# Mixup Augmentation

# Outline

Data augmentation

Mixup

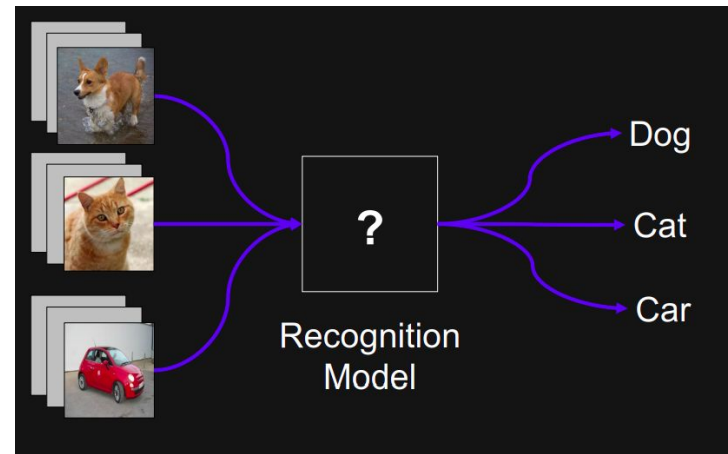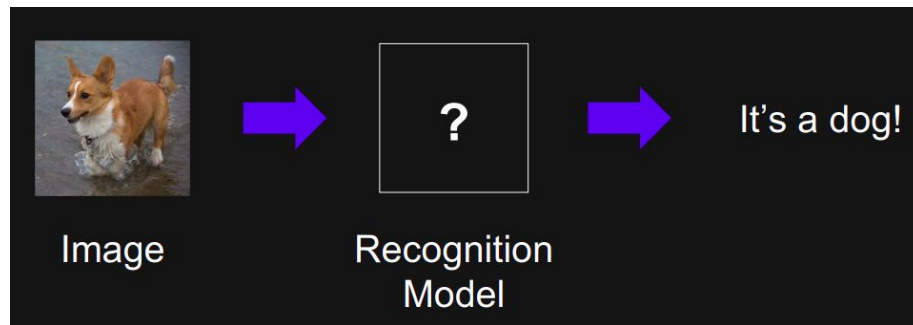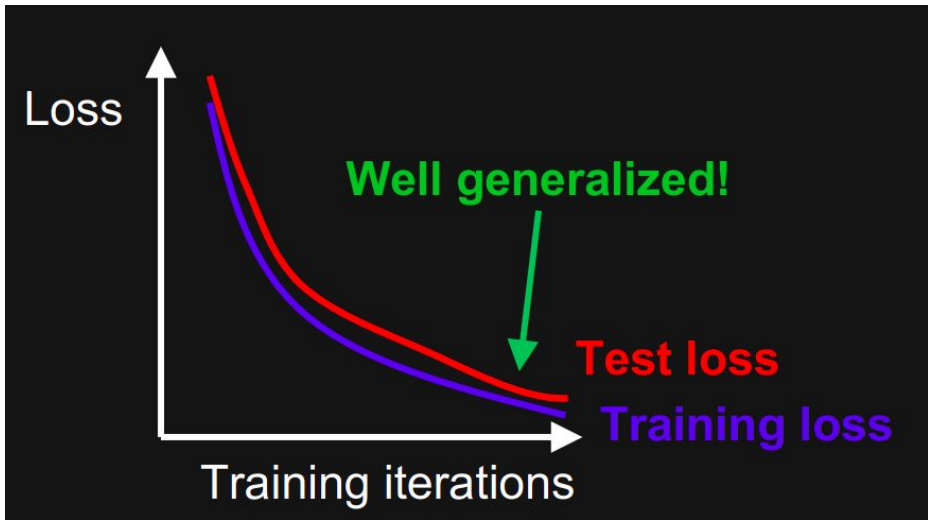Manifold Mixup

Cutmix

# Data Augmentation

# Data Augmentation

# Data Augmentation

# Data Augmentation

- **Data augmentation (DA)** is essential for ML
  - Increases the coverage of training data
  - Improves the generalizability of estimators
- An example of DA in the image domain:



goal is to maximize the performance using the same model & same dataset

# Data Augmentation

Illumination

Deformation

Occlusion

Background

Flipping

Rotating

Cropping

Original

Color distortion

Sharpening

Random erasing

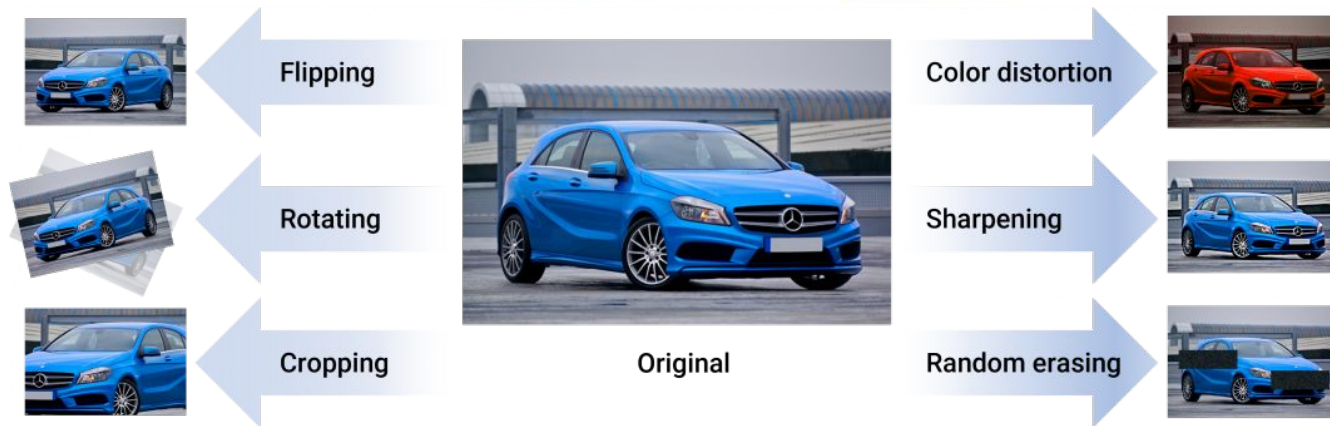# Mixup

# Mixup

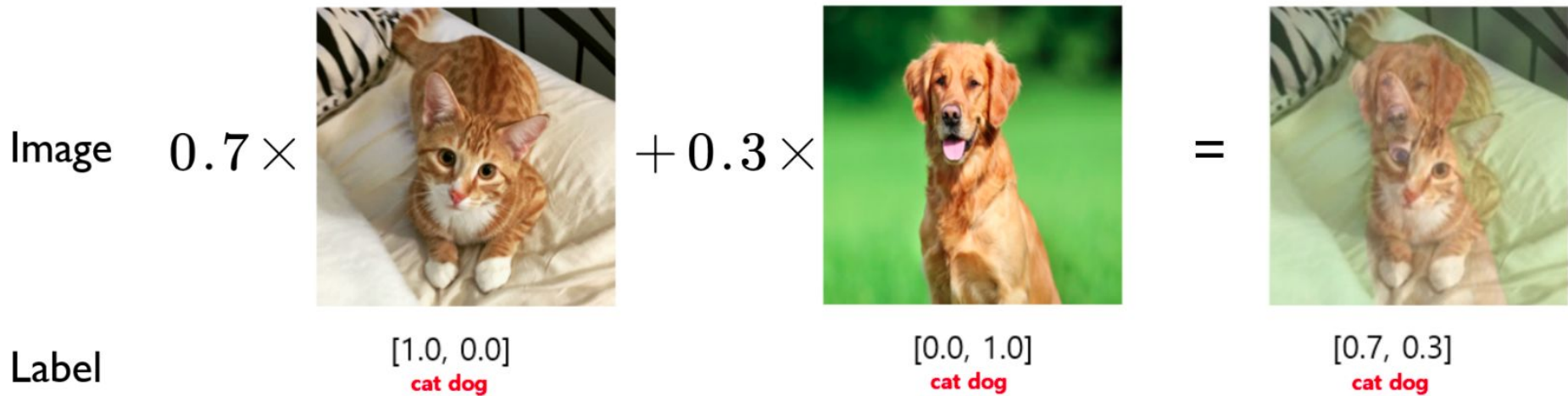**Figure:** Mixup for Image Classification

# Mixup formulation

With coefficient $\lambda \sim \text{Beta}(\alpha, \alpha)$, for $\lambda \in [0, 1], \alpha \in (0, \infty)$.
Mixup generates a virtual in-between sample,

$$\tilde{x} = \lambda x_i + (1 - \lambda)x_j,$$

$$\tilde{y} = \lambda y_i + (1 - \lambda)y_j,$$

The mixup hyper-parameter α controls the strength of interpolation between feature-target pair

# Smoother feature space

**Figure:** Illustrative sample referred from [Zhang et al., 2018]. The green and orange dots represent different classes. Blue shading indicates the probability $p(y = 1|x)$. Mixup yields a smoother decision boundary in feature space than ERM.

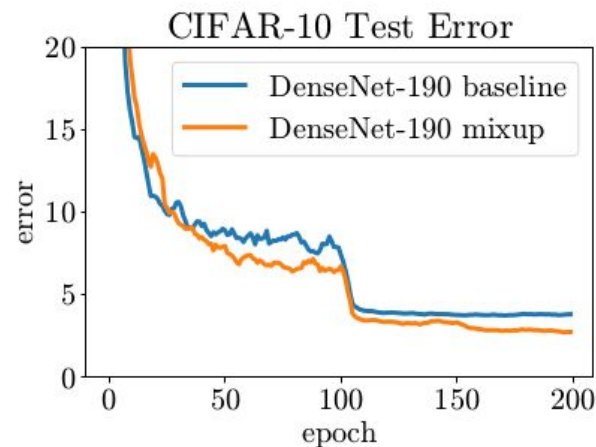| Dataset | Model | ERM | *mixup* |
|---------|-------|-----|---------|
| CIFAR-10 | PreAct ResNet-18 | 5.6 | **4.2** |
| | WideResNet-28-10 | 3.8 | **2.7** |
| | DenseNet-BC-190 | 3.7 | **2.7** |
| CIFAR-100 | PreAct ResNet-18 | 25.6 | **21.1** |
| | WideResNet-28-10 | 19.4 | **17.5** |
| | DenseNet-BC-190 | 19.0 | **16.8** |

(a) Test errors for the CIFAR experiments.



(b) Test error evolution for the best ERM and *mixup* models.

Figure 3: Test errors for ERM and *mixup* on the CIFAR experiments.

# Manifold Mixup

# Manifold Mixup

- Train on hidden states which are randomly interpolated between examples.

- Then train these interpolated hidden states to lead to lower confidence outputs.

- This also forces the model to learn representations which permit consistent interpolations.



90% blue,
10% red

vs.

10% blue,
90% red

High error!

50% red,
50% blue

vs.

50% blue,
50% red

Low error!

# Manifold Mixup

Should be
low
confidence,
since
there's no
data here.



Should be low
confidence
because it's
pretty close to
both classes.

# Manifold Mixup

- On each update, pick a random layer uniformly (including the input).

-Sample lambda ~ Beta(alpha, alpha)

-Mix between two random examples from the minibatch at that layer with rate lambda.

-Mix the labels for those two examples accordingly (soft label).



None

Input Mixup

Manifold Mixup

# Manifold Mixup



Input Space

Hidden space

Weight Decay    Noise    Dropout    Batch-Norm    Input Mixup

# Manifold Mixup

-Encourage most of the hidden space to correspond to low confidence classifications.

-Encourage real data's representations to be concentrated into local regions.

# Manifold Mixup

Table 1: Classification errors on (a) CIFAR-10 and (b) CIFAR-100. We include results from (Zhang et al. 2018)† and (Guo et al. 2016)‡. Standard deviations over five repetitions.

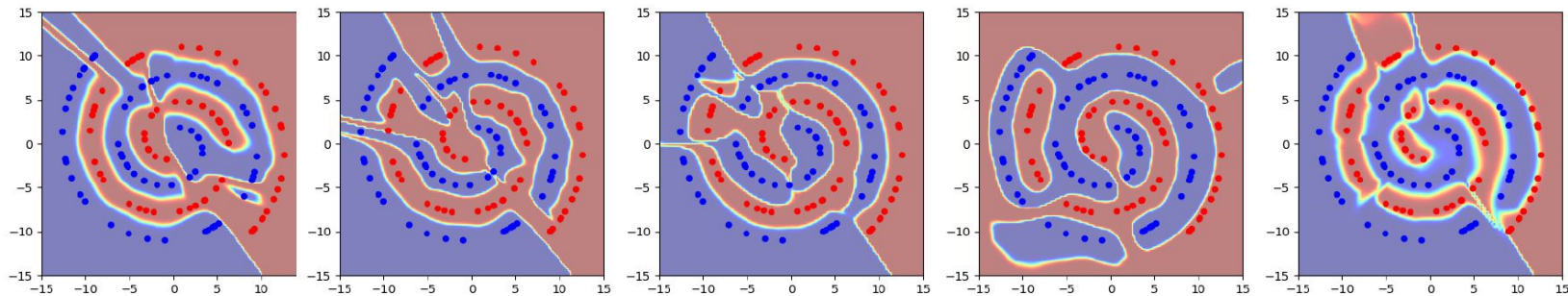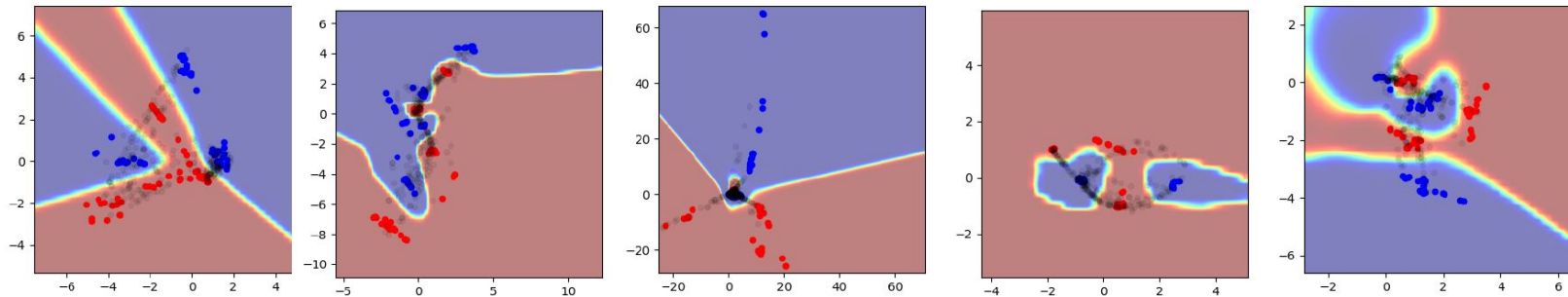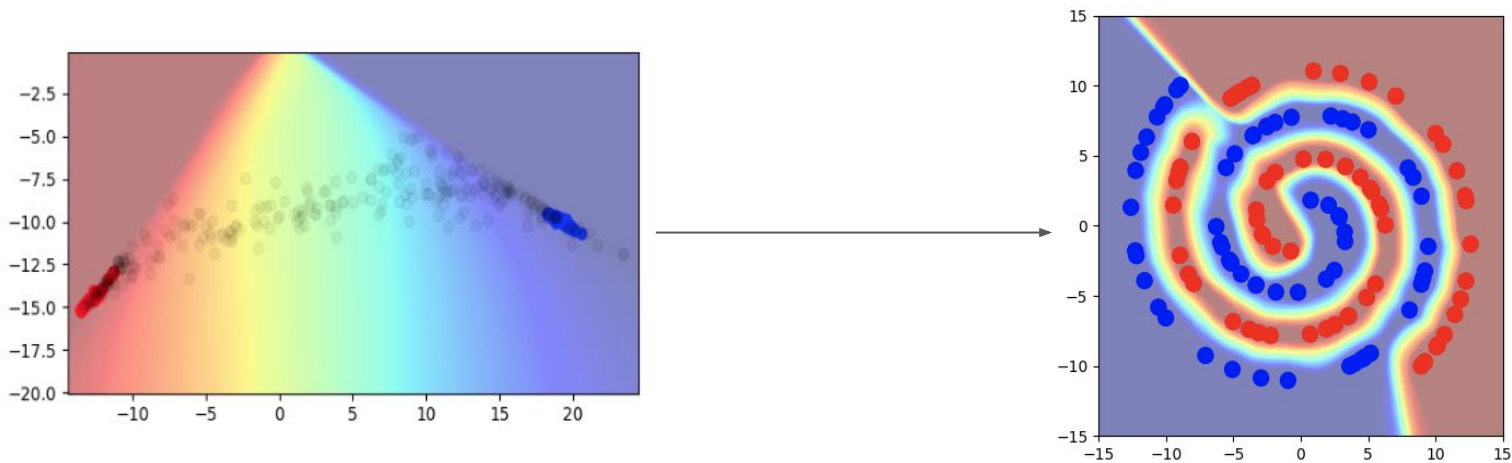| PreActResNet18 | Test Error (%) | Test NLL |
|---|---|---|
| No Mixup | $4.83 \pm 0.066$ | $0.190 \pm 0.003$ |
| AdaMix‡ | 3.52 | NA |
| Input Mixup† | 4.20 | NA |
| Input Mixup ($\alpha = 1$) | $3.82 \pm 0.048$ | $0.186 \pm 0.004$ |
| *Manifold Mixup* ($\alpha = 2$) | $\underline{2.95 \pm 0.046}$ | $\underline{0.137 \pm 0.003}$ |

| PreActResNet34 | Test Error (%) | Test NLL |
|---|---|---|
| No Mixup | $4.64 \pm 0.072$ | $0.200 \pm 0.002$ |
| Input Mixup ($\alpha = 1$) | $2.88 \pm 0.043$ | $0.176 \pm 0.002$ |
| *Manifold Mixup* ($\alpha = 2$) | $\underline{2.54 \pm 0.047}$ | $\underline{0.118 \pm 0.002}$ |

| Wide-Resnet-28-10 | Test Error (%) | Test NLL |
|---|---|---|
| No Mixup | $3.99 \pm 0.118$ | $0.162 \pm 0.004$ |
| Input Mixup ($\alpha = 1$) | $2.92 \pm 0.088$ | $0.173 \pm 0.001$ |
| *Manifold Mixup* ($\alpha = 2$) | $\underline{2.55 \pm 0.024}$ | $\underline{0.111 \pm 0.001}$ |

(a) CIFAR-10

| PreActResNet18 | Test Error (%) | Test NLL |
|---|---|---|
| No Mixup | $24.01 \pm 0.376$ | $1.189 \pm 0.002$ |
| AdaMix‡ | 20.97 | n/a |
| Input Mixup† | 21.10 | n/a |
| Input Mixup ($\alpha = 1$) | $22.11 \pm 0.424$ | $1.055 \pm 0.006$ |
| *Manifold Mixup* ($\alpha = 2$) | $\underline{20.34 \pm 0.525}$ | $\underline{0.912 \pm 0.002}$ |

| PreActResNet34 | Test Error (%) | Test NLL |
|---|---|---|
| No Mixup | $23.55 \pm 0.399$ | $1.189 \pm 0.002$ |
| Input Mixup ($\alpha = 1$) | $20.53 \pm 0.330$ | $1.039 \pm 0.045$ |
| *Manifold Mixup* ($\alpha = 2$) | $\underline{18.35 \pm 0.360}$ | $\underline{0.877 \pm 0.053}$ |

| Wide-Resnet-28-10 | Test Error (%) | Test NLL |
|---|---|---|
| No Mixup | $21.72 \pm 0.117$ | $1.023 \pm 0.004$ |
| Input Mixup ($\alpha = 1$) | $18.89 \pm 0.111$ | $0.927 \pm 0.031$ |
| *Manifold Mixup* ($\alpha = 2$) | $\underline{18.04 \pm 0.171}$ | $\underline{0.809 \pm 0.005}$ |

(b) CIFAR-100

# Manifold Mixup

Table 2: Classification errors and neg-log-likelihoods on SVHN. We run each experiment five times.

| PreActResNet18 | Test Error (%) | Test NLL |
|---|---|---|
| No Mixup | $2.89 \pm 0.224$ | $0.136 \pm 0.001$ |
| Input Mixup ($\alpha = 1$) | $2.76 \pm 0.014$ | $0.212 \pm 0.011$ |
| *Manifold Mixup* ($\alpha = 2$) | $2.27 \pm 0.011$ | $0.122 \pm 0.006$ |
| PreActResNet34 | | |
| No Mixup | $2.97 \pm 0.004$ | $0.165 \pm 0.003$ |
| Input Mixup ($\alpha = 1$) | $2.67 \pm 0.020$ | $0.199 \pm 0.009$ |
| *Manifold Mixup* ($\alpha = 2$) | $2.18 \pm 0.004$ | $0.137 \pm 0.008$ |
| Wide-Resnet-28-10 | | |
| No Mixup | $2.80 \pm 0.044$ | $0.143 \pm 0.002$ |
| Input Mixup ($\alpha = 1$) | $2.68 \pm 0.103$ | $0.184 \pm 0.022$ |
| *Manifold Mixup* ($\alpha = 2$) | $2.06 \pm 0.068$ | $0.126 \pm 0.008$ |

Table 3: Accuracy on TinyImagenet.

| PreActResNet18 | top-1 | top-5 |
|---|---|---|
| No Mixup | 55.52 | 71.04 |
| Input Mixup ($\alpha = 0.2$) | 56.47 | 71.74 |
| Input Mixup ($\alpha = 0.5$) | 55.49 | 71.62 |
| Input Mixup ($\alpha = 1.0$) | 52.65 | 70.70 |
| Input Mixup ($\alpha = 2.0$) | 44.18 | 68.26 |
| *Manifold Mixup* ($\alpha = 0.2$) | 58.70 | 73.59 |
| *Manifold Mixup* ($\alpha = 0.5$) | 57.24 | 73.48 |
| *Manifold Mixup* ($\alpha = 1.0$) | 56.83 | 73.75 |
| *Manifold Mixup* ($\alpha = 2.0$) | 48.14 | 71.69 |

# Cutmix

Cutmix

# Cutmix

| Image | ResNet-50 | Mixup [48] | Cutout [3] | CutMix |
|---|---|---|---|---|
| Label | Dog 1.0 | Dog 0.5<br>Cat 0.5 | Dog 1.0 | Dog 0.6<br>Cat 0.4 |
| ImageNet Cls (%) | 76.3 (+0.0) | 77.4 (+1.1) | 77.1 (+0.8) | **78.6 (+2.3)** |
| ImageNet Loc (%) | 46.3 (+0.0) | 45.8 (-0.5) | 46.7 (+0.4) | **47.3 (+1.0)** |
| Pascal VOC Det (mAP) | 75.6 (+0.0) | 73.9 (-1.7) | 75.1 (-0.5) | **76.7 (+1.1)** |

# Cutmix

Let $x \in \mathbb{R}^{W \times H \times C}$ and $y$ denote a training image and its label, respectively. The goal of CutMix is to generate a new training sample $(\tilde{x}, \tilde{y})$ by combining two training samples $(x_A, y_A)$ and $(x_B, y_B)$. The generated training sample $(\tilde{x}, \tilde{y})$ is used to train the model with its original loss function. We define the combining operation as

$$\tilde{x} = \mathbf{M} \odot x_A + (\mathbf{1} - \mathbf{M}) \odot x_B$$
$$\tilde{y} = \lambda y_A + (1 - \lambda) y_B, \qquad (1)$$

$$r_x \sim \text{Unif}\,(0, W), \quad r_w = W\sqrt{1 - \lambda},$$
$$r_y \sim \text{Unif}\,(0, H), \quad r_h = H\sqrt{1 - \lambda}$$



Original Samples

Input Image

CAM for 'St. Bernard'

CAM for 'Poodle'

Mixup     Cutout     CutMix

| PyramidNet-200 ($\tilde{\alpha}$=240) (# params: 26.8 M) | Top-1 Err (%) | Top-5 Err (%) |
|---|---|---|
| Baseline | 16.45 | 3.69 |
| + StochDepth [17] | 15.86 | 3.33 |
| + Label smoothing ($\epsilon$=0.1) [38] | 16.73 | 3.37 |
| + Cutout [3] | 16.53 | 3.65 |
| + Cutout + Label smoothing ($\epsilon$=0.1) | 15.61 | 3.88 |
| + DropBlock [8] | 15.73 | 3.26 |
| + DropBlock + Label smoothing ($\epsilon$=0.1) | 15.16 | 3.86 |
| + Mixup ($\alpha$=0.5) [48] | 15.78 | 4.04 |
| + Mixup ($\alpha$=1.0) [48] | 15.63 | 3.99 |
| + Manifold Mixup ($\alpha$=1.0) [42] | 16.14 | 4.07 |
| + Cutout + Mixup ($\alpha$=1.0) | 15.46 | 3.42 |
| + Cutout + Manifold Mixup ($\alpha$=1.0) | 15.09 | 3.35 |
| + ShakeDrop [46] | 15.08 | 2.72 |
| + CutMix | 14.47 | 2.97 |
| + CutMix + ShakeDrop [46] | **13.81** | **2.29** |

Table 5: Comparison of state-of-the-art regularization methods on CIFAR-100.

| Model | # Params | Top-1 Err (%) | Top-5 Err (%) |
|---|---|---|---|
| PyramidNet-110 ($\tilde{\alpha} = 64$) [11] | 1.7 M | 19.85 | 4.66 |
| PyramidNet-110 + CutMix | 1.7 M | **17.97** | **3.83** |
| ResNet-110 [12] | 1.1 M | 23.14 | 5.95 |
| ResNet-110 + CutMix | 1.1 M | **20.11** | **4.43** |

Table 6: Impact of CutMix on lighter architectures on CIFAR-100.

| PyramidNet-200 ($\tilde{\alpha}$=240) | Top-1 Error (%) |
|---|---|
| Baseline | 3.85 |
| + Cutout | 3.10 |
| + Mixup ($\alpha$=1.0) | 3.09 |
| + Manifold Mixup ($\alpha$=1.0) | 3.15 |
| + CutMix | **2.88** |

Table 7: Impact of CutMix on CIFAR-10.

| Feature/Technique | Mixup | Manifold Mixup | CutMix |
|---|---|---|---|
| Basic Concept | Combines two or more input images and their labels linearly. | Similar to Mixup, but mixes hidden representations at various layers of the network. | Cuts and pastes patches from one image onto another, mixing the labels accordingly. |
| Data Augmentation | Operates at the input level (pixel values). | Operates at both input and hidden layers within the network. | Operates at the input level with a focus on spatial regions. |
| Primary Goal | Encourages linear behavior between training examples. | Encourages learning more robust features across different tasks. | Aims at improving localization and understanding of spatial context. |
| Label Mixing | Labels are mixed in a linear fashion according to the mix ratio. | Labels are mixed based on the level at which mixing occurs. | Labels are mixed proportionally to the area of the patches involved. |
| Image Mixing | Linear interpolation of pixel values. | Interpolation of features at different network layers. | Physical combination of image patches. |
| Impact on Training | Helps in generalizing to unseen data by smoothing the decision boundary. | Promotes learning of more generalizable and robust intermediate features. | Enhances the ability of the model to localize and recognize objects within a varied context. |
| Use Cases | Generally used in image classification tasks. | Useful in tasks requiring deeper feature understanding and abstraction. | Particularly beneficial in object detection and classification tasks. |

| Technique | Advantages | Disadvantages |
|---|---|---|
| **Traditional Data Augmentation** | - Realistic modifications (rotation, flipping, scaling)<br><br>-  Simple to implement<br><br>- Improves generalization and reduces overfitting | - Limited to predefined variations<br><br>- Can be computationally inefficient<br><br>- May not represent actual data distribution |
| **Mixup** | - Enhances regularization, favoring linear behavior between training examples<br><br>- Improves generalization to unseen data<br><br>- Robustness to label noise | - Generates potentially unrealistic synthetic samples<br><br>- Complexity in interpretation of mixed images and labels |
| **Manifold Mixup** | - Encourages learning of abstract and robust features<br><br>- Versatile (applicable at multiple network layers)<br><br>- Potentially better regularization than Mixup | - More complex implementation<br><br>- Additional computational overhead<br><br>- Risk of over-smoothing decision boundaries |
| **CutMix** | - Enhances spatial context learning and object localization<br><br>- Robust to occlusion<br><br>- Balanced regularization (mix of dropout and Mixup) | - Can introduce artificial artifacts<br><br>- Complex label handling based on patch area<br><br>- Risk of feature misalignment in cut-and-paste |