

WEDNESDAY

FEB '18

7

038-327 WK-06

# [Original RESNET Paper] Deep Residual Learning for Image Recognition

How deep should we make our  
Neural Network?

- It depends on
  - a) Complexity of task at hand
  - b) Available computational capacity in the time of training
  - c) Available computational capacity in the time of inference
- If the task needs lot of parameters
  - a) Can we train very deep networks efficiently using current optimization solvers
  - b) Is training a better model as simple as adding more & more layers?  
Answer is no

February 2018

Mon	Tue	Wed	Thu	Fri	Sat	Sun	Mon	Tue	Wed	Thu	Fri	Sat	Sun	Mon	Tue	Wed	Thu	Fri	Sat	Sun
			1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18
19	20	21	22	23	24	25	26	27	28											



Why is it not OK to just add more and more layers?

THURSDAY

FEB '18

8

039-326 WK-06

- Because it causes Vanishing gradient problem
- Degradation Problem.

Degradation Problem-

- Adding more non-linear layer (adding more nodes) lead to drop in accuracy
- Our current optimization solvers are not able to approximate the identity mapping of a stack of added non-linear layers.

Residual Learning

$H(x)$  is true mapping function we want to learn.

$$F(x) = H(x) - x$$

Residual block adds explicitly identity connections throughout the network to help learning the required identity mappings

March 2018

Mon	Tue	Wed	Thu	Fri	Sat	Sun	Mon	Tue	Wed	Thu	Fri	Sat	Sun	Mon	Tue	Wed	Thu	Fri	Sat	Sun
			1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18
19	20	21	22	23	24	25	26	27	28	29	30	31								

Opportunities are usually disguised as hard work, so most people don't recognize them.

08-02

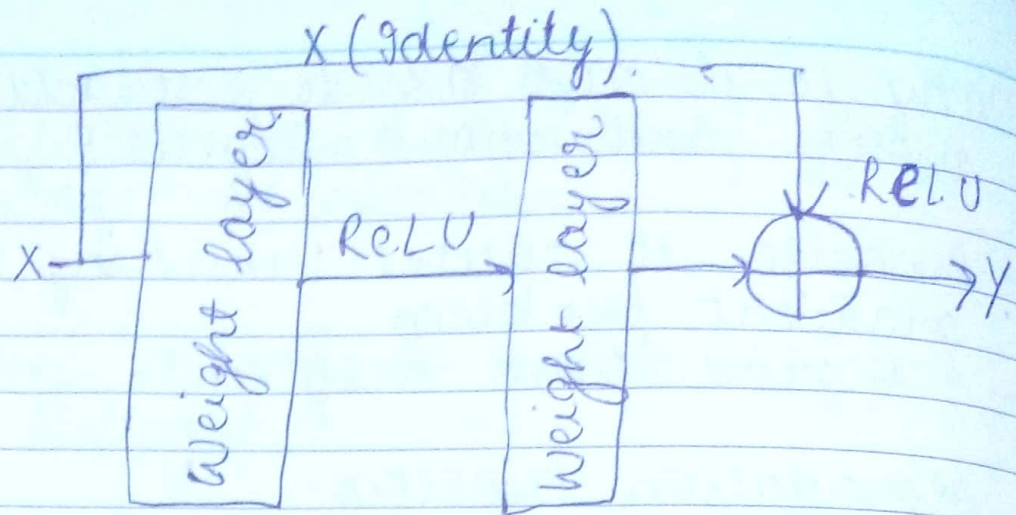


FRIDAY

FEB '18

9

040-325 WK-06



While moving from one kind of layer to different kind of layer the connection is shown by dotted connection which is linear projection for dimension mapping.

$$Y = F(X; \{W_i\}) + W_s X$$

for making two different layer dimension of same dimension, we do padding. We can also use Matrix multiplication in order to project dimension of next layer to previous layer, which adds more ~~data~~ parameters to the data.

Batchnormalization - At end of each layer, we need to normalize the output because it is in process of learning, bias and weights, the gradient can go to huge value or.

February 2018

Mon	Tue	Wed	Thu	Fri	Sat	Sun	Mon	Tue	Wed	Thu	Fri	Sat	Sun	Mon	Tue	Wed	Thu	Fri	Sat	Sun
			1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18
19	20	21	22	23	24	25	26	27	28											



go to zero so no charges happen.

Key take aways

- Residual Network eases the training process
- Addressing 'the degradation problem' in training process
- Leveraging deeper representations of neural network for image recognition tasks.
- Training neural network with more than 1200 layers.

Other than tackling the vanishing gradient problem, the architecture of ResNet encourage feature reuse making network highly parameter efficient.

We know  $l$ -th layer will have  $K * (l-1) + K_0$  feature maps

where  $K_0$  is number of channels in the input image,  $K$  is growth rate

March 2018

Mon	Tue	Wed	Thu	Fri	Sat	Sun	Mon	Tue	Wed	Thu	Fri	Sat	Sun	Mon	Tue	Wed	Thu	Fri	Sat	Sun
			1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18
19	20	21	22	23	24	25	26	27	28	29	30	31								

Now is the time to act and solve all problems and not to succumb to lame excuses later.



SUNDAY

FEB '18

11

042-323 WK-06

$1 \times 1$  convolutional bottleneck layer is also used to reduce the number of feature maps before the expensive  $3 \times 3$  convolution.

One major drawback of ResNet is the ~~of deeper~~

One major drawback of training deep neural Network is that it requires weeks for training, making it practically infeasible, hence counter-intuitive method of randomly dropping layers during training and using the full network in testing. In training time, each layer has "survival probability" and is randomly dropped. In testing time, all blocks are kept active and re-calibrated according to its survival probability during training.

$H_l$  be the output of  $l$ -th residual block,  $f_l$  be the mapping defined by the  $l$ -th block's weighted mapping,  $b_l$  be a Bernoulli random variable (0 or 1) (0 or 1, indicating whether block is active or not) during training.

February 2018

Mon	Tue	Wed	Thu	Fri	Sat	Sun	Mon	Tue	Wed	Thu	Fri	Sat	Sun	Mon	Tue	Wed	Thu	Fri	Sat	Sun
			1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18
19	20	21	22	23	24	25	26	27	28											

Truly great people usually don't feel important; they make others feel important.



$$H_l = \text{ReLU}(b_l * f_l(H_{l-1}) + \text{id}(H_{l-1}))$$

$b_l = 1$ , the block is normal residual block

$b_l = 0$ , the above formula becomes

$$H_l = \text{ReLU}(\text{id}(H_{l-1}))$$

Let  $p_l$  be survival probability of layer  $l$  during training, during test we have.

$$H_l = \text{ReLU}(p_l * f_l(H_{l-1}) + \text{id}(H_{l-1}))$$

Since earlier layer use extract low level feature, it should not be dropped too frequently, hence resulting rule becomes.

$$p_l = 1 - \frac{l}{L} (1 - p_L)$$

$L$  - denotes total number of blocks

In Dropout, it drops part of hidden unit in one layer during training whereas in ResNet, it drops entire layer

March 2018

Mon	Tue	Wed	Thu	Fri	Sat	Sun	Mon	Tue	Wed	Thu	Fri	Sat	Sun	Mon	Tue	Wed	Thu	Fri	Sat	Sun
			1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18
19	20	21	22	23	24	25	26	27	28	29	30	31								

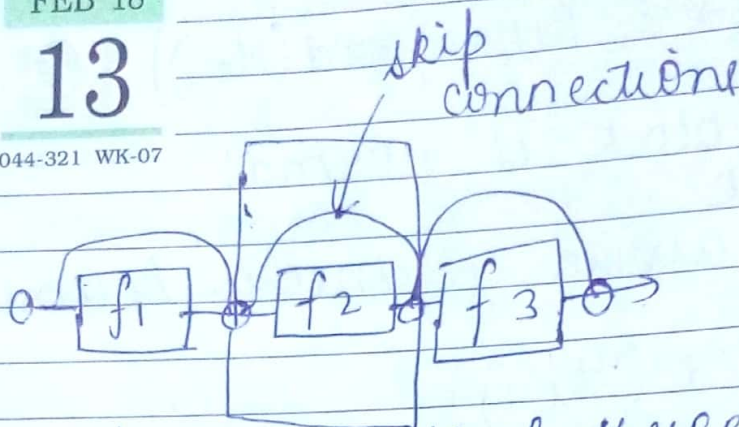


TUESDAY

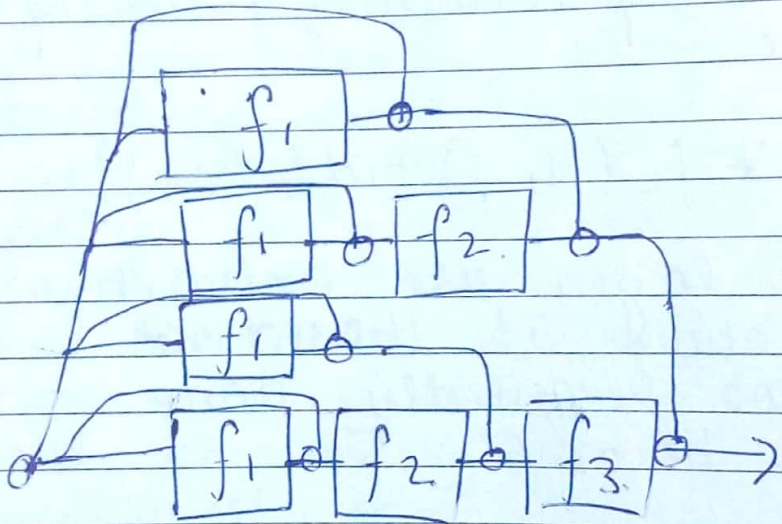
FEB '18

13

044-321 WK-07



a) conventional three block residual network



b) unravelled view of (a)

It is quite clear removing a couple of layers in ResNet architecture doesn't compromise its performance too much, the architecture has many independent effective paths

February 2018

Mon	Tue	Wed	Thu	Fri	Sat	Sun	Mon	Tue	Wed	Thu	Fri	Sat	Sun	Mon	Tue	Wed	Thu	Fri	Sat	Sun
			1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18
19	20	21	22	23	24	25	26	27	28											

TOGETHER we stand, TOGETHER we fall, TOGETHER we win, and winners take ALL.



and majority of them remain intact after we remove couple of layers.

To investigate the relationship between path length and magnitude of gradient flowing through it. To get the magnitude of gradients in path of length  $k$ , the authors first fed a batch of data to the network and randomly sample  $k$  residual blocks. When backpropagating the gradients they propagated through weight layer only for the sampled residual blocks — It shows that magnitude of gradients decreased rapidly as the path becomes longer, hence ResNet did not solve vanishing gradient problem for very long paths, and ResNet actually enables training very deep network by shortening its effective paths.

March 2018

Mon	Tue	Wed	Thu	Fri	Sat	Sun	Mon	Tue	Wed	Thu	Fri	Sat	Sun	Mon	Tue	Wed	Thu	Fri	Sat	Sun
			1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18
19	20	21	22	23	24	25	26	27	28	29	30	31								

When performing any task, stop for a moment, think of the effect it will have, and then begin.