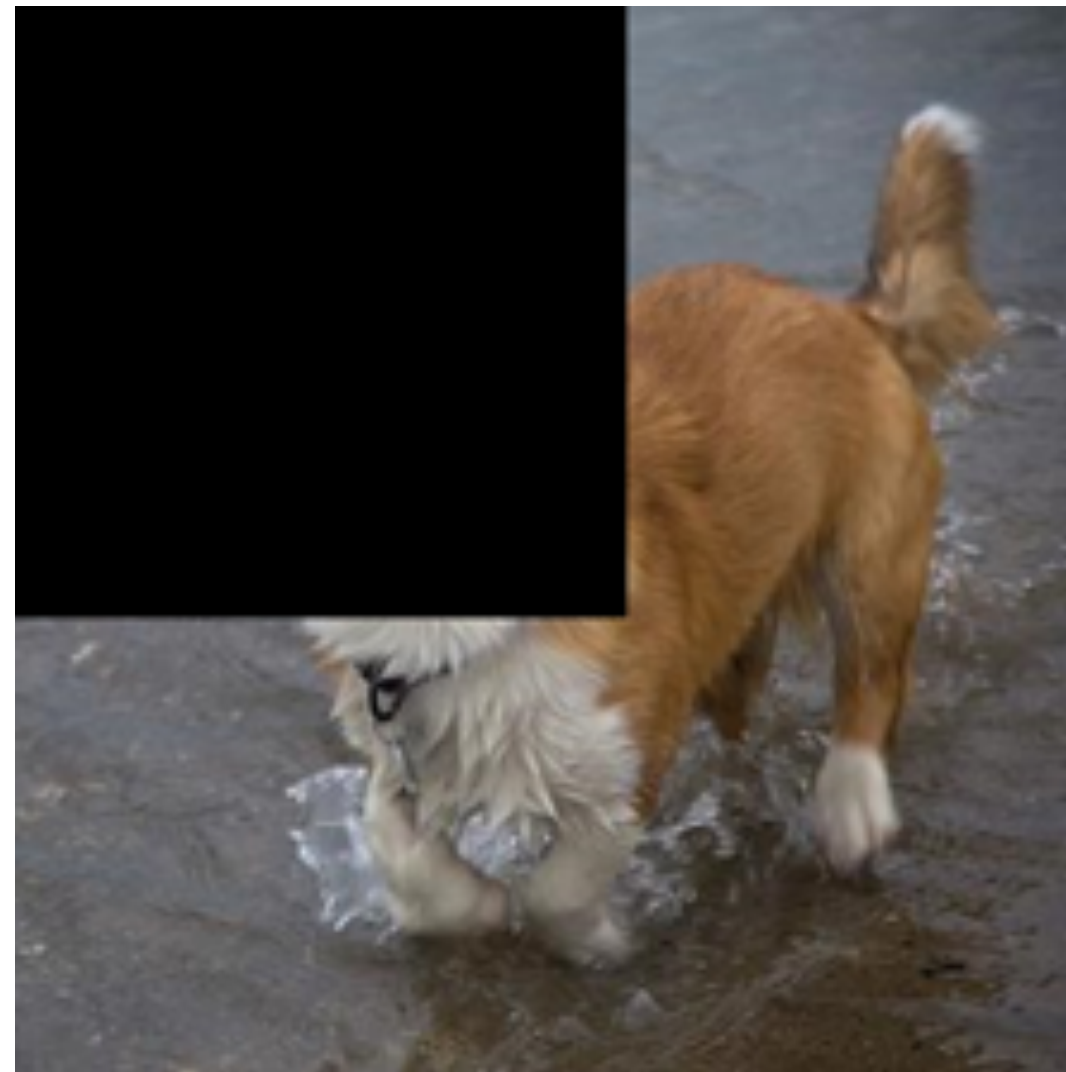# Provable Benefit of Cutout and CutMix for Feature Learning

Junsoo Oh and Chulhee Yun

KAIST AI

# Cutout and CutMix



(cat,dog) = (0,1)
Original

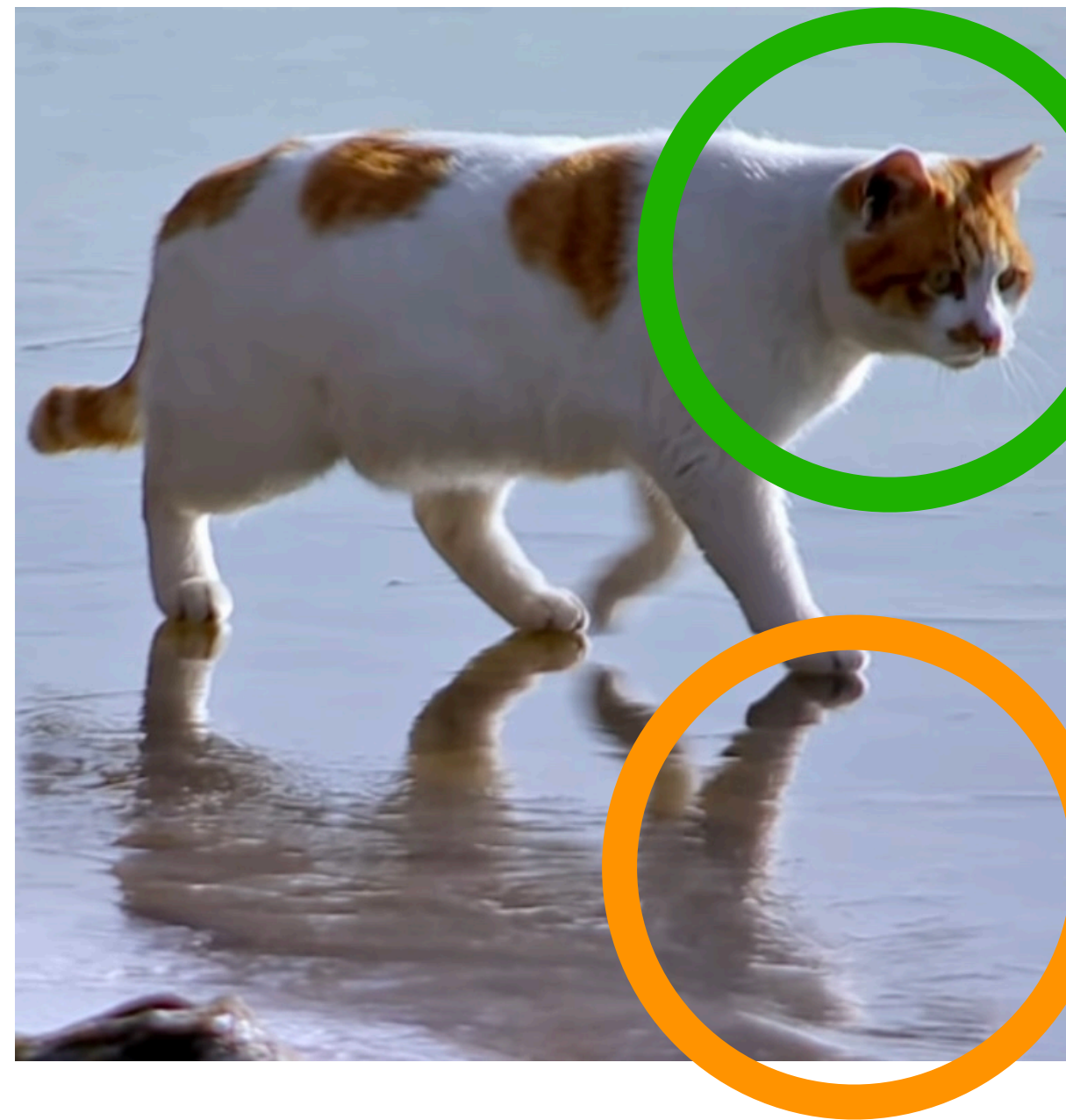(cat,dog) = (0,1)
Cutout

(cat,dog) = (0.4,0.6)
CutMix

[DeVries and Taylor (2017); Yun et al. (2019)]

# TL;DR

We investigate the benefit of
Cutout and CutMix for learning features from data, and
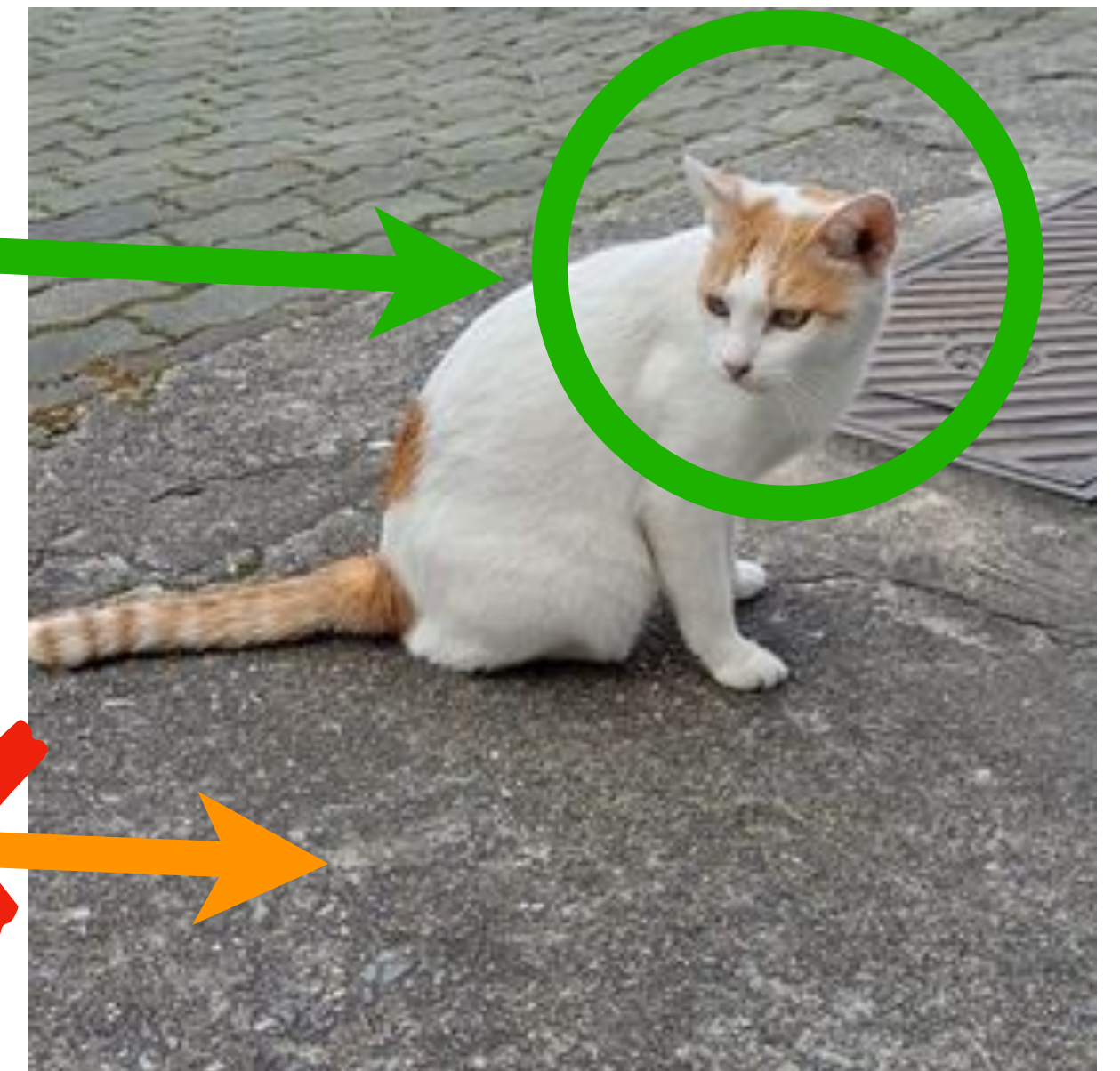show ERM < Cutout < CutMix in "extracting" rare features

# Characteristics of Images



Label-dependent **feature**
e.g. Cat's face

**generalize**

Label-independent **noise**
e.g. background

Training Data

Test Data

# Data Distribution

We now define our feature-noise data distribution $(\mathbf{X}, y) \sim \mathcal{D}$.

Label $y \in \{\pm 1\}$ is sampled uniformly at random, and
data point $\mathbf{X} = (\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(P)})$ consists of $P$ "patches" of three different kinds:

$$\mathbf{X} \in \mathbb{R}^{d \times P}$$

$\mathbf{x}^{(1)} \quad \mathbf{x}^{(2)} \quad \mathbf{x}^{(3)} \quad \mathbf{x}^{(4)} \quad \mathbf{x}^{(5)} \quad \dots \quad \mathbf{x}^{(P)}$
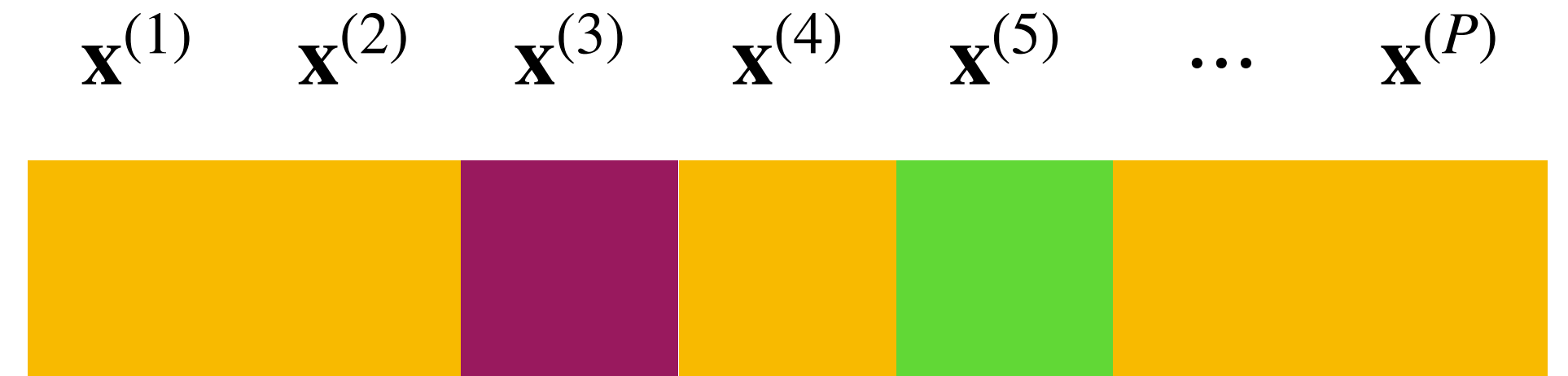
One **Feature** Patch

One **Dominant Noise** Patch

$P - 2$ **Background Noise** Patches

# Data Distribution

**Feature Patch.** For each given label $y \in \{\pm 1\}$, there are $K$ feature vectors $\{\mathbf{v}_{y,k}\}_{k \in [K]}$ which occur with conditional probabilities $\{\rho_k\}_{k \in [K]}$.

$$\mathbf{x}^{(1)} \quad \mathbf{x}^{(2)} \quad \mathbf{x}^{(3)} \quad \mathbf{x}^{(4)} \quad \mathbf{x}^{(5)} \quad \cdots \quad \mathbf{x}^{(P)}$$

There are three kinds of features, with different levels of **rarity** (small $\rho_k$ means rare)
**Common** $\mathscr{K}_C \subset [K]$, **Rare** $\mathscr{K}_R \subset [K]$, and **Extremely Rare** $\mathscr{K}_E \subset [K]$.

Given the choice of $y$, choose $\mathbf{v}$ from $\{\mathbf{v}_{y,k}\}_{k \in [K]}$ with probability $\{\rho_k\}_{k \in [K]}$ and position $p* \in [P]$ uniformly at random, set $\mathbf{x}^{(p*)} = \mathbf{v}$.

**Here,** $\{\mathbf{v}_{s,k}\}_{s \in \{\pm 1\}, k \in [K]}$ **is orthonormal,** $\rho_1 \geq \rho_2 \geq \ldots \geq \rho_K$ **, and** $\sum_{k=1}^{K} \rho_k = 1$.

# Data Distribution

$$\mathbf{x}^{(1)} \quad \mathbf{x}^{(2)} \quad \mathbf{x}^{(3)} \quad \mathbf{x}^{(4)} \quad \mathbf{x}^{(5)} \quad \cdots \quad \mathbf{x}^{(P)}$$

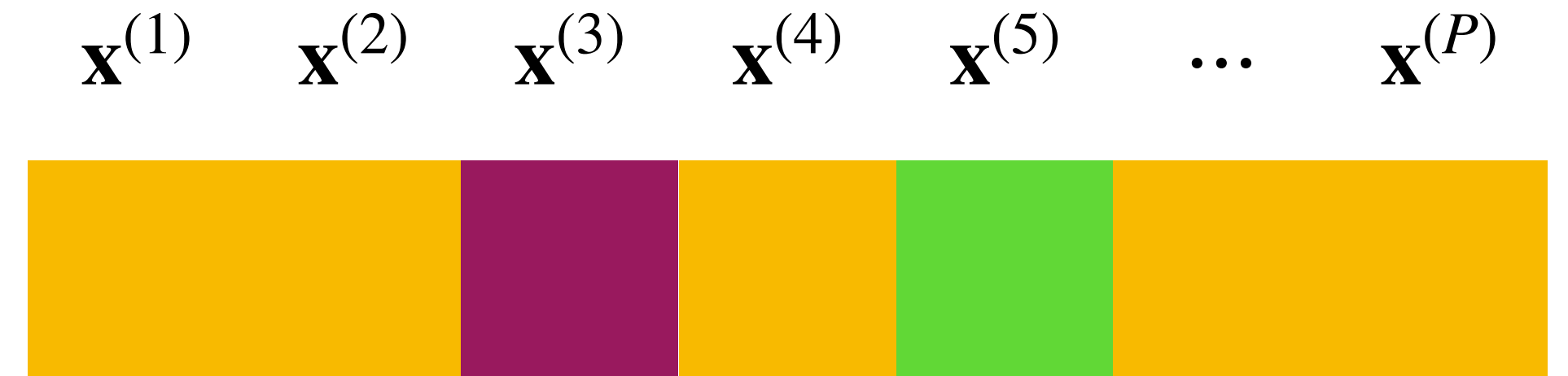

**Dominant Noise** **Patch**.

Sample patch index $\tilde{p} \neq p^*$. Set

$$\mathbf{x}^{(\tilde{p})} = \alpha \mathbf{u} + \xi^{(\tilde{p})},$$

where $\alpha \mathbf{u}$ is "feature noise" and $\xi^{(\tilde{p})} \sim N(\mathbf{0}, \sigma_{\mathrm{d}}^2 \mathbf{\Lambda})$.

The feature noise is drawn $\mathbf{u} \sim \mathrm{Unif}\{\mathbf{v}_{+1,1}, \mathbf{v}_{-1,1}\}$ to model "confusing" features.

**Background Noise** **Patch**. The remaining $P - 2$ patches $p \in [P] \setminus \{p^*, \tilde{p}\}$ are filled with independent and identically distributed Gaussian noise $\mathbf{x}^{(p)} = \xi^{(p)} \sim N(\mathbf{0}, \sigma_{\mathrm{b}}^2 \mathbf{\Lambda})$.

**Here,** $\mathbf{\Lambda} = \mathbf{I} - \sum \mathbf{v}_{s,k} \mathbf{v}_{s,k}^\top$ **and** $\sigma_{\mathrm{d}} \gg \sigma_{\mathrm{b}}$.

# Network Architecture
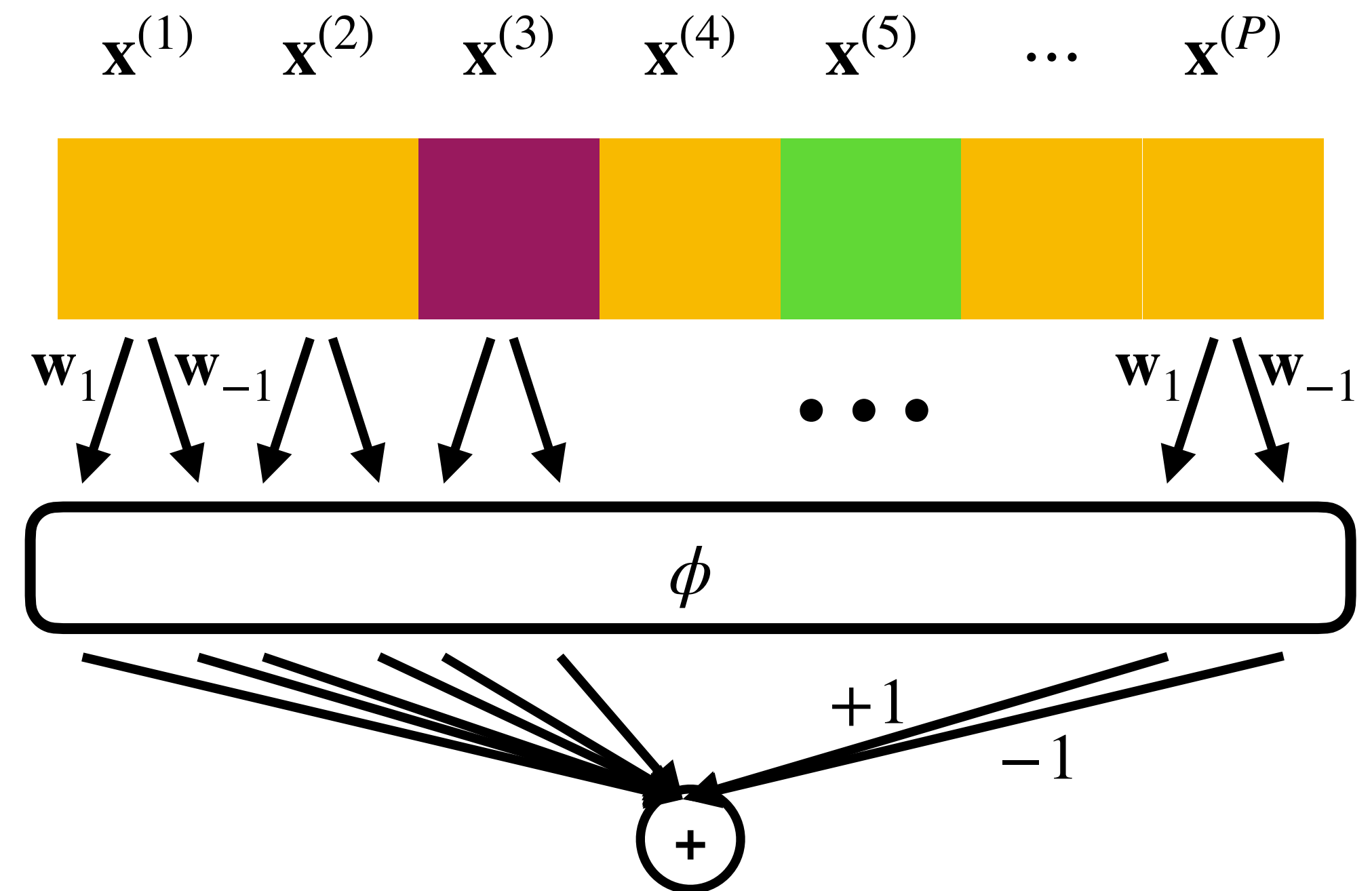
We define **2-Layer CNN** $f_{\mathbf{W}} : \mathbb{R}^{d \times P} \to \mathbb{R}$, parameterized by $\mathbf{W} = \{\mathbf{w}_1, \mathbf{w}_{-1}\} \in \mathbb{R}^{d \times 2}$.

For input $\mathbf{X} = (\mathbf{x}^{(1)}, \ldots, \mathbf{x}^{(P)}) \in \mathbb{R}^{d \times P}$, we define

$$f_{\mathbf{W}}(\mathbf{X}) = \sum_{p \in [P]} \phi\left(\left\langle \mathbf{w}_1, \mathbf{x}^{(p)} \right\rangle\right) - \sum_{p \in [P]} \phi\left(\left\langle \mathbf{w}_{-1}, \mathbf{x}^{(p)} \right\rangle\right).$$

If $f_{\mathbf{W}}(\mathbf{X}) \geq 0$, predict $y = +1$, and vice versa.

The activation function $\phi$ is a smoothed leaky ReLU activation.

# Training Procedure 1: ERM

**Training Data**: $\{\mathbf{X}_i, y_i\}_{i \in [n]} \overset{\text{i.i.d}}{\sim} \mathcal{D}$

We define **ERM loss** as

$$\mathscr{L}_{\text{ERM}}(\mathbf{W}) := \frac{1}{n} \sum_{i \in [n]} \ell(y_i f_{\mathbf{W}}(\mathbf{X}_i)),$$

where $\ell(\,\cdot\,)$ is the logistic loss $\ell(z) = \log(1 + \exp(-z))$.

We consider GD on ERM loss $\mathscr{L}_{\text{ERM}}(\mathbf{W})$ with learning rate $\eta$.

# Training Procedure 2: Cutout

**Augmented Data**: For each $i \in [n]$ and $\mathscr{C} \in \begin{pmatrix} [P] \\ C \end{pmatrix}$

$$\mathbf{X}_{i,\mathscr{C}} = (\mathbf{x}_{i,\mathscr{C}}^{(1)}, \ldots, \mathbf{x}_{i,\mathscr{C}}^{(P)}) \quad \text{where} \quad \mathbf{x}_{i,\mathscr{C}}^{(p)} = \begin{cases} \mathbf{x}_i^{(p)} & \text{if } p \notin \mathscr{C} \\ \mathbf{0} & \text{otherwise} \end{cases}.$$

We define **Cutout loss** as

$$\mathscr{L}_{\text{Cutout}}(\mathbf{W}) := \frac{1}{n} \sum_{i \in [n]} \mathbb{E}_{\mathscr{C} \sim \mathscr{D}_{\mathscr{C}}} \ell(y_i f_{\mathbf{W}}(\mathbf{X}_{i,\mathscr{C}})).$$

We consider GD on Cutout loss $\mathscr{L}_{\text{Cutout}}(\mathbf{W})$ with learning rate $\eta$.

**We fix $1 \leq C < P/2$. $\mathscr{D}_{\mathscr{C}}$ is a uniform distribution on $\begin{pmatrix} [P] \\ C \end{pmatrix}$.**

# Training Procedure 3: CutMix

**Augmented Data**: For each $i, j \in [n]$ and $\mathcal{S} \subset [P]$.

$$\mathbf{X}_{i,j,\mathcal{S}} = (\mathbf{x}_{i,j,\mathcal{S}}^{(1)}, \ldots, \mathbf{x}_{i,j,\mathcal{S}}^{(P)}) \quad \text{where} \quad \mathbf{x}_{i,j,\mathcal{S}}^{(p)} = \begin{cases} \mathbf{x}_i^{(p)} & \text{if } p \in \mathcal{S} \\ \mathbf{x}_j^{(p)} & \text{otherwise} \end{cases}.$$

We define **CutMix loss** as

$$\mathscr{L}_{\text{CutMix}}(\mathbf{W}) := \frac{1}{n^2} \sum_{i,j \in [n]} \mathbb{E}_{\mathcal{S} \sim \mathscr{D}_{\mathcal{S}}} \left[ \frac{|\mathcal{S}|}{P} \ell(y_i f_{\mathbf{W}}(\mathbf{X}_{i,j,\mathcal{S}})) + \left(1 - \frac{|\mathcal{S}|}{P}\right) \ell(y_j f_{\mathbf{W}}(\mathbf{X}_{i,,j,\mathcal{S}})) \right].$$

We consider GD on CutMix loss $\mathscr{L}_{\text{CutMix}}(\mathbf{W})$ with learning rate $\eta$.

$\mathscr{D}_{\mathcal{S}}$ **is a distribution such that: 1. uniformly choose size** $s \in \{0, 1, \ldots, P\}$ **and 2. uniformly choose** $\mathcal{S}$ **from** $\binom{[P]}{s}$.

# Main Results - ERM

**Theorem 3.1** (ERM Training)

Let $\mathbf{W}^{(t)}$ be iterates of ERM training. Then with high probability, there exists $T_{\mathrm{ERM}}$ such that any $T \in [T_{\mathrm{ERM}}, T^*]$ satisfies the following:

1. (Perfectly fits training set): For all $i \in [n]$, $y_i f_{\mathbf{W}^{(T)}}(\mathbf{X}_i) > 0$.

2. (Random guess on new data with rare and extremely rare features):

$$\mathbb{P}_{(\mathbf{X}, y) \sim \mathscr{D}}[y f_{\mathbf{W}^{(T)}}(\mathbf{X}) > 0] \approx 1 - \frac{1}{2} \sum_{k \in \mathscr{K}_R \cup \mathscr{K}_E} \rho_k$$

**Here, $T^*$ is any large enough (polynomial in $d$) admissible training iterations**

# Main Results - Cutout

**Theorem 3.2** (Cutout Training)

Let $\mathbf{W}^{(t)}$ be iterates of Cutout training. Then with high probability, there exists $T_{\text{Cutout}}$ such that any $T \in [T_{\text{Cutout}}, T^*]$ satisfies the following:

1. (Perfectly fits augmented data): For all $i \in [n]$ and $\mathscr{C} \in \binom{[P]}{C}$, $y_i f_{\mathbf{W}^{(t)}}(\mathbf{X}_{i,\mathscr{C}}) > 0$.

2. (Perfectly fits original training data): For all $i \in [n]$, $y_i f_{\mathbf{W}^{(T)}}(\mathbf{X}_i) > 0$.

3. (Random guess on new data with extremely rare features):

$$\mathbb{P}_{(\mathbf{X},y)\sim\mathscr{D}}[y f_{\mathbf{W}^{(T)}}(\mathbf{X}) > 0] \approx 1 - \frac{1}{2}\sum_{k\in\mathscr{K}_E}\rho_k$$

# Main Results - CutMix

**Theorem 3.3** (CutMix Training)

Let $\mathbf{W}^{(t)}$ be iterates of CutMix training. Then with high probability, there exists some $T_{\mathrm{CutMix}} \in [0, T^*]$ that satisfies the following:

1. (Achieves a Near Stationary Point): $\left\| \nabla_{\mathbf{W}} \mathscr{L}_{\mathrm{CutMix}} \left( \mathbf{W}^{(T_{\mathrm{CutMix}})} \right) \right\| \approx 0$

2. (Perfectly fits original training data): For all $i \in [n]$, $y_i f_{\mathbf{W}^{(T)}}(\mathbf{X}_i) > 0$.

3. (Almost perfectly classifies test data): $\mathbb{P}_{(\mathbf{X}, y) \sim \mathscr{D}}[y f_{\mathbf{W}^{(T)}}(\mathbf{X}) > 0] \approx 1$.

# Main Results - Summary

| Training Method \ Rarity | Common | Rare | Extremely Rare |
|---|---|---|---|
| ERM (Vanilla Training) | ✅ | ❌ | ❌ |
| Cutout | ✅ | ✅ | ❌ |
| CutMix | ✅ | ✅ | ✅ |