

Deep Residual Learning for Image Recognition

By
Anup Joseph





Background

In theory for CNNs deeper == better

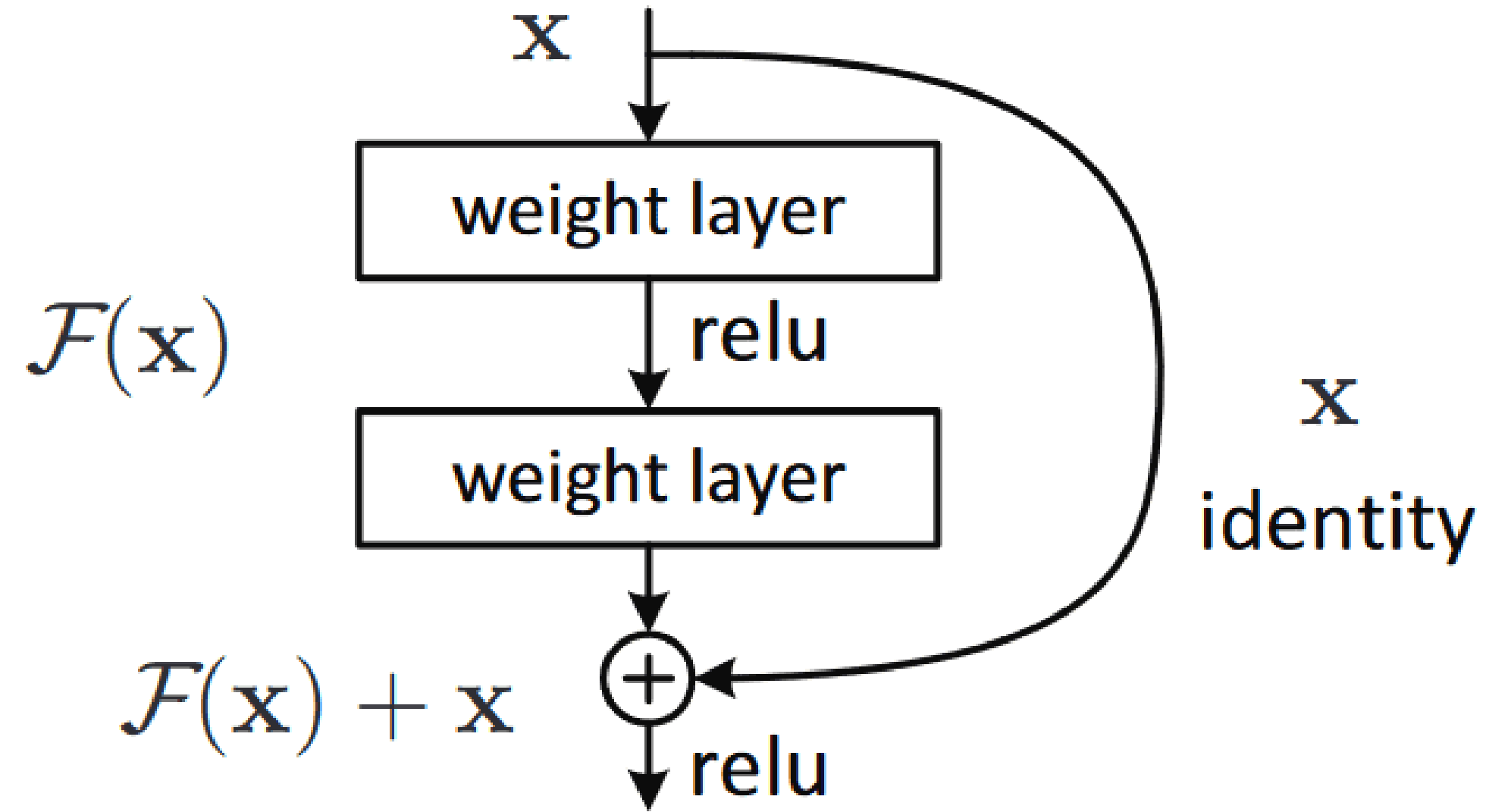
In practice however, after some depth the performance decreases

This was one of the bottlenecks for VGG. They started to lose generalization capabilities as they went very deep.

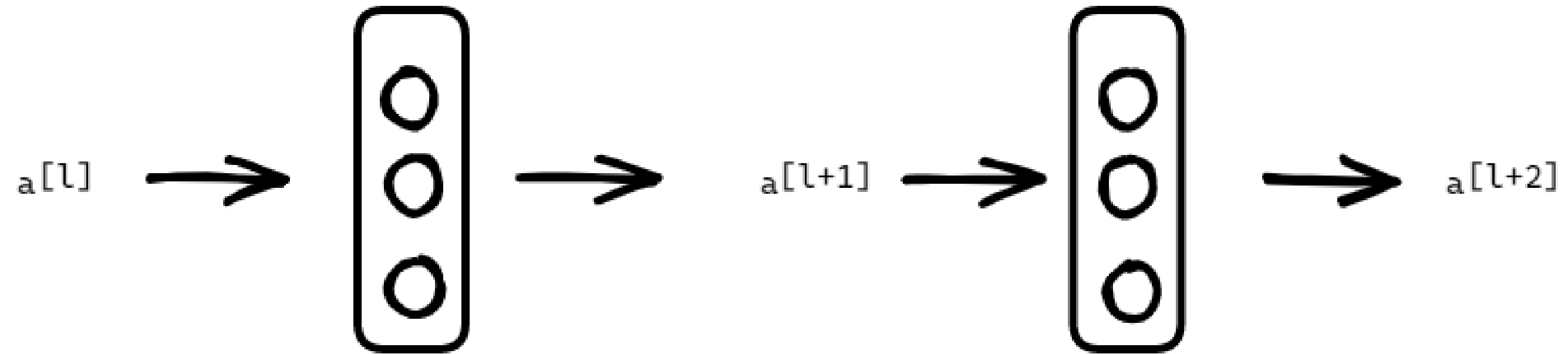
Skip Connections

Skip connections are the central idea of the Resnet paper.

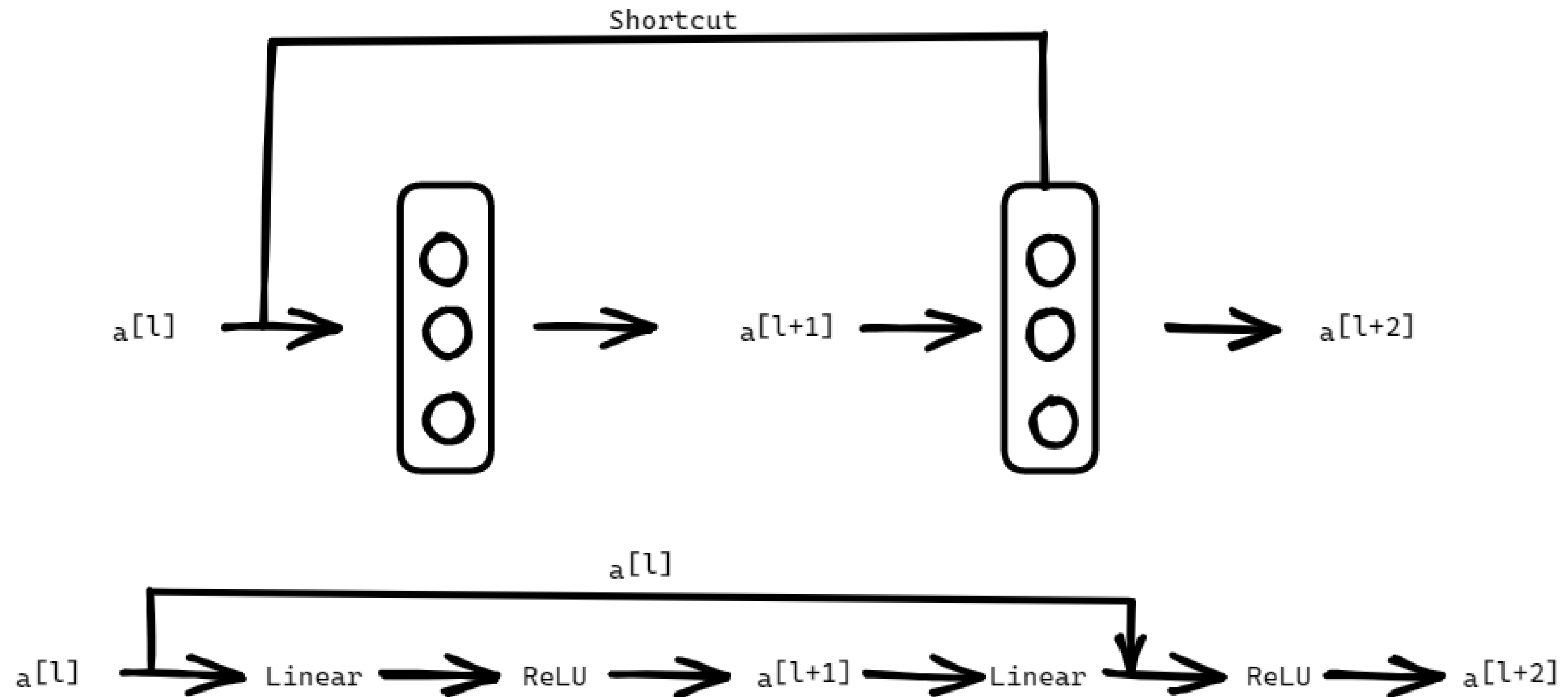
The idea is to take the output of one layer and then plug it to a layer much further in the network



Plain Neural Network



Skip connections



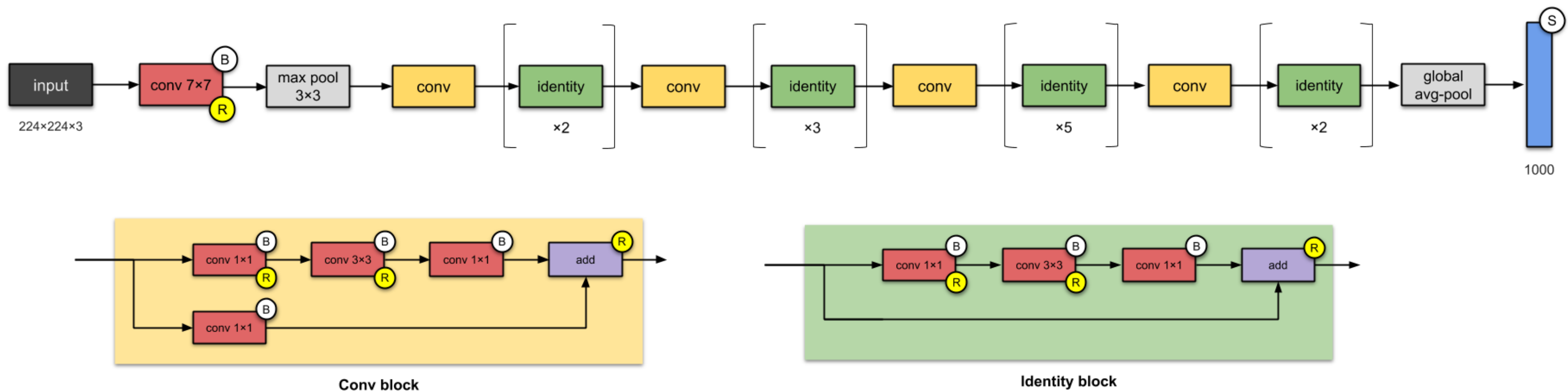
Configurations

Dataset Augmentation

- A 224x224 random crop of the image/horizontal flip is taken.
- Normalization is done by subtracting the per-pixel mean of the image from the original image.

Training Settings

- Learning Rate - 0.1 and then divided by a factor of 10 on plateau
- Iterations = 60×10^4 , weight decay - 0.0001, momentum - 0.9



Model Architecture

Resnet Models are made as a combination of "Conv" and "Identity" blocks.

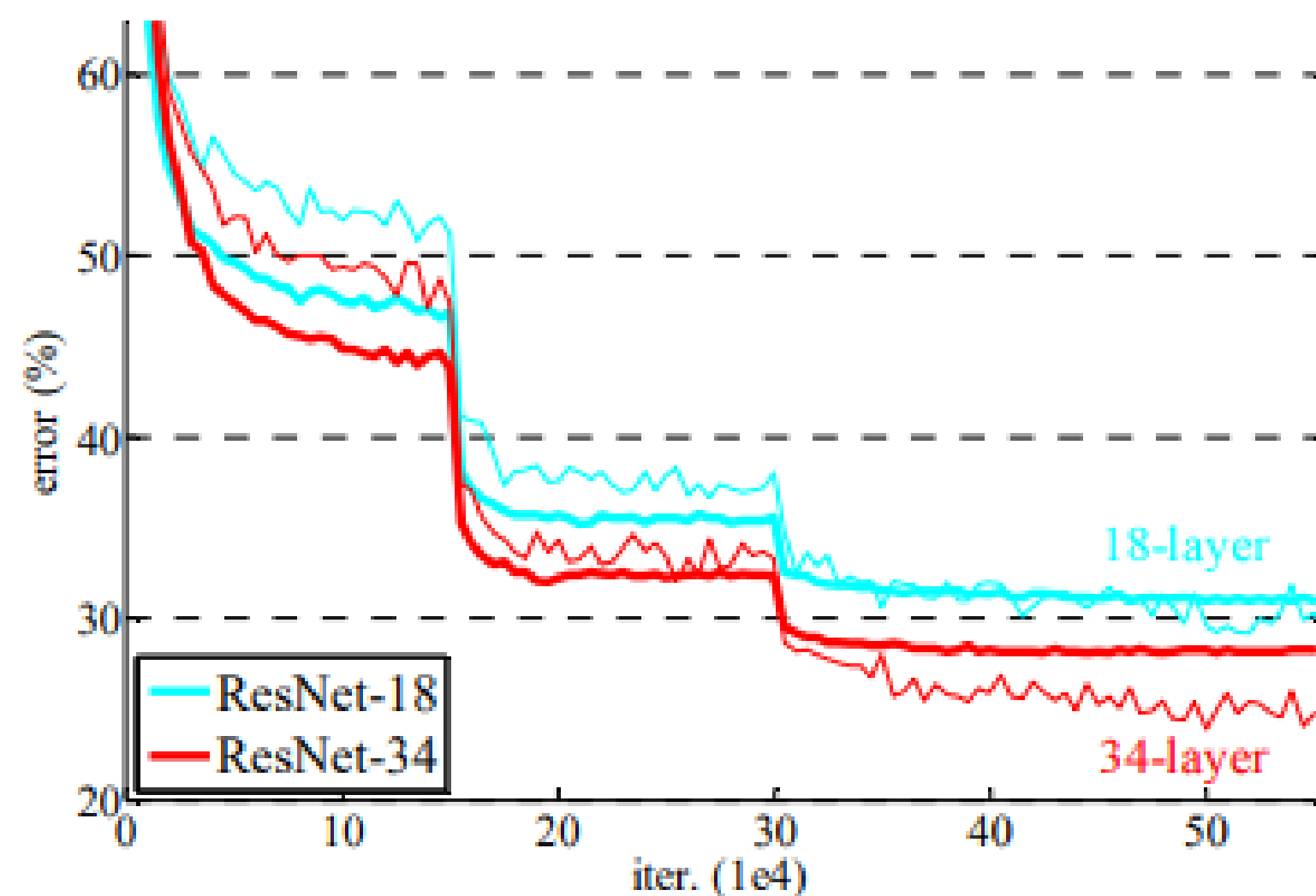
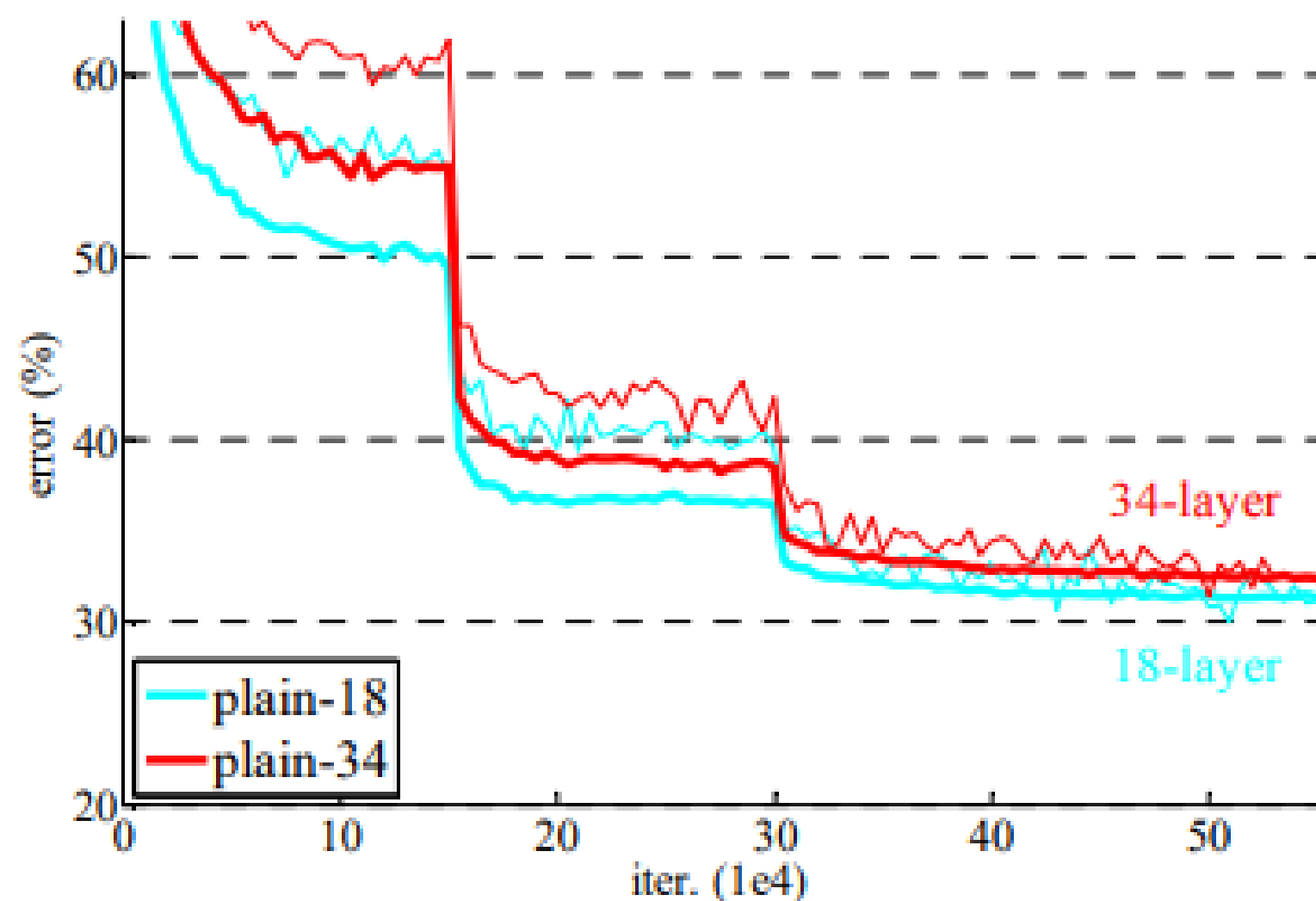


Figure 4. Training on **ImageNet**. Thin curves denote training error, and bold curves denote validation error of the center crops. Left: plain networks of 18 and 34 layers. Right: ResNets of 18 and 34 layers. In this plot, the residual networks have no extra parameter compared to their plain counterparts.

Results

Resnet was tested for 3 tasks and managed to perform exceptionally well in all of them

This the result on the ImageNet classification tasks and are the best achieved results at that point.

| method | top-1 err. | top-5 err. |
|----------------------------|--------------|-------------------|
| VGG [41] (ILSVRC'14) | - | 8.43 [†] |
| GoogLeNet [44] (ILSVRC'14) | - | 7.89 |
| VGG [41] (v5) | 24.4 | 7.1 |
| PReLU-net [13] | 21.59 | 5.71 |
| BN-inception [16] | 21.99 | 5.81 |
| ResNet-34 B | 21.84 | 5.71 |
| ResNet-34 C | 21.53 | 5.60 |
| ResNet-50 | 20.74 | 5.25 |
| ResNet-101 | 19.87 | 4.60 |
| ResNet-152 | 19.38 | 4.49 |

Table 4. Error rates (%) of **single-model** results on the ImageNet validation set (except [†] reported on the test set).



Object Detection

Object detection task on the MS-COCO dataset.
Its better than the best in class models for the tasks.

| training data | COCO train | | COCO trainval | |
|------------------------------------|------------|------------|---------------|------------|
| test data | COCO val | | COCO test-dev | |
| mAP | @.5 | @[.5, .95] | @.5 | @[.5, .95] |
| baseline Faster R-CNN (VGG-16) | 41.5 | 21.2 | | |
| baseline Faster R-CNN (ResNet-101) | 48.4 | 27.2 | | |
| +box refinement | 49.9 | 29.9 | | |
| +context | 51.1 | 30.0 | 53.3 | 32.2 |
| +multi-scale testing | 53.8 | 32.5 | 55.7 | 34.9 |
| ensemble | | | 59.0 | 37.4 |

Table 9. Object detection improvements on MS COCO using Faster R-CNN and ResNet-101.

| system | net | data | mAP | areo | bike | bird | boat | bottle | bus | car | cat | chair | cow | table | dog | horse | mbike | person | plant | sheep | sofa | train | tv |
|-------------|------------|------------|------|------|------|------|------|--------|------|------|------|-------|------|-------|------|-------|-------|--------|-------|-------|------|-------|------|
| baseline | VGG-16 | 07+12 | 73.2 | 76.5 | 79.0 | 70.9 | 65.5 | 52.1 | 83.1 | 84.7 | 86.4 | 52.0 | 81.9 | 65.7 | 84.8 | 84.6 | 77.5 | 76.7 | 38.8 | 73.6 | 73.9 | 83.0 | 72.6 |
| baseline | ResNet-101 | 07+12 | 76.4 | 79.8 | 80.7 | 76.2 | 68.3 | 55.9 | 85.1 | 85.3 | 89.8 | 56.7 | 87.8 | 69.4 | 88.3 | 88.9 | 80.9 | 78.4 | 41.7 | 78.6 | 79.8 | 85.3 | 72.0 |
| baseline+++ | ResNet-101 | COCO+07+12 | 85.6 | 90.0 | 89.6 | 87.8 | 80.8 | 76.1 | 89.9 | 89.9 | 89.6 | 75.5 | 90.0 | 80.7 | 89.6 | 90.3 | 89.1 | 88.7 | 65.4 | 88.1 | 85.6 | 89.0 | 86.8 |

Table 10. Detection results on the PASCAL VOC 2007 test set. The baseline is the Faster R-CNN system. The system “baseline+++” include box refinement, context, and multi-scale testing in Table 9.

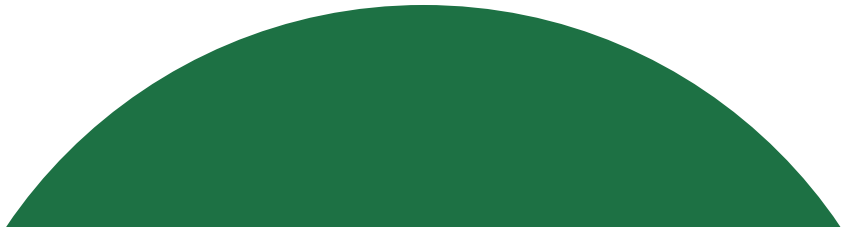


Image Localization

Using an ensemble of networks for classification, ResNet achieve a top-5 localization error of 9.0% on the test set.

This number significantly outperforms the ILSVRC 14 results showing a 64% relative reduction of error. This result won the 1st place in the ImageNet localization task in ILSVRC 2015.

| method | top-5 localization err | |
|----------------------------|------------------------|------------|
| | val | test |
| OverFeat [40] (ILSVRC'13) | 30.0 | 29.9 |
| GoogLeNet [44] (ILSVRC'14) | - | 26.7 |
| VGG [41] (ILSVRC'14) | 26.9 | 25.3 |
| ours (ILSVRC'15) | 8.9 | 9.0 |

Table 14. Comparisons of localization error (%) on the ImageNet dataset with state-of-the-art methods.



Thank
you

