



Implicit Semantic Data Augmentation for Deep Networks

Yulin Wang^{1*} Xuran Pan^{1*} Shiji Song¹ Hong Zhang² Cheng Wu¹ Gao Huang^{1†}

¹Department of Automation, Tsinghua University, Beijing, China

Beijing National Research Center for Information Science and Technology (BNRist),

²Baidu Inc., China

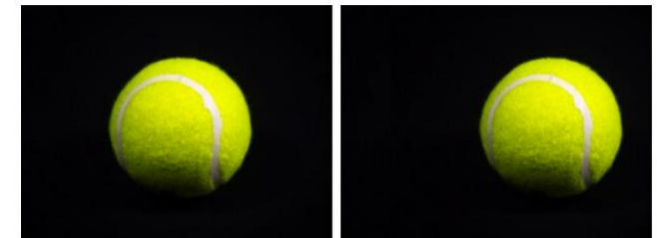
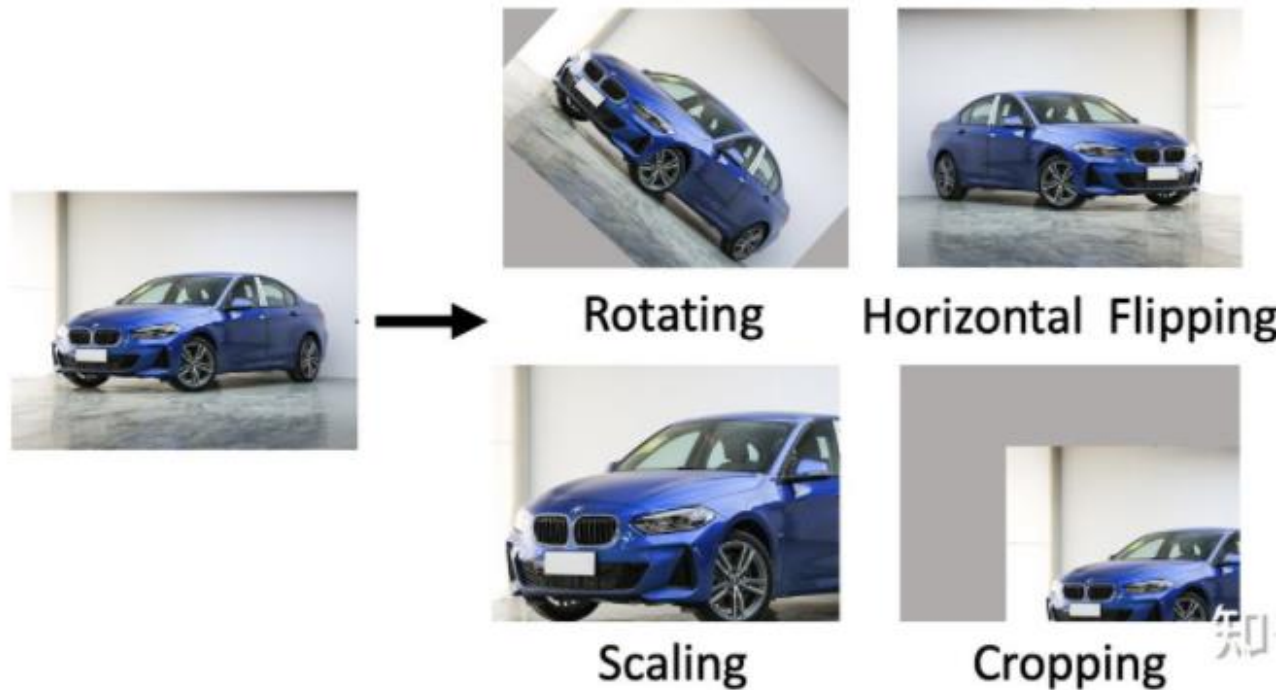
{yulin.bh, fykalviny}@gmail.com, pxr18@mails.tsinghua.edu.cn,

{shijis, wuc, gaohuang}@tsinghua.edu.cn

Data Augmentation

Data Augmentation is a widely used technique to alleviate overfitting in training deep networks in the case of a small number of data sets.

Traditional Data Augmentation



Translation

Data Augmentation

□ Generative adversarial network



- Complicated
- High cost of time and computation
- Relying on large amounts of data

Implicit Semantic Data Augmentation

This work is inspired by the observation that certain directions in the feature space correspond to meaningful semantic transformations.

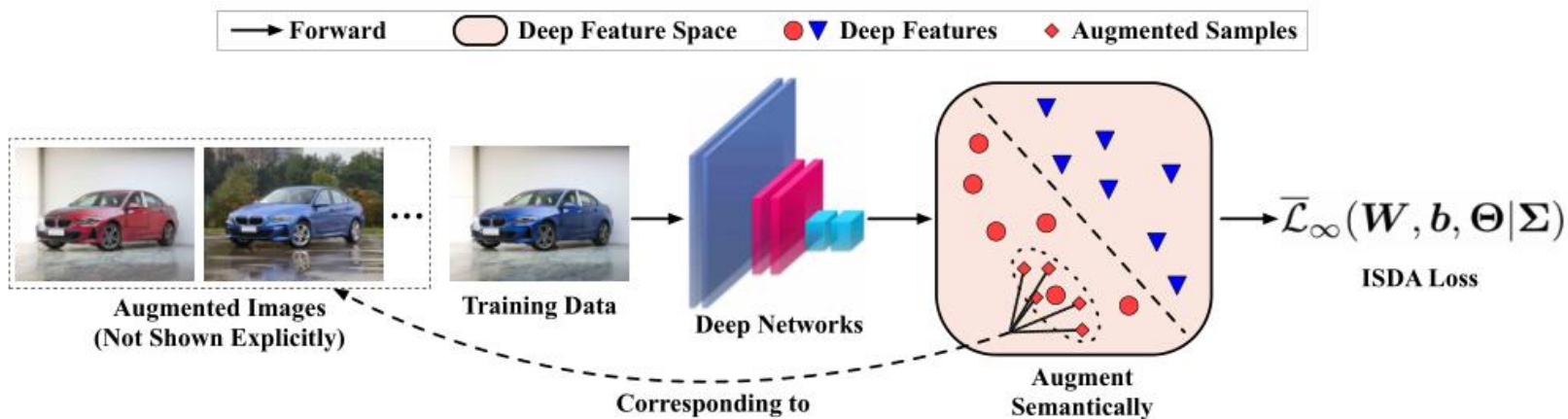
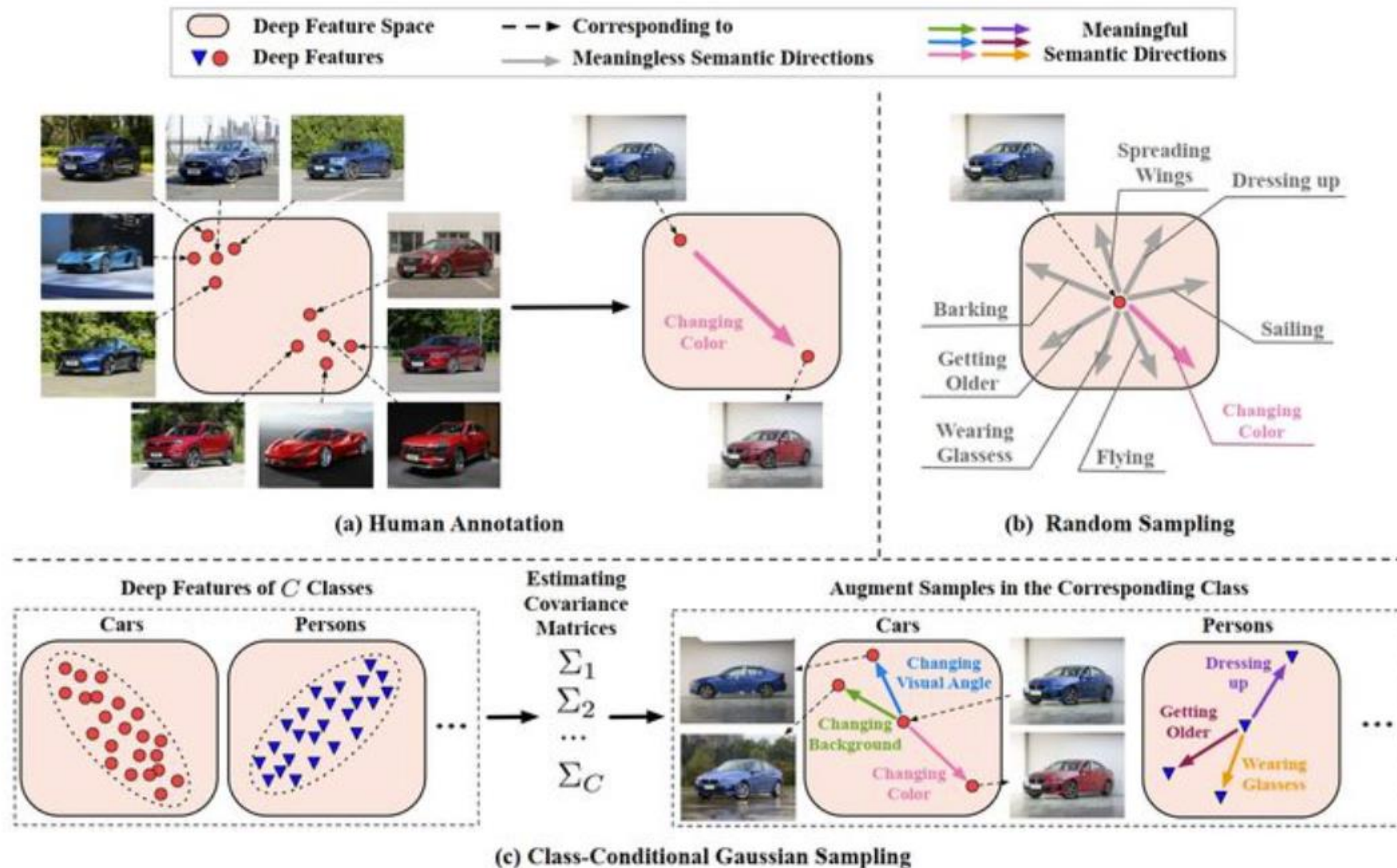


Figure 1: An **overview of ISDA**. Inspired by the observation that certain directions in the feature space correspond to meaningful semantic transformations, we augment the training data semantically by translating their features along **these semantic directions**, without involving auxiliary deep networks. The **directions** are obtained by sampling random vectors from a zero-mean normal distribution with dynamically estimated class-conditional covariance matrices. In addition, instead of performing augmentation explicitly, ISDA boils down to minimizing a closed-form upper-bound of the expected cross-entropy loss on the augmented training set, which makes our method highly efficient.

How to find meaningful semantic transformations?



✗ Manually searching for semantic directions is infeasible for large scale problems.

□ Online estimation of class-conditional covariance matrices.

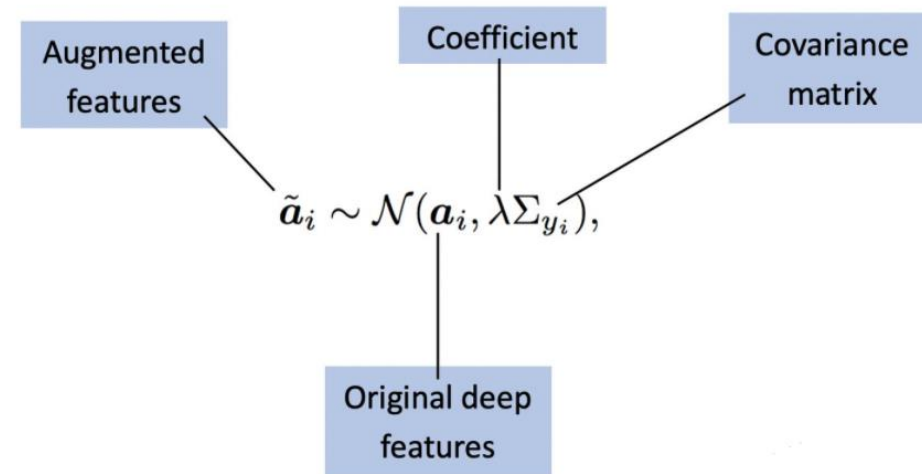
Then randomly sample vectors from a zero-mean multivariate normal distribution.

$$\mathcal{N}(0, \Sigma_{y_i})$$

covariance matrices for every class

$$\mathbf{a}_i = [a_{i1}, \dots, a_{iA}]^T = G(\mathbf{x}_i, \Theta)$$

$$\tilde{\mathbf{a}}_i \sim \mathcal{N}(\mathbf{a}_i, \lambda \Sigma_{y_i}) \quad \lambda = (t/T) \times \lambda_0$$



Dynamic estimation of covariance matrices

1. Estimation of average values of the features of j class at t step.

2. Estimation of average values of the features of j class s in t mini-batch.

$$\mu_j^{(t)} = \frac{n_j^{(t-1)} \mu_j^{(t-1)} + m_j^{(t)} \mu_j^{\prime(t)}}{n_j^{(t-1)} + m_j^{(t)}}, \quad (9)$$

$$\Sigma_j^{(t)} = \frac{n_j^{(t-1)} \Sigma_j^{(t-1)} + m_j^{(t)} \Sigma_j^{\prime(t)}}{n_j^{(t-1)} + m_j^{(t)}} + \frac{n_j^{(t-1)} m_j^{(t)} (\mu_j^{(t-1)} - \mu_j^{\prime(t)}) (\mu_j^{(t-1)} - \mu_j^{\prime(t)})^T}{(n_j^{(t-1)} + m_j^{(t)})^2}, \quad (10)$$

$$n_j^{(t)} = n_j^{(t-1)} + m_j^{(t)} \quad (11)$$

3. Total number of training samples belonging to j class in all t mini-batches

4. Total number of training samples belonging to j class in t mini-batches

Novel Loss

$$\{(\mathbf{a}_i^1, y_i), \dots, (\mathbf{a}_i^M, y_i)\}_{i=1}^N$$

cross-entropy (CE) loss:

$$\mathcal{L}_M(\mathbf{W}, \mathbf{b}, \Theta) = \frac{1}{N} \sum_{i=1}^N \frac{1}{M} \sum_{k=1}^M -\log\left(\frac{e^{\mathbf{w}_{y_i}^T \mathbf{a}_i^k + b_{y_i}}}{\sum_{j=1}^C e^{\mathbf{w}_j^T \mathbf{a}_i^k + b_j}}\right), \quad (2)$$

$M \rightarrow \infty$

Final goal

$$\mathcal{L}_\infty(\mathbf{W}, \mathbf{b}, \Theta | \Sigma) \leq \frac{1}{N} \sum_{i=1}^N -\log\left(\frac{e^{\mathbf{w}_{y_i}^T \mathbf{a}_i + b_{y_i}}}{\sum_{j=1}^C e^{\mathbf{w}_j^T \mathbf{a}_i + b_j + \frac{\lambda}{2} (\mathbf{w}_j^T - \mathbf{w}_{y_i}^T) \Sigma_{y_i} (\mathbf{w}_j - \mathbf{w}_{y_i})}}\right) \triangleq \bar{\mathcal{L}}_\infty. \quad (4)$$

Proof.

$$\mathcal{L}_\infty(\mathbf{W}, \mathbf{b}, \Theta | \Sigma) = \frac{1}{N} \sum_{i=1}^N \mathbb{E}_{\tilde{\mathbf{a}}_i} \left[\log \left(\sum_{j=1}^C e^{(\mathbf{w}_j^T - \mathbf{w}_{y_i}^T) \tilde{\mathbf{a}}_i + (b_j - b_{y_i})} \right) \right] \quad (5)$$

Jensen's inequality

$$\mathbb{E}[\log X] \leq \log \mathbb{E}[X]$$

$$\leq \frac{1}{N} \sum_{i=1}^N \log \left(\sum_{j=1}^C \mathbb{E}_{\tilde{\mathbf{a}}_i} \left[e^{(\mathbf{w}_j^T - \mathbf{w}_{y_i}^T) \tilde{\mathbf{a}}_i + (b_j - b_{y_i})} \right] \right) \quad (6)$$

$$= \frac{1}{N} \sum_{i=1}^N \log \left(\sum_{j=1}^C e^{(\mathbf{w}_j^T - \mathbf{w}_{y_i}^T) \mathbf{a}_i + (b_j - b_{y_i}) + \frac{\lambda}{2} (\mathbf{w}_j^T - \mathbf{w}_{y_i}^T) \Sigma_{y_i} (\mathbf{w}_j - \mathbf{w}_{y_i})} \right) \quad (7)$$

$$= \underline{\underline{\mathcal{L}}}_\infty. \quad (8)$$

Moment-generating function

$$\mathbb{E}[e^{tX}] = e^{t\mu + \frac{1}{2}\sigma^2 t^2}, \quad X \sim \mathcal{N}(\mu, \sigma^2)$$

Algorithm 1 The ISDA Algorithm.

- 1: **Input:** \mathcal{D} , λ_0
 - 2: Randomly initialize \mathbf{W} , \mathbf{b} and Θ
 - 3: **for** $t = 0$ **to** T **do**
 - 4: Sample a mini-batch $\{\mathbf{x}_i, y_i\}_{i=1}^B$ from \mathcal{D}
 - 5: Compute $\mathbf{a}_i = G(\mathbf{x}_i, \Theta)$
 - 6: Estimate the covariance matrices $\Sigma_1, \Sigma_2,$
 \dots, Σ_C
 - 7: Compute $\bar{\mathcal{L}}_\infty$ according to Eq. (4)
 - 8: Update \mathbf{W} , \mathbf{b} , Θ with SGD
 - 9: **end for**
 - 10: **Output:** \mathbf{W} , \mathbf{b} and Θ
-

Experiments

Table 1: Evaluation of ISDA on CIFAR with different models. The **average test error** over the last 10 epochs is calculated in each experiment, and we report mean values and standard deviations in three independent experiments. The best results are **bold-faced**.

Method	Params	CIFAR-10	CIFAR-100
ResNet-32 [4]	0.5M	7.39 ± 0.10%	31.20 ± 0.41%
ResNet-32 + ISDA	0.5M	7.09 ± 0.12%	30.27 ± 0.34%
ResNet-110 [4]	1.7M	6.76 ± 0.34%	28.67 ± 0.44%
ResNet-110 + ISDA	1.7M	6.33 ± 0.19%	27.57 ± 0.46%
SE-ResNet-110 [33]	1.7M	6.14 ± 0.17%	27.30 ± 0.03%
SE-ResNet-110 + ISDA	1.7M	5.96 ± 0.21%	26.63 ± 0.21%
Wide-ResNet-16-8 [34]	11.0M	4.25 ± 0.18%	20.24 ± 0.27%
Wide-ResNet-16-8 + ISDA	11.0M	4.04 ± 0.29%	19.91 ± 0.21%
Wide-ResNet-28-10 [34]	36.5M	3.82 ± 0.15%	18.53 ± 0.07%
Wide-ResNet-28-10 + ISDA	36.5M	3.58 ± 0.15%	17.98 ± 0.15%
ResNeXt-29, 8x64d [35]	34.4M	3.86 ± 0.14%	18.16 ± 0.13%
ResNeXt-29, 8x64d + ISDA	34.4M	3.67 ± 0.12%	17.43 ± 0.25%
DenseNet-BC-100-12 [5]	0.8M	4.90 ± 0.08%	22.61 ± 0.10%
DenseNet-BC-100-12 + ISDA	0.8M	4.54 ± 0.07%	22.10 ± 0.34%
DenseNet-BC-190-40 [5]	25.6M	3.52%	17.74%
DenseNet-BC-190-40 + ISDA	25.6M	3.24%	17.42%

Experiments

Table 2: Evaluation of ISDA with state-of-the-art *non-semantic augmentation techniques*. ‘AA’ refers to AutoAugment [32]. We report mean values and standard deviations in three independent experiments. The best results are **bold-faced**.

Dataset	Networks	Cutout [31]	Cutout + ISDA	AA [32]	AA + ISDA
CIFAR-10	Wide-ResNet-28-10 [34]	$2.99 \pm 0.06\%$	$2.83 \pm 0.04\%$	$2.65 \pm 0.07\%$	$2.56 \pm 0.01\%$
	Shake-Shake (26, 2x32d) [36]	$3.16 \pm 0.09\%$	$2.93 \pm 0.03\%$	$2.89 \pm 0.09\%$	$2.68 \pm 0.12\%$
	Shake-Shake (26, 2x112d) [36]	2.36%	2.25%	2.01%	1.82%
CIFAR-100	Wide-ResNet-28-10 [34]	$18.05 \pm 0.25\%$	$17.02 \pm 0.11\%$	$16.60 \pm 0.40\%$	$15.62 \pm 0.32\%$
	Shake-Shake (26, 2x32d) [36]	$18.92 \pm 0.21\%$	$18.17 \pm 0.08\%$	$17.50 \pm 0.19\%$	$17.21 \pm 0.33\%$
	Shake-Shake (26, 2x112d) [36]	$17.34 \pm 0.28\%$	$16.24 \pm 0.20\%$	$15.21 \pm 0.20\%$	$13.87 \pm 0.26\%$

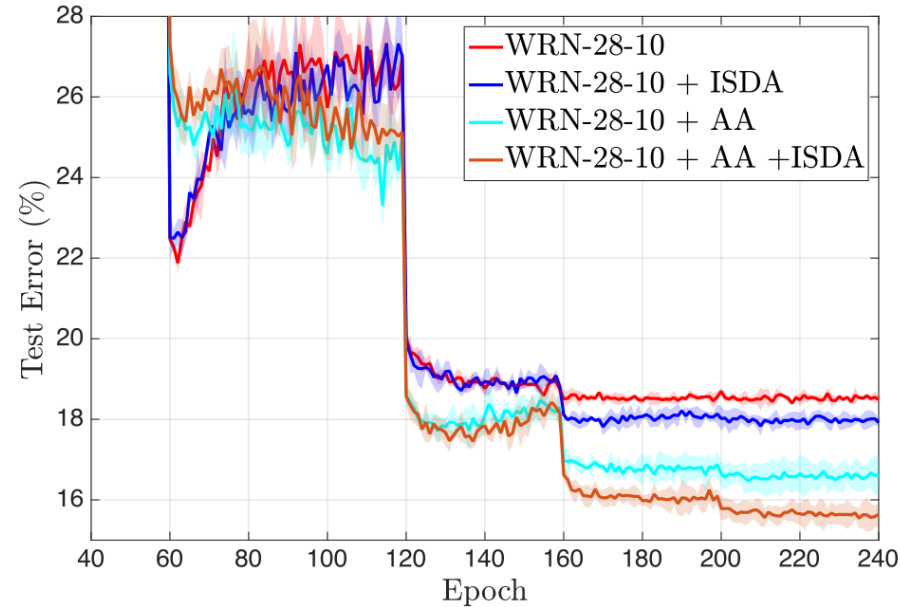


Figure 2: Curves of test errors on CIFAR-100 with Wide-ResNet (WRN).

Experiments

Table 3: Comparisons with the state-of-the-art methods. We report mean values and standard deviations of the test **error** in three independent experiments. Best results are **bold-faced**.

Method	ResNet-110		Wide-ResNet-28-10	
	CIFAR-10	CIFAR-100	CIFAR-10	CIFAR-100
Large Margin [18]	6.46±0.20%	28.00±0.09%	3.69±0.10%	18.48±0.05%
Disturb Label [38]	6.61±0.04%	28.46±0.32%	3.91±0.10%	18.56±0.22%
Focal Loss [17]	6.68±0.22%	28.28±0.32%	3.62±0.07%	18.22±0.08%
Center Loss [22]	6.38±0.20%	27.85±0.10%	3.76±0.05%	18.50±0.25%
L_q Loss [16]	6.69±0.07%	28.78±0.35%	3.78±0.08%	18.43±0.37%
WGAN [39]	6.63±0.23%	-	3.81±0.08%	-
CGAN [40]	6.56±0.14%	28.25±0.36%	3.84±0.07%	18.79±0.08%
ACGAN [41]	6.32±0.12%	28.48±0.44%	3.81±0.11%	18.54±0.05%
infoGAN [42]	6.59±0.12%	27.64±0.14%	3.81±0.05%	18.44±0.10%
Basic	6.76±0.34%	28.67±0.44%	-	-
Basic + Dropout	6.23±0.11%	27.11±0.06%	3.82±0.15%	18.53±0.07%
ISDA	6.33±0.19%	27.57±0.46%	-	-
ISDA + Dropout	5.98±0.20%	26.35±0.30%	3.58±0.15%	17.98±0.15%

For generator-based semantic augmentation methods

Wide-ResNet-28-10 and ResNeXt-29, 8x64d, our method outperforms the competitive baselines by nearly 0.7%. Compared to ResNets, DenseNets generally suffer less from overfitting due to their architecture design, thus appear to benefit less from our algorithm.

Experiments

Table 4: Evaluation of ISDA on ImageNet.

Method	Top-1	Top-5
ResNet-50 [4]	23.58%	6.92%
ResNet-50 + ISDA	23.30%	6.82%
ResNet-152 [4]	21.65%	6.01%
ResNet-152 + ISDA	21.20%	5.67%
ResNeXt-50, 32x4d [35]	22.42%	6.42%
ResNeXt-50, 32x4d + ISDA	21.88%	6.23%

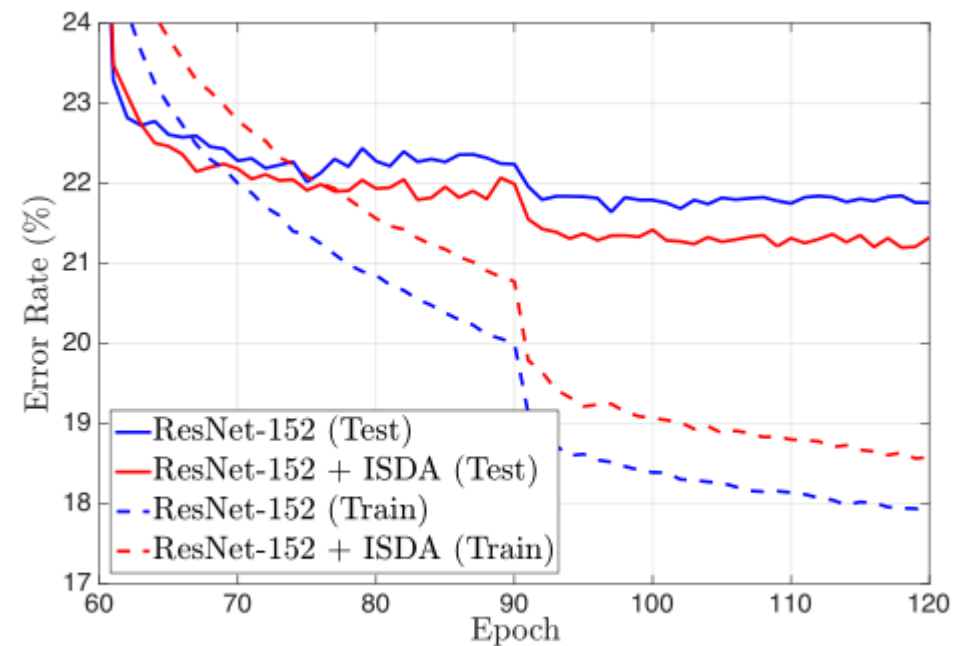


Figure 3: Training and test errors on ImageNet.

Table 5: The ablation study for ISDA.

Setting	CIFAR-10	CIFAR-100
Basic	$3.82 \pm 0.15\%$	$18.58 \pm 0.10\%$
Identity matrix	$3.63 \pm 0.12\%$	$18.53 \pm 0.02\%$
Diagonal matrix	$3.70 \pm 0.15\%$	$18.23 \pm 0.02\%$
Single covariance matrix	$3.67 \pm 0.07\%$	$18.29 \pm 0.13\%$
Constant λ_0	$3.69 \pm 0.08\%$	$18.33 \pm 0.16\%$
ISDA	$3.58 \pm 0.15\%$	$17.98 \pm 0.15\%$

Experiments

Visualization Results

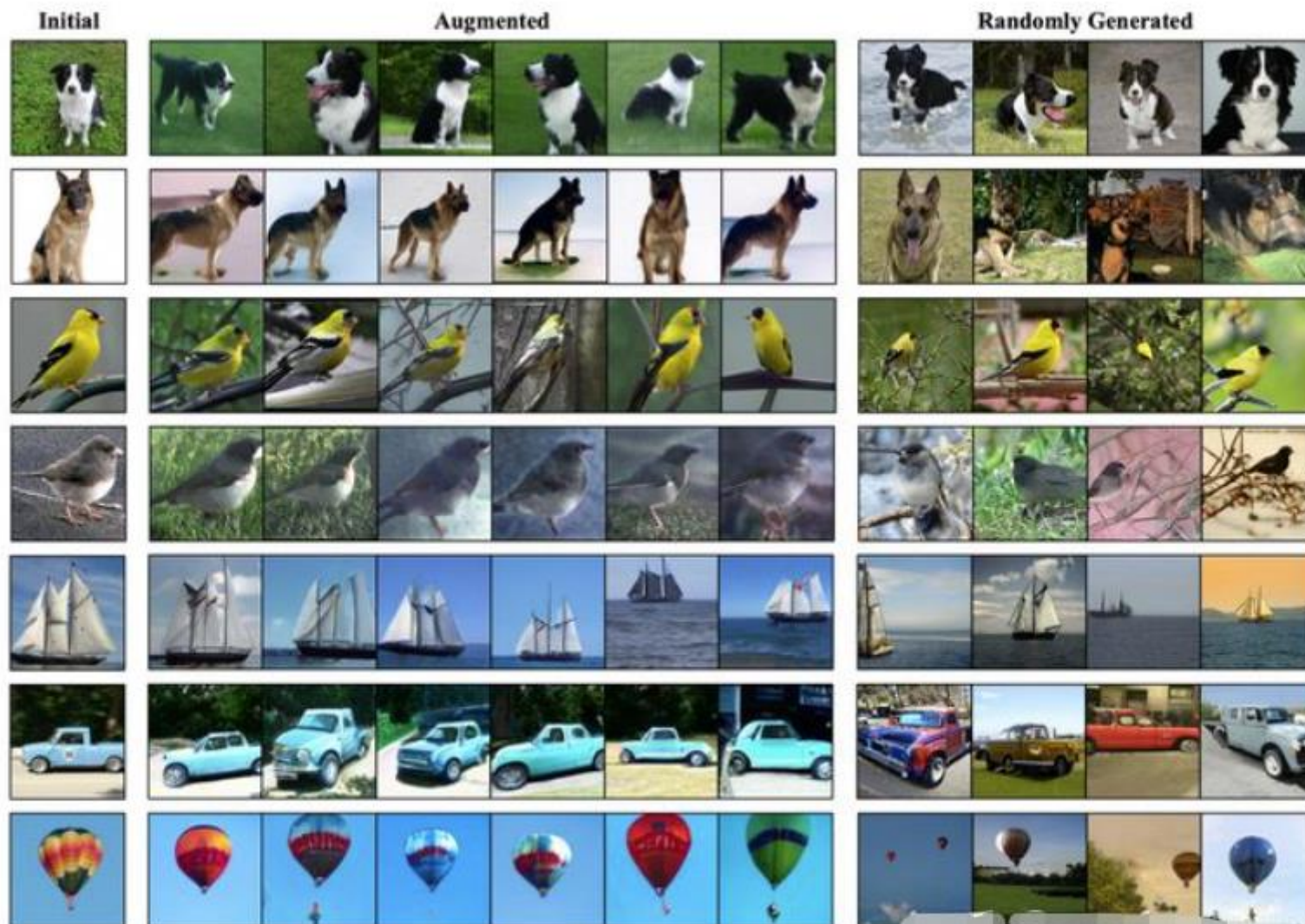


Fig. 7. Visualization of the semantically augmented images on ImageNet. ISDA is able to alter the semantics of images that are unrelated to the class identity, like backgrounds, actions of animals, visual angles, etc. We also present the randomly generated images of the same class.

Thanks

