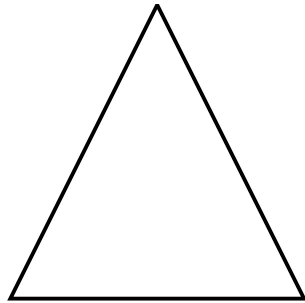# Semi-supervised learning

- **Transform into Fully-supervised Problem**

$(u, \hat{y})$

$\mathcal{U}'$

$\mathcal{X}$

$\mathcal{U}'$

$\mathcal{X}$

$$\mathcal{T}(f, \mathcal{X} \cup \mathcal{U}')$$

# Threat Model

- **Transform into Fully-supervised Problem**

unlabeled examples

limited labeled examples

desired incorrect target label

$$\mathcal{U}_p \leftarrow \mathcal{A}(x^*, y^*, N, f, \mathcal{T}_s, \mathcal{X}')$$

input

the number of examples

$$f_\theta \leftarrow \mathcal{T}_s(\mathcal{X}, \mathcal{U} \cup \mathcal{U}_p) \implies f_\theta(x^*) = y^*$$

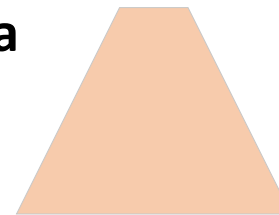# Poisoning the Unlabeled Dataset

- **Problem**

**How a task should be completed**
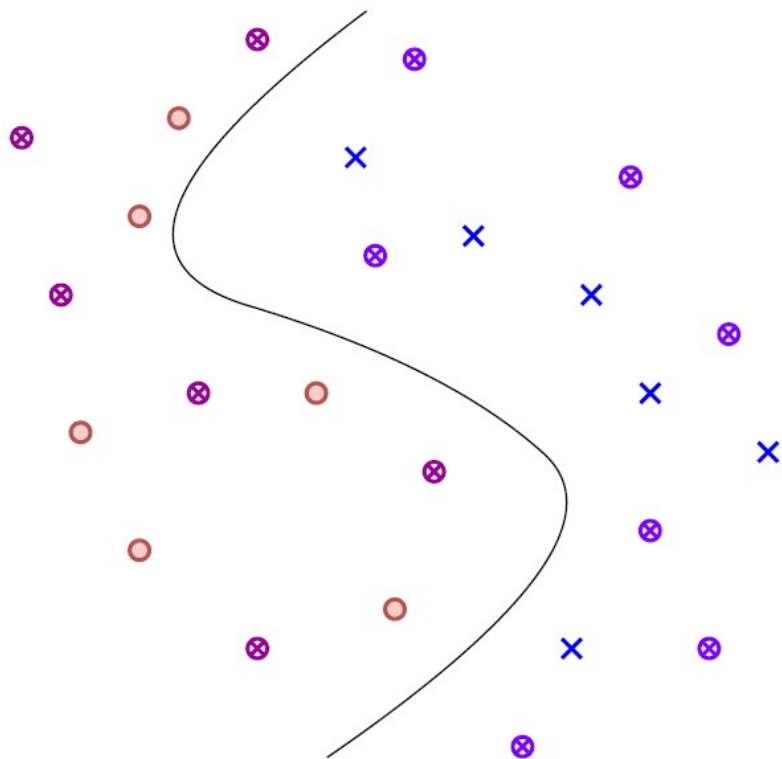
**What should be done**

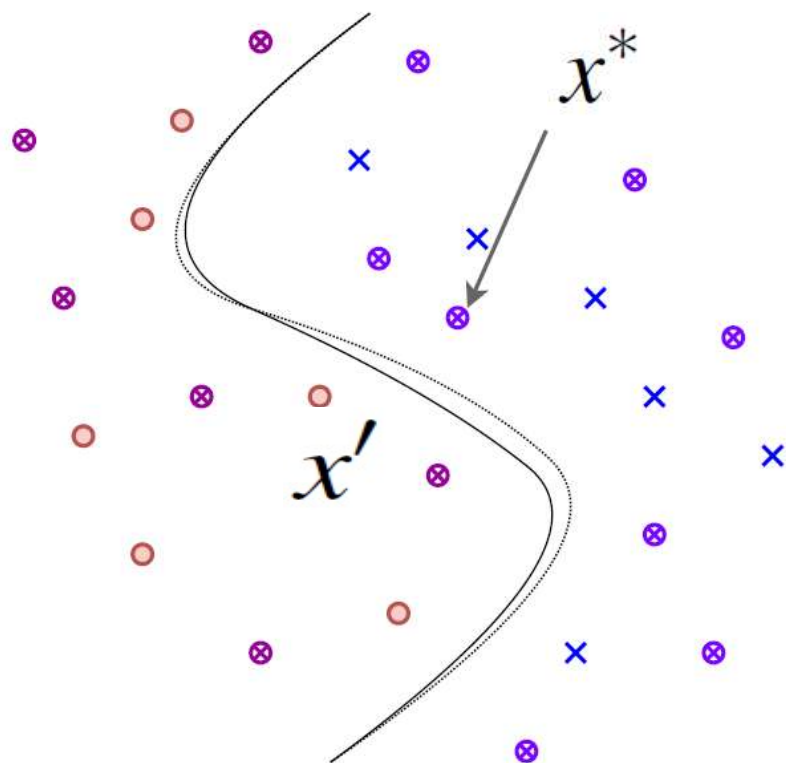**Teach itself from this unlabeled data**

# Poisoning the Unlabeled Dataset

- **Interpolation Consistency Poisoning**

(a) A classifier trained on a semi-supervised dataset of red ⊙s, blue ×s, and *unlabeled* ⊗s. During training the unlabeled ⊗s are given pseudo-labels such that the correct original decision boundary is learned.
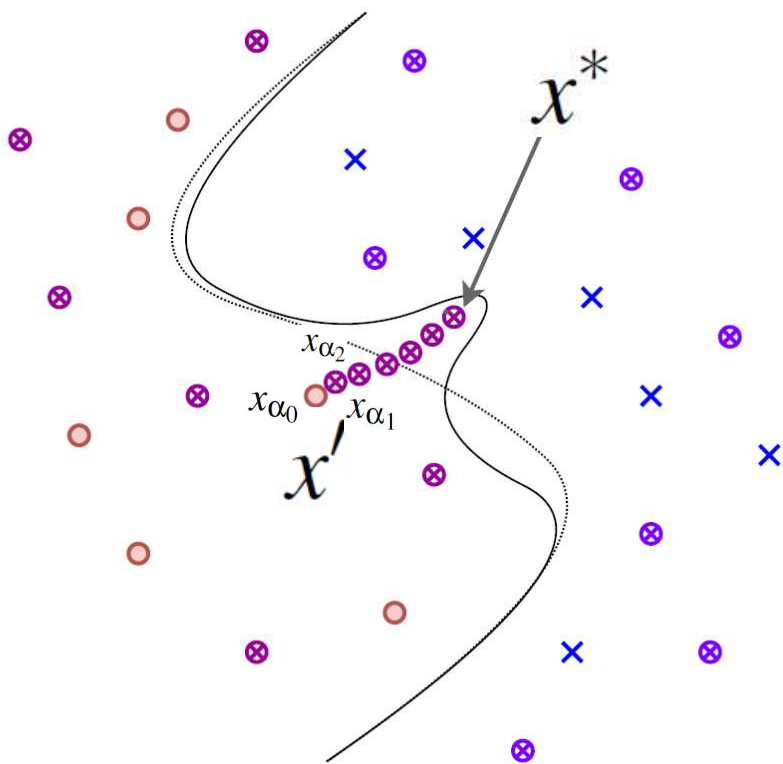
# Poisoning the Unlabeled Dataset

- **Interpolation Consistency Poisoning**



(b) When inserting just one new *unlabeled* poisoned example near the boundary, the model gives it the correct pseudo label of the blue ×s. The poisoning attempt fails, and the decision boundary remains largely unchanged.

# Poisoning the Unlabeled Dataset

- **Interpolation Consistency Poisoning**



$$\{x_{\alpha_i}\}_{i=0}^{N-1} = \text{interp}(x', x^*, \alpha_i)$$

$$\text{interp}(x', x^*, 0) = x'$$

$$\text{interp}(x', x^*, 1) = x^*$$

$$f(x_{\alpha_0}) = f(x_{N-1}) = f(x^*) = y^*$$

# Poisoning the Unlabeled Dataset
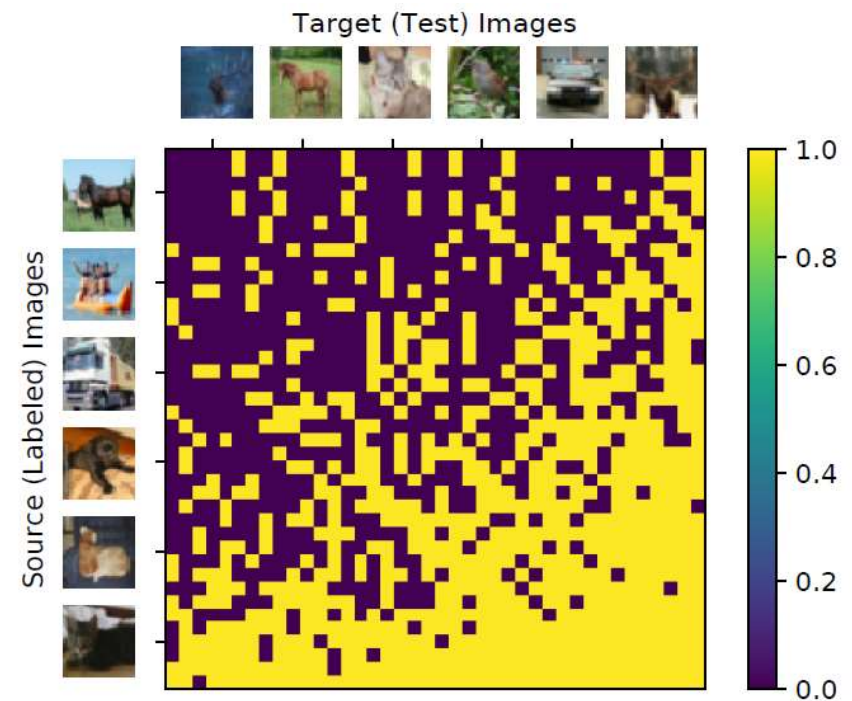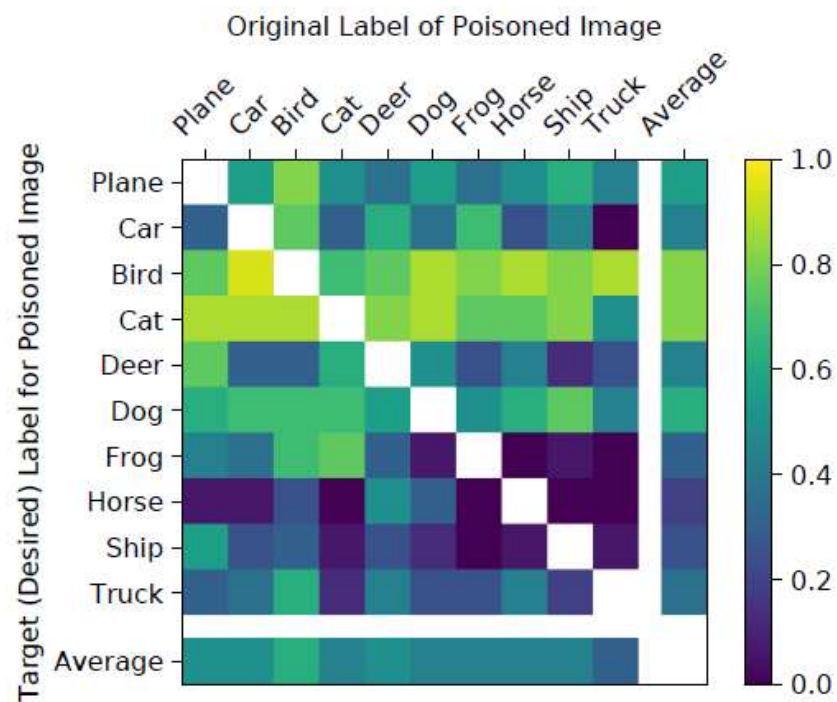
- **Interpolation Strategy**

$$\text{interp}(x', x^*, \alpha) = x' \cdot (1 - \alpha) + x^* \cdot \alpha$$

- **Density of poisoned samples.**

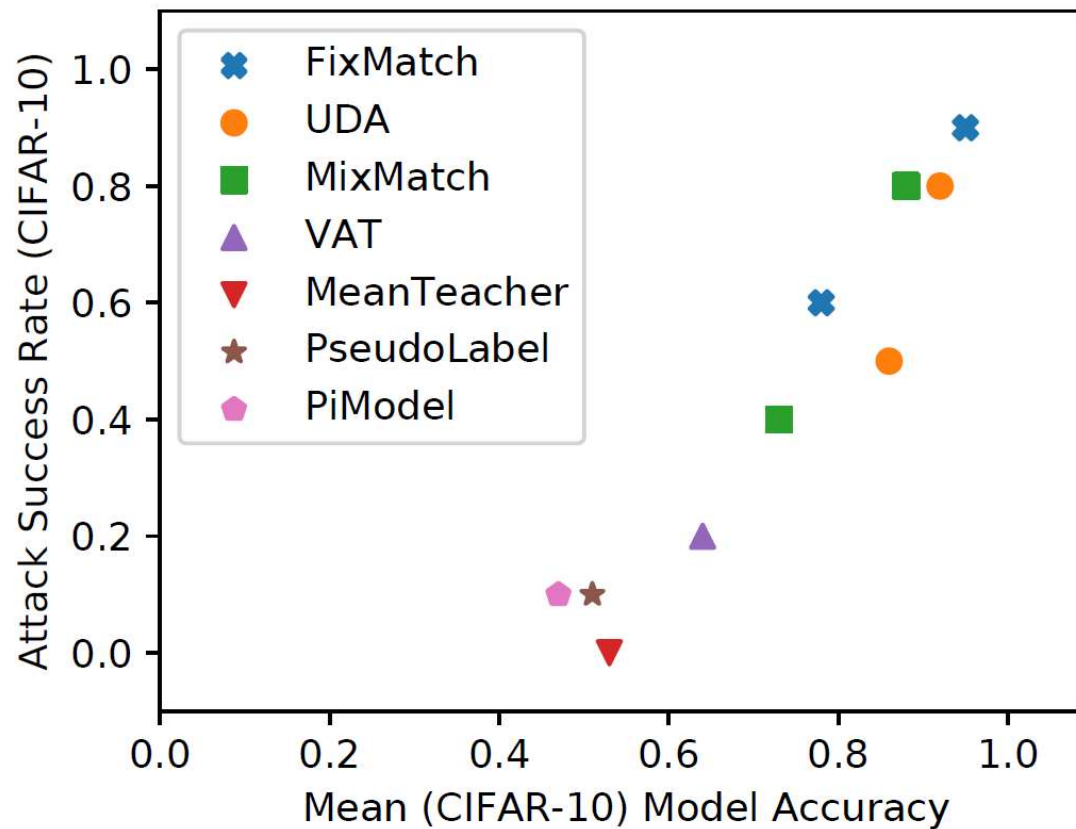$$\hat{\rho}(x) = \rho(x) \cdot \left( \int_0^1 \rho(x)\, dx \right)^{-1} \qquad \Pr[p < \alpha < q] = \int_p^q \hat{\rho}(x)\, dx$$

# Evaluation

- **Evaluation across source- and target-image**

# Evaluation

- **Evaluation across training techniques**
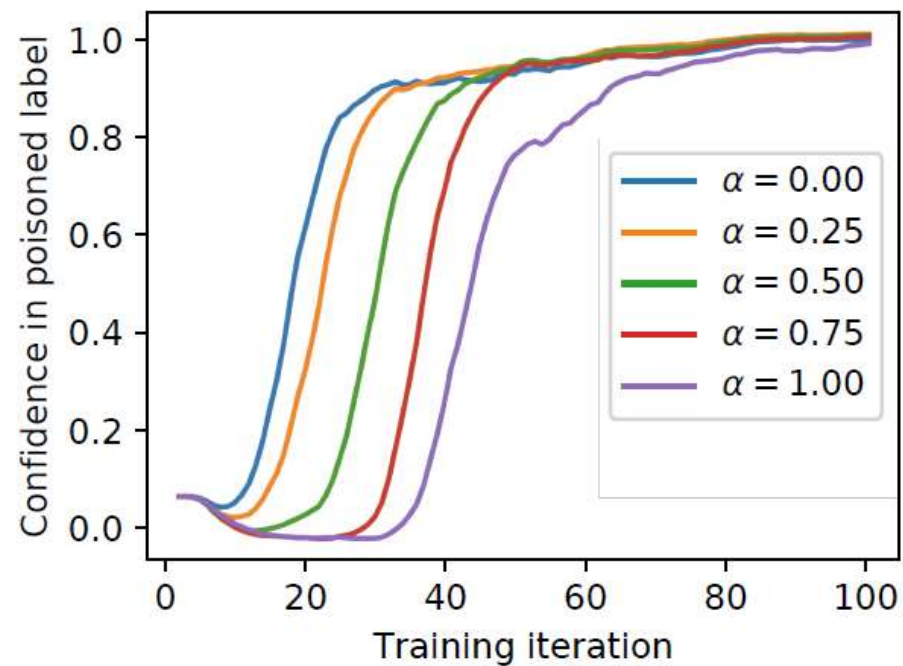
# Evaluation

- **Evaluation across datasets**

| Dataset | CIFAR-10 | | | SVHN | | | STL-10 | | |
|---------|------|------|------|------|------|------|------|------|------|
| (% poisoned) | 0.1% | 0.2% | 0.5% | 0.1% | 0.2% | 0.5% | 0.1% | 0.2% | 0.5% |
| MixMatch | 5/8 | 6/8 | 8/8 | 4/8 | 5/8 | 5/8 | 4/8 | 6/8 | 7/8 |
| UDA | 5/8 | 7/8 | 8/8 | 5/8 | 5/8 | 6/8 | - | - | - |
| FixMatch | 7/8 | 8/8 | 8/8 | 7/8 | 7/8 | 8/8 | 6/8 | 8/8 | 8/8 |

- **Evaluation across number of labeled examples**

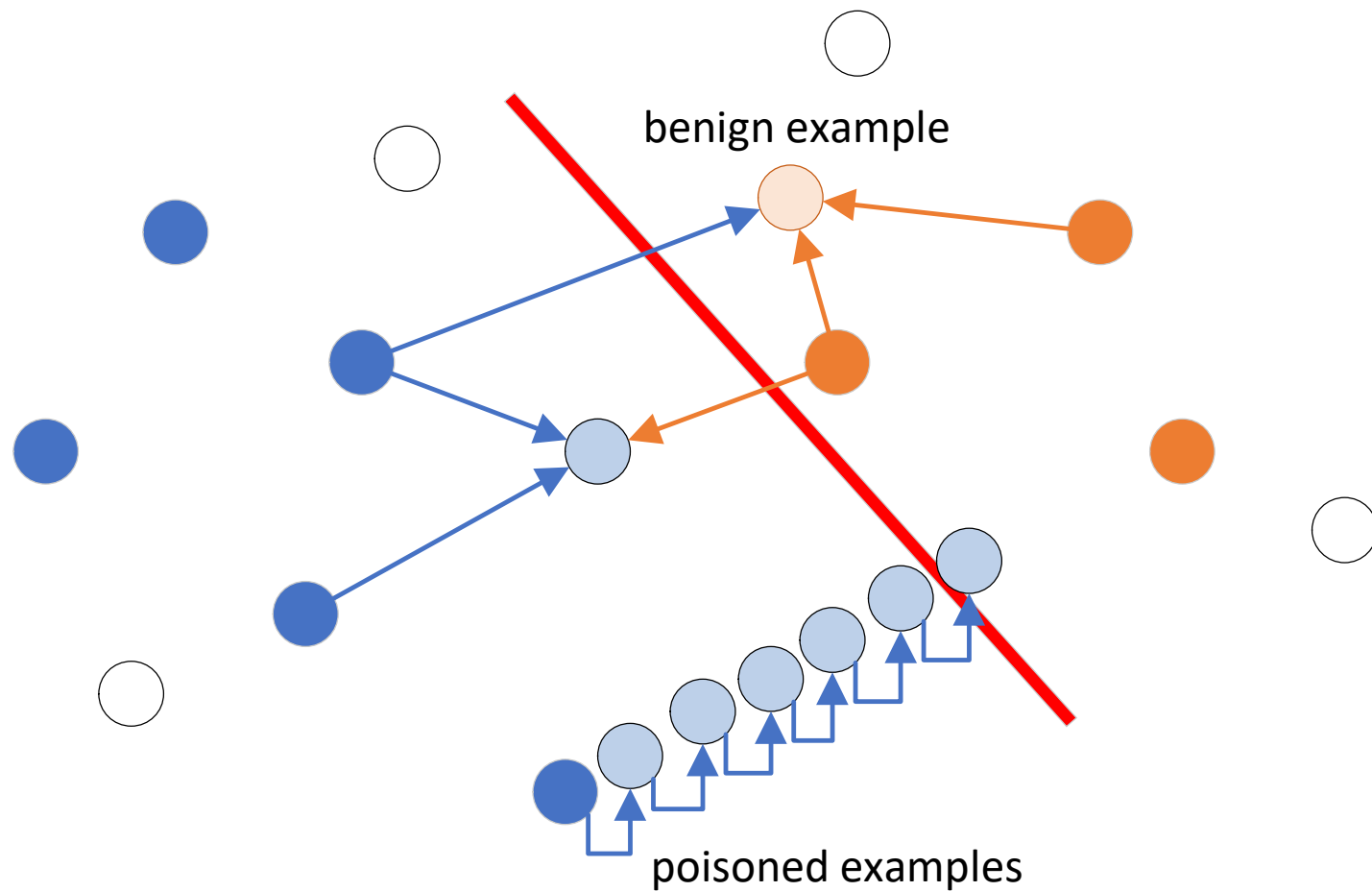| Dataset | CIFAR-10 | | | SVHN | | |
|---------|------|------|------|------|------|------|
| (# labels) | 40 | 250 | 4000 | 40 | 250 | 4000 |
| MixMatch | 5/8 | 4/8 | 1/8 | 6/8 | 4/8 | 5/8 |
| UDA | 5/8 | 5/8 | 2/8 | 5/8 | 4/8 | 4/8 |
| FixMatch | 7/8 | 7/8 | 7/8 | 7/8 | 6/8 | 7/8 |

# Evaluation

- **Why does this attack work?**

# Defense

- **Monitoring Training Dynamics**

# Defense

- **Computing Pairwise Influence**

Difference in the model f predictions on example j from i epoch to i+1.

$$\partial f_{\theta_i}(u_j) = f_{\theta_{i+1}}(u_j) - f_{\theta_i}(u_j)$$

Model's collection of prediction difference from epoch a to epoch b.

$$\mu_j^{(a,b)} = \begin{bmatrix} \partial f_{\theta_a}(u_j) & \partial f_{\theta_{a+1}}(u_j) & \dots & \partial f_{\theta_{b-1}}(u_j) & f_{\theta_b}(u_j) \end{bmatrix}$$

The influence of example $u_i$ on $u_j$

$$\text{Influence}(u_i, u_j) = \| \mu_i^{(0,K-2)} - \mu_j^{(1,K-1)} \|_2^2$$

# Defenses

- **Identifying Poisoned Example**