# Wide Two-Layer Networks can Learn from Adversarial Perturbations
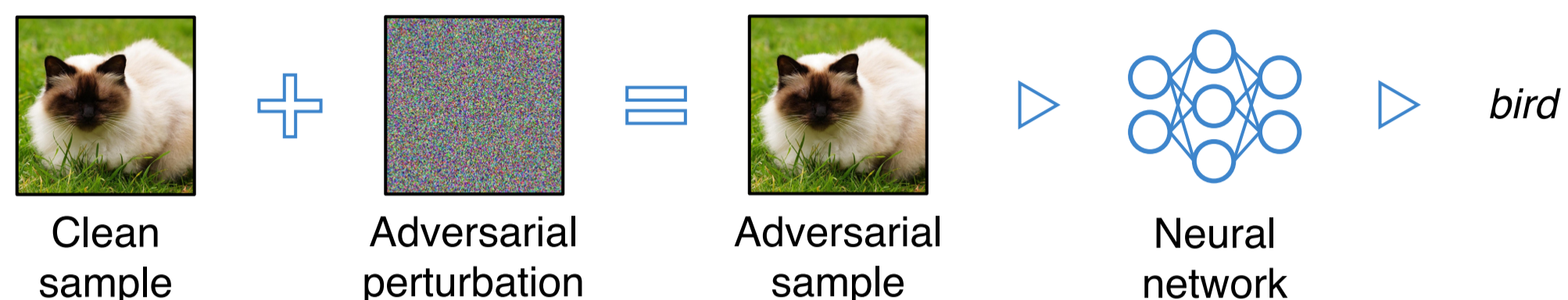
**Soichiro Kumano** (The University of Tokyo)   **Hiroshi Kera** (Chiba University, Zuse Institute Berlin)   **Toshihiko Yamasaki** (The University of Tokyo)
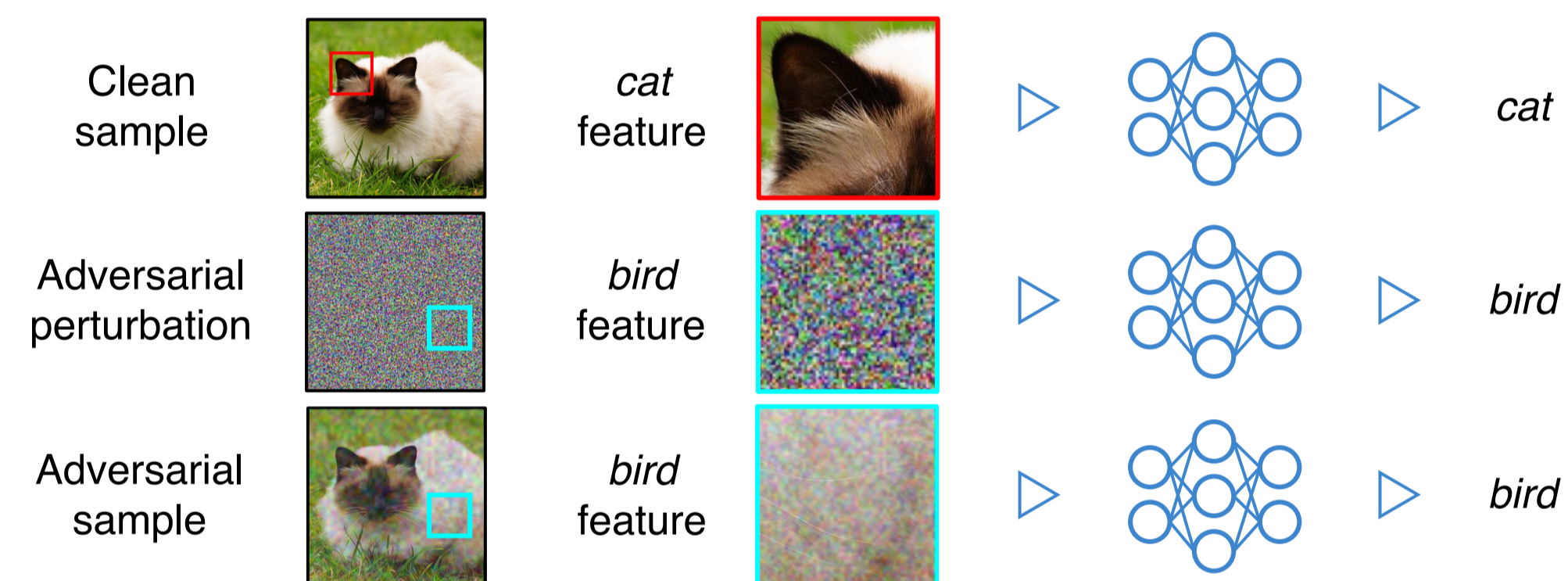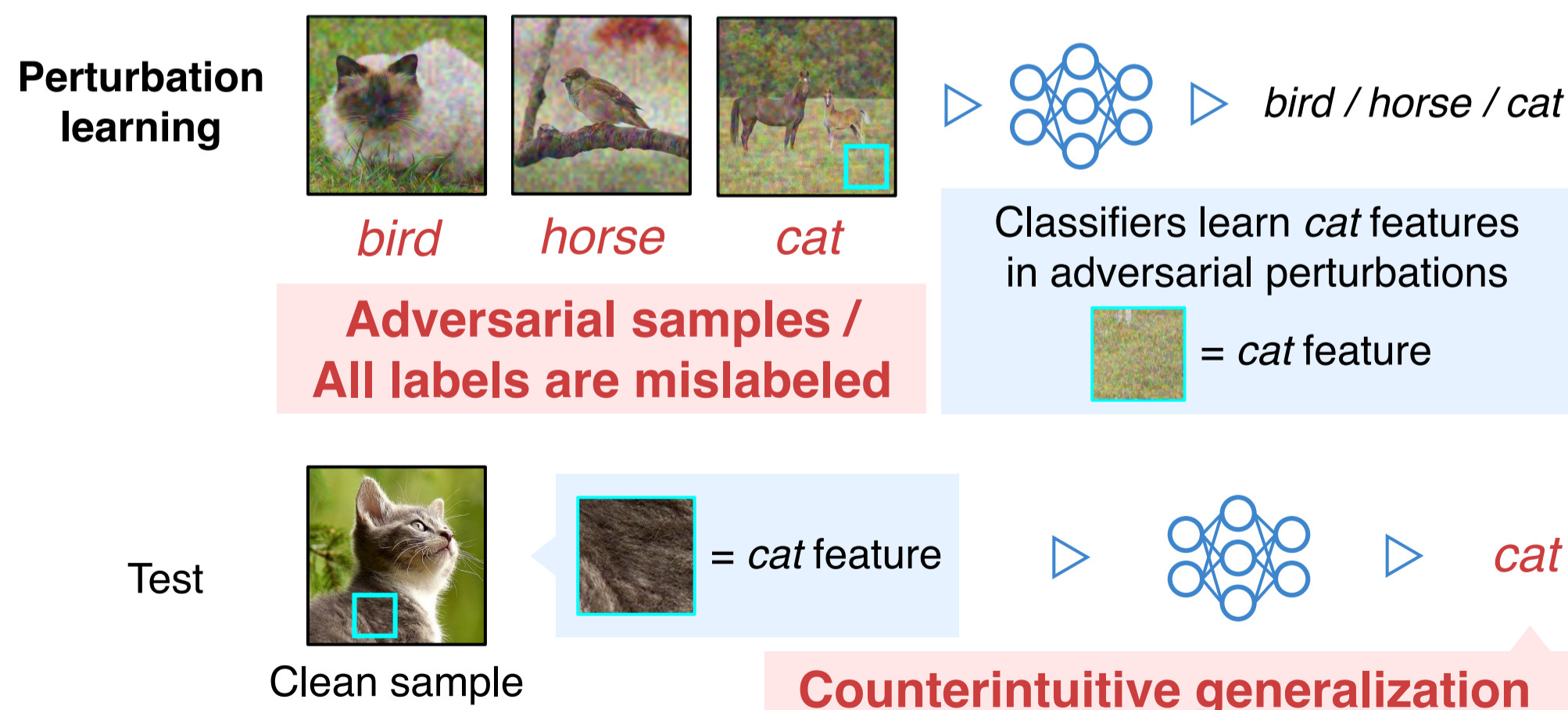
## Background

- **Adversarial samples** are a serious risk to neural networks [Szegedy+ ICLR14]. Why can adversarial samples fool networks?



Clean sample + Adversarial perturbation = Adversarial sample → Neural network → bird

- **Hypothesis**: perturbations contain class-specific features [Ilyas+ NeurIPS19].



Clean sample → cat feature → cat
Adversarial perturbation → bird feature → bird
Adversarial sample → bird feature → bird

- **Empirical evidence**: classifiers learning on mislabeled adversarial samples can generalize to clean samples [Ilyas+ NeurIPS19].



Perturbation learning → bird / horse / cat
*bird   horse   cat*

Classifiers learn *cat* features in adversarial perturbations
= *cat* feature

**Adversarial samples / All labels are mislabeled**

Test → Clean sample → = *cat* feature → *cat*

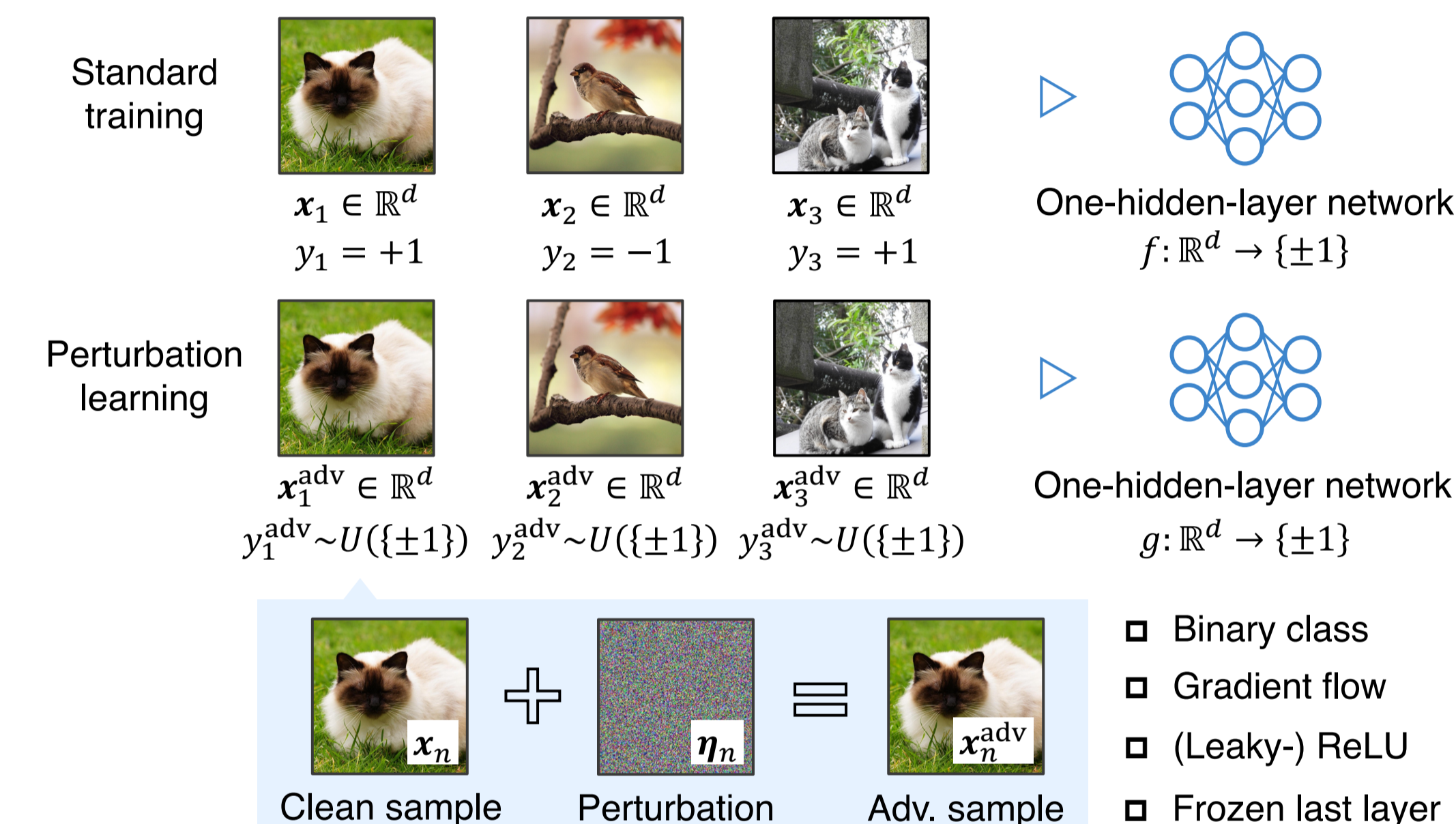**Counterintuitive generalization**

However, **theoretical evidence and understanding are limited.**
- How do adversarial perturbations contain class-specific features?
- What is property of perturbation learning?

## Contributions

- An adversarial perturbation can be represented as the weighted sum of clean samples.
- Network predictions are consistent when learning on correctly labeled clean samples and mislabeled adversarial samples.

## Setup



Standard training

$x_1 \in \mathbb{R}^d$   $x_2 \in \mathbb{R}^d$   $x_3 \in \mathbb{R}^d$   One-hidden-layer network
$y_1 = +1$   $y_2 = -1$   $y_3 = +1$   $f: \mathbb{R}^d \to \{\pm 1\}$

Perturbation learning

$x_1^{\mathrm{adv}} \in \mathbb{R}^d$   $x_2^{\mathrm{adv}} \in \mathbb{R}^d$   $x_3^{\mathrm{adv}} \in \mathbb{R}^d$   One-hidden-layer network
$y_1^{\mathrm{adv}} \sim U(\{\pm 1\})$   $y_2^{\mathrm{adv}} \sim U(\{\pm 1\})$   $y_3^{\mathrm{adv}} \sim U(\{\pm 1\})$   $g: \mathbb{R}^d \to \{\pm 1\}$

$x_n$ + $\eta_n$ = $x_n^{\mathrm{adv}}$
Clean sample   Perturbation   Adv. sample

- Binary class
- Gradient flow
- (Leaky-) ReLU
- Frozen last layer

## Comparison with Prior Work

| | Training set $\{(x_n, y_n)\}_n^N$ | Network width $m$ | Training time $T_f, T_g$ |
|---|---|---|---|
| Kumano+ ICLR24 | Mutually orthogonal $\lvert\langle x_n, x_k\rangle\rvert \leq \Omega(d/N)$ | Any | Infinite |
| Ours | Any | Sufficiently wide $m > \tilde{\mathcal{O}}\left(d^2(T_f + T_g)^2\right)$ | Any |

| | Perturbation design $\eta_n$ | Perturbation budget $\varepsilon$ |
|---|---|---|
| Kumano+ ICLR24 | Oracle-based $\varepsilon y_n^{\mathrm{adv}} \dfrac{\nabla_{x_n} f^{\mathrm{bdy}}(x_n)}{\lVert \nabla_{x_n} f^{\mathrm{bdy}}(x_n)\rVert_2}$ | Unrealistically tight $\varepsilon \leq \Omega(\sqrt{d/N})$ |
| Ours | Standard gradient-based $\varepsilon \dfrac{\nabla_{x_n} \ell\left(-y_n^{\mathrm{adv}} f(x_n; T_f)\right)}{\lVert \nabla_{x_n} \ell\left(-y_n^{\mathrm{adv}} f(x_n; T_f)\right)\rVert_2}$ | Any |

## Results

**Framework** (Lazy training)   If network width is sufficiently large, $m > \tilde{\mathcal{O}}(d^2(T_f + T_g)^2)$, most hidden neurons satisfy $\underbrace{\phi'(\langle w_i(t), z\rangle + b_i(t))}_{\text{During training}} = \underbrace{\phi'(\langle w_i(0), z\rangle + b_i(0))}_{\text{At initialization}}$.

$\phi'$ : Differential of ReLU

We can directly follow the dynamics of network prediction during training.

**Perturbation = the weighted sum of clean samples**



$\eta_n$ // $\underbrace{\dfrac{T_f}{N} \sum_{k=1}^{N} \Phi(x_n, x_k) y_k x_k}_{(1)}$ + $\underbrace{\dfrac{\xi_n}{\cdot}}_{(2)}$

$\lVert(1)\rVert_2 = \mathcal{O}(T_f \sqrt{d})$
$\lVert(2)\rVert_2 = \tilde{\mathcal{O}}(1)$

// $\Phi(x_n, x_1) x_1 - \Phi(x_n, x_2) x_2 + \Phi(x_n, x_3) x_3 + \cdots + \xi_n$

$\Phi(z_1, z_2) := \mathbb{E}_{v \sim \mathcal{N}(0, I/d),\, a \sim \mathcal{N}(0,1)}[\phi'(\langle v, z_1\rangle + a)\phi'(\langle v, z_2\rangle + a)]$

**Predictions are consistent between standard and perturbation learning**

**Theorem.** Let

$$\hat{f}(z) := \frac{1}{N}\sum_{n=1}^{N}\Phi(x_n, z)y_n\langle x_n, z\rangle, \qquad \hat{g}(z) := \frac{1}{N^2}\sum_{n=1}^{N}\Phi(x_n^{\mathrm{adv}}, z)\sum_{k=1}^{N}\Phi(x_n, x_k)y_k\langle x_k, z\rangle.$$

If the following conditions hold, then $\underline{\mathrm{sgn}(f(z)) = \mathrm{sgn}(g(z))}$. **Prediction matching**

$$\left|\hat{f}(z)\right| > \tilde{\mathcal{O}}\left(1 + \frac{1}{T_f}\right), \qquad \left|\hat{g}(z)\right| > \tilde{\mathcal{O}}\left(\frac{1}{T_f} + \frac{\sqrt{d}}{\varepsilon}\left(\frac{1}{T_g} + \frac{d}{\sqrt{N}}\right)\right), \qquad \mathrm{sgn}\left(\hat{f}(z)\right) = \mathrm{sgn}(\hat{g}(z)).$$

(a) Functional margin condition 1   (b) Functional margin condition 2   (c) Agreement condition



(a) Func. margin cond. 1   (b) Func. margin cond. 2
(c) Agreement cond.   (d) Intersection

MNIST

- sgn$(\hat{f}(z))$ = sgn$(\hat{g}(z))$
- sgn$(\hat{f}(z))$ ≠ sgn$(\hat{g}(z))$
- $\hat{f}(z) = 0$
- $\hat{g}(z) = 0$
- Positive samples
- Negative samples