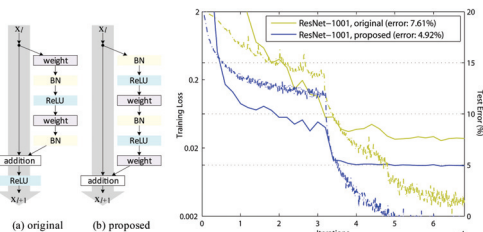


Identity Mappings in Deep Residual Networks

Kaiming He, Xiangyu Zhang, Shaoqing Ren, Jian Sun
Microsoft Research Asia (MSRA)

Highlights

- Novel pre-activation residual structure
- Improved results using **1001**-layer ResNet on CIFAR-10 and **200**-layer on ImageNet



Results on CIFAR

CIFAR-10	error (%)	CIFAR-100	error (%)
NIN [13]	8.81	NIN [13]	35.68
DSN [14]	8.22	DSN [14]	34.57
FitNet [15]	8.39	FitNet [15]	35.04
Highway [7]	7.72	Highway [7]	32.39
ELU [12]	6.55	ELU [12]	24.28
FitResNet, LSUV [16]	5.84	FitResNet, LSUV [16]	27.66
ResNet-110 [1] (1.7M)	6.61	ResNet-164 [1] (1.7M)	25.16
ResNet-1202 [1] (19.4M)	7.93	ResNet-1001 [1] (10.2M)	27.82
ResNet-164 [ours] (1.7M)	5.46	ResNet-164 [ours] (1.7M)	24.33
ResNet-1001 [ours] (10.2M)	4.92 (4.89±0.14)	ResNet-1001 [ours] (10.2M)	22.71 (22.68±0.22)
ResNet-1001 [ours] (10.2M) ²	4.62 (4.69±0.20)		

Results on ImageNet

method	train crop size	test crop size	top-1 (%)	top-5 (%)
ResNet-152, original Residual Unit [1]	224×224	224×224	23.0	6.7
ResNet-152, original Residual Unit [1]	224×224	320×320	21.3	5.5
ResNet-152, proposed Residual Unit	224×224	320×320	21.1	5.5
ResNet-200, original Residual Unit [1]	224×224	320×320	21.8	6.0
ResNet-200, proposed Residual Unit	224×224	320×320	20.7	5.3
Inception v3 [17]	299×299	299×299	21.2	5.6

Importance of Identity Skip Connections

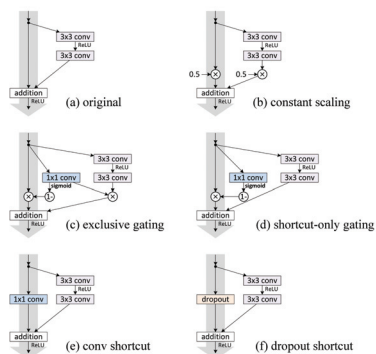
- ResNet with Identity mapping

$$x_L = x_i + \sum_{l=i+1}^{L-1} \mathcal{F}(x_l, \mathcal{W}_l), \quad \frac{\partial \mathcal{E}}{\partial x_i} = \frac{\partial \mathcal{E}}{\partial x_L} \frac{\partial x_L}{\partial x_i} = \frac{\partial \mathcal{E}}{\partial x_i} \left(1 + \frac{\partial}{\partial x_i} \sum_{l=i+1}^{L-1} \mathcal{F}(x_l, \mathcal{W}_l) \right)$$

- What if we break the identity shortcut?

$$x_L = \prod_{i=1}^{L-1} \lambda_i x_i + \sum_{i=1}^{L-1} \tilde{\mathcal{F}}(x_i, \mathcal{W}_i), \quad \frac{\partial \mathcal{E}}{\partial x_i} = \frac{\partial \mathcal{E}}{\partial x_L} \left(\prod_{l=i+1}^{L-1} \lambda_l + \frac{\partial}{\partial x_i} \sum_{l=i+1}^{L-1} \tilde{\mathcal{F}}(x_l, \mathcal{W}_l) \right)$$

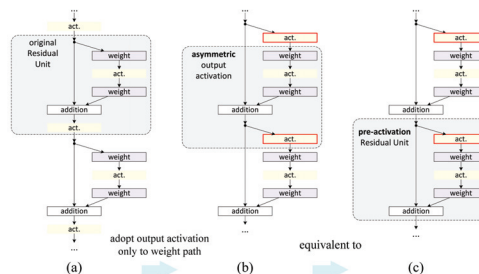
- Various types of shortcut connections



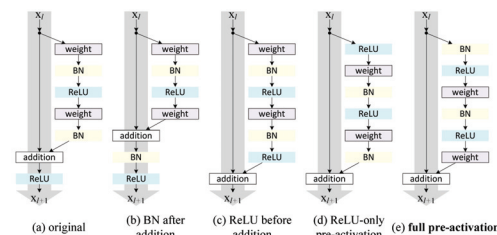
case	Fig.	on shortcut	on \mathcal{F}	error (%)	remark
original [1]	Fig. 2(a)	1	1	6.61	
constant scaling	Fig. 2(b)	0	1	fail	This is a plain net
		0.5	1	fail	
exclusive gating	Fig. 2(c)	$1 - g(x)$	$g(x)$	fail	init $b_g=0$ to -5
		$1 - g(x)$	$g(x)$	8.70	init $b_g=-6$
shortcut-only gating	Fig. 2(d)	$1 - g(x)$	1	12.86	init $b_g=0$
		$1 - g(x)$	1	6.91	init $b_g=-6$
1×1 conv shortcut	Fig. 2(e)	1×1 conv	1	12.22	
dropout shortcut	Fig. 2(f)	dropout 0.5	1	fail	

Usage of Activation Function

- Post-activation to pre-activation



- Various usage of activation function



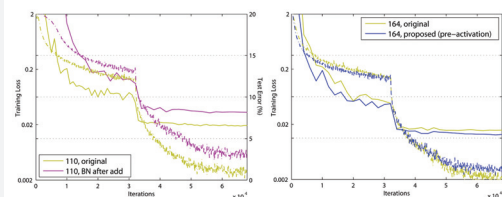
case	Fig.	ResNet-110	ResNet-164
original Residual Unit [1]	Fig. 4(a)	6.61	5.93
BN after addition	Fig. 4(b)	8.17	6.50
ReLU before addition	Fig. 4(c)	7.84	6.14
ReLU-only pre-activation	Fig. 4(d)	6.71	5.91
full pre-activation	Fig. 4(e)	6.37	5.46

Analysis of Pre-activation Structure

- Ease of optimization

dataset	network	baseline unit	pre-activation unit
CIFAR-10	ResNet-110(1layer)	9.90	8.91
	ResNet-110	6.61	6.37
	ResNet-164	5.93	5.46
CIFAR-100	ResNet-1001	7.61	4.92
	ResNet-164	25.16	24.33
	ResNet-1001	27.82	22.71

- Reducing overfitting



Code available:

- Deep Residual Networks with 1K Layers: <https://github.com/KaimingHe/resnet-1k-layers>

[1] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. "Deep Residual Learning for Image Recognition". IEEE Conference on Computer Vision and Pattern Recognition (CVPR), 2016