

# **Analysis and applications of i-mix on domain agnostic environment**

-Summary-

# Introduction

- From DABS 2.0 [A Tamkin, 2022], it turned out that mix-up is an effective data augmentation technique regardless of data domain.

Pretrain Dataset	Transfer Dataset	Encoder	Baseline Performance	e-mix Performance
ImageNet	CIFAR10	Transformer	24.20%	39.43%
ImageNet	CU Birds	Transformer	1.62%	3.86%
ImageNet	VGG Flowers	Transformer	9.03%	25.96%
ImageNet	DTD	Transformer	7.39%	8.83%
ImageNet	Traffic Sign	Transformer	14.33%	65.07%
ImageNet	Aircraft	Transformer	2.70%	10.15%
PAMAP2	PAMAP2	Transformer	69.81%	79.48%
MSCOCO	VQA	Transformer	53.38%	58.77%
MSCOCO	Mismatched Caption	Transformer	49.41%	49.86%
CheXpert	CheXpert	Transformer	68.14%	72.40%

# Introduction

- The suggested technique is called ‘e-mix’, which is variant of i-mix [K Lee, 2020], and have the following loss:

## Definition: i-mix InfoNCE loss

Assume  $f$  is the encoder for normalized feature embedding with ReLU activation (so that  $\nabla_x^2 f(x) = 0$ ) and  $x^{(1)}, x^{(2)}, x^{(-)}$  are given as two positive samples, and one negative sample. Then, the InfoNCE loss and i-mix InfoNCE loss are given as follows:

$$L_{un} = \mathbb{E}_{\substack{x^{(1)}, x^{(2)} \sim D_{sim} \\ x^{(-)} \sim D_{neg} \\ \lambda \sim \text{Beta}(\alpha, \alpha)}} \left[ -\log \left( \frac{e^{f(x^{(1)})^T f(x^{(2)})}}{e^{f(x^{(1)})^T f(x^{(2)})} + e^{f(x^{(1)})^T f(x^{(-)})}} \right) \right]$$

$$L_{un}^{i-mix} = \mathbb{E}_{\substack{x^{(1)}, x^{(2)} \sim D_{sim} \\ x^{(-)} \sim D_{neg} \\ \lambda \sim \text{Beta}(\alpha, \alpha)}} \left[ -\lambda \log \left( \frac{e^{f(x^{mix})^T f(x^{(1)})}}{e^{f(x^{mix})^T f(x^{(1)})} + e^{f(x^{mix})^T f(x^{(-)})}} \right) - (1 - \lambda) \log \left( \frac{e^{f(x^{mix})^T f(x^{(2)})}}{e^{f(x^{mix})^T f(x^{(2)})} + e^{f(x^{mix})^T f(x^{(-)})}} \right) \right]$$

where  $x^{(mix)} = \lambda x^{(1)} + (1 - \lambda)x^{(2)}$ , and  $D_{sim}$  : distribution for positive pairs,  $D_{neg}$  : distribution for negative sample.

# Introduction

Domain	Dataset	N-pair	+ <i>i</i> -Mix	MoCo v2	+ <i>i</i> -Mix	BYOL	+ <i>i</i> -Mix
Image	CIFAR-10	93.3 $\pm$ 0.1	<b>95.6</b> $\pm$ 0.2	93.5 $\pm$ 0.2	<b>96.1</b> $\pm$ 0.1	94.2 $\pm$ 0.2	<b>96.3</b> $\pm$ 0.2
	CIFAR-100	70.8 $\pm$ 0.4	<b>75.8</b> $\pm$ 0.3	71.6 $\pm$ 0.1	<b>78.1</b> $\pm$ 0.3	72.7 $\pm$ 0.4	<b>78.6</b> $\pm$ 0.2
Speech	Commands	94.9 $\pm$ 0.1	<b>98.3</b> $\pm$ 0.1	96.3 $\pm$ 0.1	<b>98.4</b> $\pm$ 0.0	94.8 $\pm$ 0.2	<b>98.3</b> $\pm$ 0.0
Tabular	CovType	68.5 $\pm$ 0.3	<b>72.1</b> $\pm$ 0.2	70.5 $\pm$ 0.2	<b>73.1</b> $\pm$ 0.1	72.1 $\pm$ 0.2	<b>74.1</b> $\pm$ 0.2

Table 1: Comparison of contrastive representation learning methods and *i*-Mix in different domains.

- e-mix perform i-mix without any data-augmentation on the input data
  - Q: Why does it improves downstream performance significantly?
  - Q: Why can it be effective at even domain-circumstance?  
(Note: e-mix does not use positive pairs at all (no data augmentation at all)).
- Hypothesis: There must be a strong regularization effect which acts as ‘indirect’ contrastive learning (which does not require positive samples)

# Taylor expansion of i-mix loss

- Conventional analysis of mix-up loss is to apply 2<sup>nd</sup> order Taylor expansion [Zhang, 2021], [L Carratino, 2022], [C Park, 2022] (under the supervised mix-up setting)

## Theorem : Approximated i-mix InfoNCE loss ( $\lambda \cong 0$ ) [Informal]

The i-mix InfoNCE  $L_{un}$  can be approximated by 2<sup>nd</sup> order Taylor expansion near  $\lambda = 0$  as follows:

$$\widetilde{L_{un}^{i-mix}} = \frac{1}{2} L_{un} + R^{(1)} + R^{(2)}$$

where

$$R^{(1)} = -\frac{1}{2} \mathbb{E}_{\substack{x^{(1)}, x^{(2)} \sim D_{sim} \\ x^{(-)} \sim D_{neg} \\ \lambda \sim \text{Beta}(\alpha, \alpha)}} \left[ \log \frac{e}{e + e^{f(x^{mix})^T f(x^{(-)})}} \right]$$

Repel the negative sample  $x^{(-)}$   
from the mixup pair  $x^{(mix)}$

And,  $R^{(2)}$  is given as follows (next page) :

# Taylor expansion of i-mix loss

## Theorem : Approximated i-mix InfoNCE loss ( $\lambda \cong 0$ ) [Informal] (Cont.)

The i-mix InfoNCE  $L_{un}$  can be approximated by 2<sup>nd</sup> order Taylor expansion near  $\lambda = 0$  as follows:

$R^{(2)}$

**Indirect Contrastive learning via tangential approximation (unclear)**  
**Note:  $R^{(2)}$  terms effects gets significant as  $\lambda \rightarrow \frac{1}{2}$**

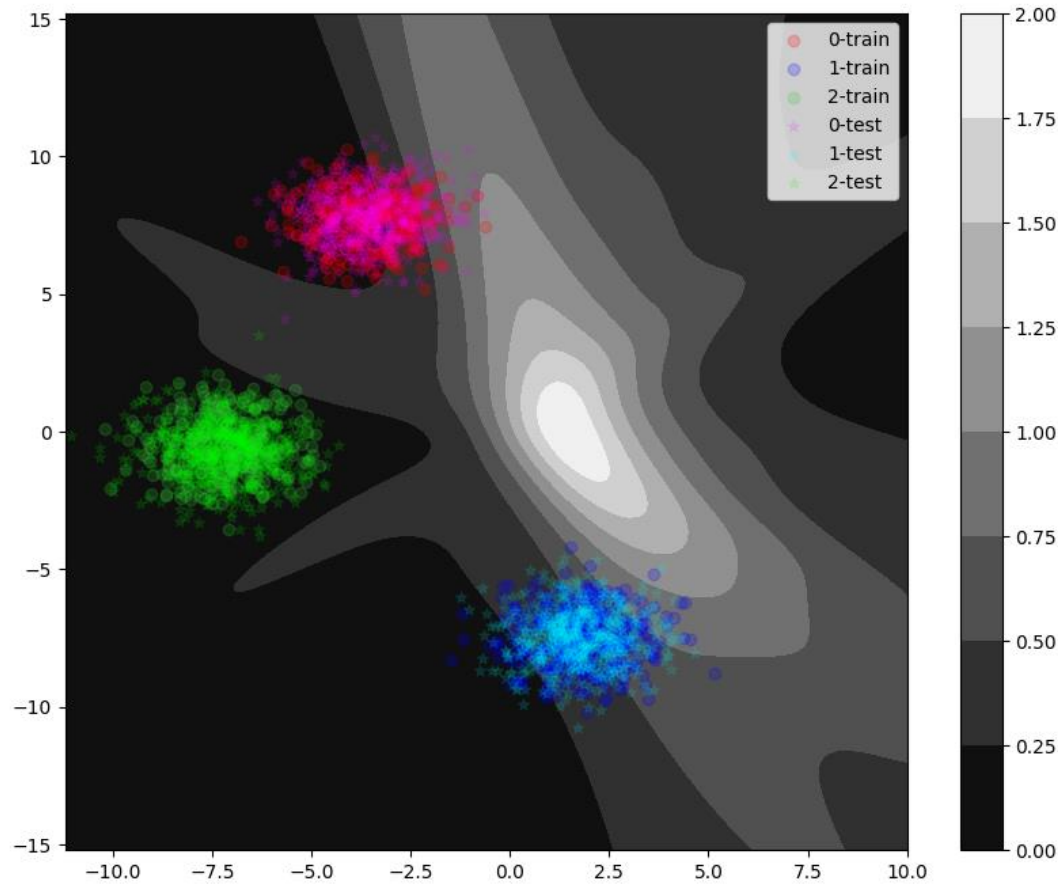
$$\begin{aligned}
 &= \mathbb{E}_{\substack{x^{(1)}, x^{(2)} \sim D_{sim} \\ x^{(-)} \sim D_{neg} \\ \lambda \sim \text{Beta}(\alpha, \alpha)}} \left[ c_1 \cdot \frac{e^{f(x^{(2)})^T f(x^{(-)})}}{e^{f(x^{(2)})^T f(x^{(1)})} + e^{f(x^{(2)})^T f(x^{(-)})}} \cdot (x^{(1)} - x^{(2)})^T \nabla f(x^{(2)})^T [f(x^{(1)}) - f(x^{(-)})] + c_2 \cdot \frac{e^{f(x^{(2)})^T f(x^{(-)})}}{e + e^{f(x^{(2)})^T f(x^{(-)})}} \right. \\
 &\quad \left. \cdot (x^{(1)} - x^{(2)})^T \nabla f(x^{(2)})^T [f(x^{(2)}) - f(x^{(-)})] + c_3 \cdot \frac{e^{f(x^{(2)})^T f(x^{(-)})}}{(e + e^{f(x^{(2)})^T f(x^{(-)})})^2} \cdot \left[ (x^{(1)} - x^{(2)})^T \nabla f(x^{(2)})^T [f(x^{(2)}) - f(x^{(-)})] \right]^2 \right]
 \end{aligned}$$

# Taylor expansion of i-mix loss

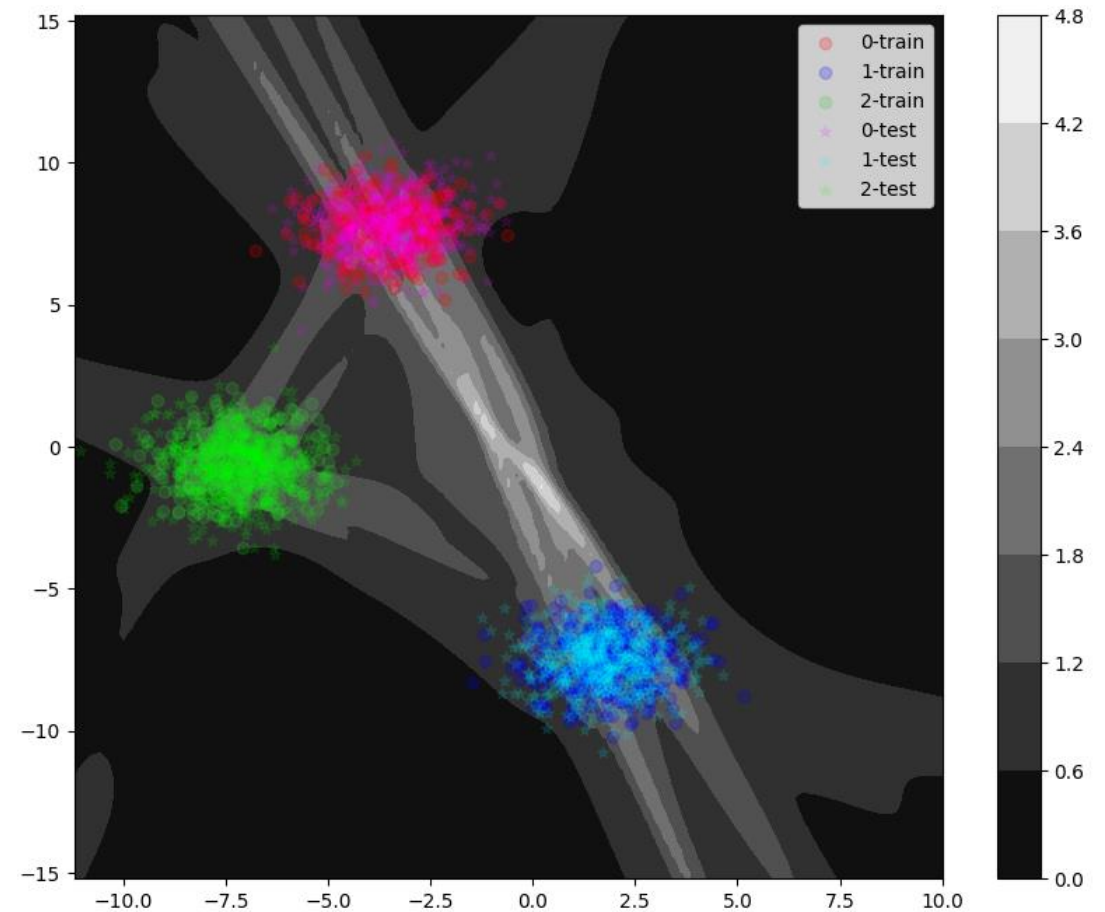
- Inside  $R^{(2)}$ , there are terms related with indirect contrastive loss
  - Attraction between ‘tangential approximation of  $x^{(1)}$  w.r.t  $x^{(2)}$ ’ and ‘ $f(x^{(1)})$ ’, and repulsion between ‘tangential approximation of  $x^{(1)}$  w.r.t  $x^{(2)}$ ’ and ‘ $f(x^{(-)})$ ’:
$$(x^{(1)} - x^{(2)})^T \nabla f(x^{(2)})^T [f(x^{(1)}) - f(x^{(-)})]$$
  - Attraction between ‘tangential approximation of  $x^{(1)}$  w.r.t  $x^{(2)}$ ’ and ‘ $f(x^{(2)})$ ’, and repulsion between ‘tangential approximation of  $x^{(1)}$  w.r.t  $x^{(2)}$ ’ and ‘ $f(x^{(-)})$ ’:
$$(x^{(1)} - x^{(2)})^T \nabla f(x^{(2)})^T [f(x^{(2)}) - f(x^{(-)})]$$
- Unsolved Q: Why does it help? (why the input gradient control can be effective?)
  - Unclear, but it is speculated that they control the input gradients to make a regularized contrasted feature space. (which might be helpful for downstream performance)

# Taylor expansion of i-mix loss

- Empirically, it turned out that the  $\|\nabla f(x)\|_2$  is get more stronger at the convex hull of data set:



Input gradient norm trained under  $L_{un}$

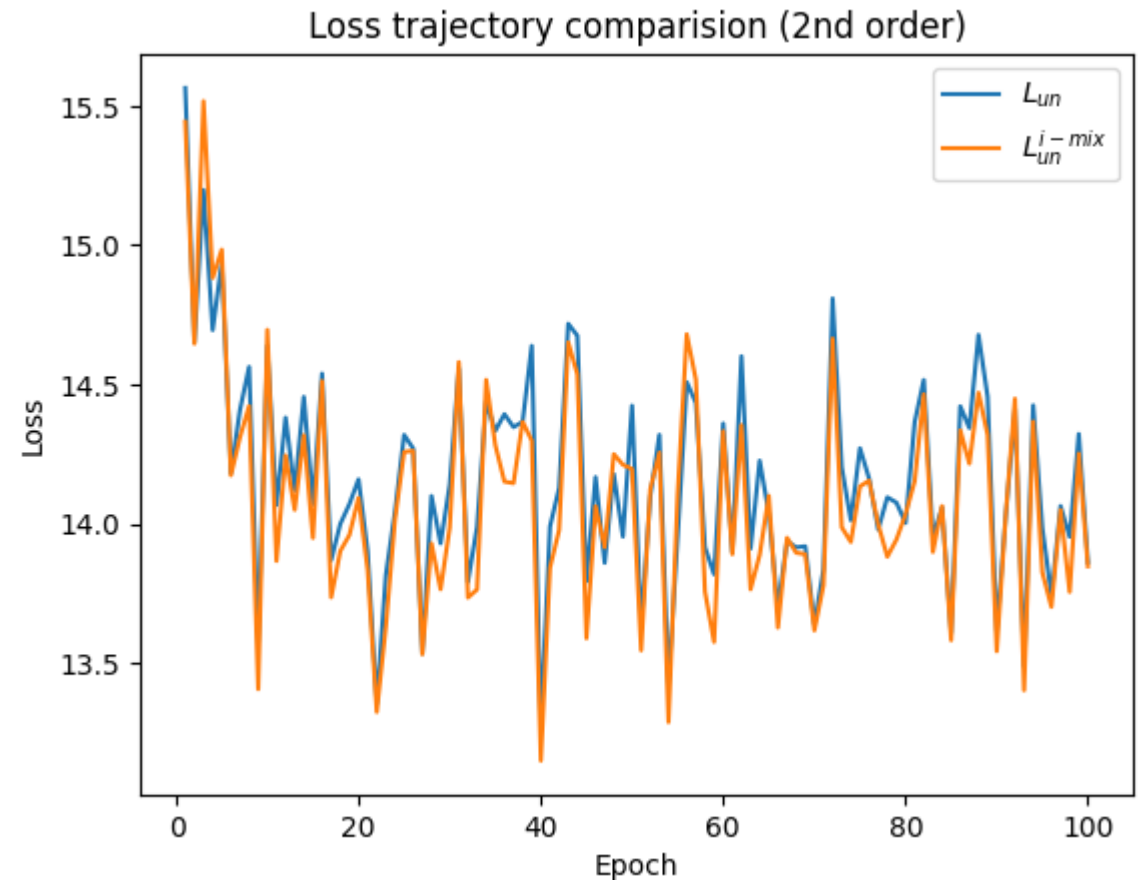
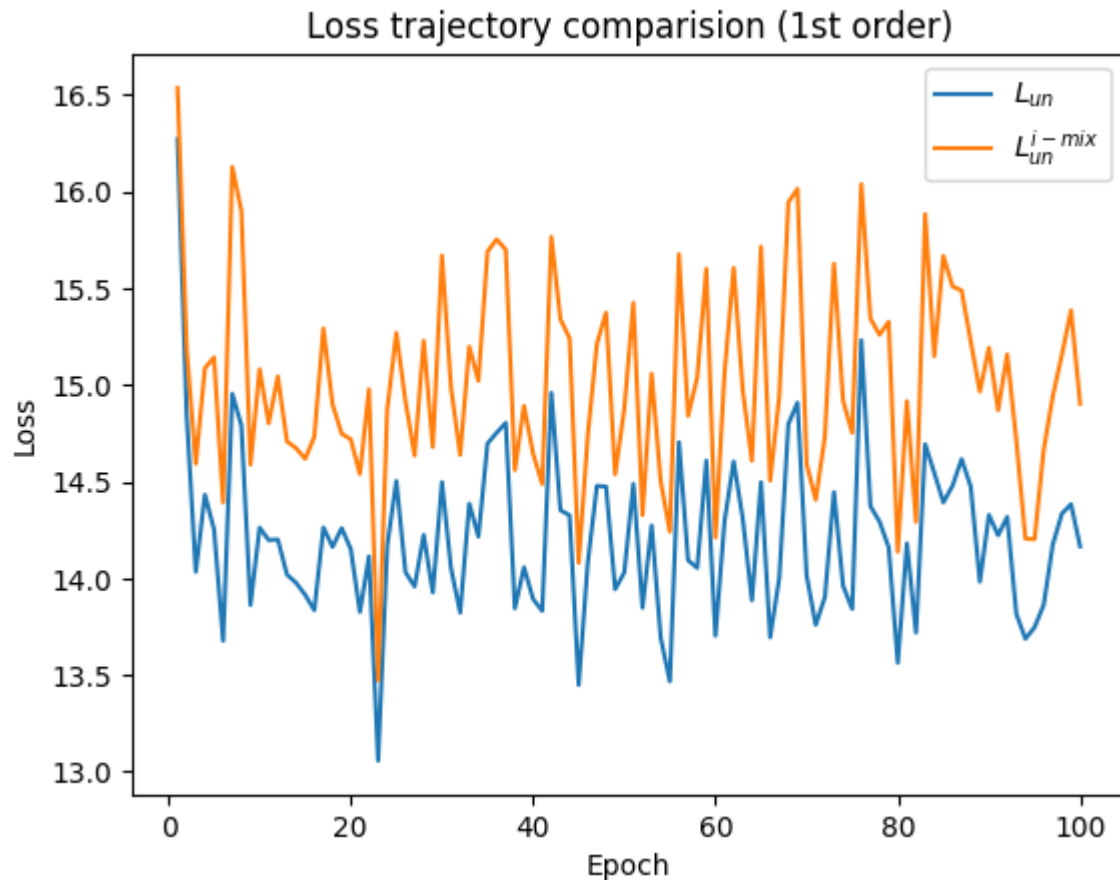


Input gradient norm trained under  $L_{un}^{i-mix}$



# Taylor expansion of i-mix loss

- Then, does the approximated mix-up loss is accurate compared to true mix-up loss?
  - Under the 3-layer ReLU Network + 2-dimensional blob dataset, the loss trajectory is given as follows (as did in [C. Park, 2022], [Zhang, 2021]):  $(\lambda \sim \text{Beta}(1,1))$



# Downstream performance by i-mix

- How about the downstream performance guarantee under i-mix loss?
- Using the framework from [Arora, 2019], [Verma, 2021], the following claim holds:

**Lemma 1: Minimizing the empirical loss  $\rightarrow$  minimizing the upper bound of true loss (similarly as in [Arora, 2019])**

With probability at least  $1 - \delta$  over the training set  $\mathcal{S}$ , for all  $f \in \mathcal{F} := \{f: \mathcal{X} \rightarrow \mathbb{R}^d \mid \|f\|_B < R \text{ for some } R > 0\}$

$$L_{un}^{i-mix}(\hat{f}) \leq \widehat{L}_{un}^{i-mix}(f) + O\left(R \cdot \frac{\mathcal{R}_S(\mathcal{F})}{M} + R^2 \sqrt{\frac{\log \frac{1}{\delta}}{M}}\right)$$

where  $\mathcal{R}_S(\mathcal{F}) := \mathbb{E}_{\sigma \sim \{\pm 1\}^{3dM}} \left[ \sup_{f \in \mathcal{F}} \langle \sigma, f|_{\mathcal{S}} \rangle \right]$  (= Rademacher complexity), and  $M = |\mathcal{S}|$  (= sample size),  $\hat{f} \in \operatorname{argmin}_{f \in \mathcal{F}} \widehat{L}_{un}(f)$

# Downstream performance by i-mix

**Theorem: Weak guarantee of downstream performance under i-mix training [similarly as in [Verma, 2021]]**

Let  $\bar{\rho}(y) = P(y' \neq y | y) > 0$ ,  $\rho(y) = P(y' = y | y)$ , then, under the binary class contrastive learning and linear evaluation setting, the following holds:

$$l_{un}^{i-mix} \cong \frac{1}{2} \mathbb{E}_{\substack{(x^{(1)}, y) \sim D \\ x^{(2)} \sim D_{1-y}}} \left[ \rho(y) l_{sup} \left( f(x^{(1)}, y)^T \tilde{w} \right) \right] + \frac{1}{2} \mathbb{E}_y [(1 - \bar{\rho}(y)) \mathcal{E}_y] + R^{(1)} + R^{(2)}$$

where

$$\mathcal{E}_y = \mathbb{E}_{\substack{x^{(1)}, x^{(2)} \sim D_y \\ \lambda \sim \text{Beta}(\alpha, \alpha)}} \left[ \log \left( 1 + \exp \left( - \frac{f(x^{(1)})^T}{\|f(x^{(1)})\|} \left( \frac{f(x^{(2)})^T}{\|f(x^{(2)})\|} - \frac{f(x^{(-)})^T}{\|f(x^{(-)})\|} \right) \right) \right) \right]$$

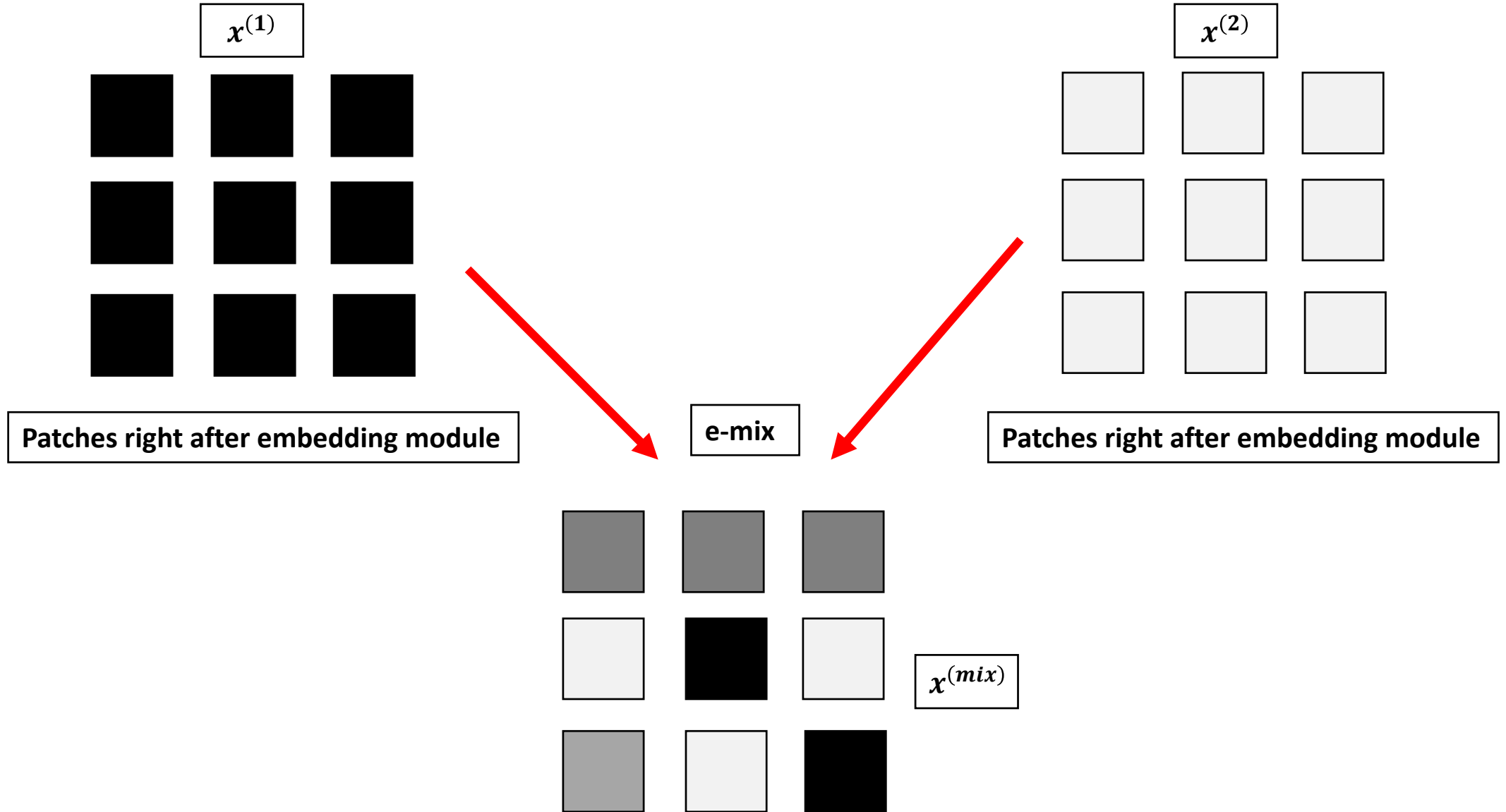
$$\tilde{w} = \frac{1}{\|h(x^{(1)})\|} \left( \frac{1}{\|h(\pi_{y,1}(x^{(2)}, x^{(-)}))\|} h(\pi_{y,1}(x^{(2)}, x^{(-)})) - \frac{1}{\|h(\pi_{y,0}(x^{(2)}, x^{(-)}))\|} h(\pi_{y,0}(x^{(2)}, x^{(-)})) \right)$$
$$\pi_{y,y'}(x^{(2)}, x^{(-)}) = P(y = y') \cdot x^{(2)} + P(y \neq y') \cdot x^{(-)}$$

# Application of i-mix on domain agnostic environment

- In DABS 2.0, they used ViT structure to cover domain agnostic data.
  - But, it turned out that the suggested domain-agnostic algorithms (ex: e-mix, ShED, ContPred) tend to show improvement in certain types of domain:
    - e-mix : strong at continuous data (spectrogram, vision, continuous tabular [sensor])
    - SheD (Shuffled embedding detection) / ContPred (mask certain portion of random patches) : strong at discontinuous data (text, vision, discontinuous tabular)
- But, How about exploiting the ViT structure for mix-up?
  - Use patch mix-up in discrete unit => may be helpful for both continuous / discontinuous data.

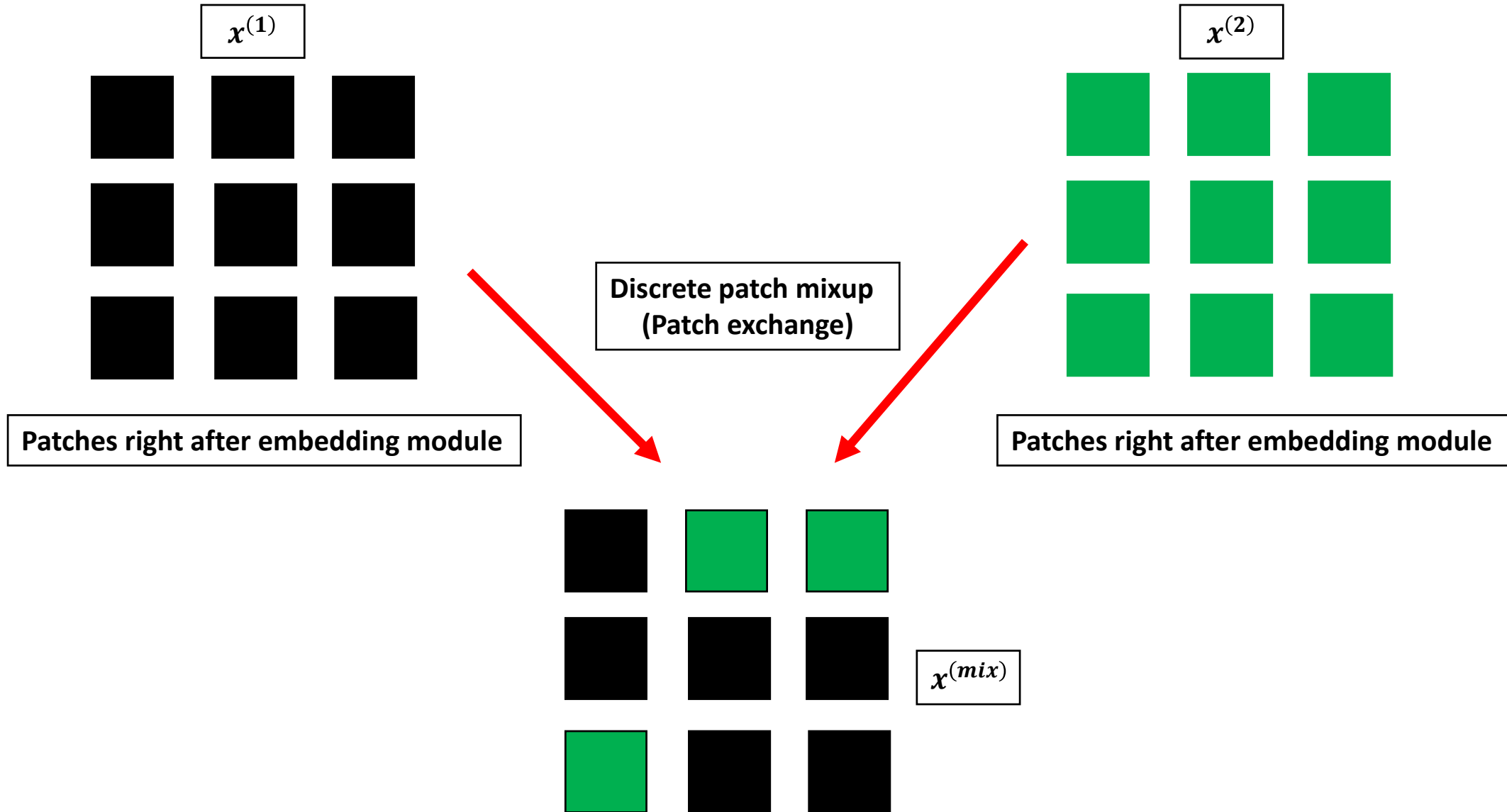
# Application of i-mix on domain agnostic environment

- e-mix (from DABS 2.0)



# Application of i-mix on domain agnostic environment

- Discrete patch mixup:



# Application of i-mix on domain agnostic environment

- Performance comparison ('dmix' = discrete patch mixup)
  - Online evaluation test accuracy (not linear evaluation, but can be a substitute)

Dataset	CIFAR-10	CIFAR-100	PAMAP2	Euroset	CUBirds	DTD	Traffic Sign	MNLI (matched)
emix	38.4	10.1	81.8	<b>89.2</b>	<b>1.54</b>	7.71	55.4	32.6
ShED	36.4	No data	65.9	55.7	No data	6.12	28.4	<b>38.1</b>
dmix	<b>48.7</b>	<b>19.7</b>	<b>82.8</b>	87.4	1.45	<b>8.24</b>	<b>63.0</b>	32.6

- Clearly shows good performance compared to emix, but unclear on the NLP domain
  - We can further generalize dmix by combining ContPred algorithm (called 'Contdmix')
    - Then, on continuous data : shED < emix < contdmix < dmix
    - On discontinuous data : emix < dmix <= contdmix < shED