



# AutoMix: Unveiling the Power of Mixup for Stronger Classifiers

ECCV 2022 Oral

PRESENTER: JIAHANG ZHANG

2023/02/12

---

## ● Outline

---

1 / **Authors**

2 / **Background**

3 / **Method**

4 / **Experiments**

5 / **Discussion**



---

## ● Outline

---

1 / **Authors**

2 / **Background**

3 / **Method**

4 / **Experiments**

5 / **Discussion**



---

## ● Outline

---

1 / Authors

2 / **Background**

3 / Method

4 / Experiments

5 / Discussion



---

## ● Background

---

### ■ Mix up & Cut Mix



Original Image  
Dog



Mix up  
Dog & Cat



Cut out  
Dog

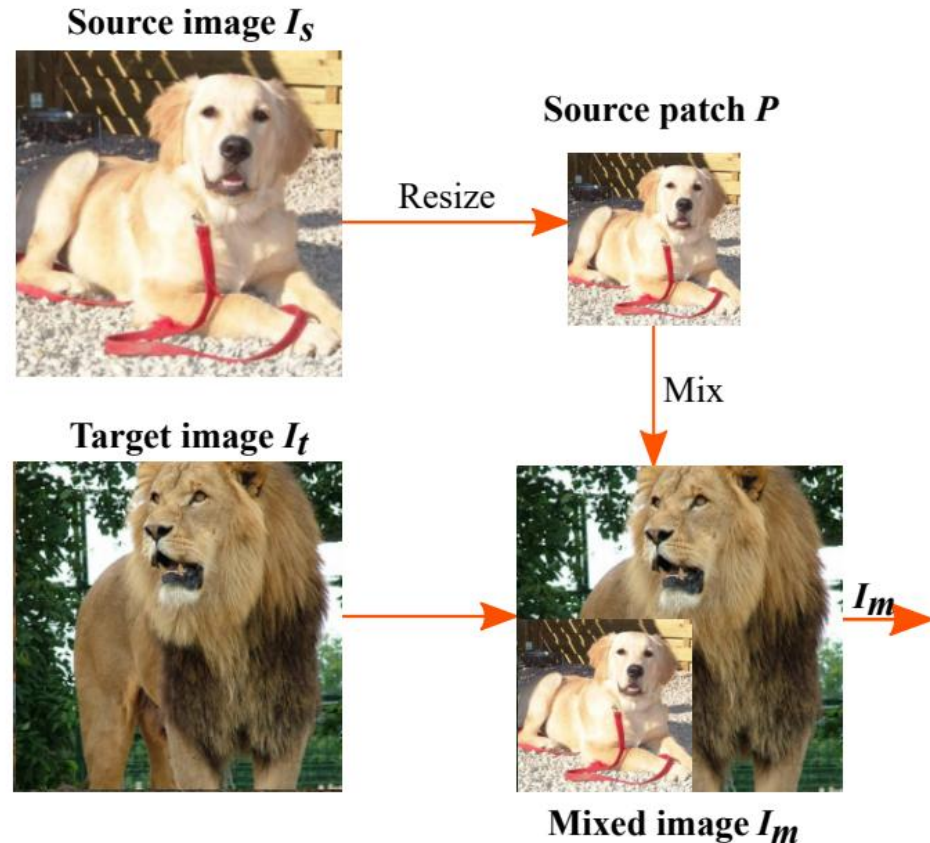


Cut Mix  
Dog & Cat

$$\tilde{y} = \lambda y_A + (1 - \lambda) y_B$$

# ● Background

## ■ Resize Mix



## ResizeMix: Mixing Data with Preserved Object Information and True Labels

Jie Qin<sup>1,2\*</sup>, Jiemin Fang<sup>3,4\*</sup>, Qian Zhang<sup>5</sup>, Wenyu Liu<sup>4</sup>, Xingang Wang<sup>2†</sup>, Xinggang Wang<sup>4</sup>

<sup>1</sup>School of Artificial Intelligence, University of Chinese Academy of Sciences

<sup>2</sup>Institute of Automation, Chinese Academy of Sciences

<sup>3</sup>Institute of Artificial Intelligence, Huazhong University of Science and Technology

<sup>4</sup>School of EIC, Huazhong University of Science and Technology <sup>5</sup>Horizon Robotics

$$l_m = \lambda l_s + (1 - \lambda) l_t$$

## ● Background

### ■ Puzzle Mix

- Maximally utilize the saliency information

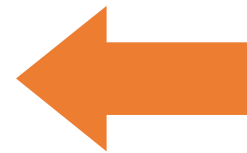
Puzzle Mix (*full*)



Input1



Input2



# ● Background

## FMix: Enhancing Mixed Sample Data Augmentation

Ethan Harris, Antonia Marcu, Matthew Painter, Mahesan Niranjan, Adam Prügel-Bennett, Jonathon Hare

### ■ FMix





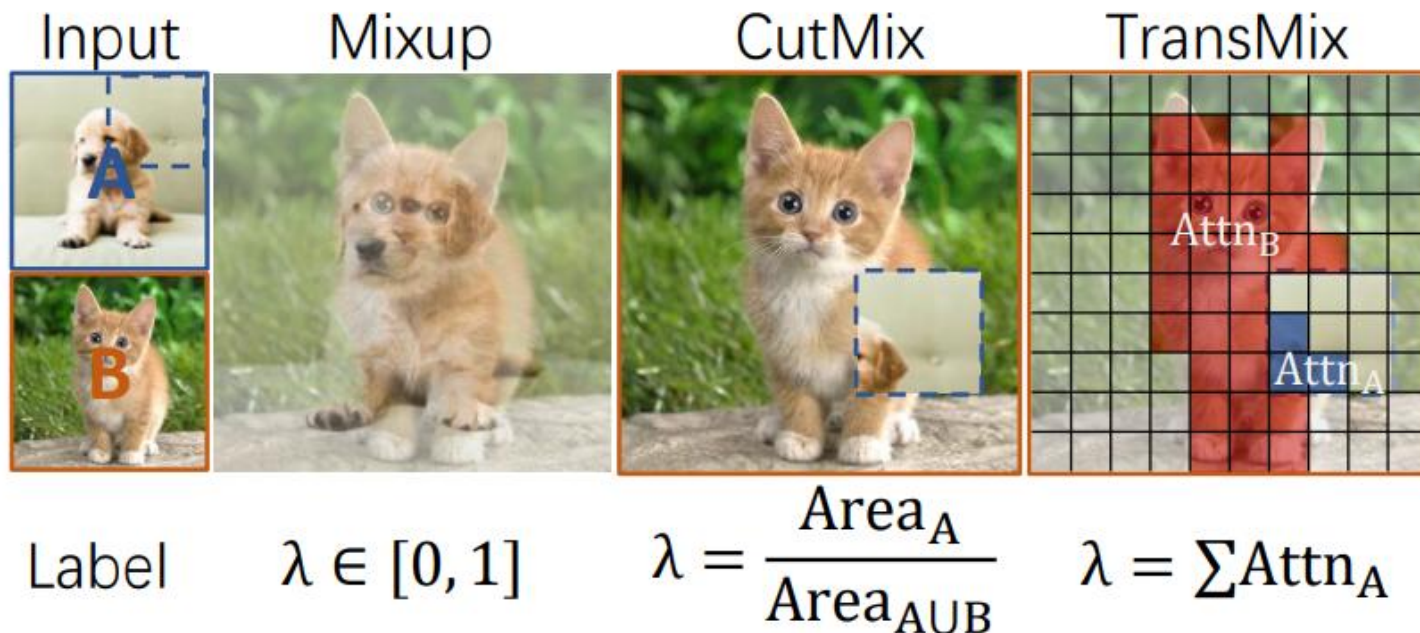
# ● Background

## TransMix: Attend to Mix for Vision Transformers

Jie-Neng Chen<sup>1\*</sup> Shuyang Sun<sup>2\*</sup> Ju He<sup>1</sup> Philip Torr<sup>2</sup> Alan Yuille<sup>1</sup> Song Bai<sup>3</sup>  
<sup>1</sup>Johns Hopkins University <sup>2</sup>University of Oxford <sup>3</sup>ByteDance Inc.

### ■ TransMix

#### ■ Attention-guided for ViTs



---

## ● Background

---

- **Why Mix effective?**
  - **Data-dependent regularization**
  - **Label smoothing**

---

## ● Outline

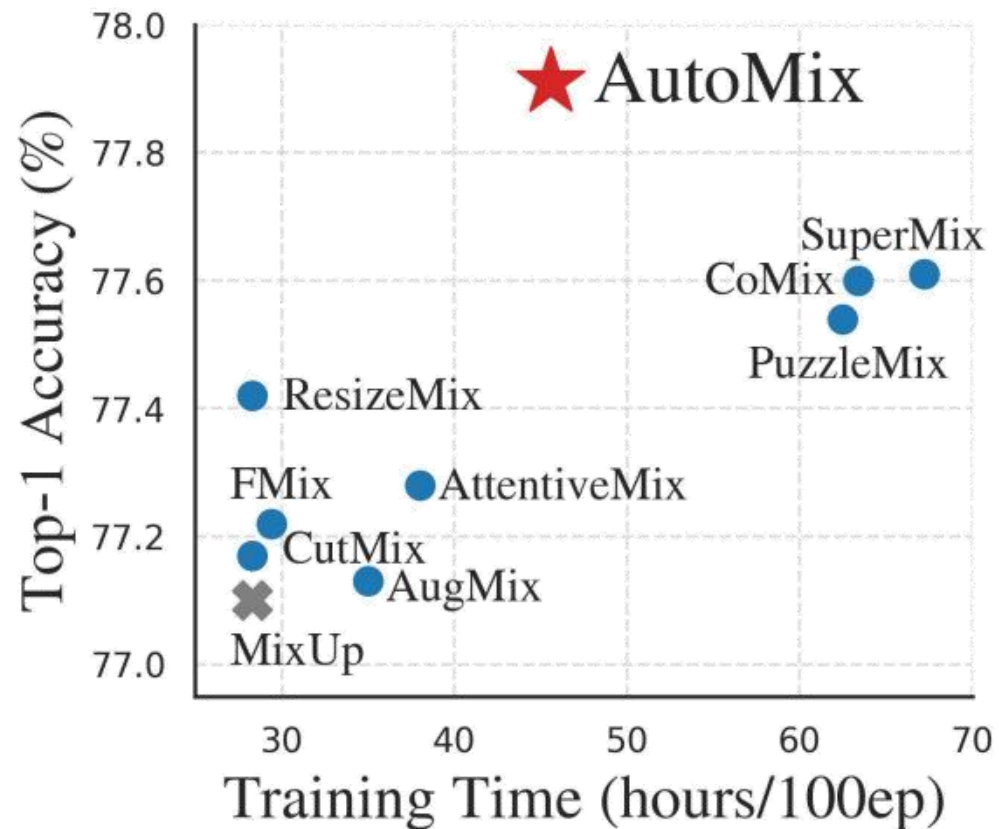
---

- 1 / Authors
- 2 / Background
- 3 / **Method**
- 4 / Experiments
- 5 / Discussion



## ● Method

### ■ Motivation:



---

## ● Method

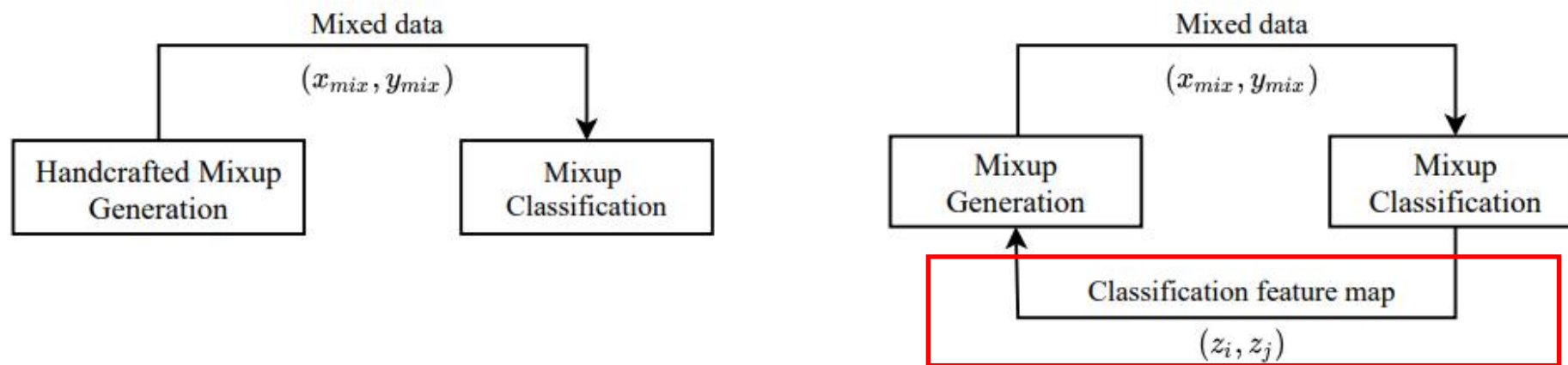
---

### ■ Motivation:

- How to design an **accurate mixing policy** to benefit the mixup classification objective? → *Label Mismatch Issue*
- How to solve generation-classification optimization problems **efficiently**?

## ● Method

### ■ Overview



**Fig. 2.** The difference between AutoMix and offline approaches. **Left:** Offline mixup methods, where a fixed mixup policy generates mixed samples for the classifier to learn from. **Right:** AutoMix, where the mixup policy is trained with the feature map.

---

## ● Method

---

- Mixup policy for input  $h_\phi$
- Mixup policy for label  $g$

$$g(y_i, y_j, \lambda) = \lambda y_i + (1 - \lambda) y_j \quad (1)$$

- Optimization objective:

$$\min_{\theta, \phi} \ell_{MCE} \left( f_\theta (h_\phi(x_i, x_j, \lambda)), g(y_i, y_j, \lambda) \right). \quad (2)$$

---

## ● Method

---

- Generate the pixel-level mask and obtain the mixed data

$$h_\phi(x_i, x_j, \lambda) = \mathcal{M}_\phi(z_{i,\lambda}^l, z_{j,1-\lambda}^l) \odot x_i + (1 - \mathcal{M}_\phi(z_{i,\lambda}^l, z_{j,1-\lambda}^l)) \odot x_j \quad (3)$$

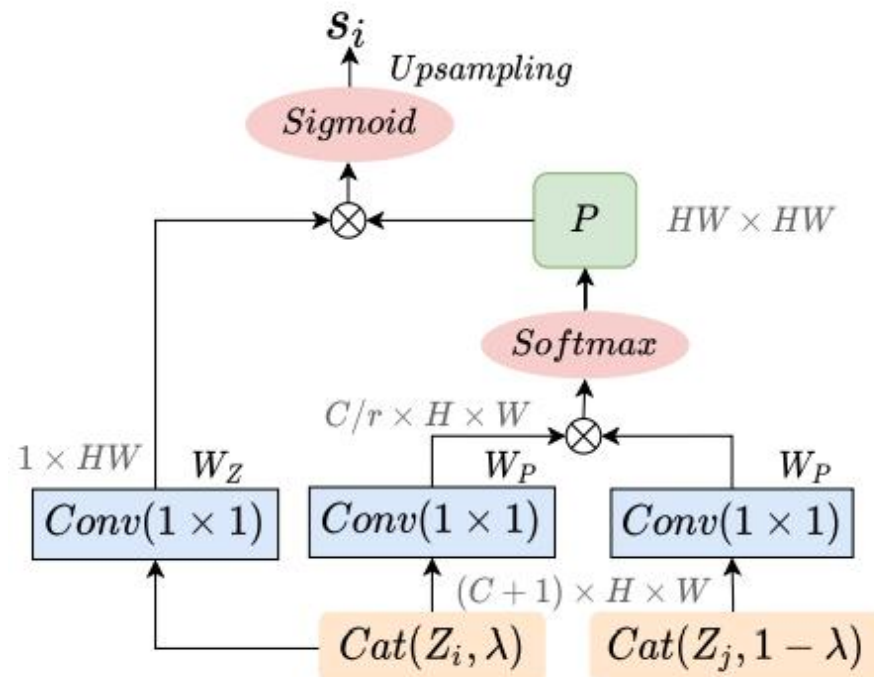
$$z_\lambda^l = \text{concat}(z, \lambda)$$



## ● Method

- Generate the pixel-level mask and obtain the mixed data

$$h_\phi(x_i, x_j, \lambda) = \mathcal{M}_\phi(z_{i,\lambda}^l, z_{j,1-\lambda}^l) \odot x_i + (1 - \mathcal{M}_\phi(z_{i,\lambda}^l, z_{j,1-\lambda}^l)) \odot x_j \quad (3)$$



---

## ● Method

---

- Generate the pixel-level mask and obtain the mixed data

$$h_\phi(x_i, x_j, \lambda) = \mathcal{M}_\phi(z_{i,\lambda}^l, z_{j,1-\lambda}^l) \odot x_i + (1 - \mathcal{M}_\phi(z_{i,\lambda}^l, z_{j,1-\lambda}^l)) \odot x_j \quad (3)$$

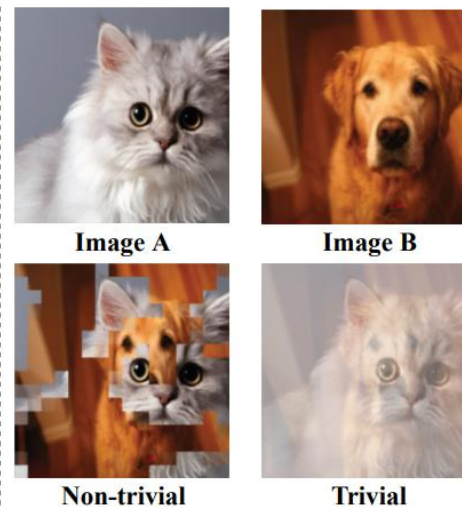
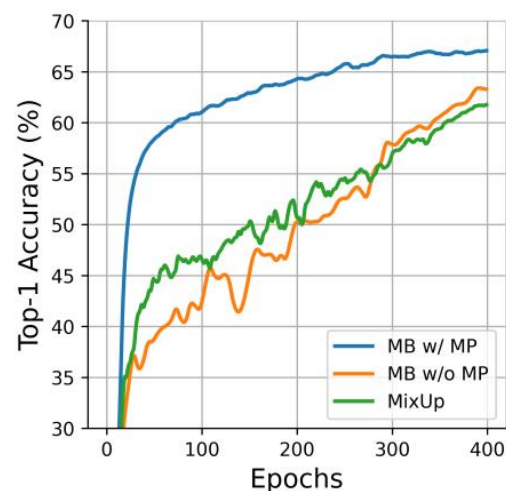
- Objective for  $\mathcal{M}_\phi$

$$\ell_\lambda = \gamma \max \left( \left\| \lambda - \frac{1}{HW} \sum_{h,w} s_{i,h,w} \right\| - \epsilon, 0 \right) \quad (4)$$

## ● Method

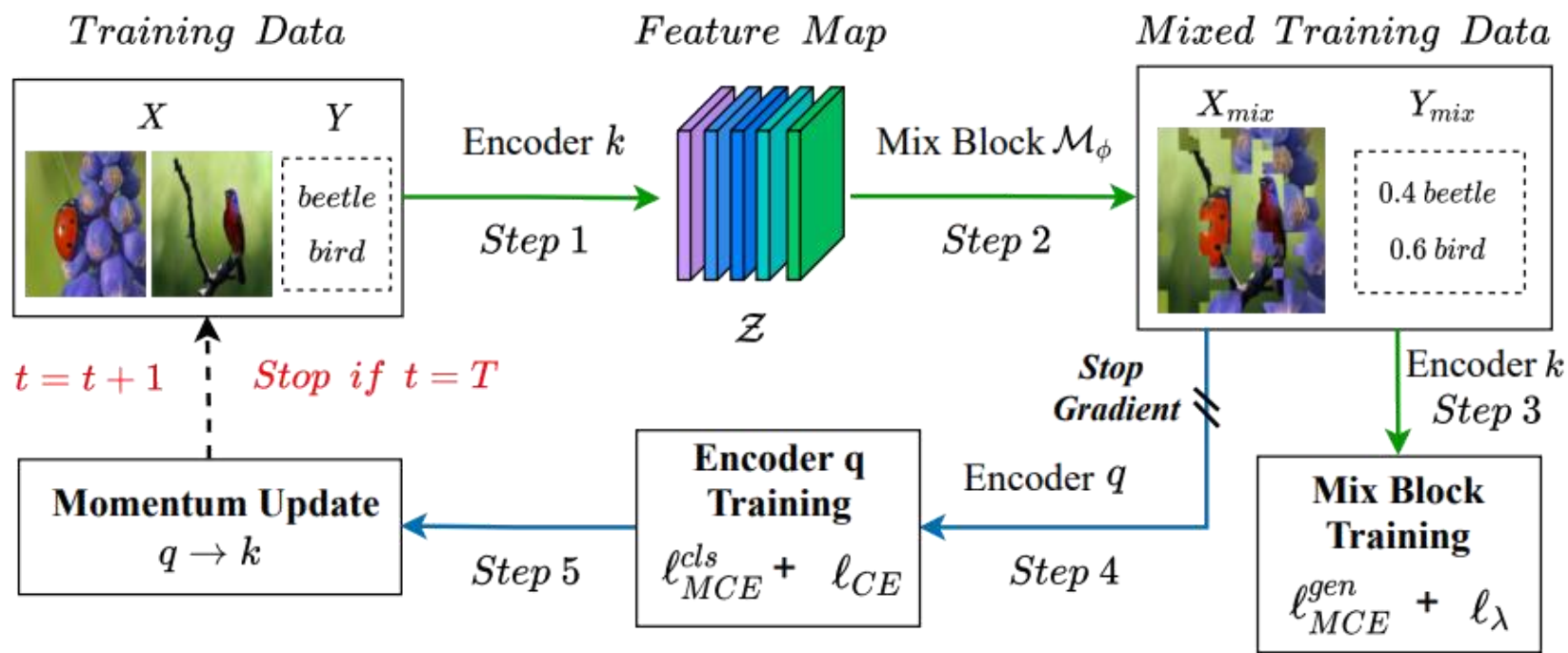
- A trick for training  $\mathcal{M}_\phi$ 
  - Direct optimizing the two sub-tasks leads to gradient entanglement problem

$$\nabla_\phi \mathcal{L}_{MCE}^{cls} \propto \nabla_\phi h_\phi(x_i, x_j, \lambda) \odot f'_\theta(h_\phi(x_i, x_j, \lambda)). \quad (5)$$



## ● Method

- A trick for training  $\mathcal{M}_\phi$ 
  - To this end, **Momentum Pipeline** are proposed for bi-level optimization
  - Inspired by the **self-supervised learning**



---

## ● Outline

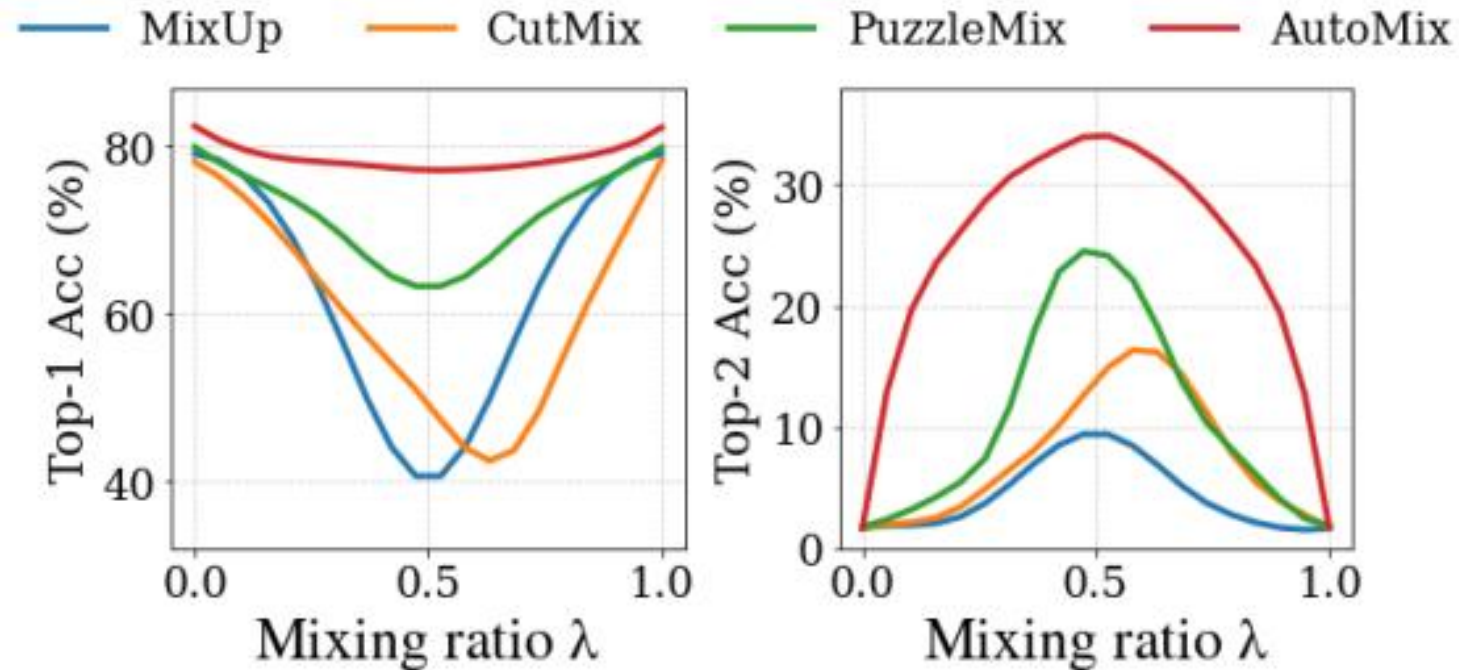
---

- 1 / Authors
- 2 / Background
- 3 / Method
- 4 / Experiments**
- 5 / Discussion



# ● Experiments

## ■ Quality of the proposed method



# ● Experiments

## ■ Classification

**Table 2.** Top-1 accuracy (%) $\uparrow$  of image classification based on ResNet variants on ImageNet-1k using PyTorch-style 100-epoch and 300-epoch training procedures.

Methods	PyTorch 100 epochs					PyTorch 300 epochs			
	R-18	R-34	R-50	R-101	RX-101	R-18	R-34	R-50	R-101
Vanilla	70.04	73.85	76.83	78.18	78.71	<b>71.83</b>	75.29	77.35	78.91
MixUp	69.98	73.97	77.12	78.97	79.98	71.72	75.73	78.44	80.60
CutMix	68.95	73.58	77.17	78.96	80.42	71.01	75.16	78.69	80.59
ManifoldMix	69.98	73.98	77.01	79.02	79.93	71.73	75.44	78.21	80.64
SaliencyMix	69.16	73.56	77.14	79.32	80.27	70.21	75.01	78.46	80.45
FMix*	69.96	74.08	77.19	79.09	80.06	70.30	75.12	78.51	80.20
PuzzleMix	<b>70.12</b>	<b>74.26</b>	<b>77.54</b>	<b>79.43</b>	80.53	71.64	<b>75.84</b>	78.86	<b>80.67</b>
ResizeMix*	69.50	73.88	77.42	79.27	<b>80.55</b>	71.32	75.64	<b>78.91</b>	80.52
<b>AutoMix</b>	<b>70.50</b>	<b>74.52</b>	<b>77.91</b>	<b>79.87</b>	<b>80.89</b>	<b>72.05</b>	<b>76.10</b>	<b>79.25</b>	<b>80.98</b>
Gain	<b>+0.38</b>	<b>+0.26</b>	<b>+0.37</b>	<b>+0.44</b>	<b>+0.34</b>	<b>+0.22</b>	<b>+0.26</b>	<b>+0.34</b>	<b>+0.31</b>



# ● Experiments

## ■ Classification

**Table 2.** Top-1 accuracy (%) $\uparrow$  of image classification based on ResNet variants on ImageNet-1k using PyTorch-style 100-epoch and 300-epoch training procedures.

Methods	PyTorch 100 epochs					PyTorch 300 epochs			
	R-18	R-34	R-50	R-101	RX-101	R-18	R-34	R-50	R-101
Vanilla	70.04	73.85	76.83	78.18	78.71	<b>71.83</b>	75.29	77.35	78.91
MixUp	69.98	73.97	77.12	78.97	79.98	71.72	75.73	78.44	80.60
CutMix	68.95	73.58	77.17	78.96	80.42	71.01	75.16	78.69	80.59
ManifoldMix	69.98	73.98	77.01	79.02	79.93	71.73	75.44	78.21	80.64
SaliencyMix	69.16	73.56	77.14	79.32	80.27	70.21	75.01	78.46	80.45
FMix*	69.96	74.08	77.19	79.09	80.06	70.30	75.12	78.51	80.20
PuzzleMix	<b>70.12</b>	<b>74.26</b>	<b>77.54</b>	<b>79.43</b>	80.53	71.64	<b>75.84</b>	78.86	<b>80.67</b>
ResizeMix*	69.50	73.88	77.42	79.27	<b>80.55</b>	71.32	75.64	<b>78.91</b>	80.52
<b>AutoMix</b>	<b>70.50</b>	<b>74.52</b>	<b>77.91</b>	<b>79.87</b>	<b>80.89</b>	<b>72.05</b>	<b>76.10</b>	<b>79.25</b>	<b>80.98</b>
Gain	<b>+0.38</b>	<b>+0.26</b>	<b>+0.37</b>	<b>+0.44</b>	<b>+0.34</b>	<b>+0.22</b>	<b>+0.26</b>	<b>+0.34</b>	<b>+0.31</b>



# ● Experiments

## ■ Classification

**Table 4.** Top-1 accuracy (%) $\uparrow$  on ImageNet-1k based on ViTs and ConvNeXt using DeiT training settings.

Methods	DeiT-S	Swin-T	ConvNeXt-T
DeiT	79.80	81.28	<b>82.10</b>
MixUp	79.65	81.01	80.88
CutMix	79.78	81.20	81.57
AttentiveMix	77.63	77.27	78.19
SaliencyMix	79.88	81.37	81.33
FMix*	77.37	79.60	81.04
PuzzleMix	80.45	<b>81.47</b>	81.48
ResizeMix*	78.61	81.36	81.64
TransMix <sup>†</sup>	<b>80.70</b>	<b>81.80</b>	-
<b>AutoMix</b>	<b>80.78</b>	<b>81.80</b>	<b>82.28</b>
Gain	<b>+0.08</b>	+0.00	<b>+0.18</b>

---

## ● Outline

---

- 1 / Authors
- 2 / Background
- 3 / Method
- 4 / Experiments
- 5 / Discussion



---

## ● Discussion

---

- Review the proposed *AutoMix* framework
  - Optimize both the **mixed sample generation task** and the **mixup classification task** in a **momentum training pipeline**.
  - **Without adding cost to inference**, AutoMix can generate various masks adaptively.

**Thanks!**



**Table 5.** Top-1 accuracy (%) $\uparrow$  of various algorithms based on ResNet variants on fine-grained and scenic classification datasets.

Method	CUB-200		FGVC-Aircraft		iNat2017		iNat2018		Places205	
	R-18	RX-50	R-18	RX-50	R-50	RX-101	R-50	RX-101	R-18	R-50
Vanilla	77.68	83.01	80.23	85.10	60.23	63.70	62.53	66.94	59.63	63.10
MixUp	78.39	84.58	79.52	85.18	61.22	66.27	62.69	67.56	59.33	63.01
CutMix	78.40	85.68	78.84	84.55	62.34	67.59	63.91	69.75	59.21	63.75
ManifoldMix	<b>79.76</b>	<b>86.38</b>	80.68	<b>86.60</b>	61.47	66.08	63.46	69.30	59.46	63.23
SaliencyMix	77.95	83.29	80.02	84.31	62.51	67.20	64.27	70.01	59.50	63.33
FMix*	77.28	84.06	79.36	86.23	61.90	66.64	63.71	69.46	59.51	63.63
PuzzleMix	78.63	84.51	<b>80.76</b>	86.23	<b>62.66</b>	<b>67.72</b>	<b>64.36</b>	<b>70.12</b>	59.62	<b>63.91</b>
ResizeMix*	78.50	84.77	78.10	84.08	62.29	66.82	64.12	69.30	<b>59.66</b>	63.88
<b>AutoMix</b>	<b>79.87</b>	<b>86.56</b>	<b>81.37</b>	<b>86.72</b>	<b>63.08</b>	<b>68.03</b>	<b>64.73</b>	<b>70.49</b>	<b>59.74</b>	<b>64.06</b>
Gain	<b>+0.11</b>	<b>+0.18</b>	<b>+0.61</b>	<b>+0.12</b>	<b>+0.42</b>	<b>+0.31</b>	<b>+0.37</b>	<b>+0.37</b>	<b>+0.08</b>	<b>+0.15</b>

**Table 6.** Top-1 accuracy (%) $\uparrow$  and FGSM error (%) $\downarrow$  on CIFAR-100 based on ResNeXt-50 (32x4d) trained 400 epochs. **Table 7.** Transfer learning of object detection task with Faster-RCNN on Pascal VOC and COCO datasets.

	Clean	Corruption	FGSM	Methods	VOC	COCO		
	Acc(%) $\uparrow$	Acc(%) $\uparrow$	Error(%) $\downarrow$		mAP	mAP	AP <sub>50</sub> <sup>bb</sup>	AP <sub>75</sub> <sup>bb</sup>
Vanilla	80.24	51.71	63.92	Vanilla	81.0	38.1	59.1	41.8
MixUp	82.44	<b>58.10</b>	<b>56.60</b>	Mixup	80.7	37.9	59.0	41.7
CutMix	81.09	49.32	76.84	CutMix	81.9	38.2	59.3	42.0
AugMix	81.18	66.54	55.59	PuzzleMix	81.9	38.3	59.3	42.1
PuzzleMix	<b>82.76</b>	57.82	63.71	ResizeMix	<b>82.1</b>	<b>38.4</b>	<b>59.4</b>	<b>42.1</b>
<b>AutoMix</b>	<b>83.13</b>	<b>58.35</b>	<b>55.34</b>	<b>AutoMix</b>	<b>82.4</b>	<b>38.6</b>	<b>59.5</b>	<b>42.2</b>

**Table 6.** Top-1 accuracy (%) $\uparrow$  and FGSM error (%) $\downarrow$  on CIFAR-100 based on ResNeXt-50 (32x4d) trained 400 epochs.

	Clean Acc(%) $\uparrow$	Corruption Acc(%) $\uparrow$	FGSM Error(%) $\downarrow$
Vanilla	80.24	51.71	63.92
MixUp	82.44	<b>58.10</b>	<b>56.60</b>
CutMix	81.09	49.32	76.84
AugMix	81.18	66.54	55.59
PuzzleMix	<b>82.76</b>	57.82	63.71
<b>AutoMix</b>	<b>83.13</b>	<b>58.35</b>	<b>55.34</b>



**Table 9.** Ablation of **Table 10.** Ablation of the proposed momentum pipeline (MP) modules in MixBlock. and the cross-entropy loss  $l_{CE}$  (CE) based on ResNet-18.

module	Tiny-ImageNet	
	R-18	RX-50
(random grids)	64.40	66.83
+cross attention	66.87	69.76
+ $\lambda$ embedding	67.15	70.41
+ $\ell_\lambda$	<b>67.33</b>	<b>70.72</b>

modules	CIFAR-100			Tiny-ImageNet			ImageNet-1k		
	MixUp	CutMix	$\mathcal{M}_\phi$	MixUp	CutMix	$\mathcal{M}_\phi$	MixUp	CutMix	$\mathcal{M}_\phi$
(none)	79.12	78.17	79.46	63.39	64.40	64.84	69.98	68.95	70.04
+MP(m=0)	-	-	81.75	-	-	67.05	-	-	70.41
+MP	<b>80.82</b>	79.57	81.93	66.02	<b>65.72</b>	67.19	<b>70.13</b>	70.02	70.45
+MP+CE	80.41	<b>79.64</b>	<b>82.04</b>	<b>66.10</b>	65.05	<b>67.33</b>	70.10	<b>70.04</b>	<b>70.50</b>