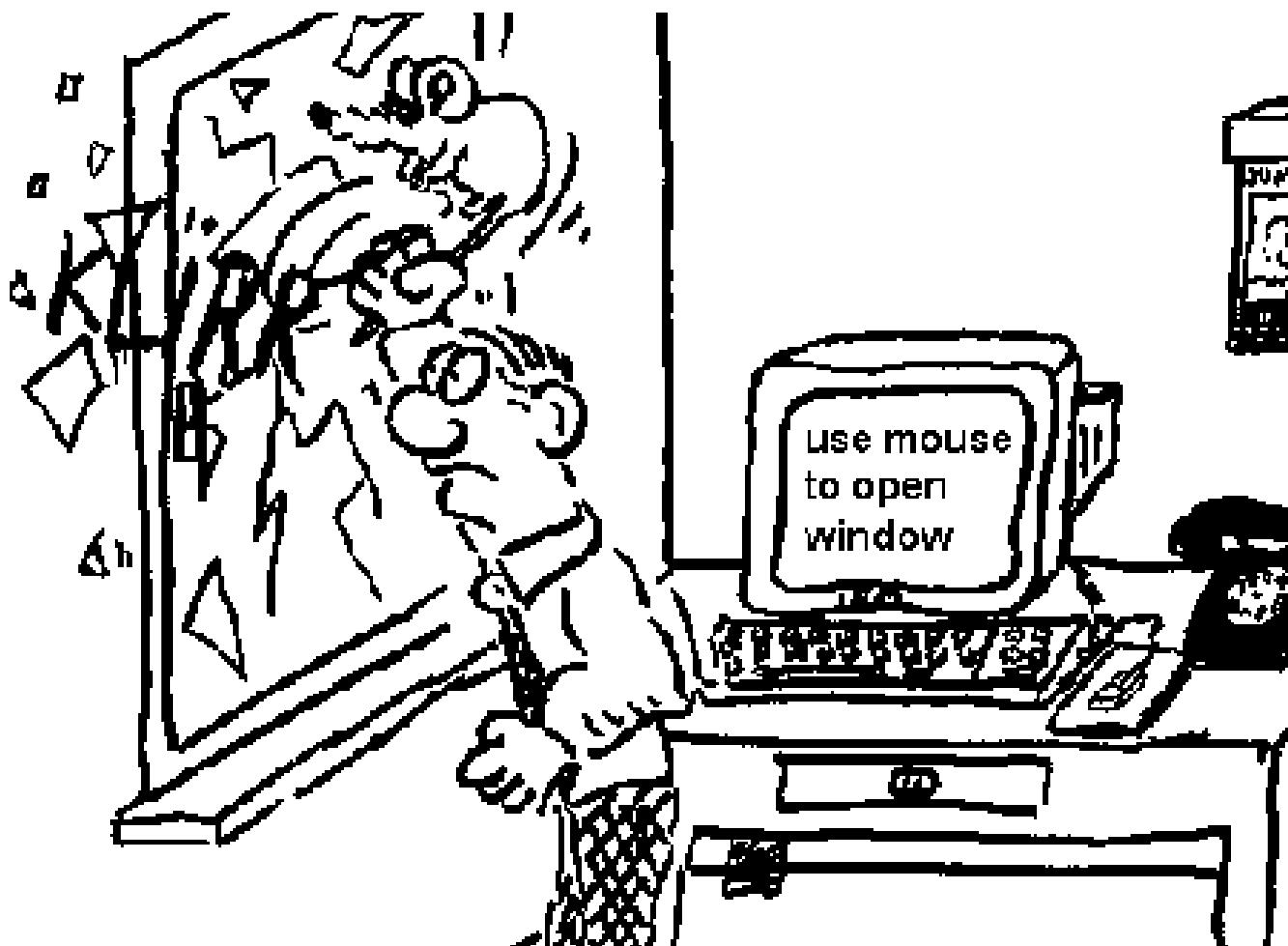


# Terminology

Sungjoon Choi  
(sungjoon.choi@cplab.snu.ac.kr)

# Terminologies are Important

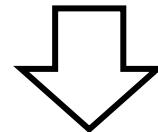
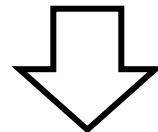
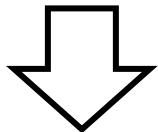
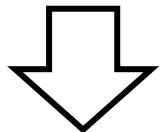


# Input Data



Dies ist ein Blindtext. An ihm lässt sich vieles über die Schrift ablesen, in der er gesetzt ist. Auf den ersten Blick wird der Grauwert der Schriftfläche sichtbar. Dann kann man darüber rüfen, wie gut die Schrift zu einer Leserart passt und wie sie auf den Leser wirkt.

# Output / Class / Label



Cat

[1 0 0 0]

Dog

[0 1 0 0]

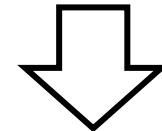
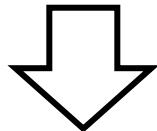
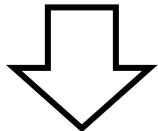
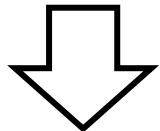
Cow

[0 0 1 0]

Horse

[0 0 0 1]

# One-hot coding



Cat

[1 0 0 0]

Dog

[0 1 0 0]

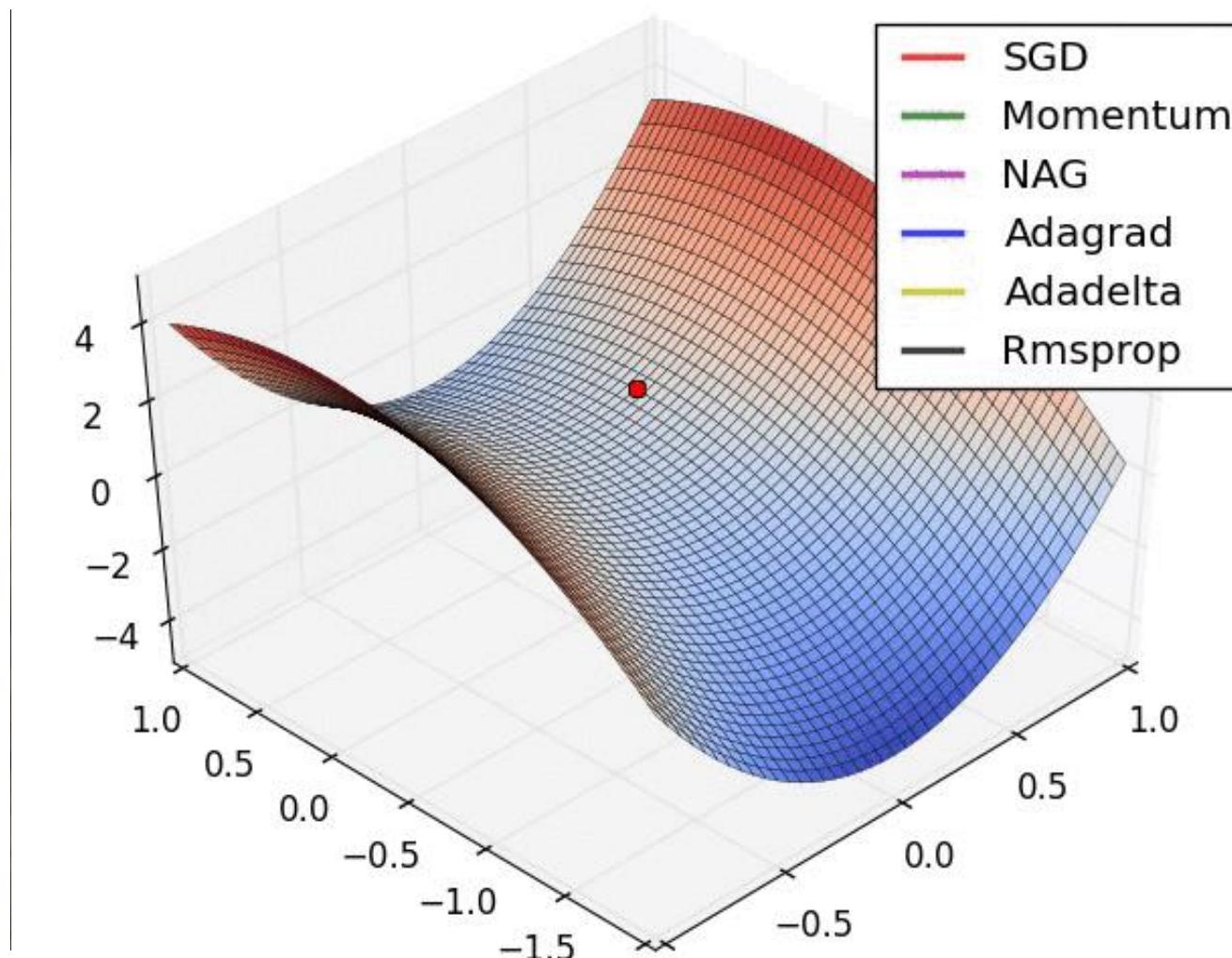
Cow

[0 0 1 0]

Horse

[0 0 0 1]

# Training / Learning



# Dataset

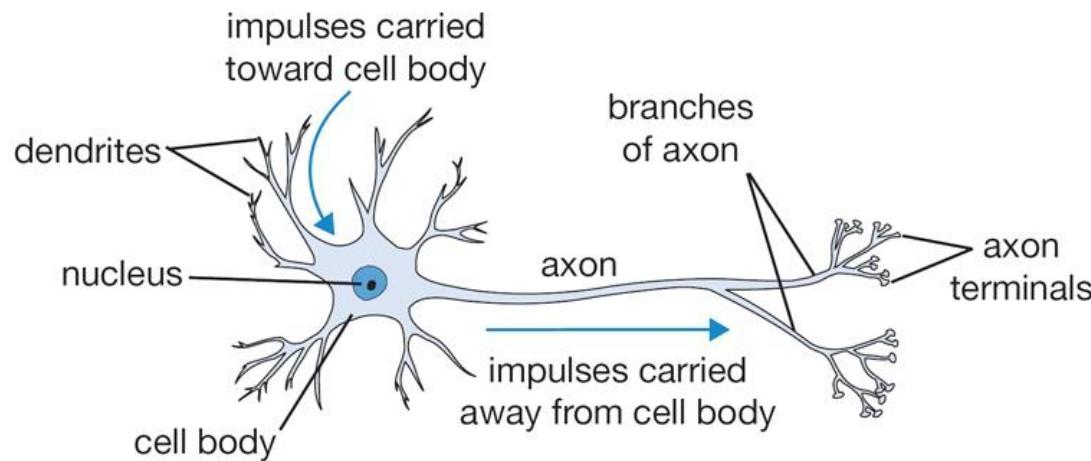


Training Data

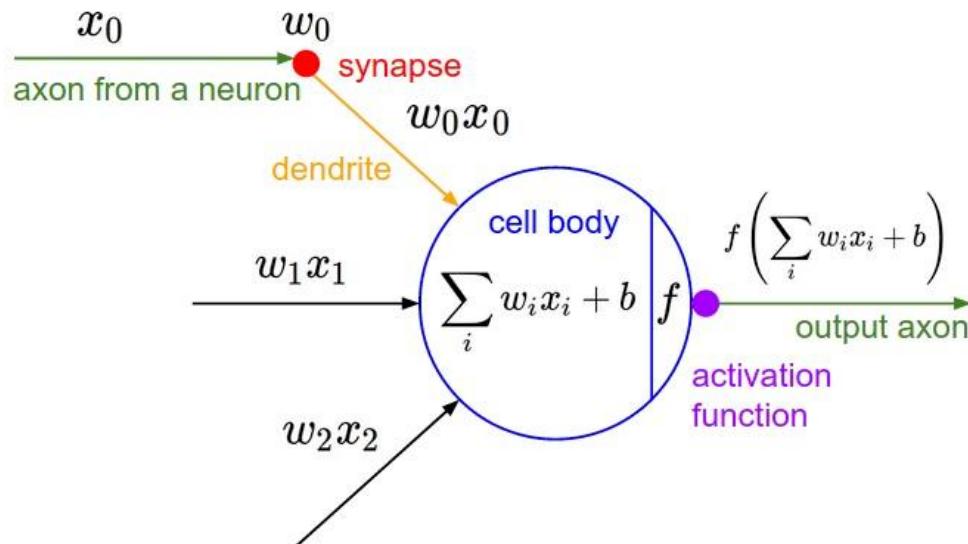
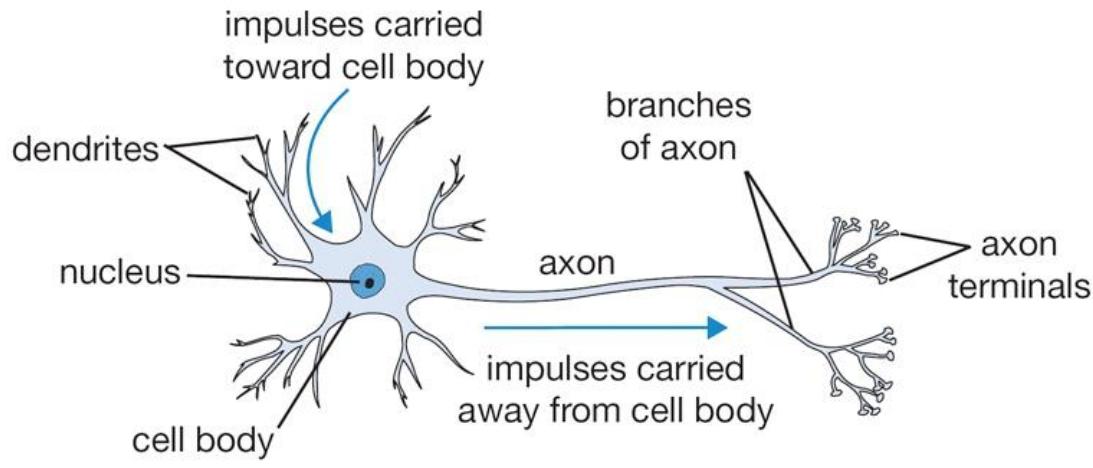
Validation Data      (Cross-validation)

Test Data

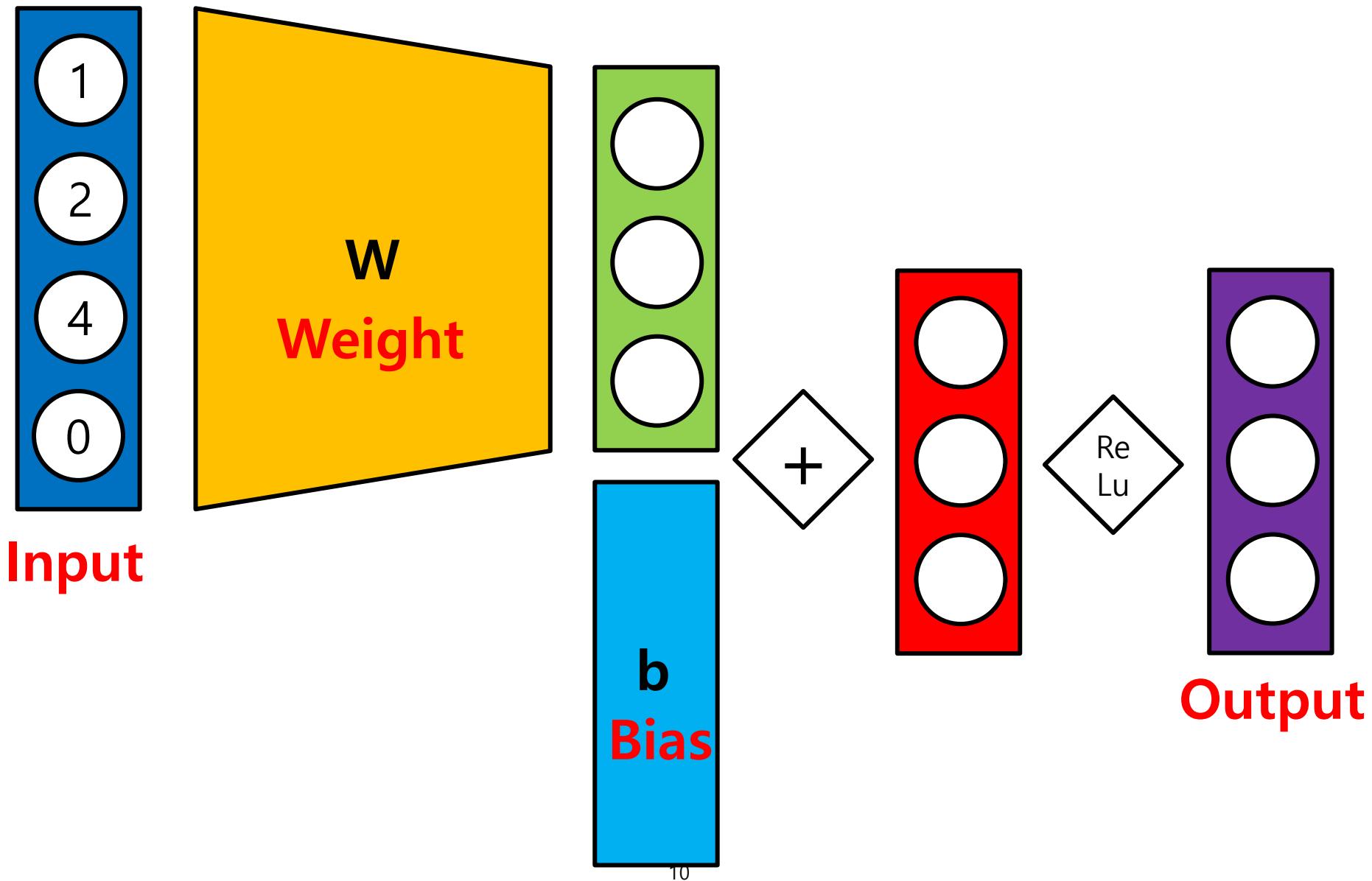
# Neuron



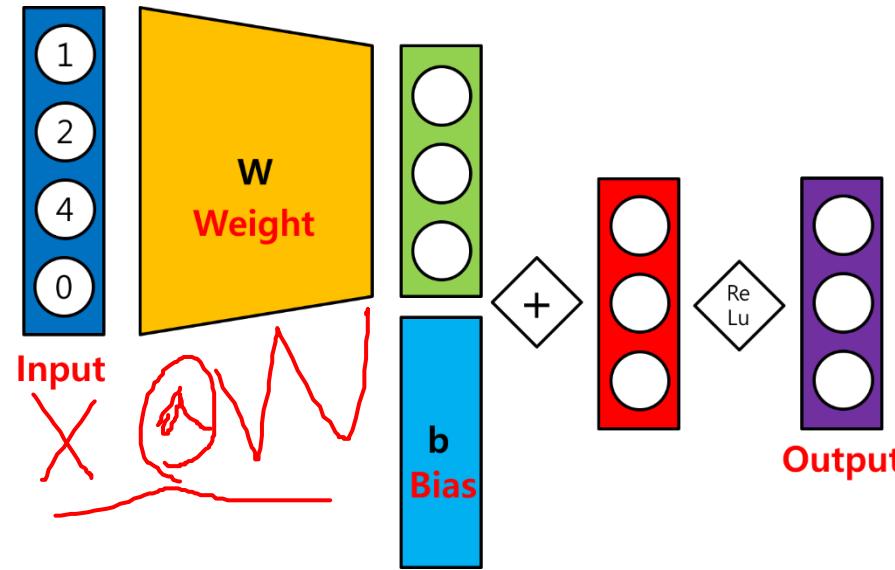
# Perceptron



# Basic single layer network



# Basic single layer network



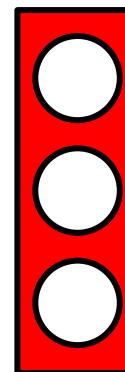
Input: [1 2 4 0] (1/4)

w:  $\begin{bmatrix} 1 & 0 & 1 \\ 3 & 0 & 2 \\ -2 & -2 & 0 \\ 1 & 1 & 0 \end{bmatrix}$  (4x3)

Bias: [2 1 3] (1x3)



[-1 -8 5]

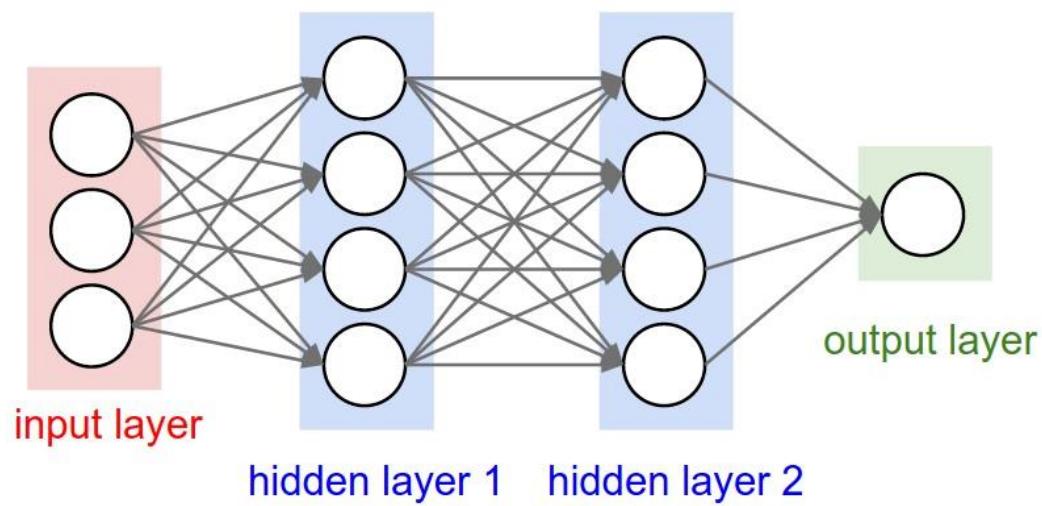
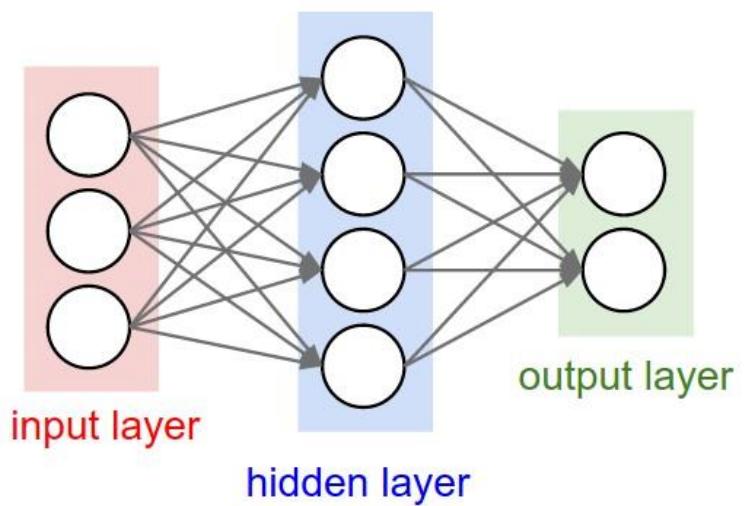


[1 -7 8]

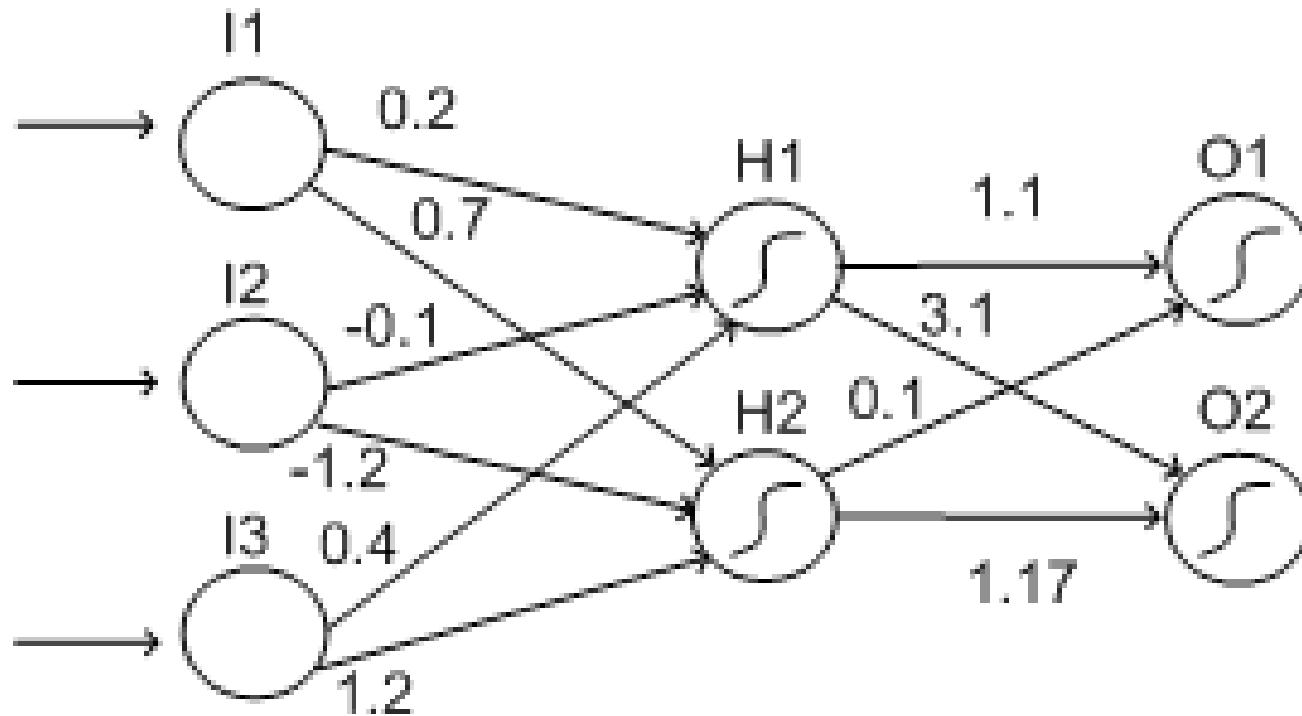


[1 0 8]

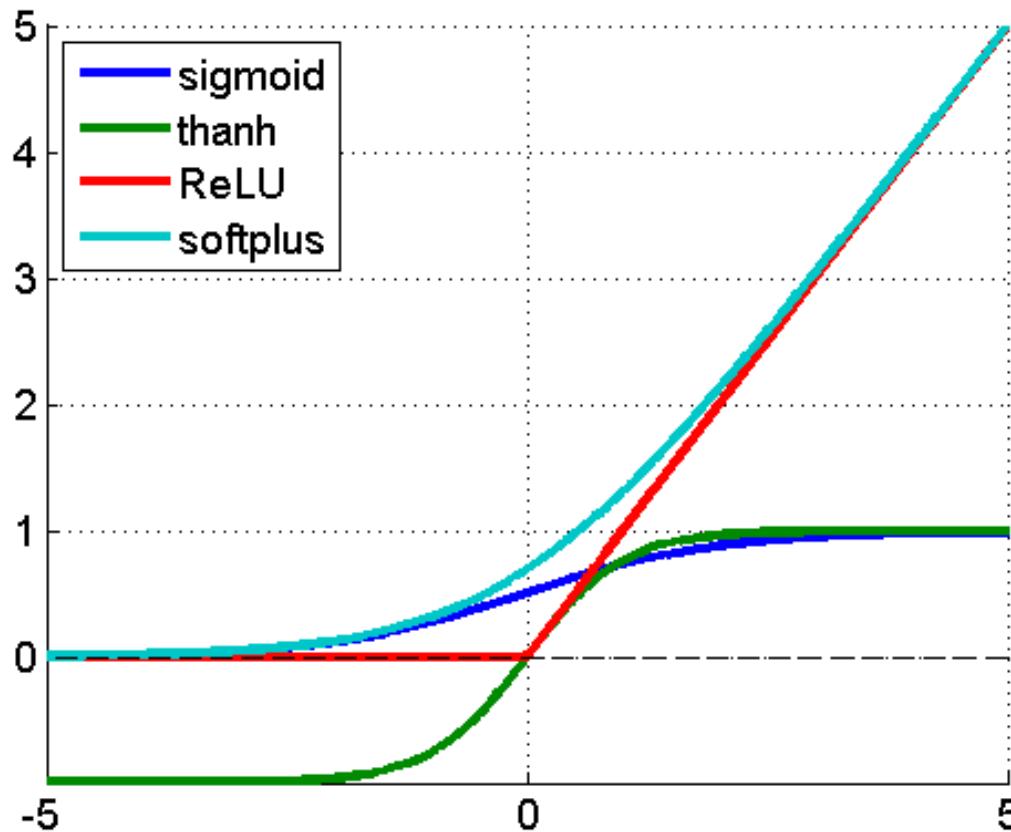
# (Neural) Network



# Multi-Layer Perceptron



# Activation Function



# Epoch / Batch size / Iteration

One **epoch**: one forward and backward pass of **all training data**

**Batch size**: the number of training examples in **one forward and backward pass**

One **iteration**: number of passes

If we have 55,000 training data, and the batch size is 1,000. Then, we need 55 iterations to complete 1 epoch.

# Cost function

This is just a function that we want to minimize.

There is **no guarantee** that it will bring us to the best solution!

It should be **differentiable**!

# MNIST

(Mixed National Institute of Standards and Technology)

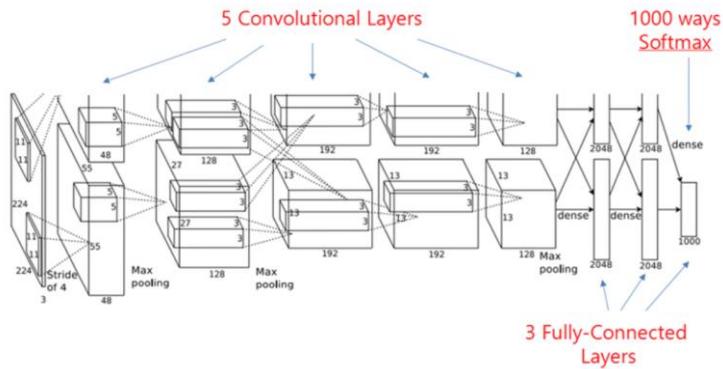


# **Modern CNNs**

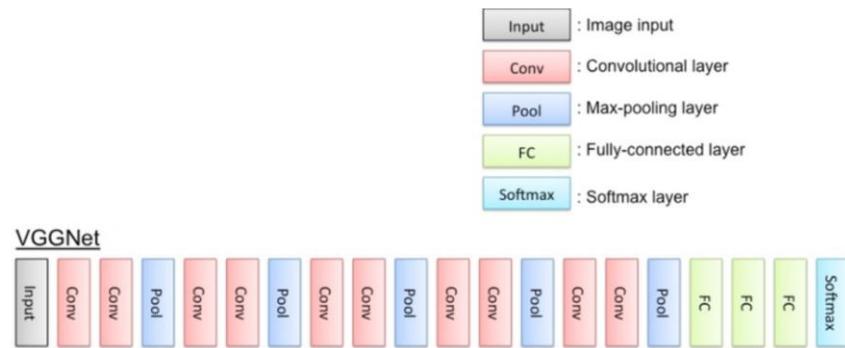
**Sungjoon Choi**  
**(sungjoon.choi@cplab.snu.ac.kr)**

# CNN Architectures

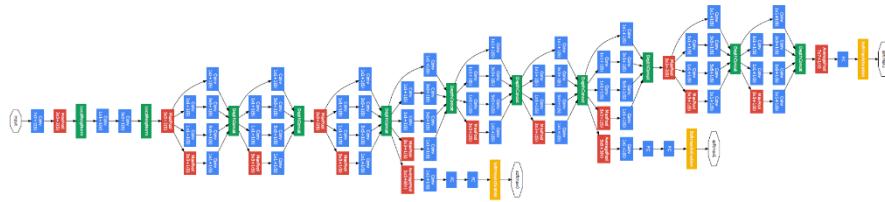
## AlexNet



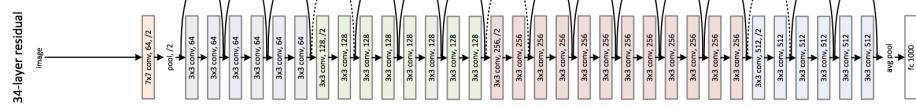
## VGG



## GoogLeNet

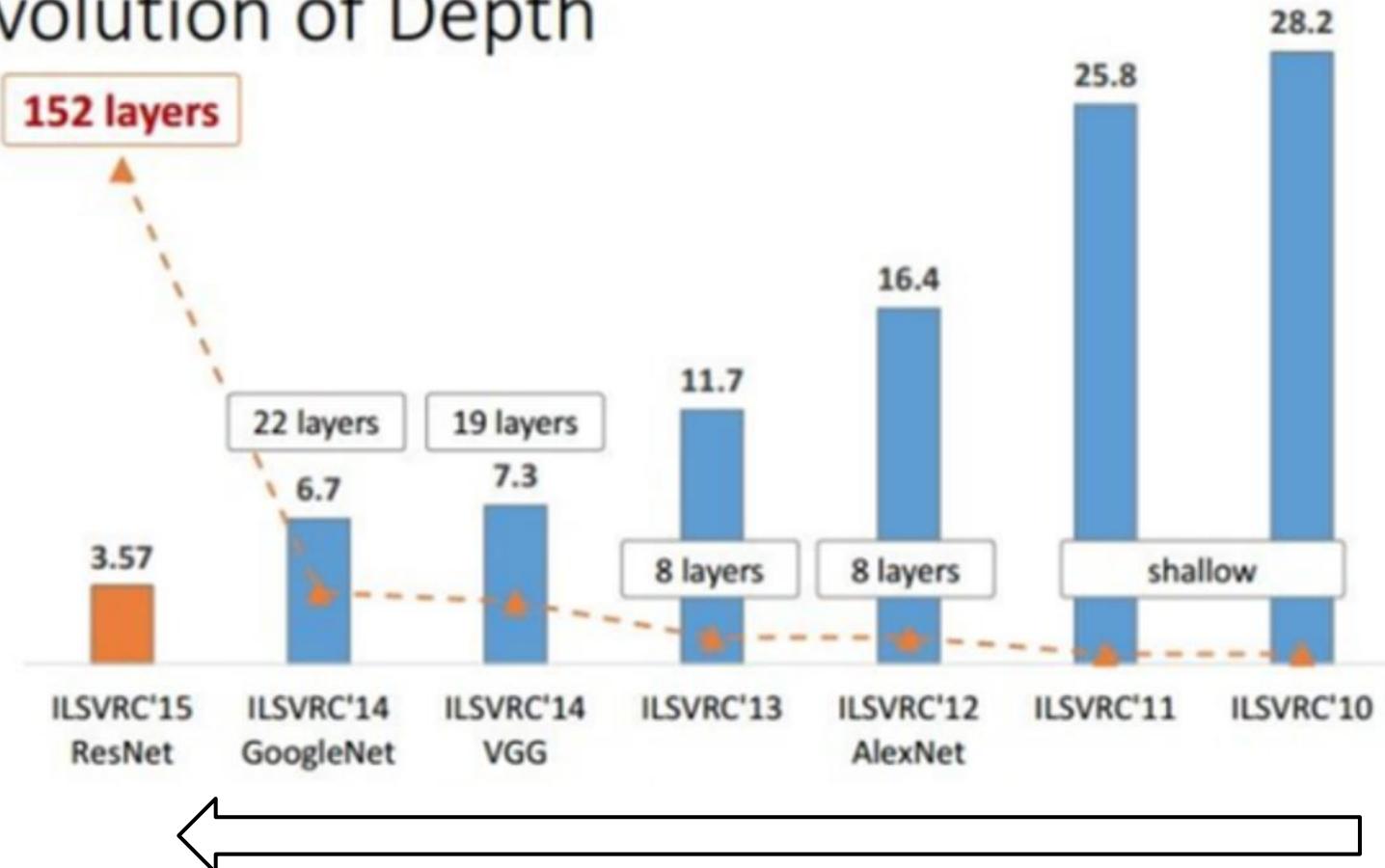


## ResNet



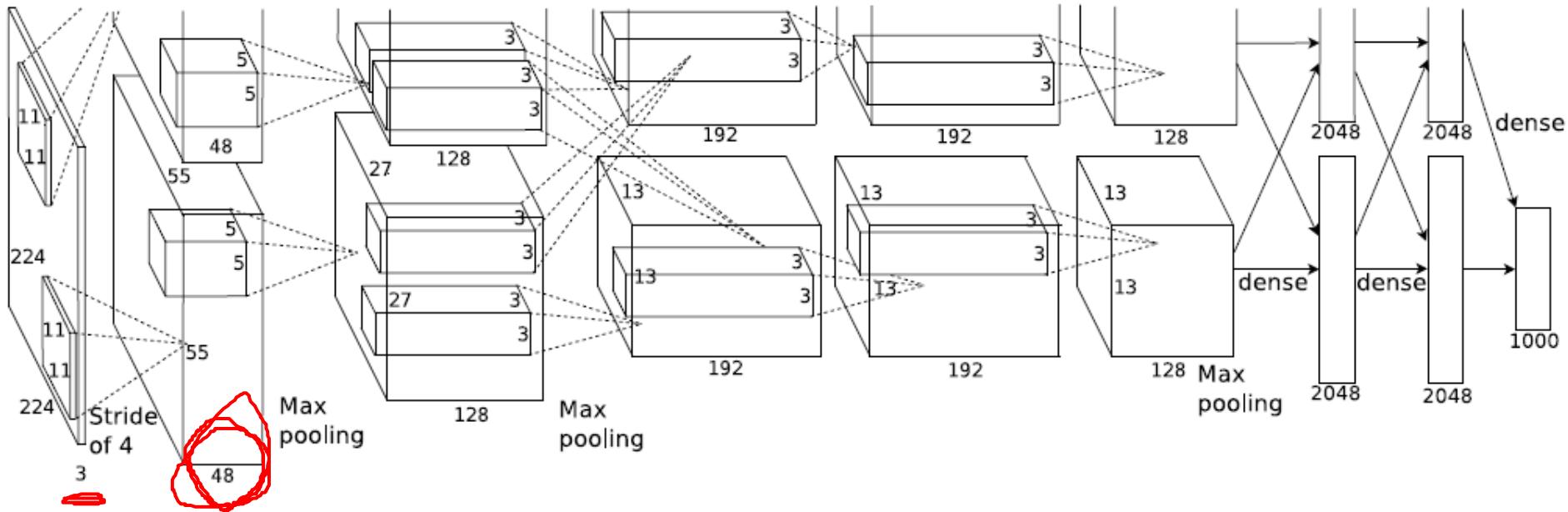
# Top-5 Classification Error

## Revolution of Depth



# AlexNet

# AlexNet

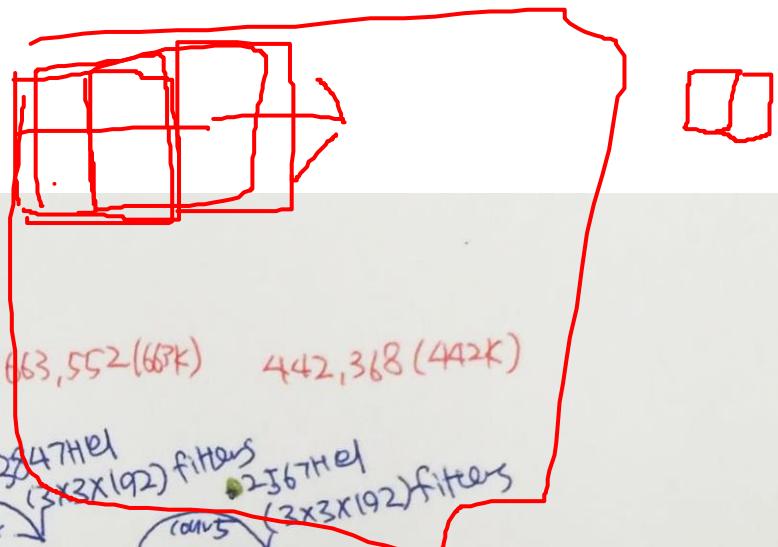


What is the number of parameters?

Why are layers divided into two parts?

# AlexNet

$(2 \times 2)$



#convparam: 34,848 (34K) 307,200 (307K) 884,736 (884K) 663,552 (663K) 442,368 (442K)

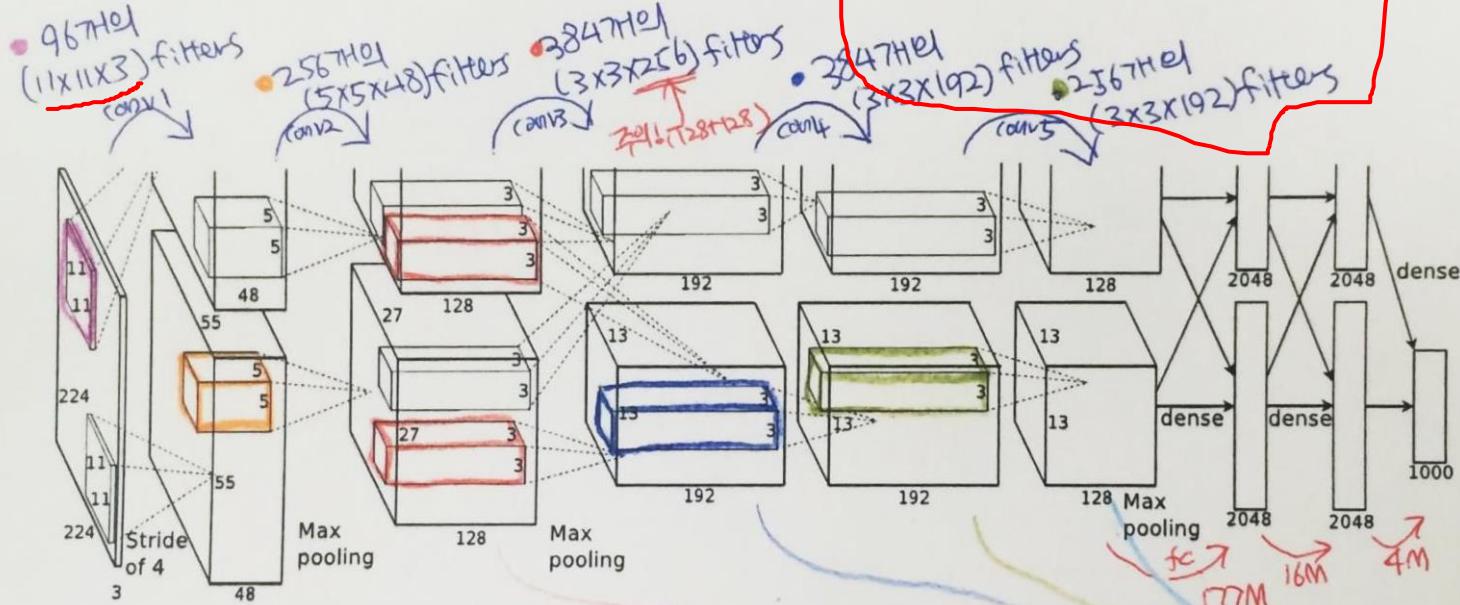
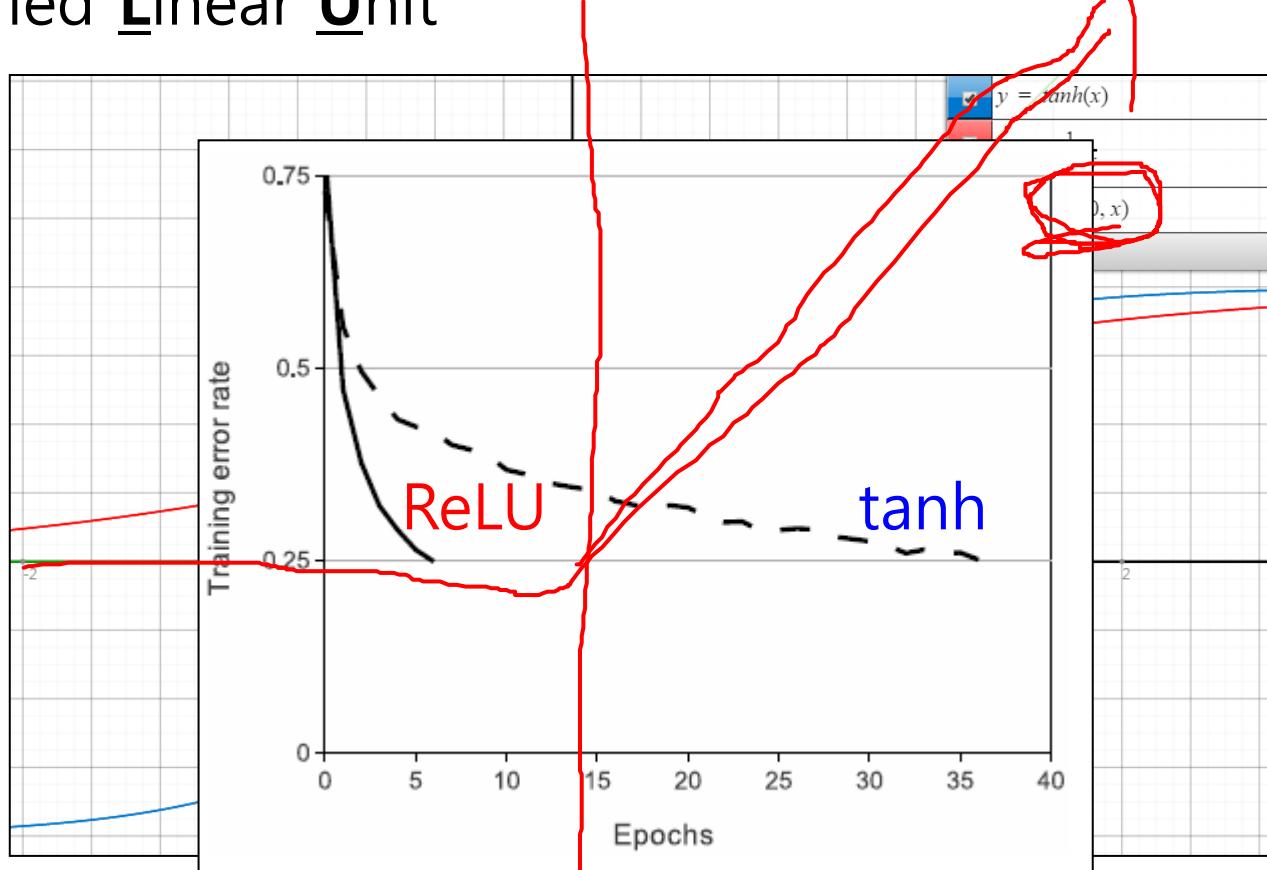


Figure 2: An illustration of the architecture of our CNN, explicitly showing the delineation of responsibilities between the two GPUs. One GPU runs the layer-parts at the top of the figure while the other runs the layer-parts at the bottom. The GPUs communicate only at certain layers. The network's input is 150,528-dimensional, and the number of neurons in the network's remaining layers is given by 253,440–186,624–64,896–64,896–43,264–4096–4096–1000.

# ReLU

## Rectified Linear Unit



Faster Convergence!

# LRN

## Local Response Normalization

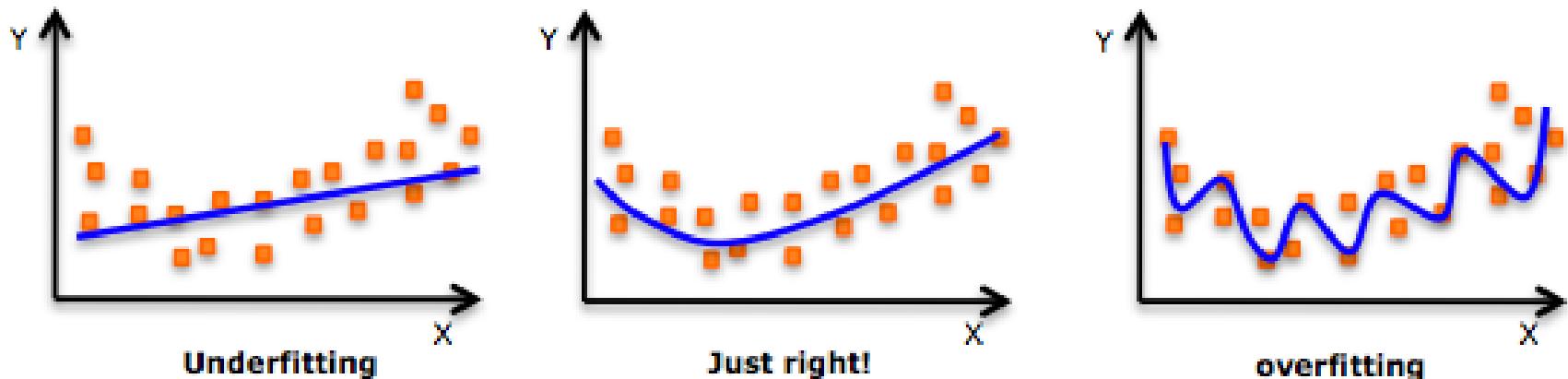


The response-normalized activity is given by:

$$b_{x,y}^i = \frac{a_{x,y}^i}{\left( k + \alpha \sum_{j=\max(0,i-n/2)}^{\min(N-1,i+n/2)} (a_{x,y}^j)^2 \right)^\beta}$$

It implements a form of **lateral inhibition** inspired by real neurons.

# Regularization in AlexNet



Main objective is to reduce **overfitting**.

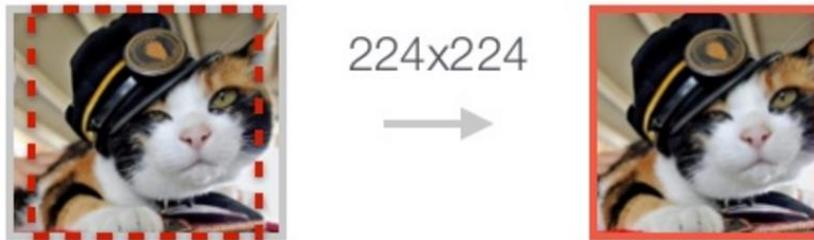
More details will be handled in next week.

In the **AlexNet**, two regularization methods are used.

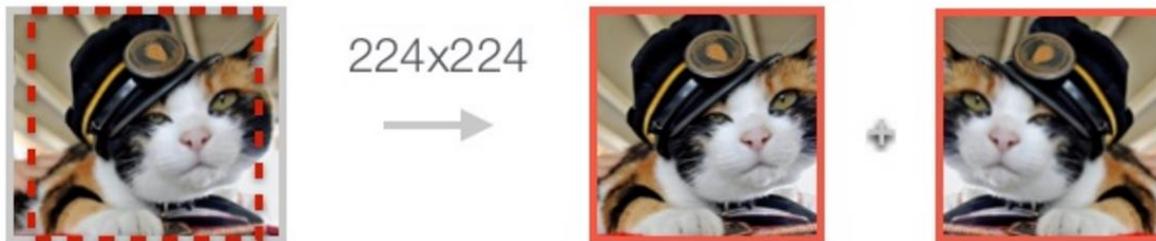
- Data augmentation
- Dropout

# Data augmentation

a. No augmentation (= 1 image)



b. Flip augmentation (= 2 images)



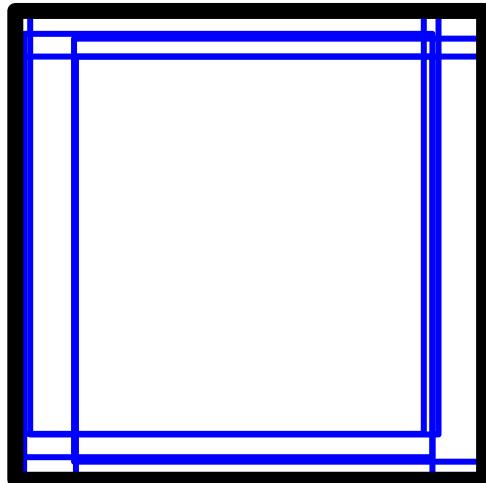
c. Crop+Flip augmentation (= 10 images)



# Data augmentation in AlexNet

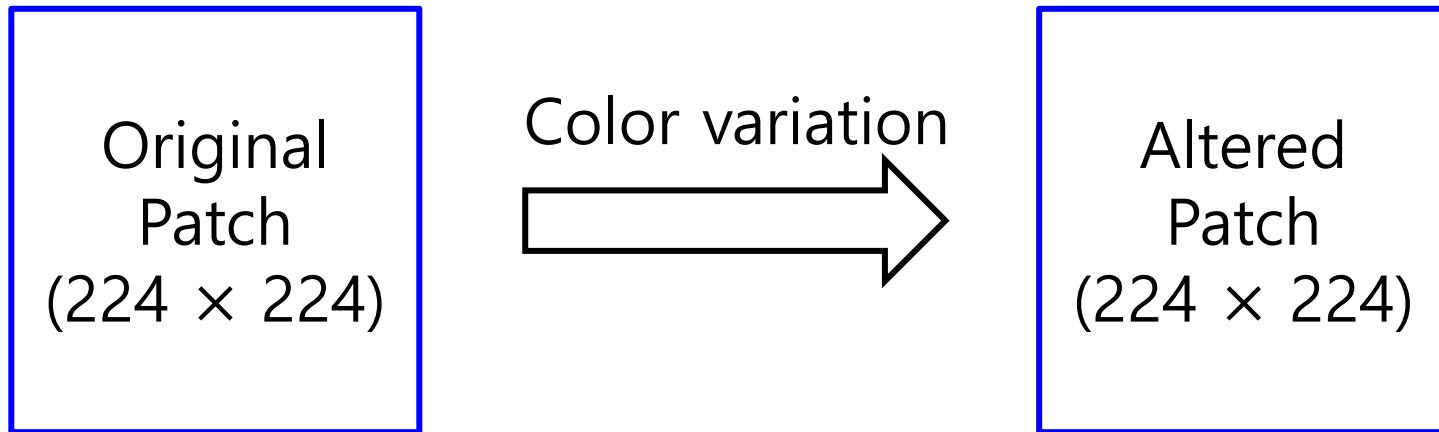
Original  
Image  
 $(256 \times 256)$

Smaller  
Patch  
 $(224 \times 224)$



This increases the size of the training set by a **factor of 2048** ( $32 * 32 * 2$ ).  
Two comes from horizontal reflections.

# Data augmentation in AlexNet



To each RGB image pixel, following quantity is added:

$$[p_1, p_2, p_3][\alpha_1 \lambda_1, \alpha_2 \lambda_2, \alpha_3 \lambda_3]^T$$

where  $p_i$  and  $\lambda_i$  are  $i$ th eigenvector and eigenvalue of  $3 \times 3$  covariance matrix of RGB pixel values.

Probabilistically, not a single patch will be same at the training phase! (a **factor of infinity!**)

# Dropout



Original dropout [1] sets the output of each hidden neuron with certain probability.

In this paper, they simply multiply the outputs by 0.5.

[1] G.E. Hinton, N. Srivastava, A. Krizhevsky, I. Sutskever, and R.R. Salakhutdinov. Improving neural networks by preventing co-adaptation of feature detectors. arXiv, 2012.

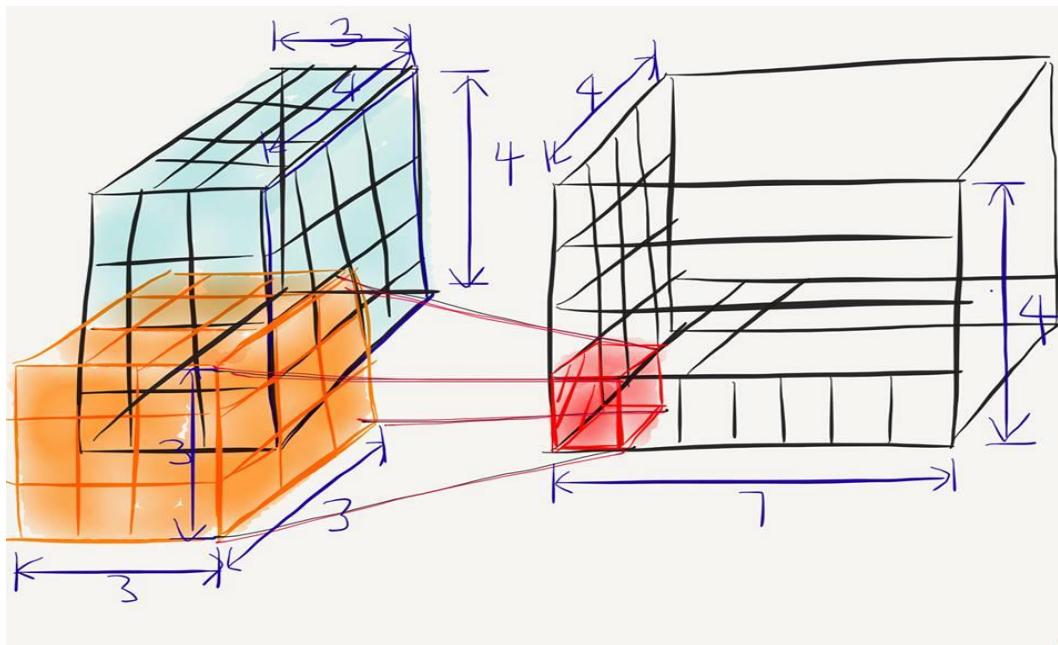
# Must remember!

# Conv2D

```
tf.nn.conv2d(input, filter, strides, padding,  
use_cudnn_on_gpu=None, name=None)
```

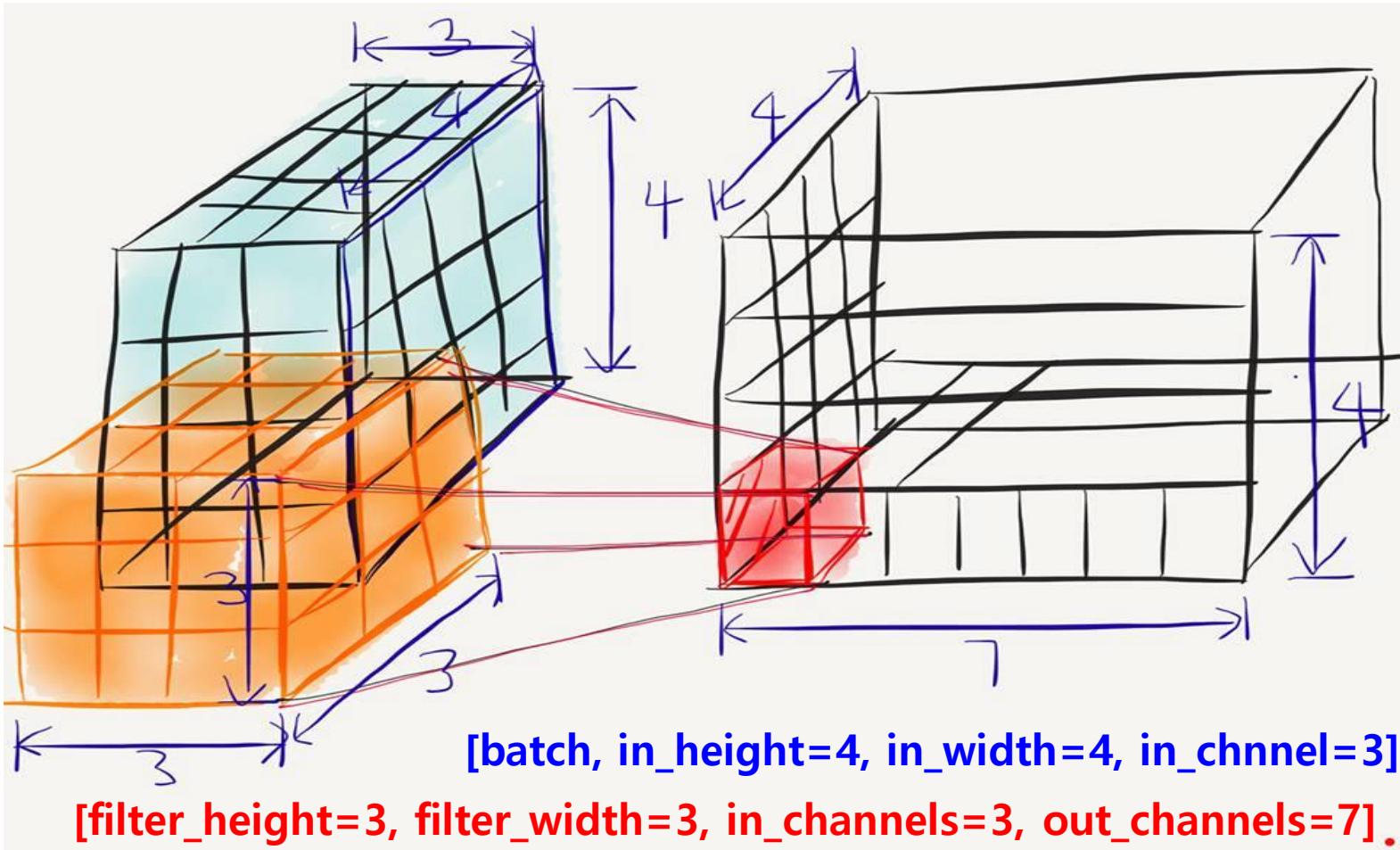
Computes a 2-D convolution given 4-D input and filter tensors.

Given an input tensor of shape **[batch, in\_height, in\_width, in\_channel]** and a filter / kernel tensor of shape **[filter\_height, filter\_width, in\_channels, out\_channels]** this op performs the following:



**[batch, in\_height=4,  
in\_width=4, in\_channel=3]**  
**[filter\_height=3, filter\_width=3,  
in\_channels=3, out\_channels=7]**

# Conv2D



What is the **number of parameters** in this convolution layer?

$$\rightarrow 189 = 3 * 3 * 3 * 7$$

**VGG**

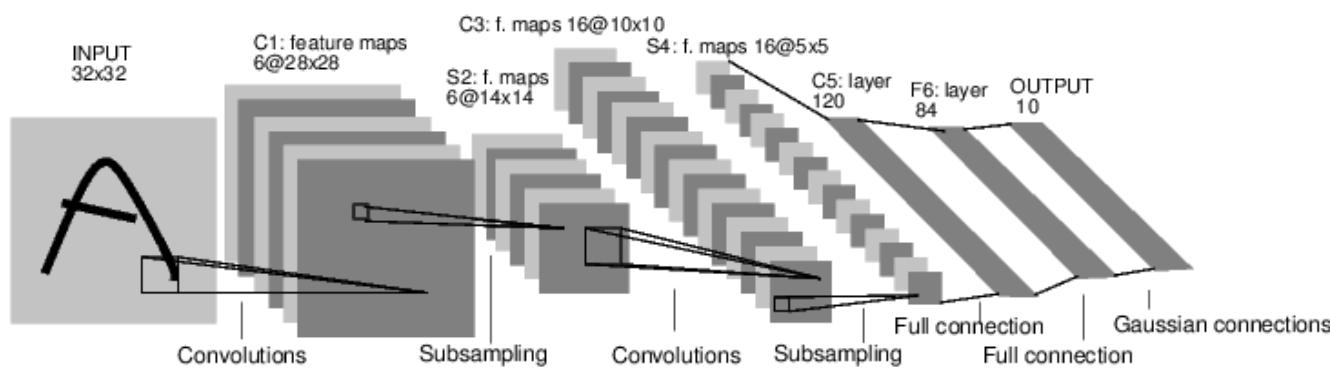
# VGG?

ConvNet Configuration					
A	A-LRN	B	C	D	E
11 weight layers	11 weight layers	13 weight layers	16 weight layers	16 weight layers	19 weight layers
input (224 × 224 RGB image)					
conv3-64	conv3-64 <b>LRN</b>	conv3-64 <b>conv3-64</b>	conv3-64 conv3-64	conv3-64 conv3-64	conv3-64 conv3-64
maxpool					
conv3-128	conv3-128	conv3-128 <b>conv3-128</b>	conv3-128 conv3-128	conv3-128 conv3-128	conv3-128 conv3-128
maxpool					
conv3-256 conv3-256	conv3-256 conv3-256	conv3-256 conv3-256	conv3-256 conv3-256 <b>conv1-256</b>	conv3-256 conv3-256 <b>conv3-256</b>	conv3-256 conv3-256 conv3-256 <b>conv3-256</b>
maxpool					
conv3-512 conv3-512	conv3-512 conv3-512	conv3-512 conv3-512	conv3-512 conv3-512 <b>conv1-512</b>	conv3-512 conv3-512 <b>conv3-512</b>	conv3-512 conv3-512 conv3-512 <b>conv3-512</b>
maxpool					
conv3-512 conv3-512	conv3-512 conv3-512	conv3-512 conv3-512	conv3-512 conv3-512 <b>conv1-512</b>	conv3-512 conv3-512 <b>conv3-512</b>	conv3-512 conv3-512 conv3-512 <b>conv3-512</b>
maxpool					
FC-4096					
FC-4096					
FC-1000					
soft-max					

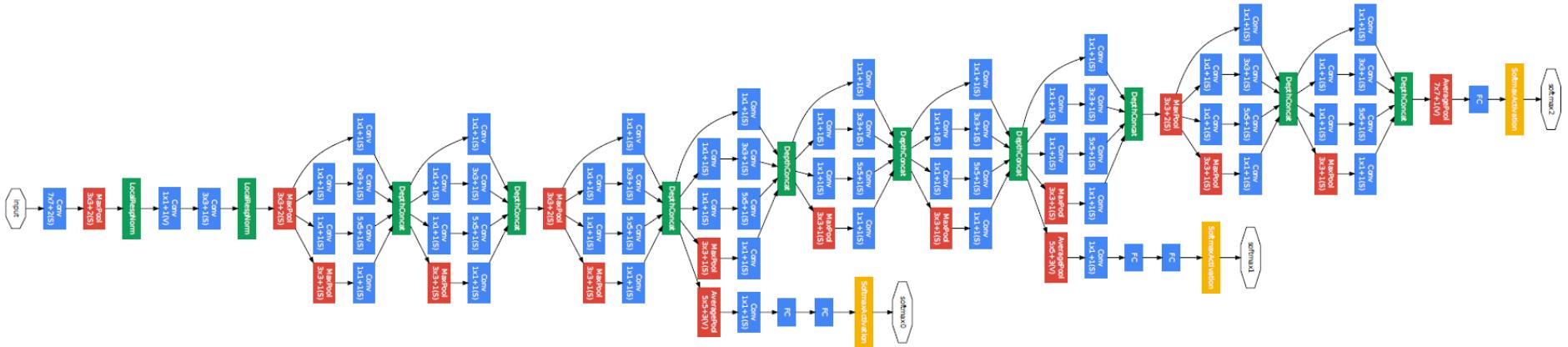
# GoogLeNet

# GoogLeNet

Google



# GoogLeNet

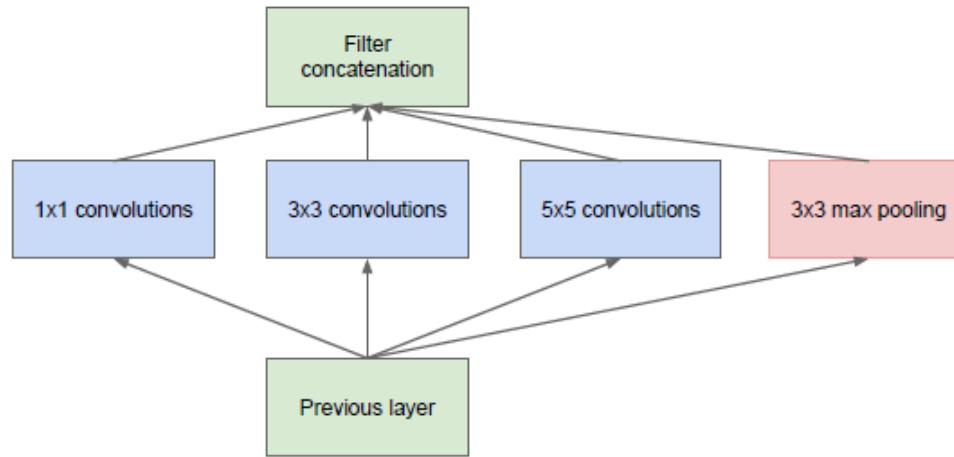


22 Layers Deep Network

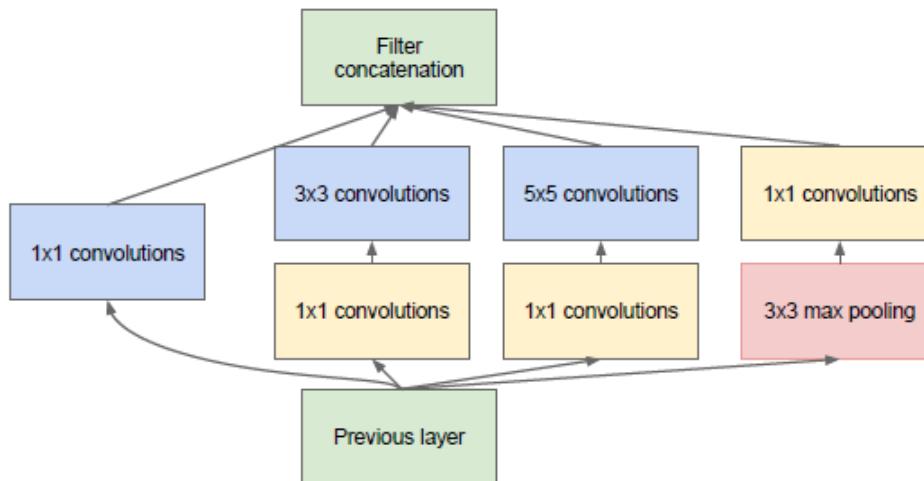
Efficiently utilized computing resources, "Inception Module"

Significantly outperforms previous methods on ILSVRC 2014

# Inception module

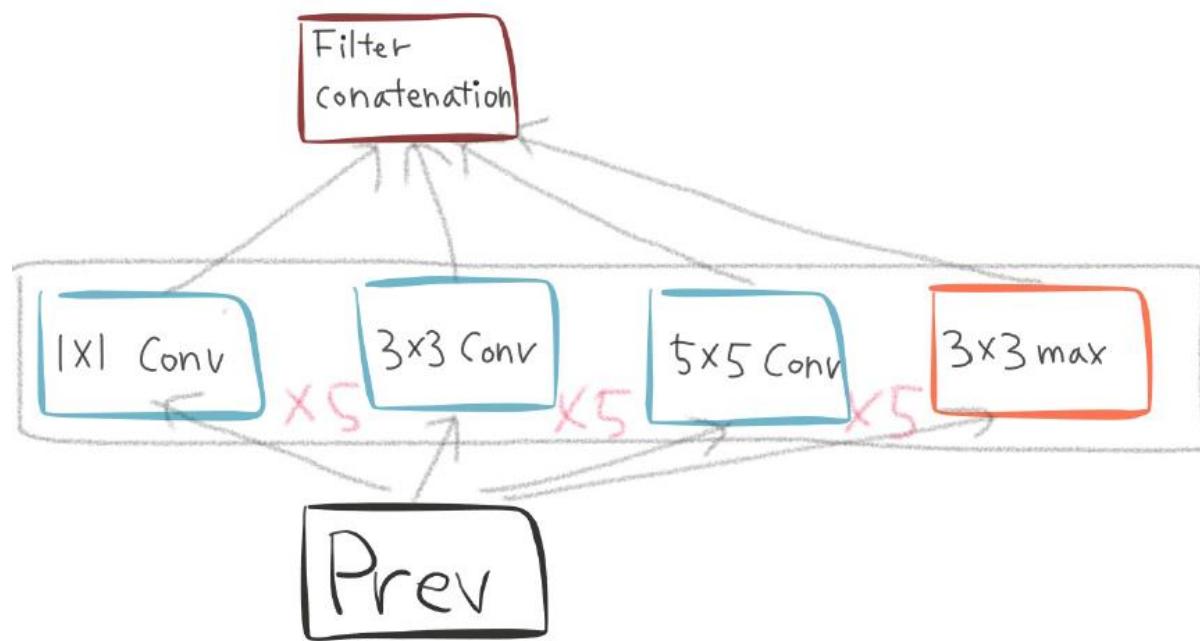


(a) Inception module, naïve version

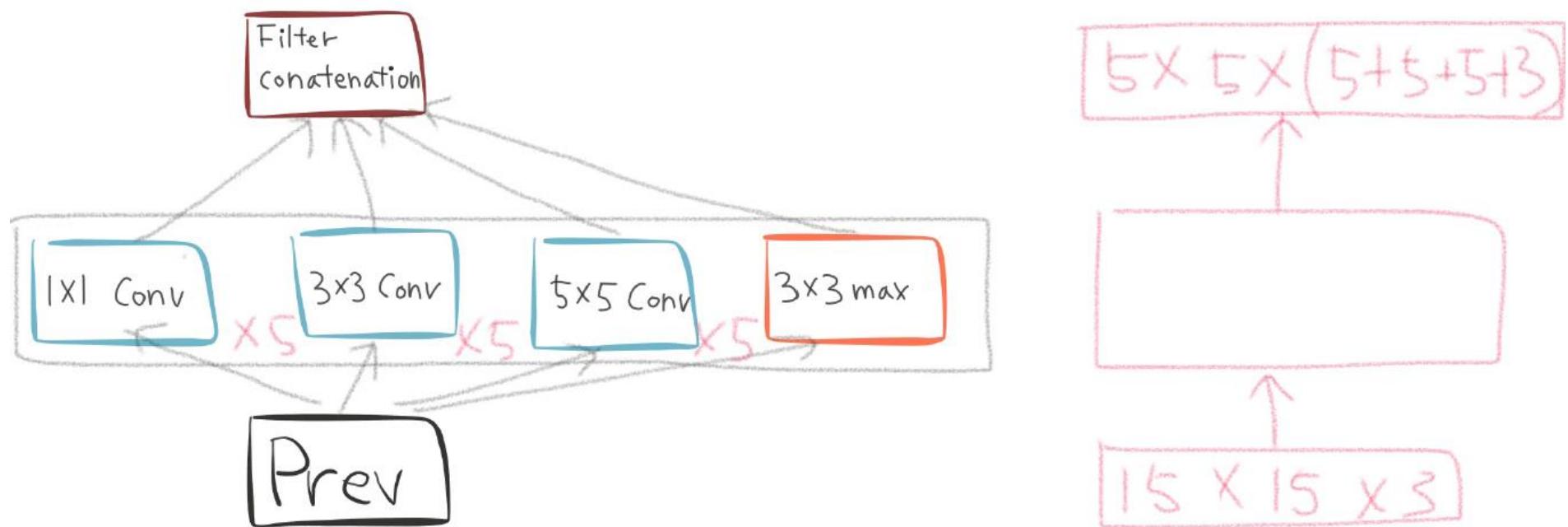


(b) Inception module with dimensionality reduction

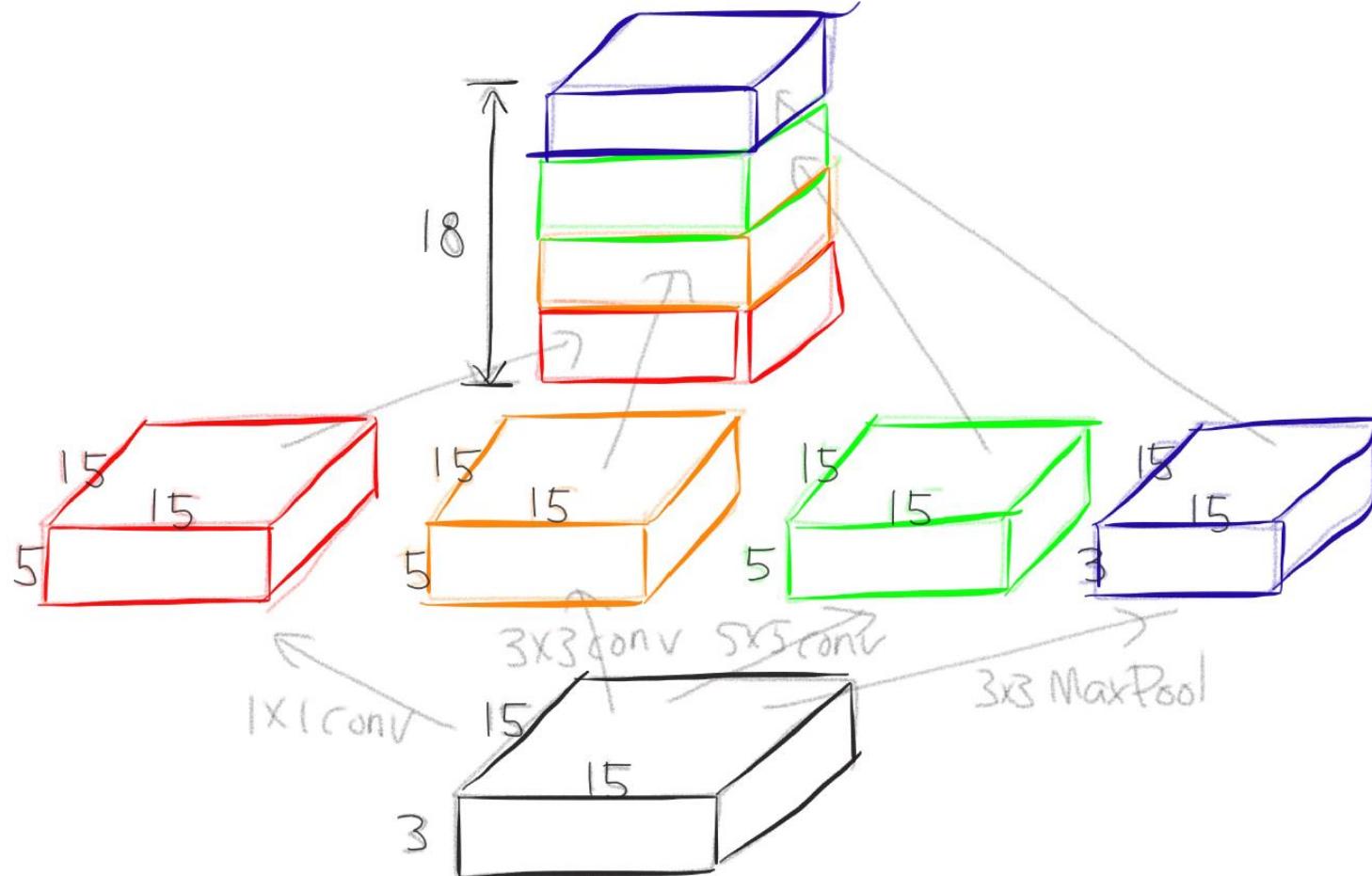
# Naïve inception module



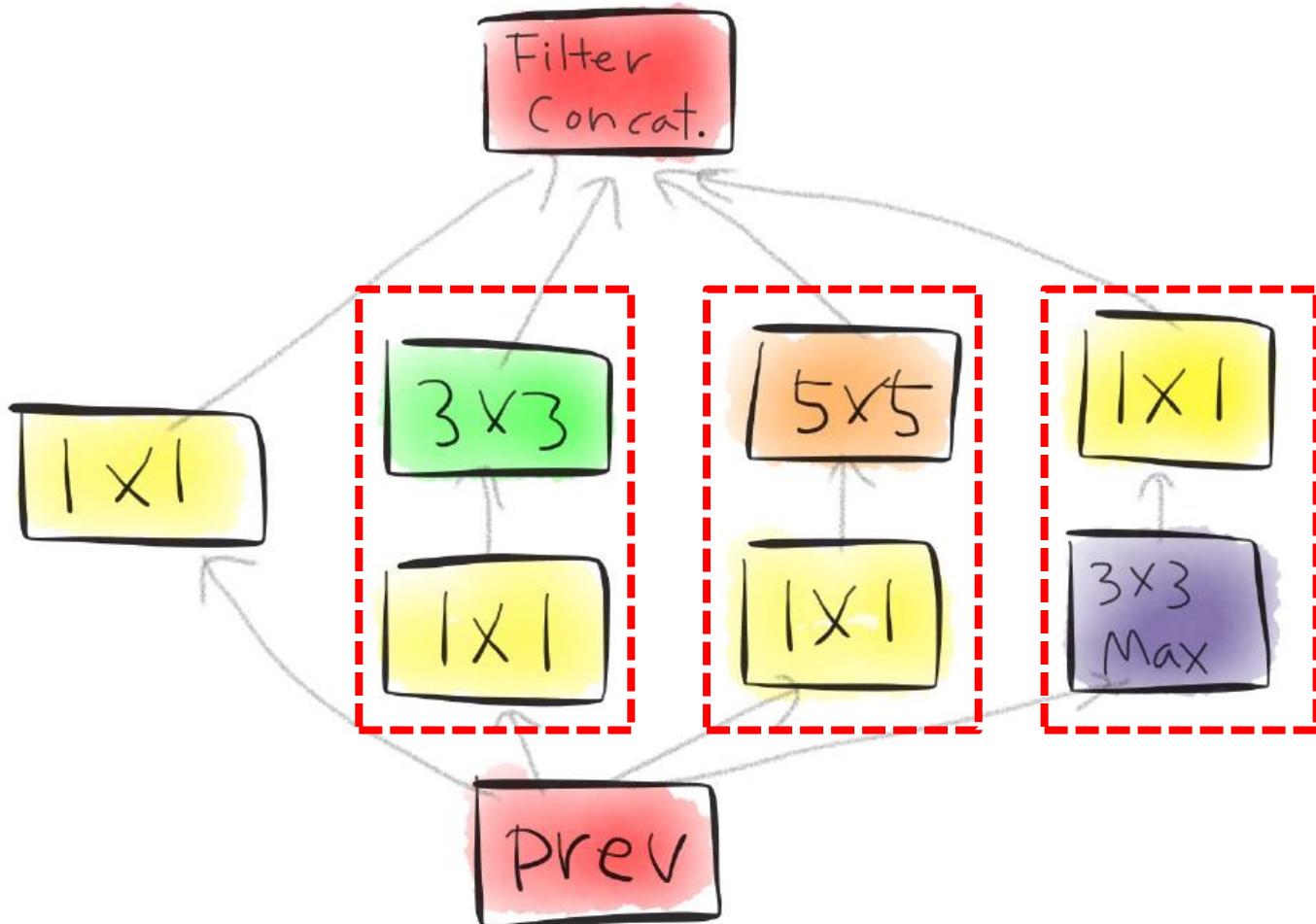
# Naïve inception module



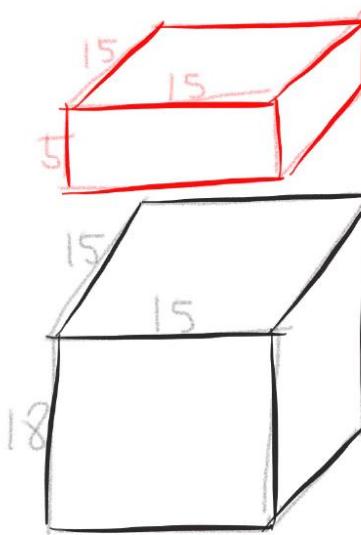
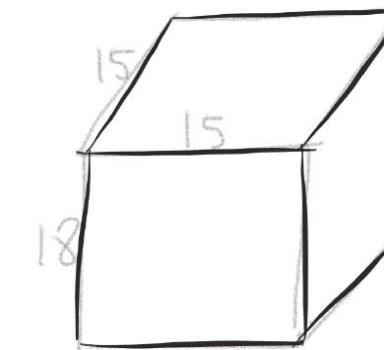
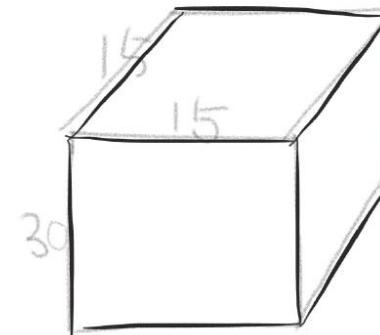
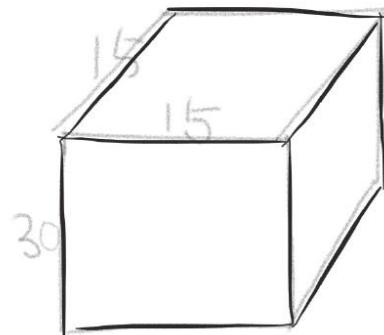
# Naïve inception module



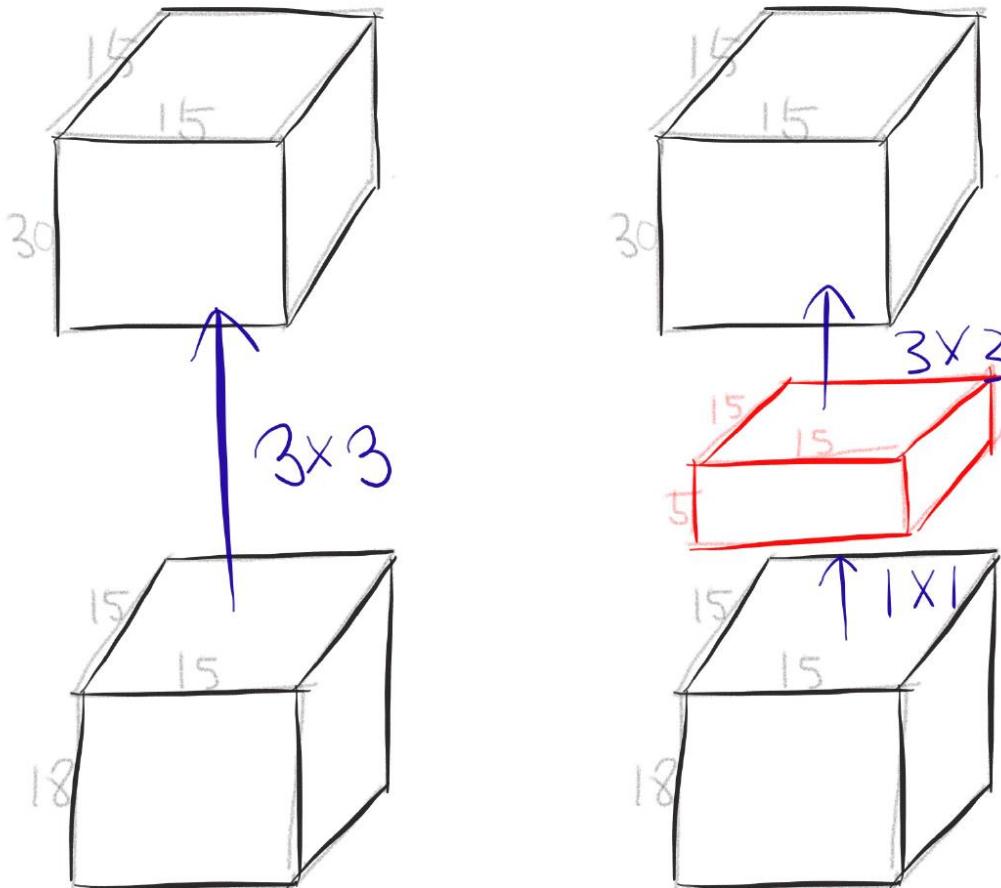
# Actual inception module



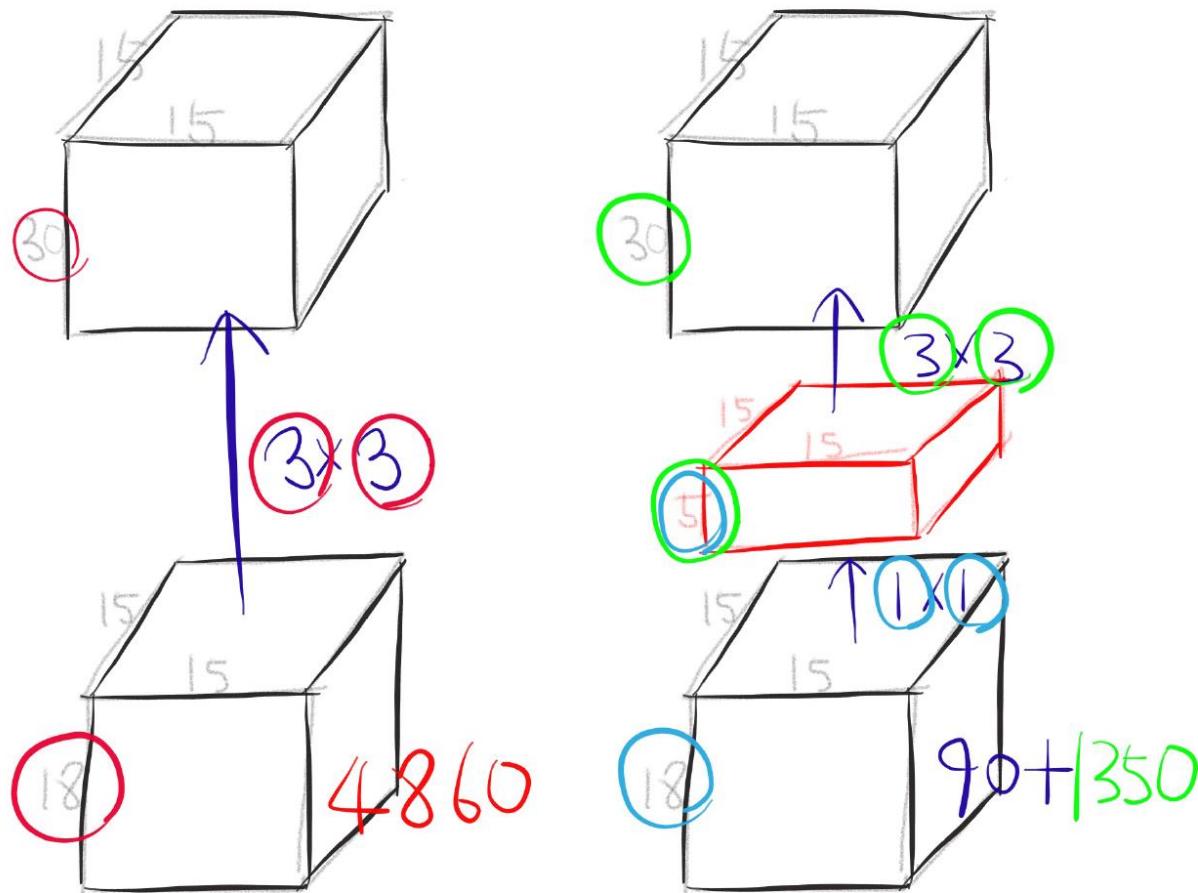
# One by one convolution



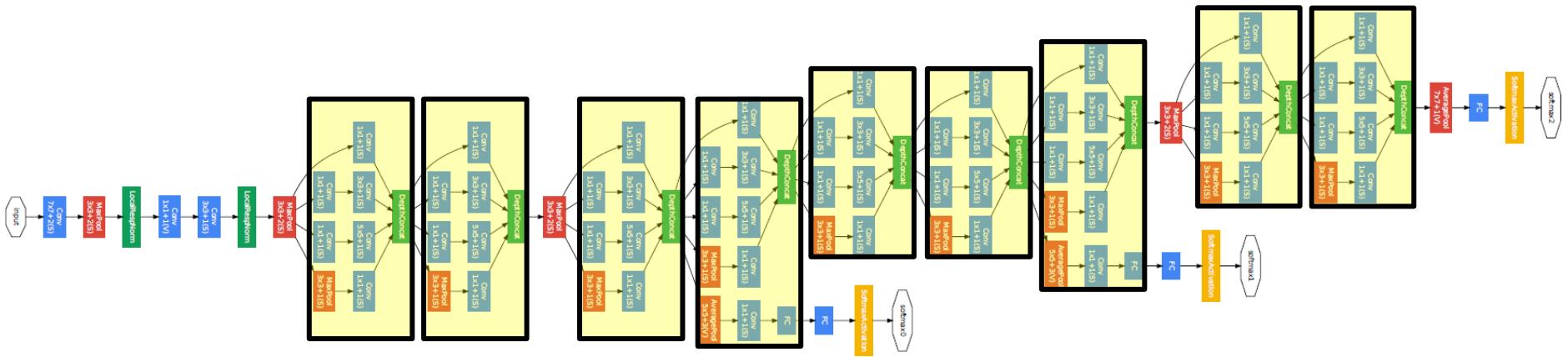
# One by one convolution



# One by one convolution



# GoogLeNet



Network in Network!

# Conclusion

Very clever idea of using **one by one convolution** for dimension reduction!

Other than that...



# Inception-v4, Inception- ResNet and the Impact of Residual Connections on Learning

# Inception v4

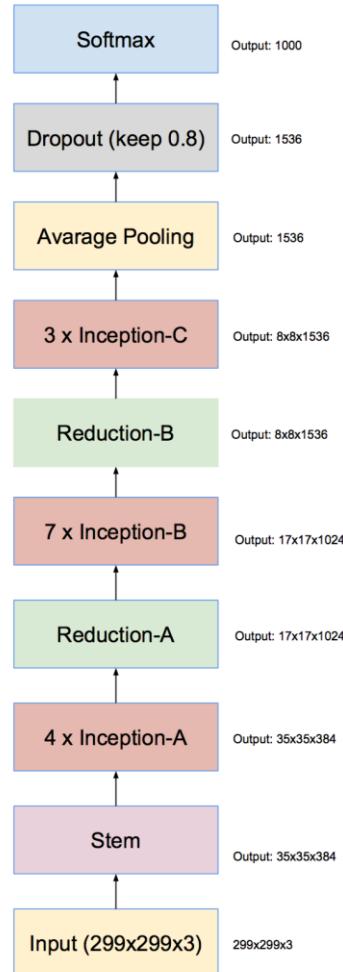


Figure 9. The overall schema of the Inception-v4 network. For the detailed modules, please refer to Figures 3, 4, 5, 6, 7 and 8 for the detailed structure of the various components.

# Inception v4

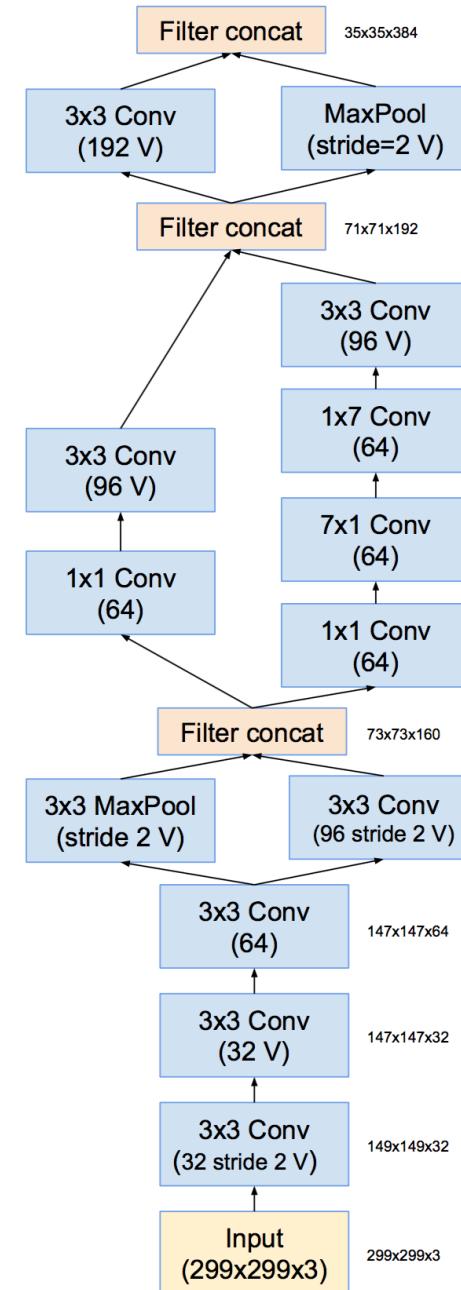
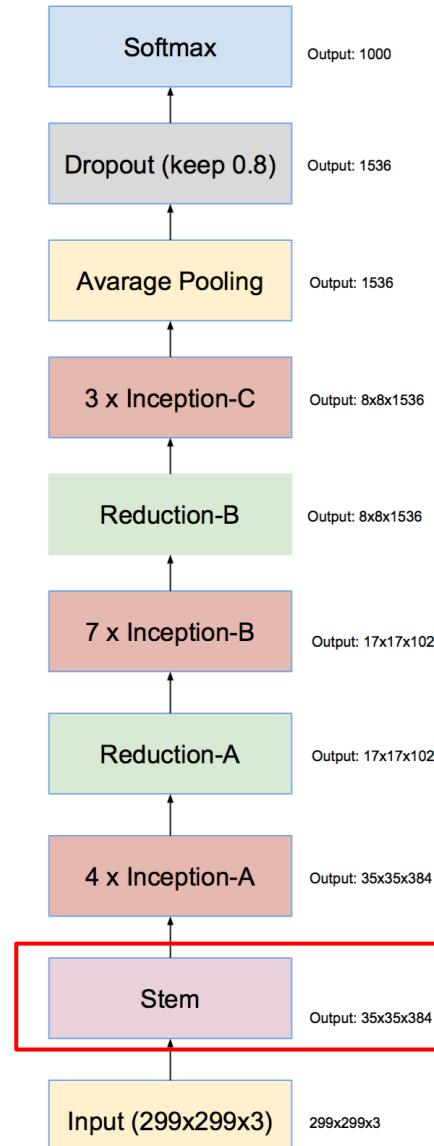


Figure 3. The schema for stem of the pure Inception-v4 and Inception-ResNet-v2 networks. This is the input part of those networks. Cf. Figures 9 and 15.

# Inception v4

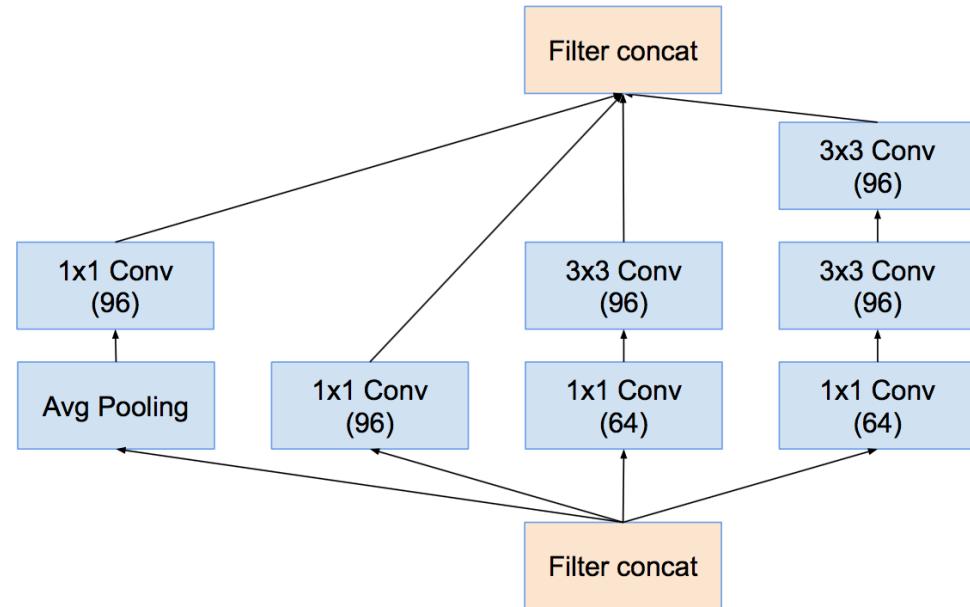
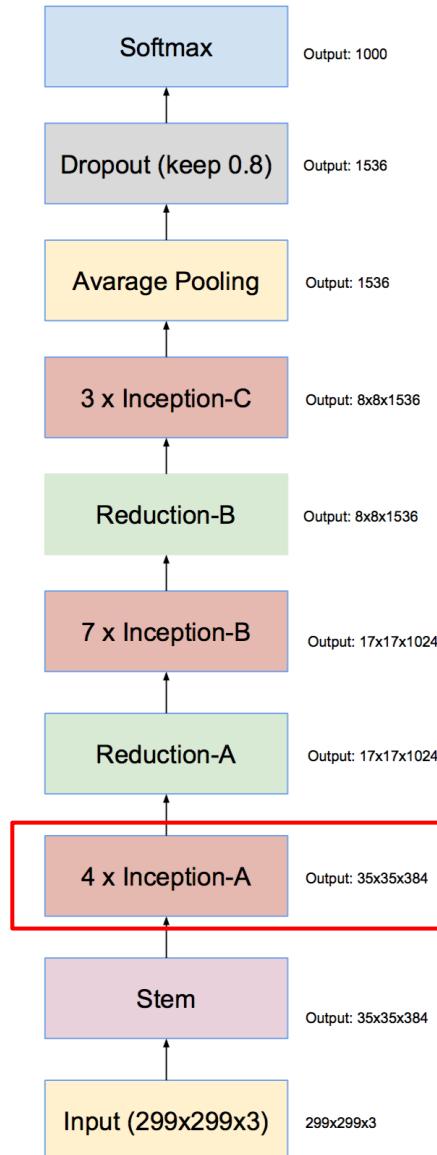


Figure 4. The schema for  $35 \times 35$  grid modules of the pure Inception-v4 network. This is the Inception-A block of Figure 9.

# Inception v4

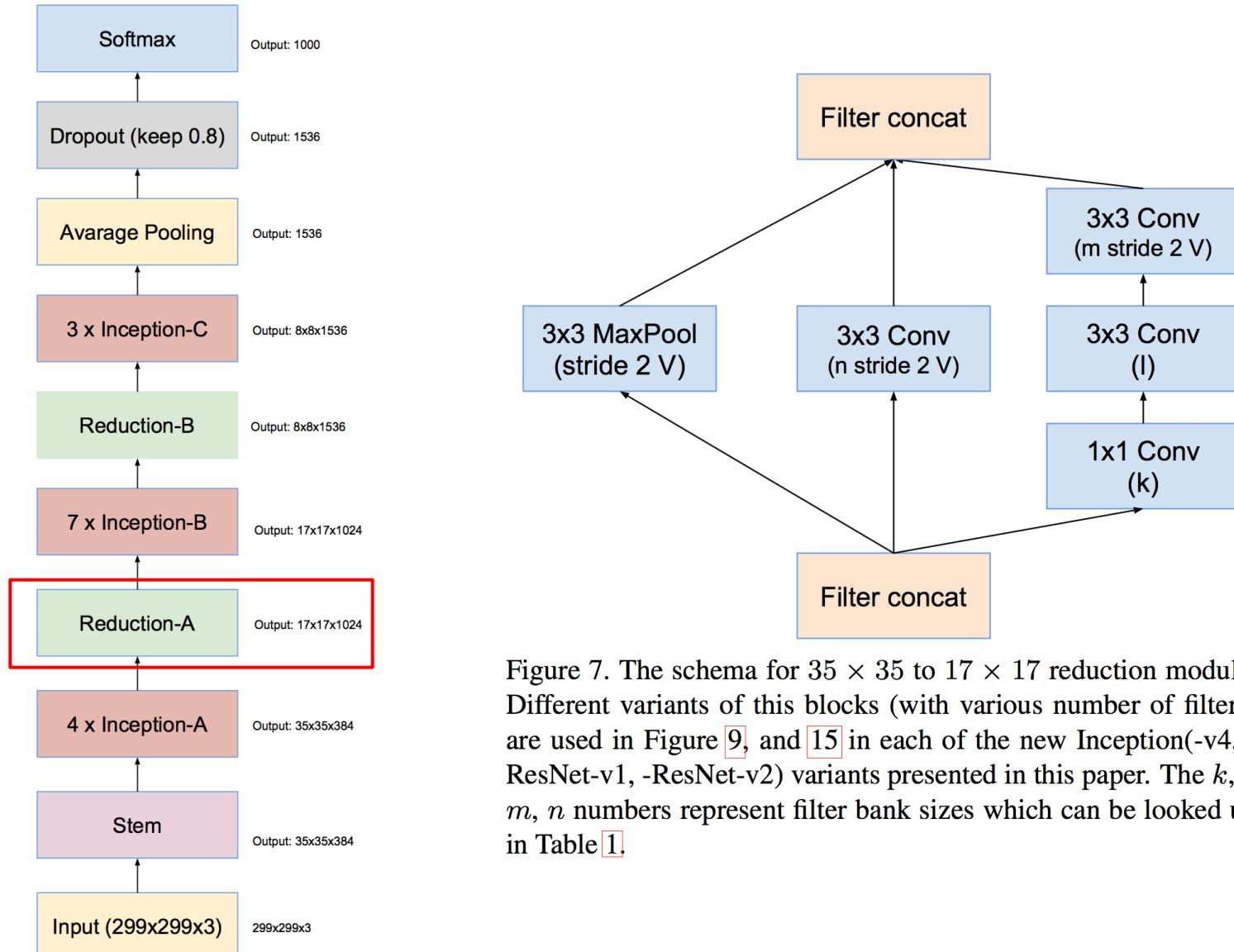


Figure 7. The schema for  $35 \times 35$  to  $17 \times 17$  reduction module. Different variants of this blocks (with various number of filters) are used in Figure 9, and 15 in each of the new Inception(-v4, -ResNet-v1, -ResNet-v2) variants presented in this paper. The  $k, l, m, n$  numbers represent filter bank sizes which can be looked up in Table 1.

# Inception v4

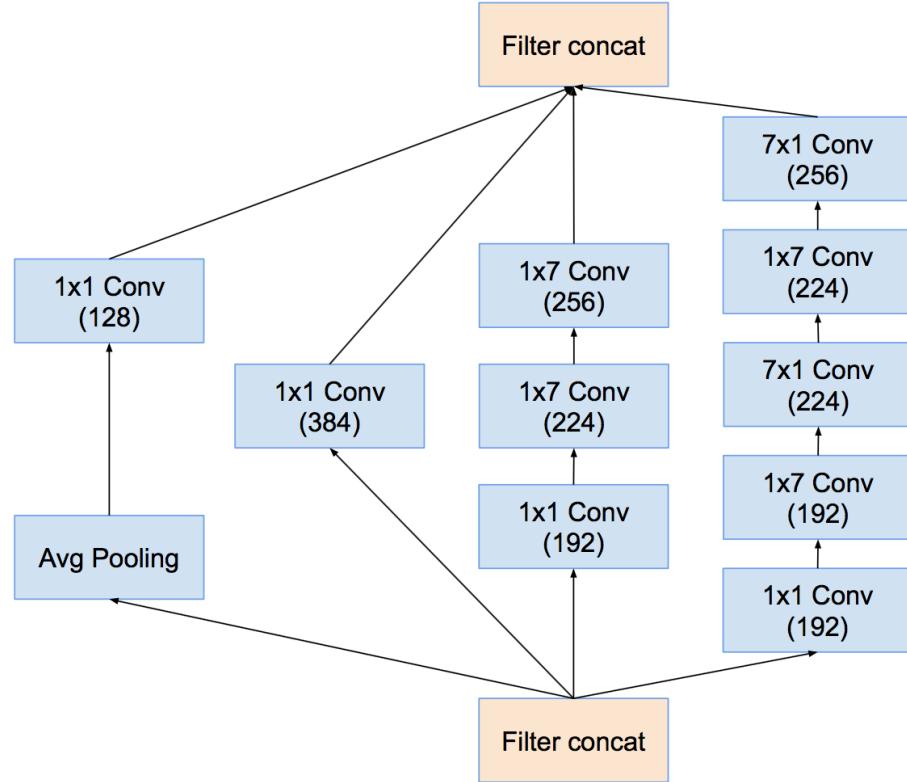
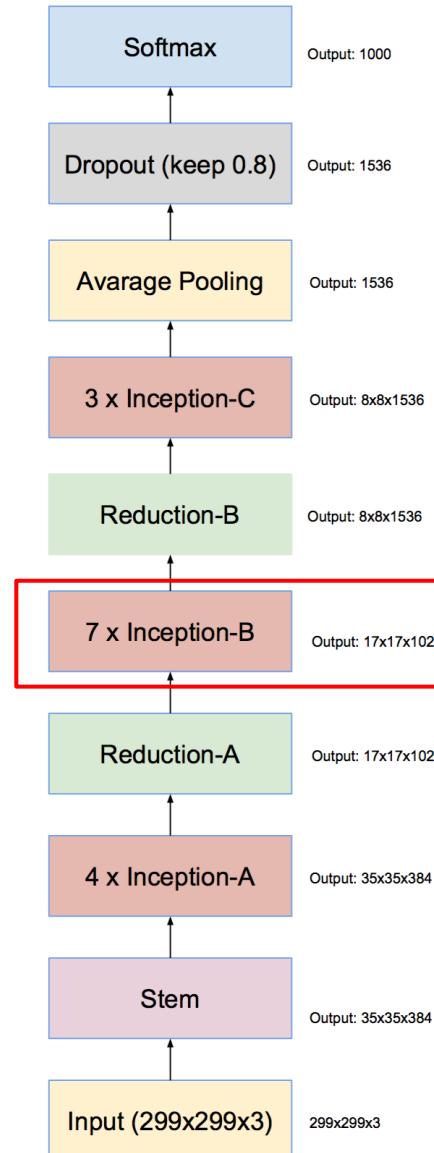


Figure 5. The schema for  $17 \times 17$  grid modules of the pure Inception-v4 network. This is the Inception-B block of Figure 9.

# Inception v4

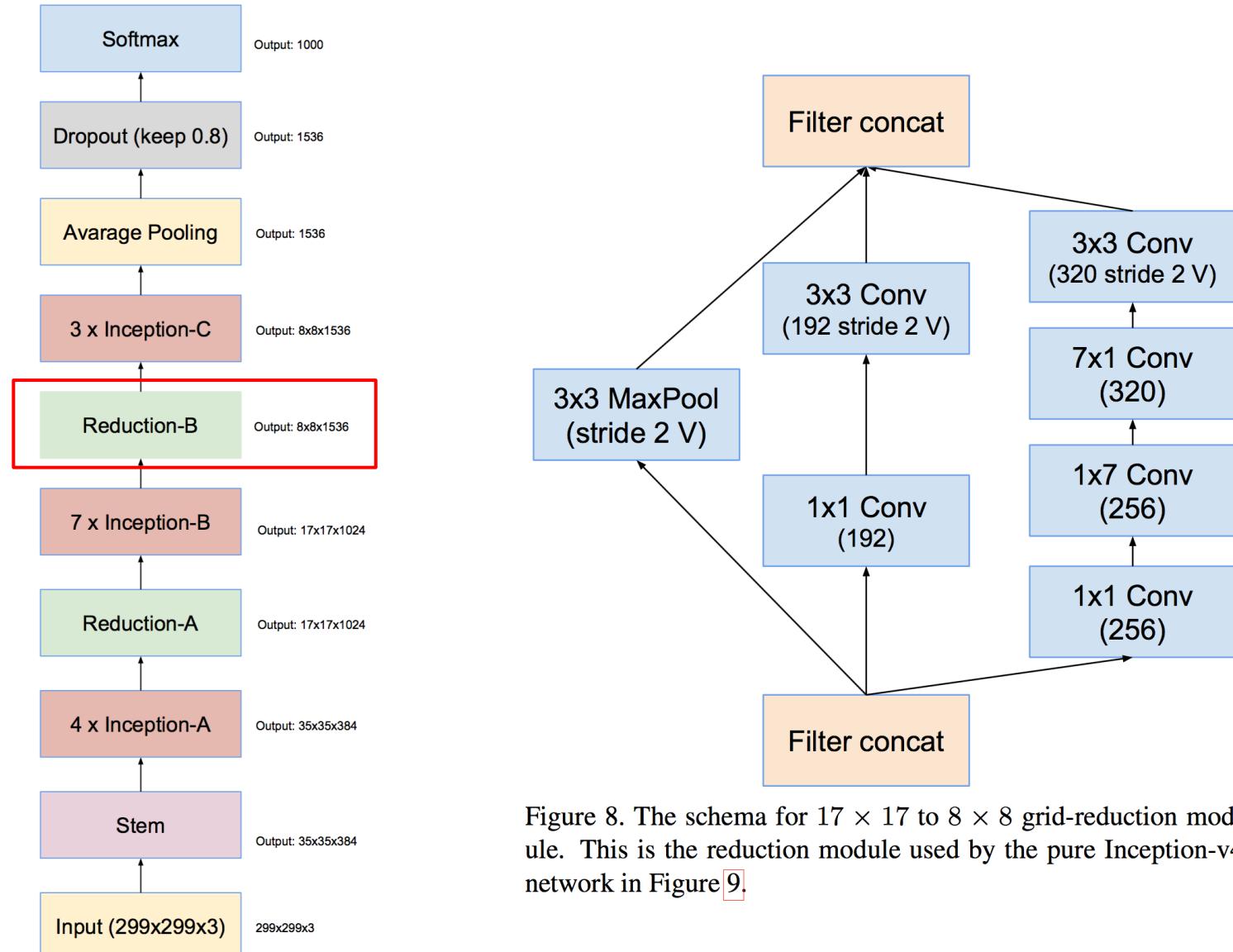


Figure 8. The schema for  $17 \times 17$  to  $8 \times 8$  grid-reduction module. This is the reduction module used by the pure Inception-v4 network in Figure 9.

# Inception v4

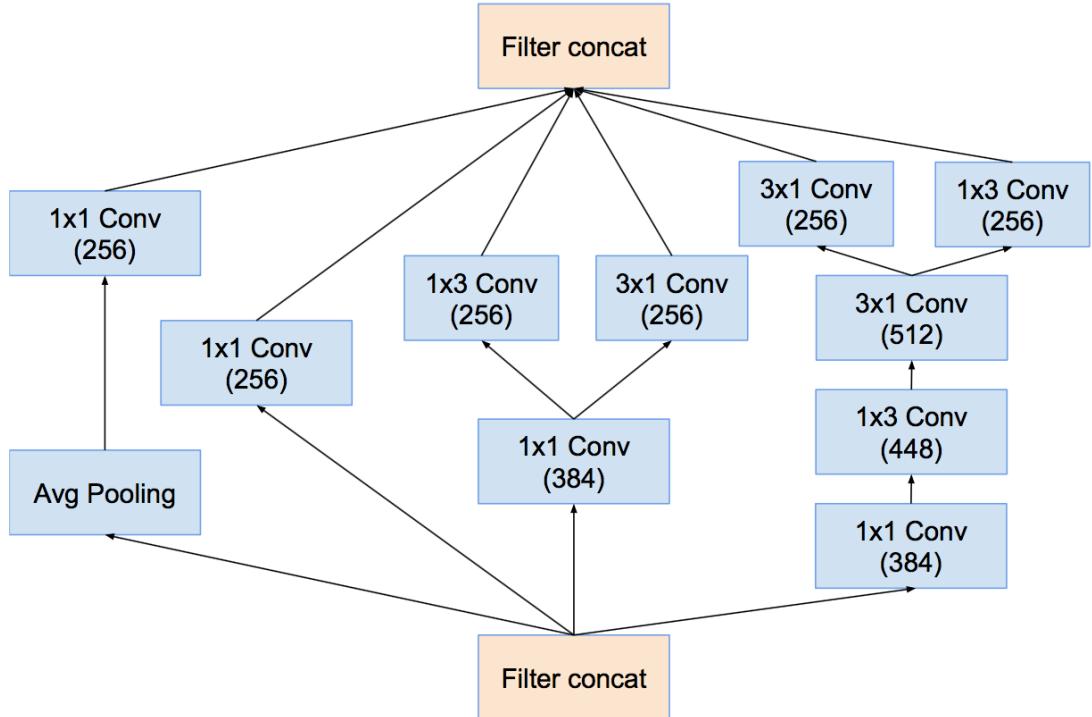
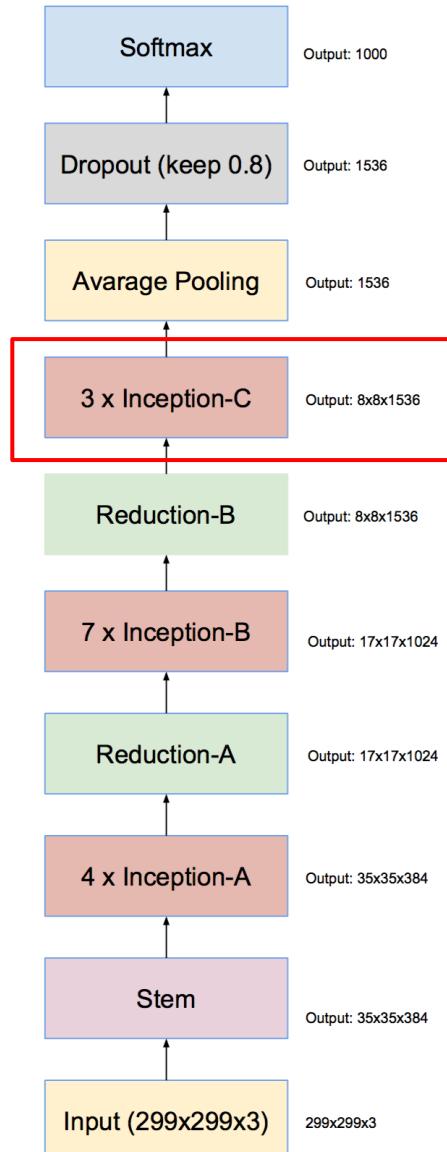


Figure 6. The schema for  $8 \times 8$  grid modules of the pure Inception-v4 network. This is the Inception-C block of Figure 9.

# Inception-ResNet-v1

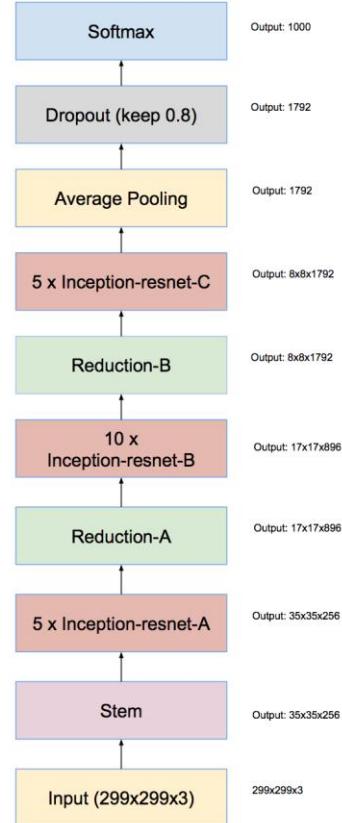


Figure 15. Schema for Inception-ResNet-v1 and Inception-ResNet-v2 networks. This schema applies to both networks but the underlying components differ. Inception-ResNet-v1 uses the blocks as described in Figures 14, 10, 7, 11, 12 and 13. Inception-ResNet-v2 uses the blocks as described in Figures 3, 16, 7, 17, 18 and 19. The output sizes in the diagram refer to the activation vector tensor shapes of Inception-ResNet-v1.

# Inception-ResNet-v1

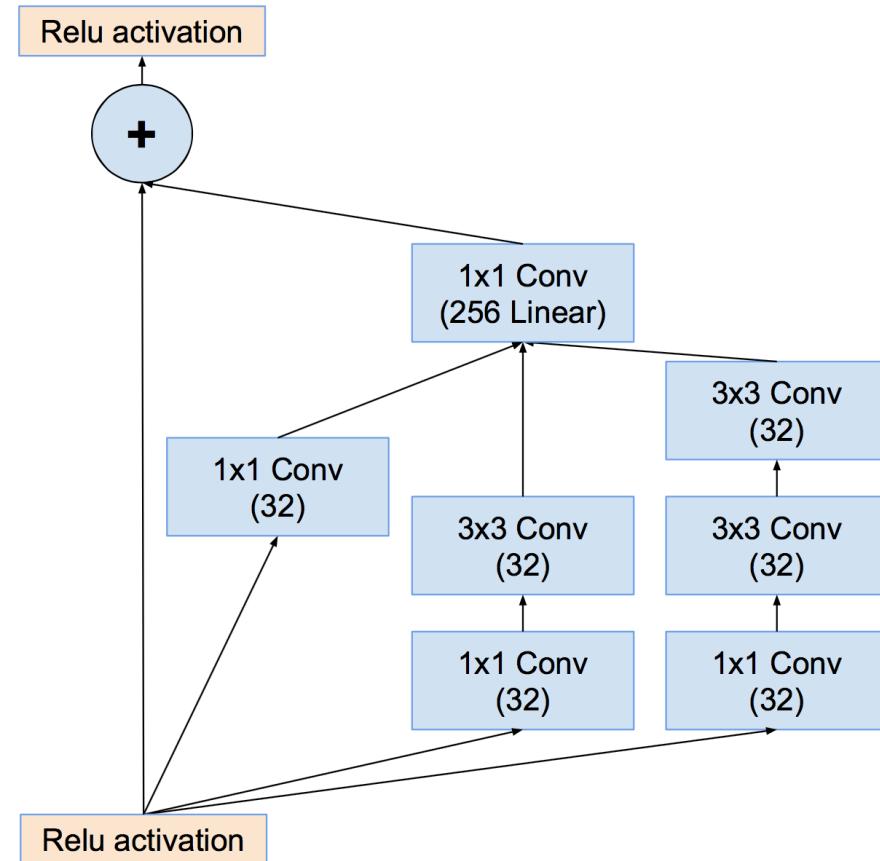
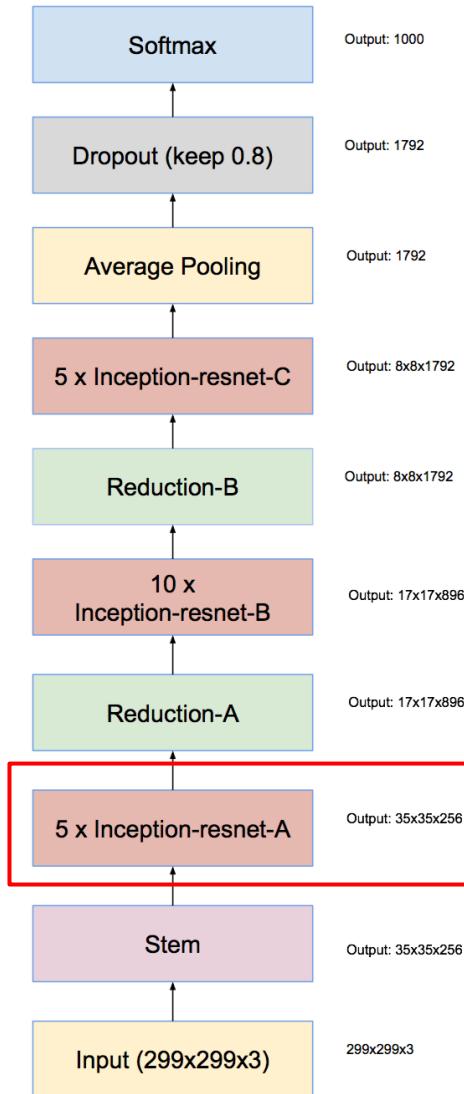


Figure 10. The schema for  $35 \times 35$  grid (Inception-ResNet-A) module of Inception-ResNet-v1 network.

# Inception-ResNet-v1

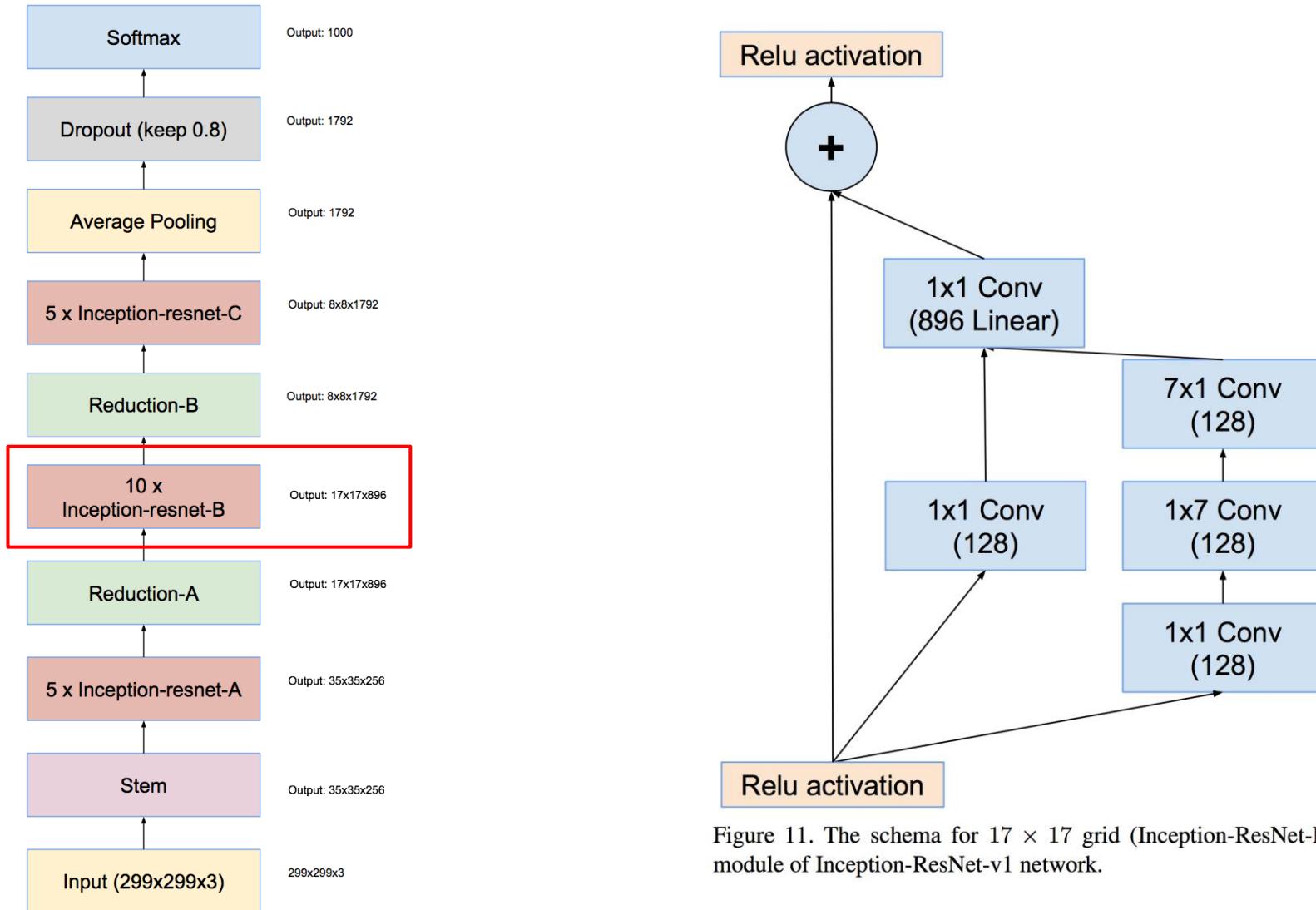


Figure 11. The schema for  $17 \times 17$  grid (Inception-ResNet-B) module of Inception-ResNet-v1 network.

# Inception-ResNet-v1

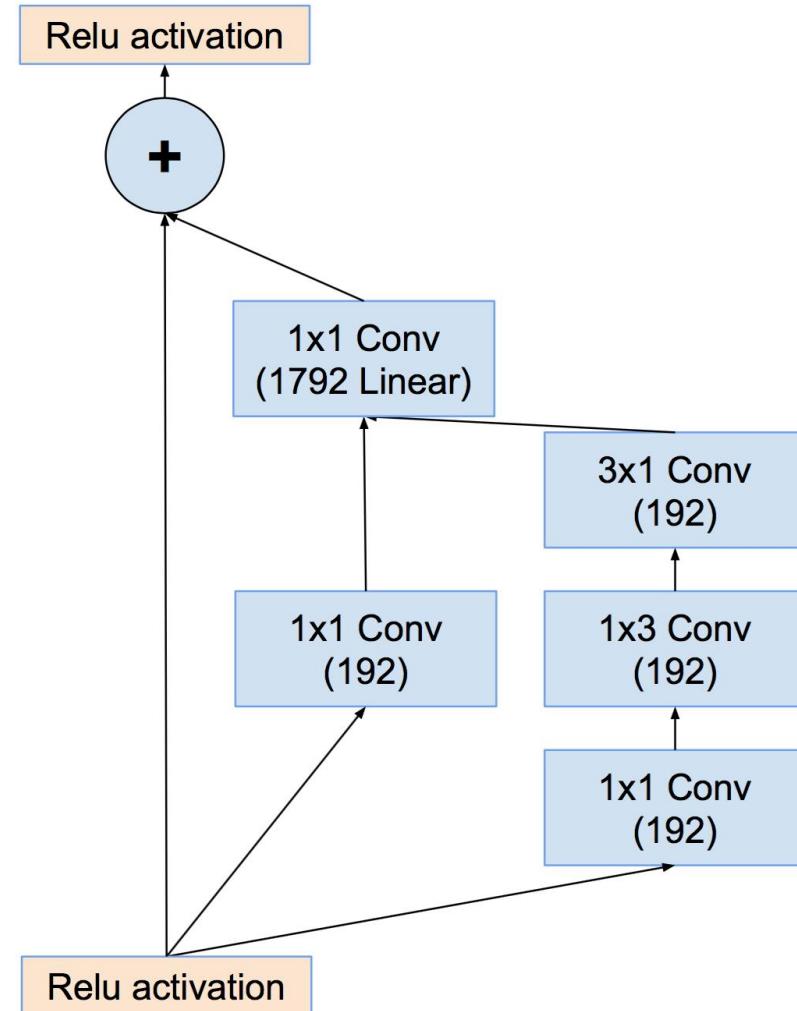
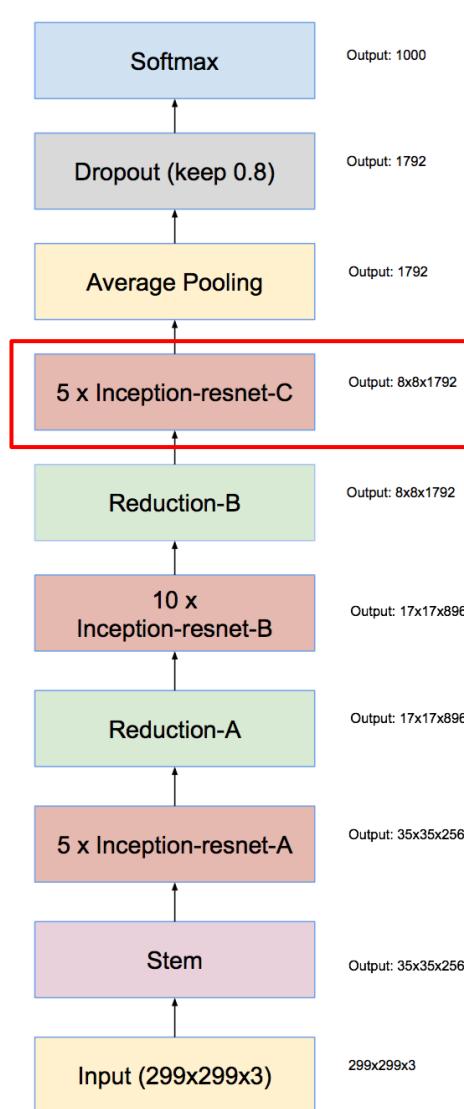


Figure 13. The schema for  $8 \times 8$  grid (Inception-ResNet-C) module of Inception-ResNet-v1 network.

# ResNet

# Deep residual networks

152 layers network

1<sup>st</sup> place on ILSVRC 2015 classification task

1<sup>st</sup> place on ImageNet detection

1<sup>st</sup> place on ImageNet localization

1<sup>st</sup> place on COCO detection

1<sup>st</sup> place on COCO segmentation

# Deeper Network?

Is deeper network always better?

What about vanishing/exploding gradients?

Better initialization methods / batch normalization / ReLU

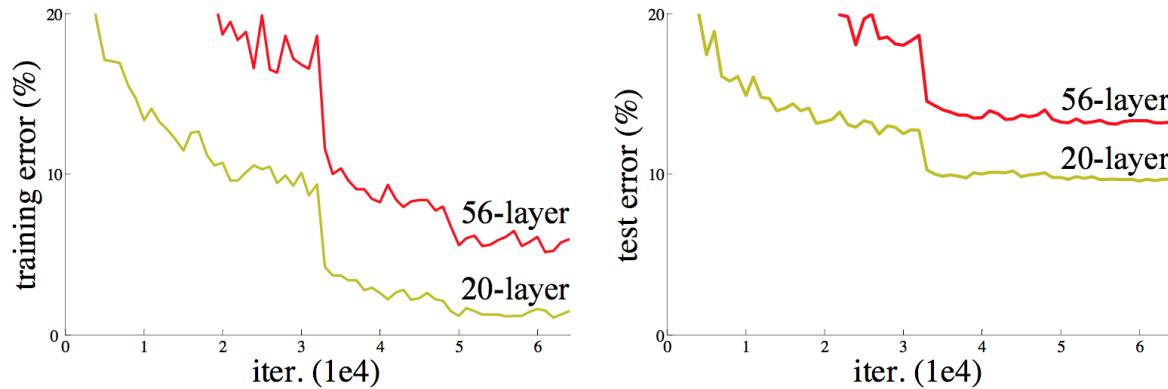
Any other problems?

Overfitting?

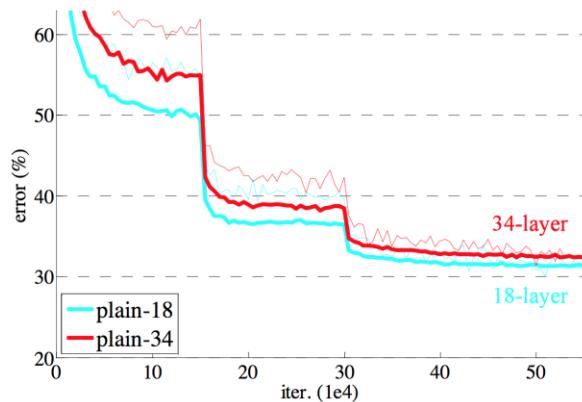
Degradation problem: more depth but lower performance

# Degradation problem

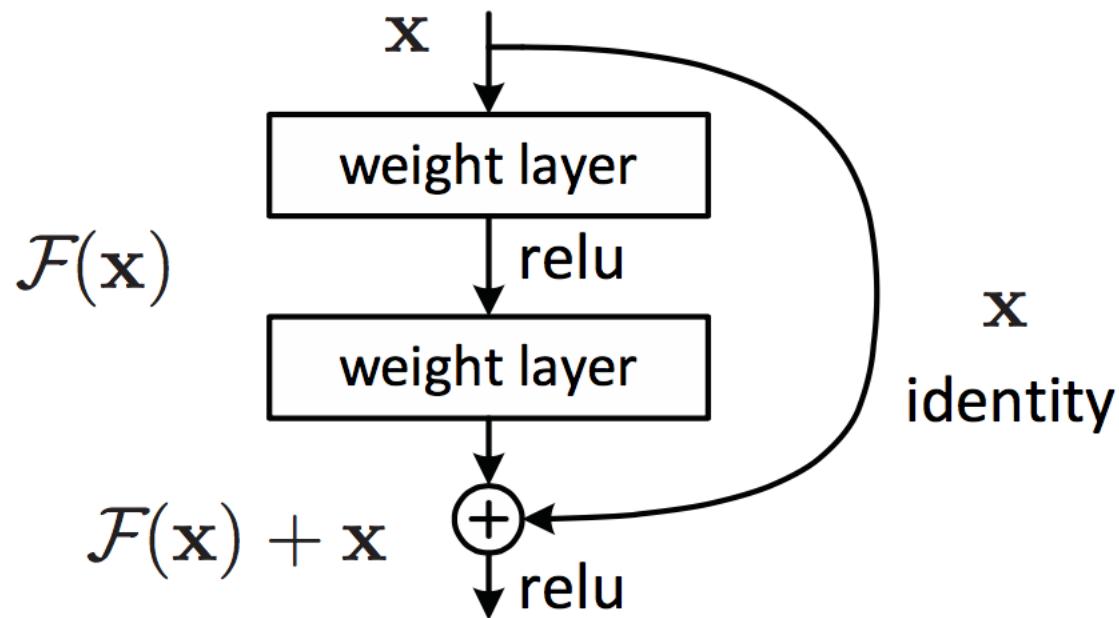
CiFAR 100 Dataset



ImageNet



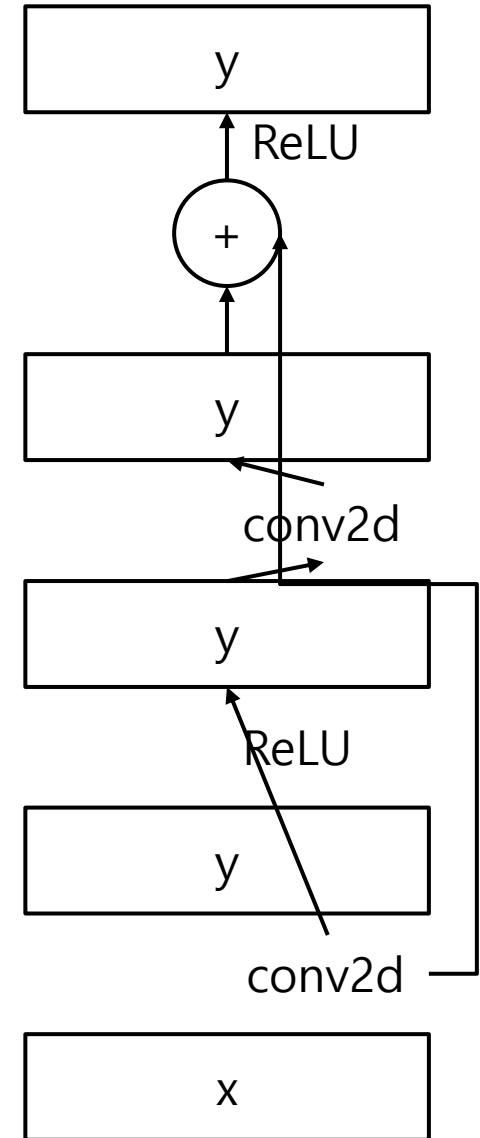
# Residual learning building block



# Residual learning building block

```
def residual_block(x, n_in, n_out, subsample, phase_train, scope='res_block'):
    with tf.variable_scope(scope):
        if subsample:
            y = conv2d(x, n_in, n_out, 3, 2, 'SAME', False, scope='conv_1')
            shortcut = conv2d(x, n_in, n_out, 3, 2, 'SAME',
                               False, scope='shortcut')
        else:
            y = conv2d(x, n_in, n_out, 3, 1, 'SAME', False, scope='conv_1')
            shortcut = tf.identity(x, name='shortcut')
            y = batch_norm(y, n_out, phase_train, scope='bn_1')
            y = tf.nn.relu(y, name='relu_1')
            y = conv2d(y, n_out, n_out, 3, 1, 'SAME', True, scope='conv_2')
            y = batch_norm(y, n_out, phase_train, scope='bn_2')
            y = y + shortcut
            y = tf.nn.relu(y, name='relu_2')
    return y
```

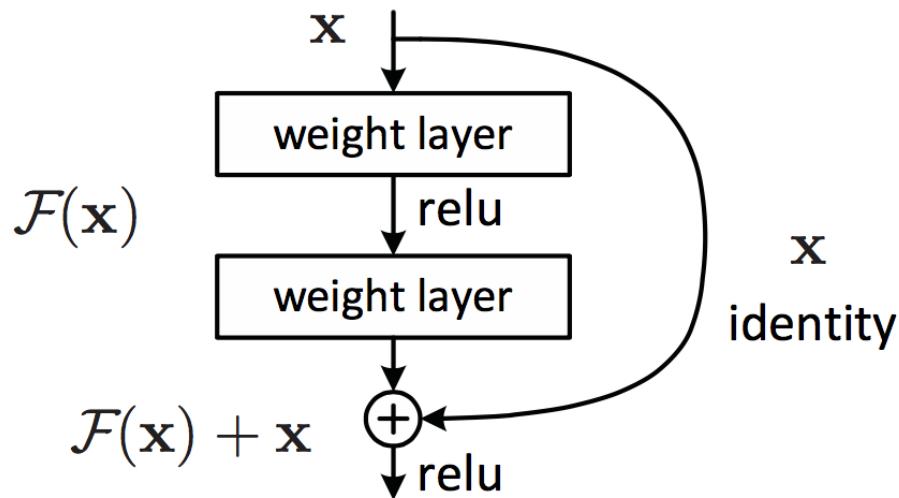
```
def conv2d(x, n_in, n_out, k, s, p='SAME', bias=False, scope='conv'):
    with tf.variable_scope(scope):
        kernel = tf.Variable(
            tf.truncated_normal([k, k, n_in, n_out],
                               stddev=math.sqrt(2/(k*k*n_in))),
            name='weight')
        tf.add_to_collection('weights', kernel)
        conv = tf.nn.conv2d(x, kernel, [1,s,s,1], padding=p)
        if bias:
            bias = tf.get_variable('bias', [n_out], initializer=tf.constant_initializer(0.0))
            tf.add_to_collection('biases', bias)
            conv = tf.nn.bias_add(conv, bias)
    return conv
```



# Why residual?

We **hypothesize** that it is easier to optimize the residual mapping than to optimize the original, unreferenced mapping.

Shortcut connections are used.



# Why residual?

“The extremely deep residual nets are **easy** to optimize. ”

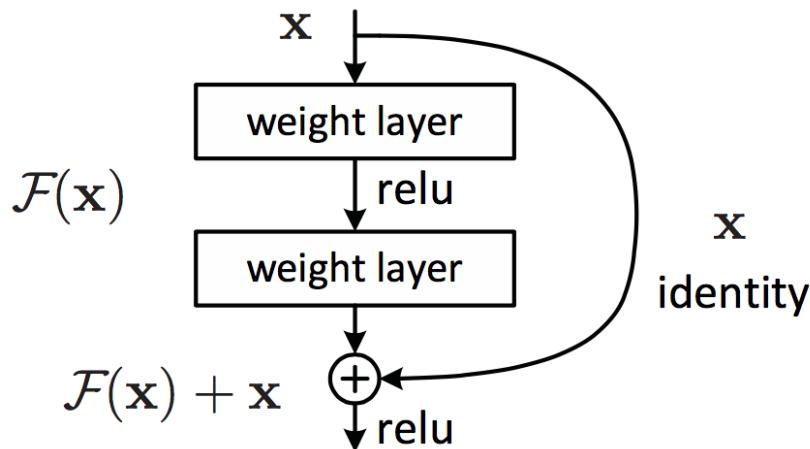
“The deep residual nets can **easily** enjoy accuracy gains from greatly increased depth, producing results substantially better than previous networks. ”



# Residual mapping

Basic residual mapping (same dim.)

$$\mathbf{y} = \mathcal{F}(\mathbf{x}, \{W_i\}) + \mathbf{x}.$$



$$\mathcal{F} = W_2 \sigma(W_1 \mathbf{x})$$

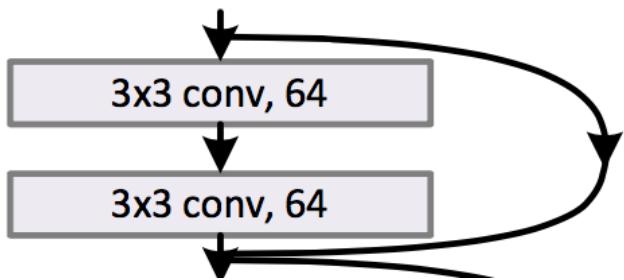
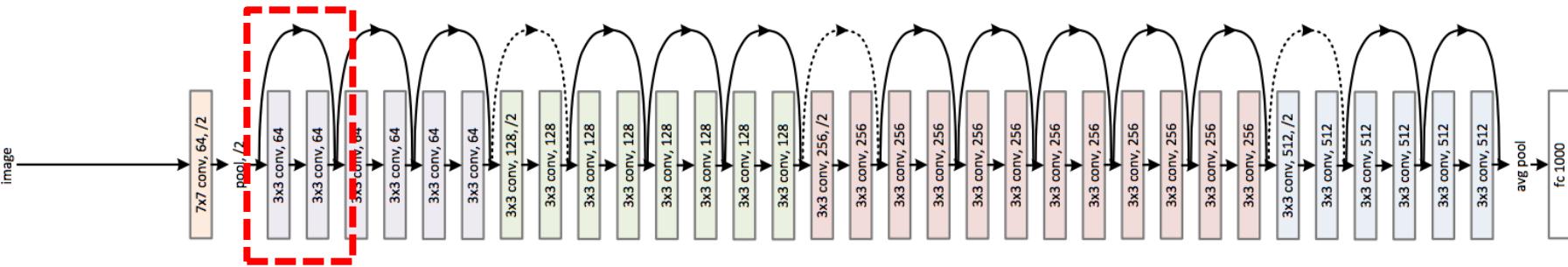
Basic residual mapping (different dim.)

$$\mathbf{y} = \mathcal{F}(\mathbf{x}, \{W_i\}) + W_s \mathbf{x}.$$

"But we will show by experiments that the identity mapping is sufficient for addressing the degradation problem and is economical, and thus  $W$  is only used when matching dimensions."

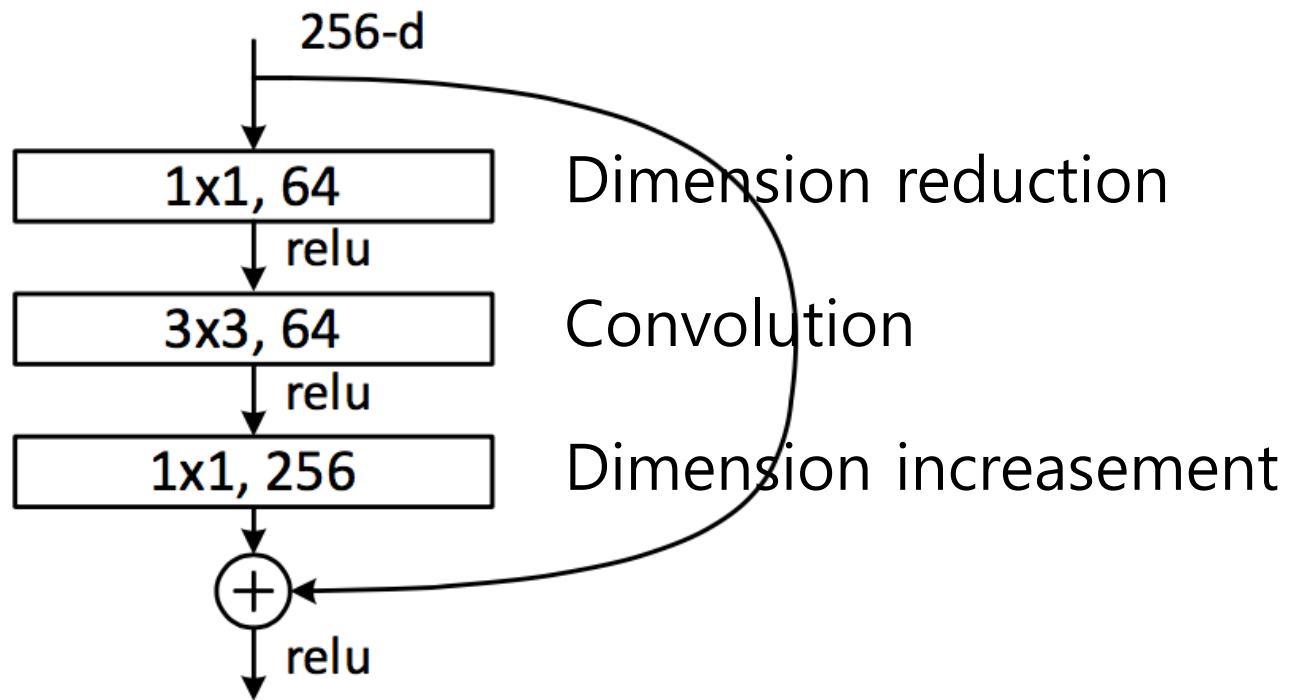
# Deep residual network

34-layer residual

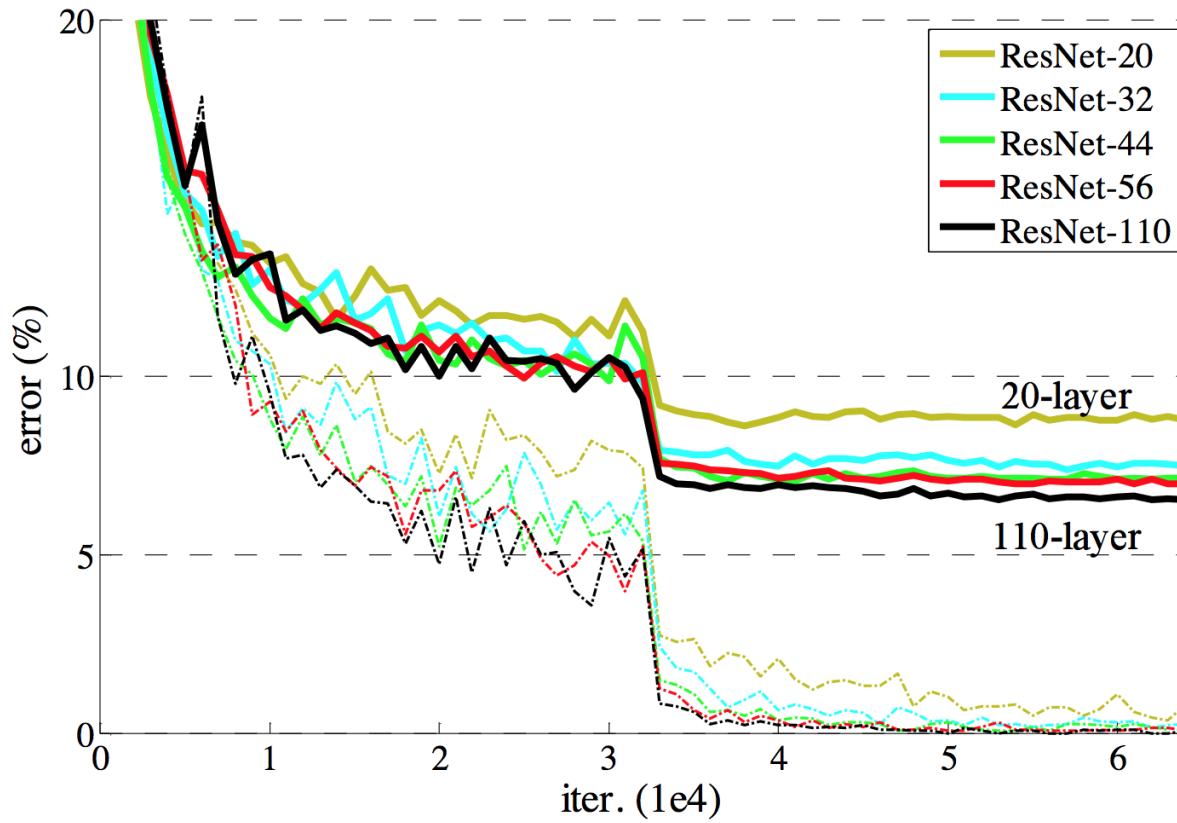


```
def residual_block(x, n_in, n_out, subsample, phase_train, scope='res_block'):
    with tf.variable_scope(scope):
        if subsample:
            y = conv2d(x, n_in, n_out, 3, 2, 'SAME', False, scope='conv_1')
            shortcut = conv2d(x, n_in, n_out, 3, 2, 'SAME',
                              False, scope='shortcut')
        else:
            y = conv2d(x, n_in, n_out, 3, 1, 'SAME', False, scope='conv_1')
            shortcut = tf.identity(x, name='shortcut')
        y = batch_norm(y, n_out, phase_train, scope='bn_1')
        y = tf.nn.relu(y, name='relu_1')
        y = conv2d(y, n_out, n_out, 3, 1, 'SAME', True, scope='conv_2')
        y = batch_norm(y, n_out, phase_train, scope='bn_2')
        y = y + shortcut
        y = tf.nn.relu(y, name='relu_2')
    return y
```

# Deeper bottle architecture



# Experimental results



# Experimental results

