

# Bag of Tricks for Image Classification with Convolutional Neural Networks

Tong He, Zhi Zhang, Hang Zhang, Zhongyue Zhang, Junyuan Xie, Mu Li

Amazon Web Services

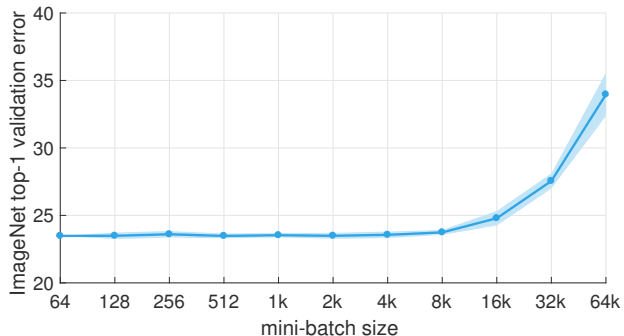
December 14, 2018

## 4% increase in ResNet50's accuracy

Model	FLOPs	top-1	top-5
ResNet-50 [9]	3.9 G	75.3	92.2
ResNeXt-50 [27]	4.2 G	77.8	-
SE-ResNet-50 [12]	3.9 G	76.71	93.38
SE-ResNeXt-50 [12]	4.3 G	78.90	94.51
DenseNet-201 [13]	4.3 G	77.42	93.66
ResNet-50 + tricks (ours)	4.3 G	<b>79.29</b>	<b>94.63</b>

Table 1: **Computational costs and validation accuracy of various models.** ResNet, trained with our “tricks”, is able to outperform newer and improved architectures trained with standard pipeline.

## Improving training efficiency: Large batch training



## Improving training efficiency: Large batch training

1. Scale learning rate with batch size (Linear Scaling).
2. Use fixed, small learning rate for first few epochs (Warmup).
3. Initialize  $\gamma$  to zero for Batch Normalization layer.
4. No bias decay.

# Improving training efficiency

Model	Efficient			Baseline		
	Time/epoch	Top-1	Top-5	Time/epoch	Top-1	Top-5
ResNet-50	<b>4.4 min</b>	<b>76.21</b>	<b>92.97</b>	13.3 min	75.87	92.70
Inception-V3	<b>8 min</b>	<b>77.50</b>	<b>93.60</b>	19.8 min	77.32	93.43
MobileNet	<b>3.7 min</b>	<b>71.90</b>	<b>90.47</b>	6.2 min	69.03	88.71

Table 3: Comparison of the training time and validation accuracy for ResNet-50 between the baseline (BS=256 with FP32) and a more hardware efficient setting (BS=1024 with FP16).

# Architecture Modifications

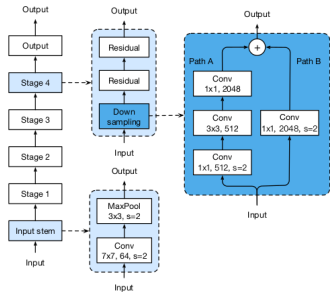


Figure 1: The architecture of ResNet-50. The convolution kernel size, output channel size and stride size (default is 1) are illustrated, similar for pooling layers.

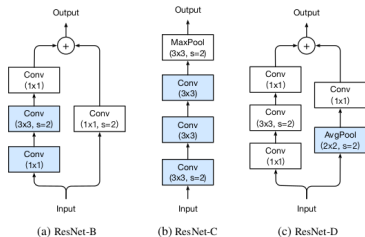


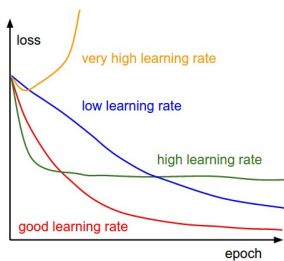
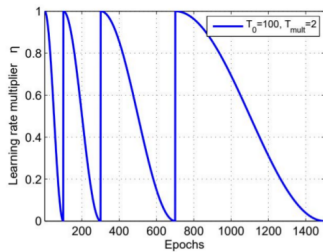
Figure 2: Three ResNet tweaks. ResNet-B modifies the downsampling block of Resnet. ResNet-C further modifies the input stem. On top of that, ResNet-D again modifies the downsampling block.

## ResNet-D is 1% more accurate

Model	#params	FLOPs	Top-1	Top-5
ResNet-50	25 M	<b>3.8 G</b>	76.21	92.97
ResNet-50-B	25 M	4.1 G	76.66	93.28
ResNet-50-C	25 M	4.3 G	76.87	93.48
ResNet-50-D	25 M	4.3 G	<b>77.16</b>	<b>93.52</b>

Table 5: Compare ResNet-50 with three model tweaks on model size, FLOPs and ImageNet validation accuracy.

# Training Refinements: Cosine LR decay

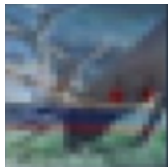




# Training Refinements: Label Smoothing

If the label is  $[0, 0, 1]$   
Use  $[0.05, 0.05, 0.9]$

## Training Refinements: Mixup



Linear interpolation between deer and ship classes

Label = [..., 0.4, ..., 0.6, ...]

# Image classification

Refinements	ResNet-50-D		Inception-V3		MobileNet	
	Top-1	Top-5	Top-1	Top-5	Top-1	Top-5
Efficient	77.16	93.52	77.50	93.60	71.90	90.53
+ cosine decay	77.91	93.81	78.19	94.06	72.83	91.00
+ label smoothing	78.31	94.09	78.40	94.13	72.93	91.14
+ distill w/o mixup	78.67	94.36	78.26	94.01	71.97	90.89
+ mixup w/o distill	79.15	94.58	<b>78.77</b>	<b>94.39</b>	<b>73.28</b>	<b>91.30</b>
+ distill w/ mixup	<b>79.29</b>	<b>94.63</b>	78.34	94.16	72.51	91.02

Table 6: The validation accuracies on ImageNet for stacking training refinements one by one. The baseline models are obtained from Section 3.

# Image detection

Refinement	Top-1	mAP
B-standard	76.14	77.54
D-efficient	77.16	78.30
+ cosine	77.91	79.23
+ smooth	78.34	80.71
+ distill w/o mixup	78.67	80.96
+ mixup w/o distill	79.16	81.10
+ distill w/ mixup	79.29	<b>81.33</b>

Table 8: Faster-RCNN performance with various pre-trained base networks evaluated on Pascal VOC.

# Image segmentation

Refinement	Top-1	PixAcc	mIoU
B-standard	76.14	78.08	37.05
D-efficient	77.16	78.88	38.88
+ cosine	77.91	<b>79.25</b>	<b>39.33</b>
+ smooth	78.34	78.64	38.75
+ distill w/o mixup	78.67	78.97	38.90
+ mixup w/o distill	79.16	78.47	37.99
+ mixup w/ distill	79.29	78.72	38.40

Table 9: FCN performance with various base networks evaluated on ADE20K.