



Software



# Convolutional Neural Net Architectures

Materials from

- Intel Deep Learning <https://www.intel.com/content/www/us/en/developer/learn/course-deep-learning.html>
- Introduction to Neural Networks <https://www.deeplearning.ai>

# Legal Notices and Disclaimers

This presentation is for informational purposes only. INTEL MAKES NO WARRANTIES, EXPRESS OR IMPLIED, IN THIS SUMMARY.

Intel technologies' features and benefits depend on system configuration and may require enabled hardware, software or service activation. Performance varies depending on system configuration. Check with your system manufacturer or retailer or learn more at [intel.com](https://www.intel.com).

This sample source code is released under the [Intel Sample Source Code License Agreement](#).

Intel and the Intel logo are trademarks of Intel Corporation in the U.S. and/or other countries.

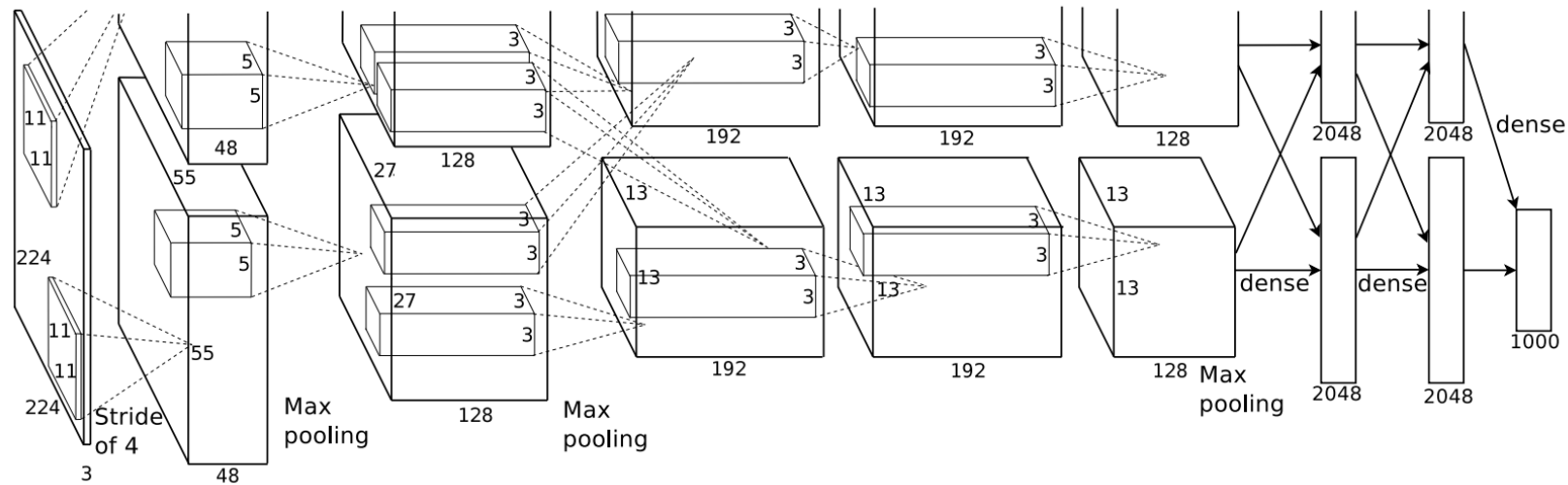
\*Other names and brands may be claimed as the property of others.

Copyright © 2017, Intel Corporation. All rights reserved.

# AlexNet

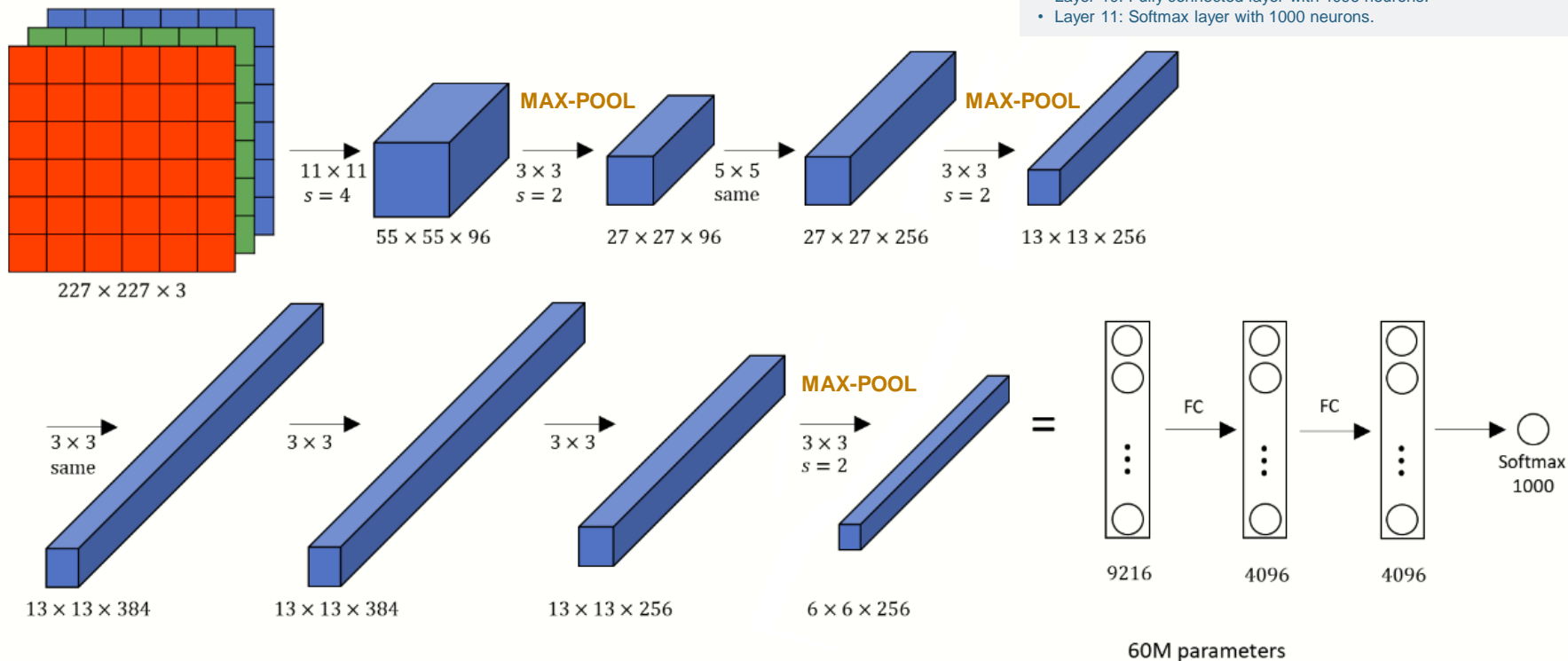
- Created in 2012 for the ImageNet Large Scale Visual Recognition Challenge (ILSVRC)
- Task: predict the correct label from among 1000 classes
- Dataset: around 1.2 million images
- Considered the “flash point” for modern deep learning
- Demolished the competition.
- Top 5 error rate of 15.3%
- Next best: 26.2%

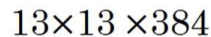
# AlexNet - Model Diagram



# AlexNet - Model Diagram

- Layer 1: Convolutional layer with 96 filters of size 11x11 and a stride of 4.
- Layer 2: Max-pooling layer with a pooling size of 3x3 and a stride of 2.
- Layer 3: Convolutional layer with 256 filters of size 5x5 and a stride of 1.
- Layer 4: Max-pooling layer with a pooling size of 3x3 and a stride of 2.
- Layer 5: Convolutional layer with 384 filters of size 3x3 and a stride of 1.
- Layer 6: Convolutional layer with 384 filters of size 3x3 and a stride of 1.
- Layer 7: Convolutional layer with 256 filters of size 3x3 and a stride of 1.
- Layer 8: Max-pooling layer with a pooling size of 3x3 and a stride of 2.
- Layer 9: Fully connected layer with 4096 neurons.
- Layer 10: Fully connected layer with 4096 neurons.
- Layer 11: Softmax layer with 1000 neurons.





- ReLU

- Multiple GPUs.

## nom parameters

# AlexNet - Details

- They performed *data augmentation* for training
  - Includes Cropping, horizontal flipping, and other manipulations
- Basic Template:
  - Convolutions with ReLUs
  - Sometimes add maxpool after convolutional layer
  - Fully connected layers at the end before a softmax classifier

# VGG

- Simplify Network Structure
- Achieved at error rate of 7.32%, Top 5 Recall in 2014
- Avoid Manual Choices of Convolution Size
- Very Deep Network with 3x3 Convolutions
- These “effectively” give rise to larger convolutions

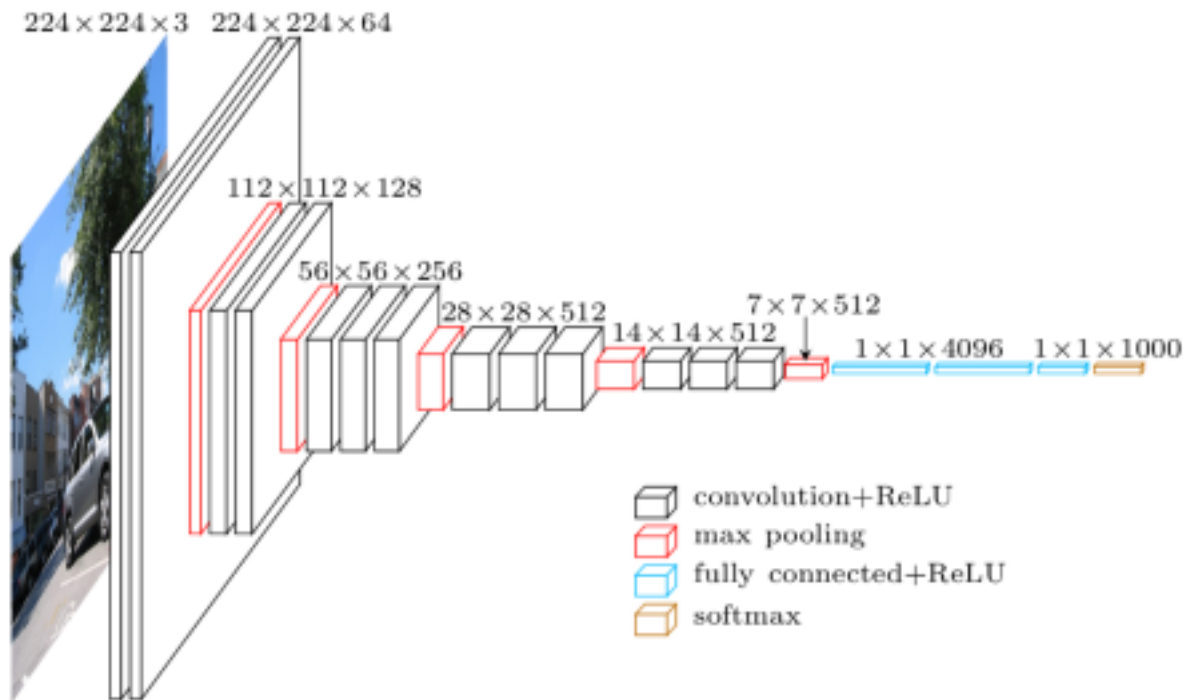
Reference:

*Very Deep Convolutional Networks for Large-Scale Image Recognition*

Karen Simonyan and Andrew Zisserman, 2014



# VGG-16 Diagram

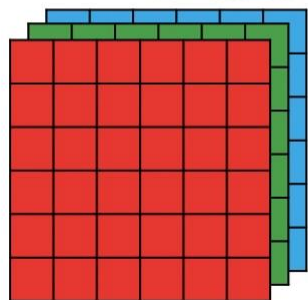


# VGG - 16

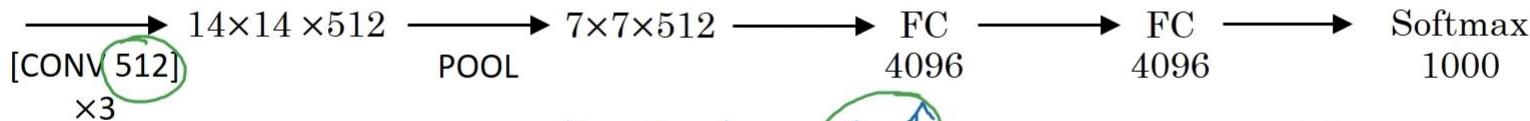
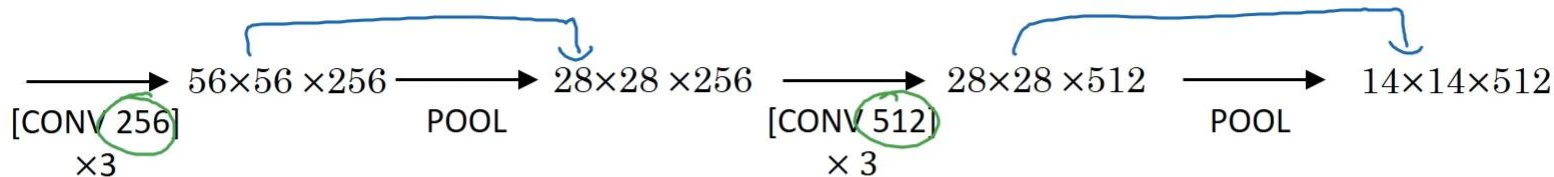
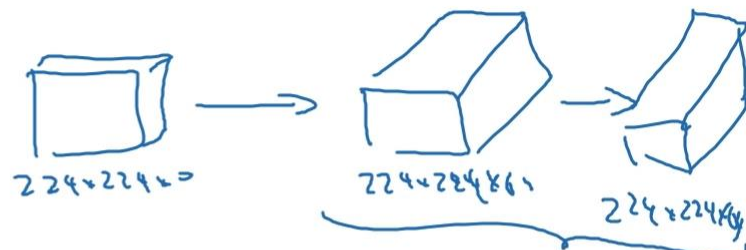
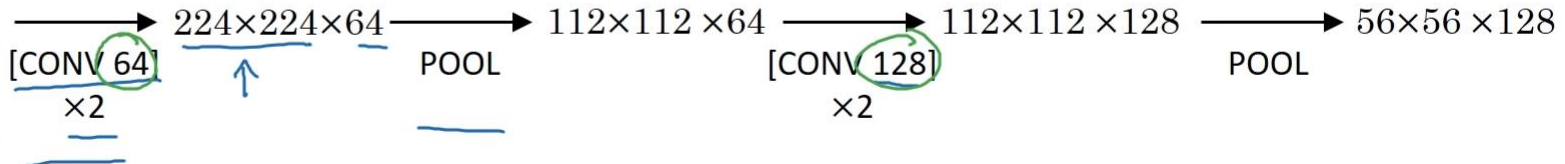
CONV = 3x3 filter, s = 1, same

VGG-19

MAX-POOL = 2x2, s = 2



224x224 x 3

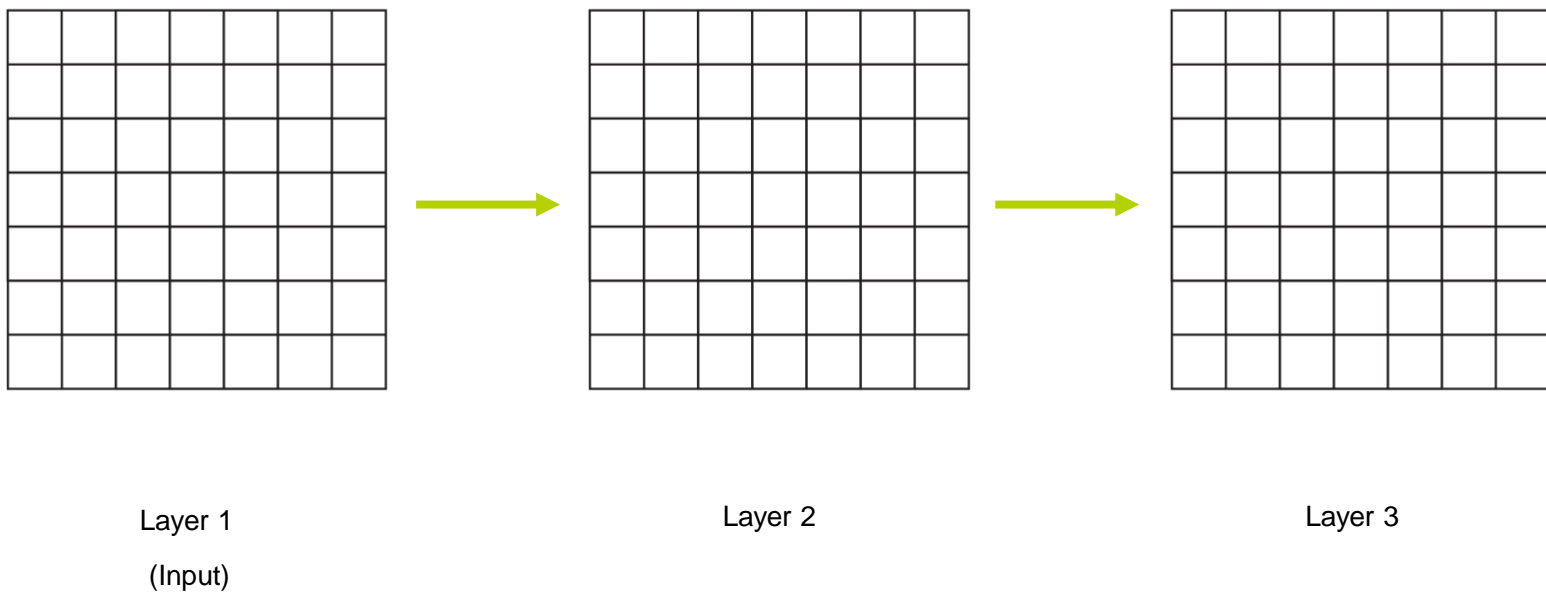


$n_H, n_W \downarrow$

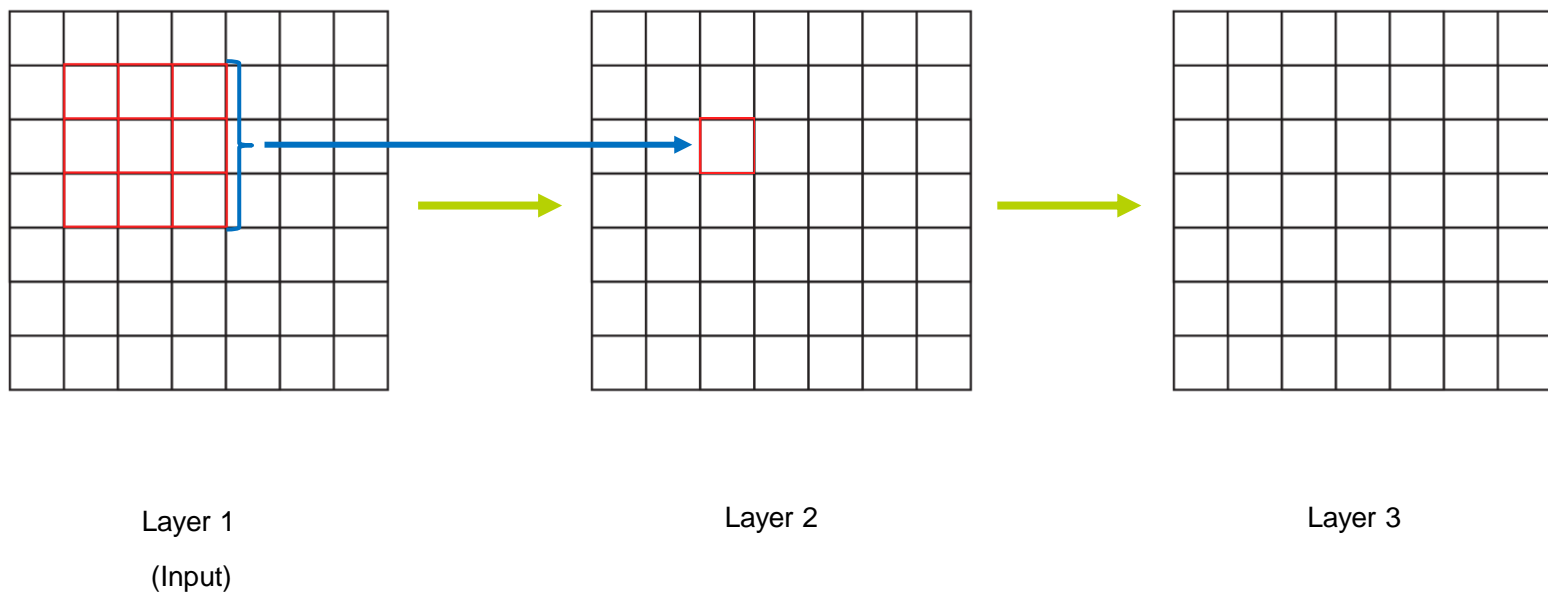
$n_C \uparrow$

~138M

# VGG

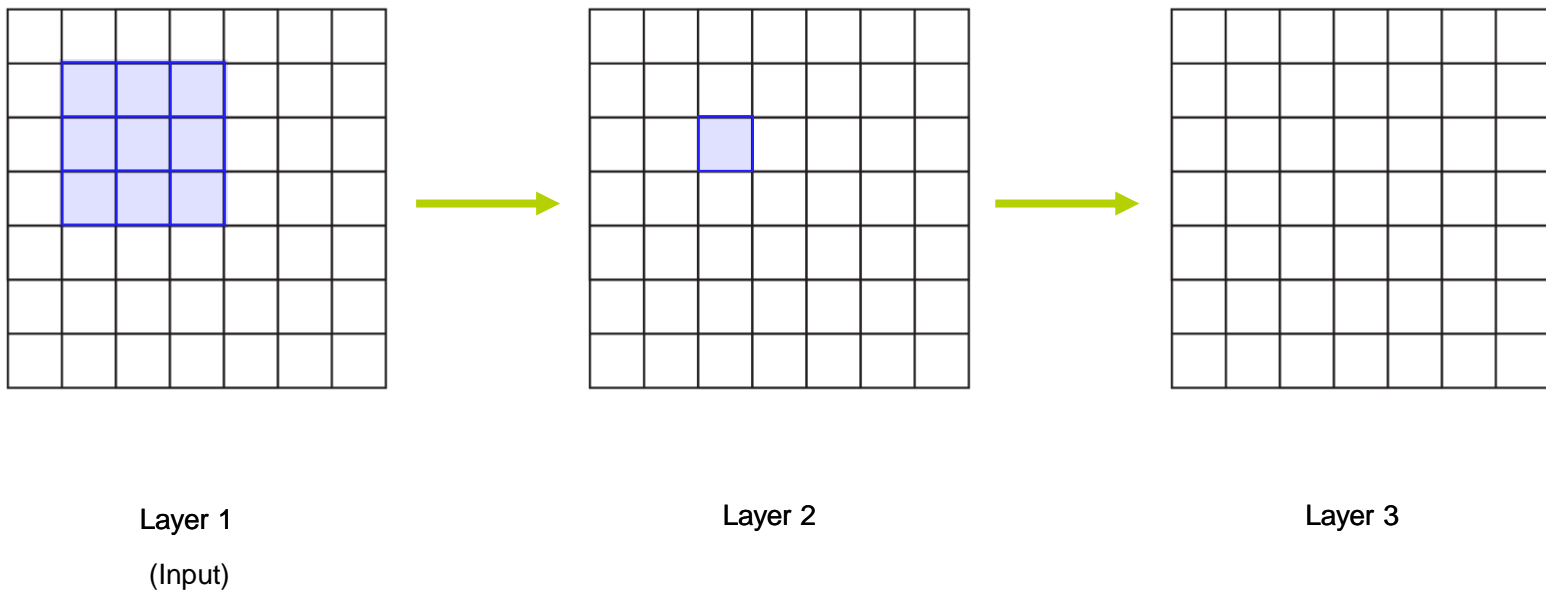


# VGG



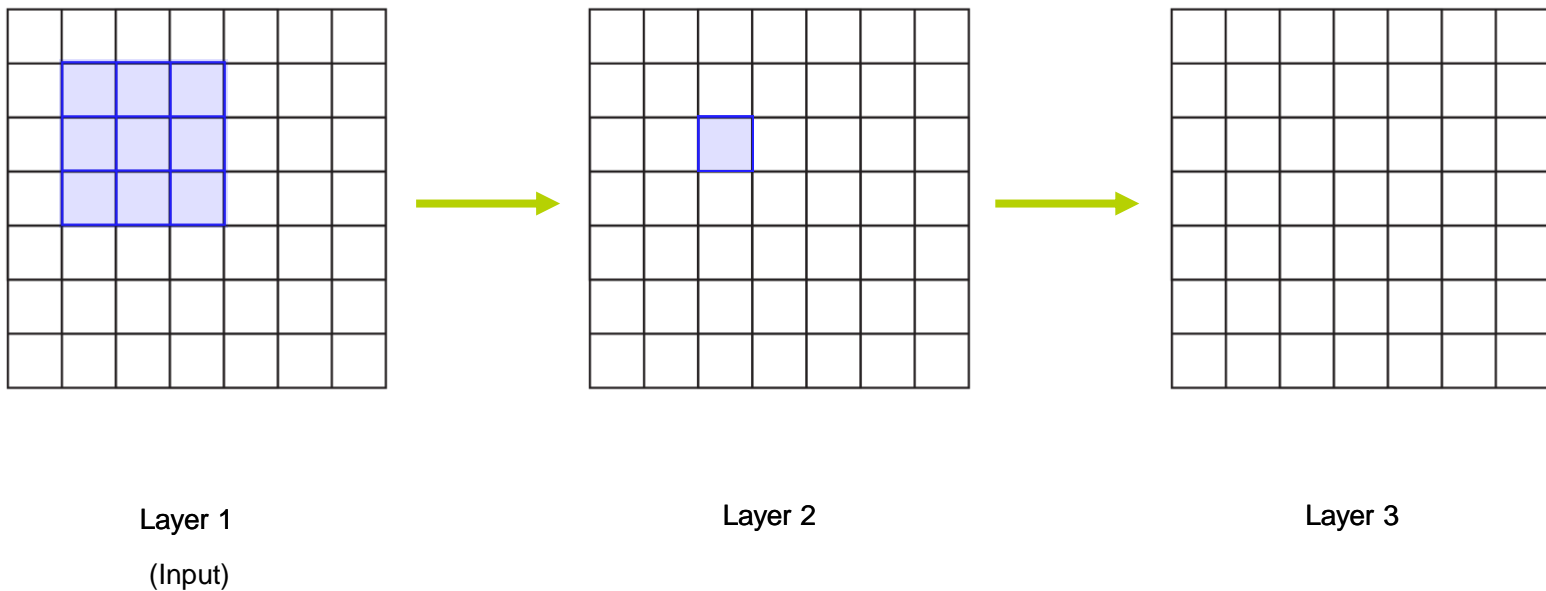
# VGG

We can say that the “receptive field” of Layer 2 is 3x3  
Each output has been influenced by a 3x3 patch of inputs



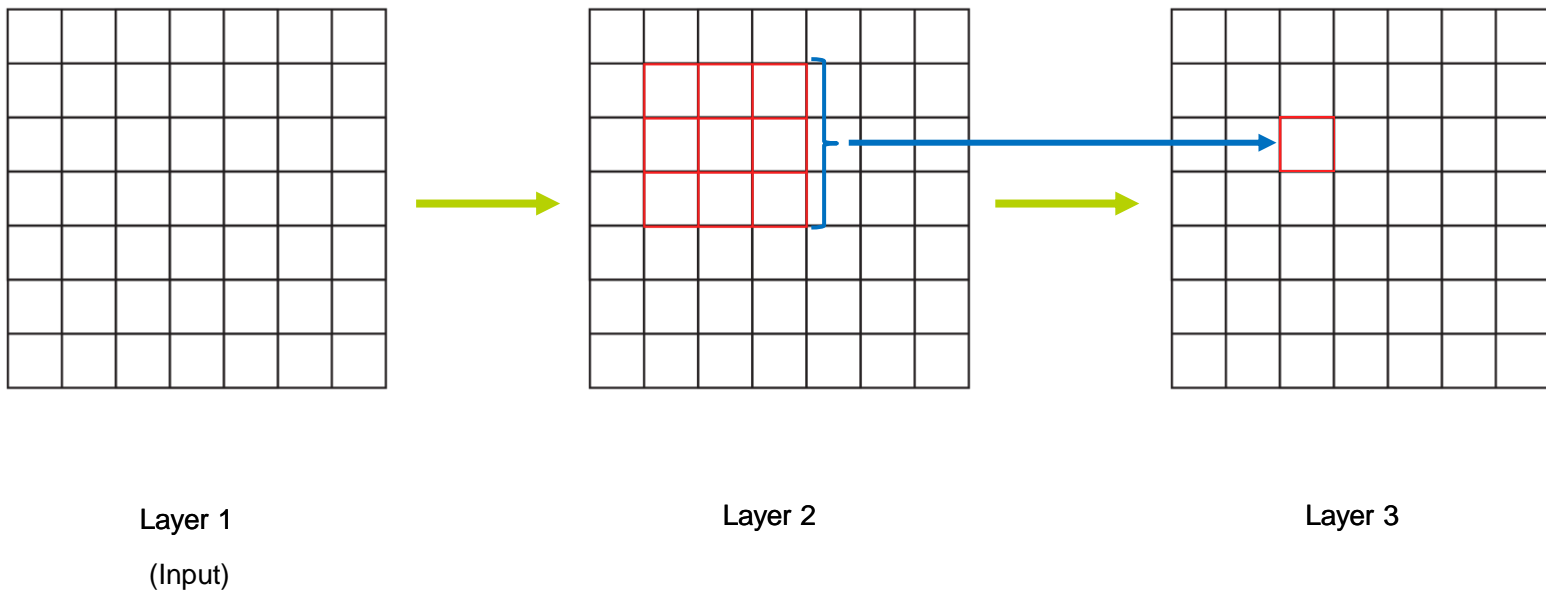
# VGG

What about on Layer 3?

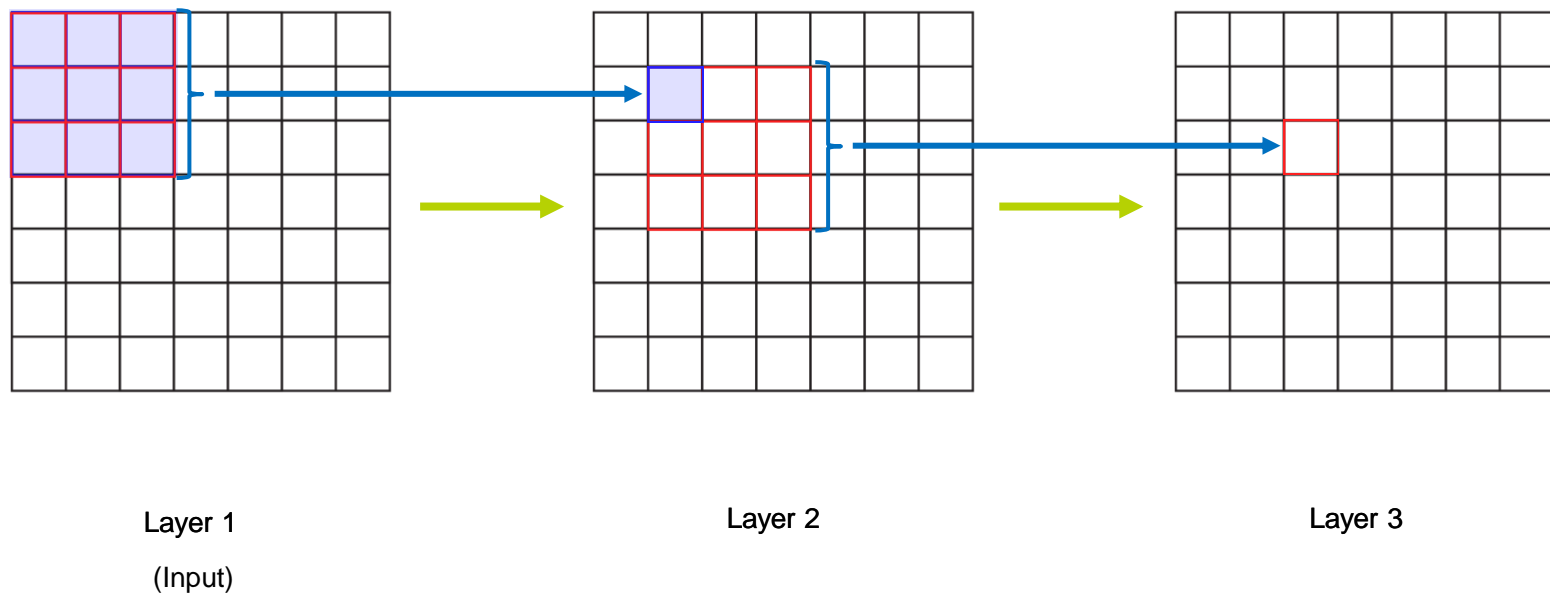


# VGG

This output on Layer 3 uses a 3x3 patch from Layer 2  
How much from Layer 1 does it use?

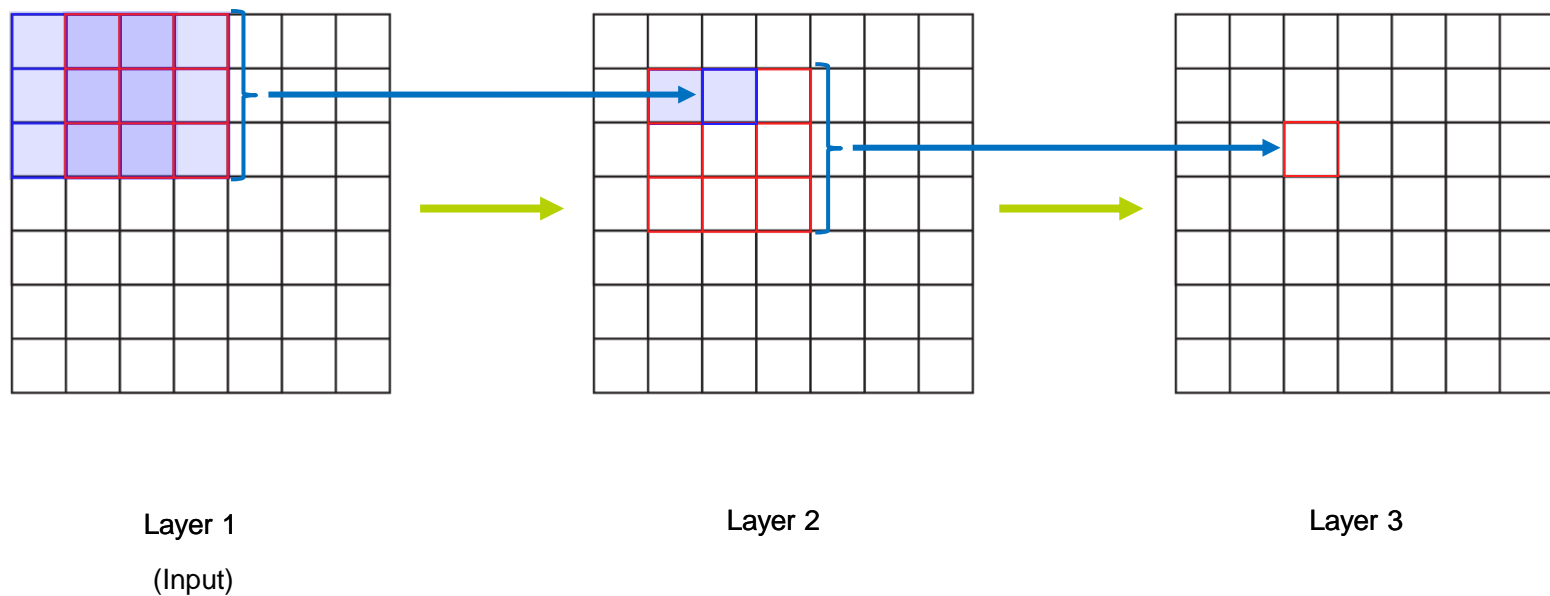


# VGG

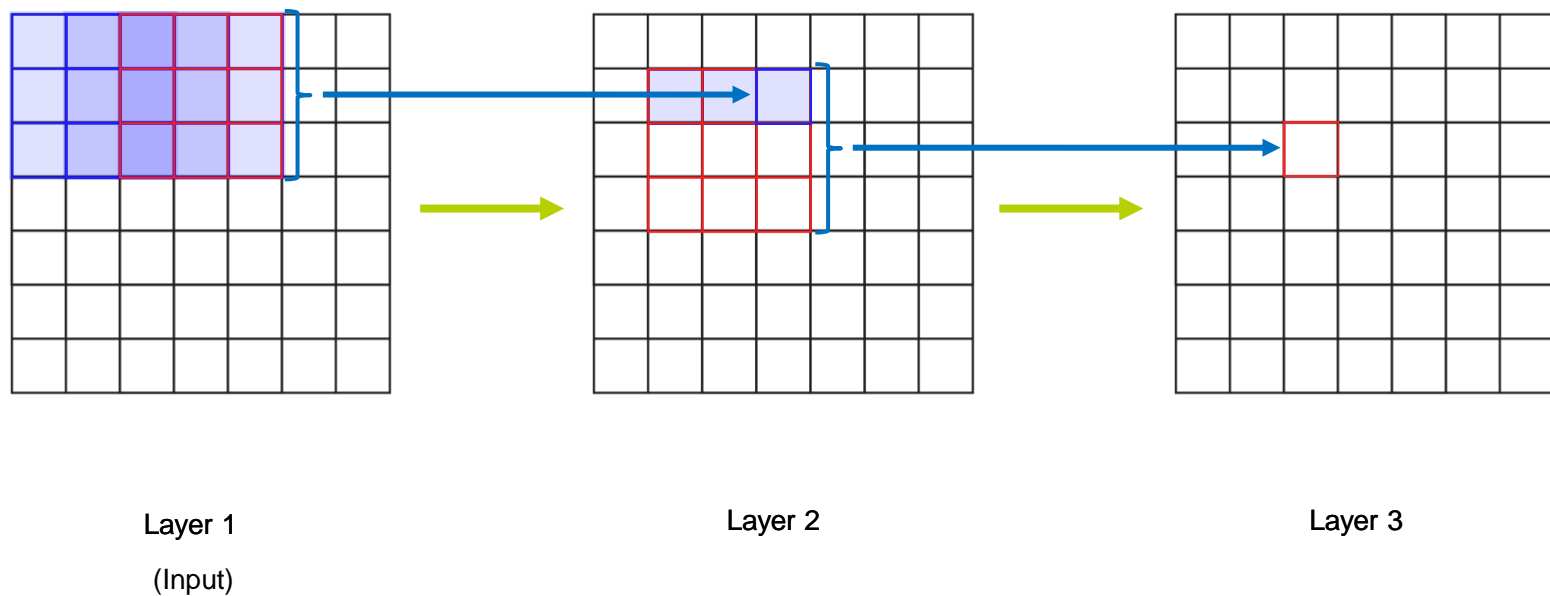




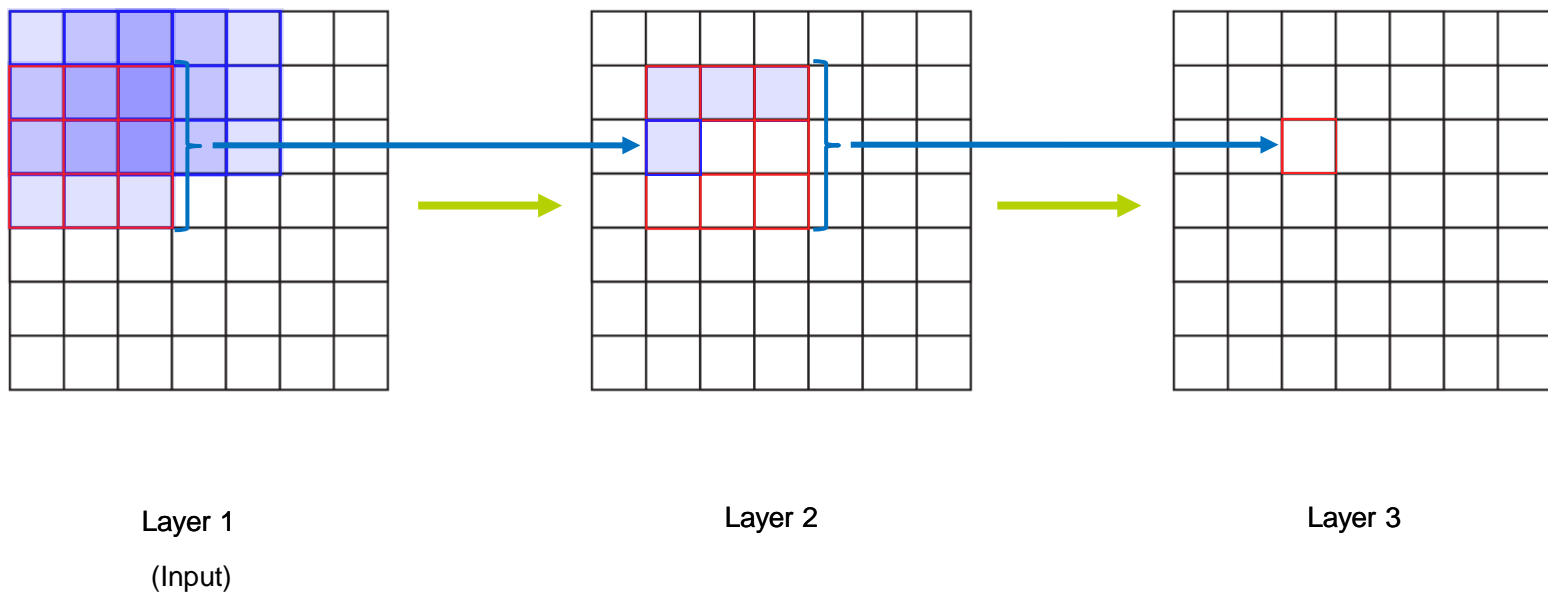
# VGG



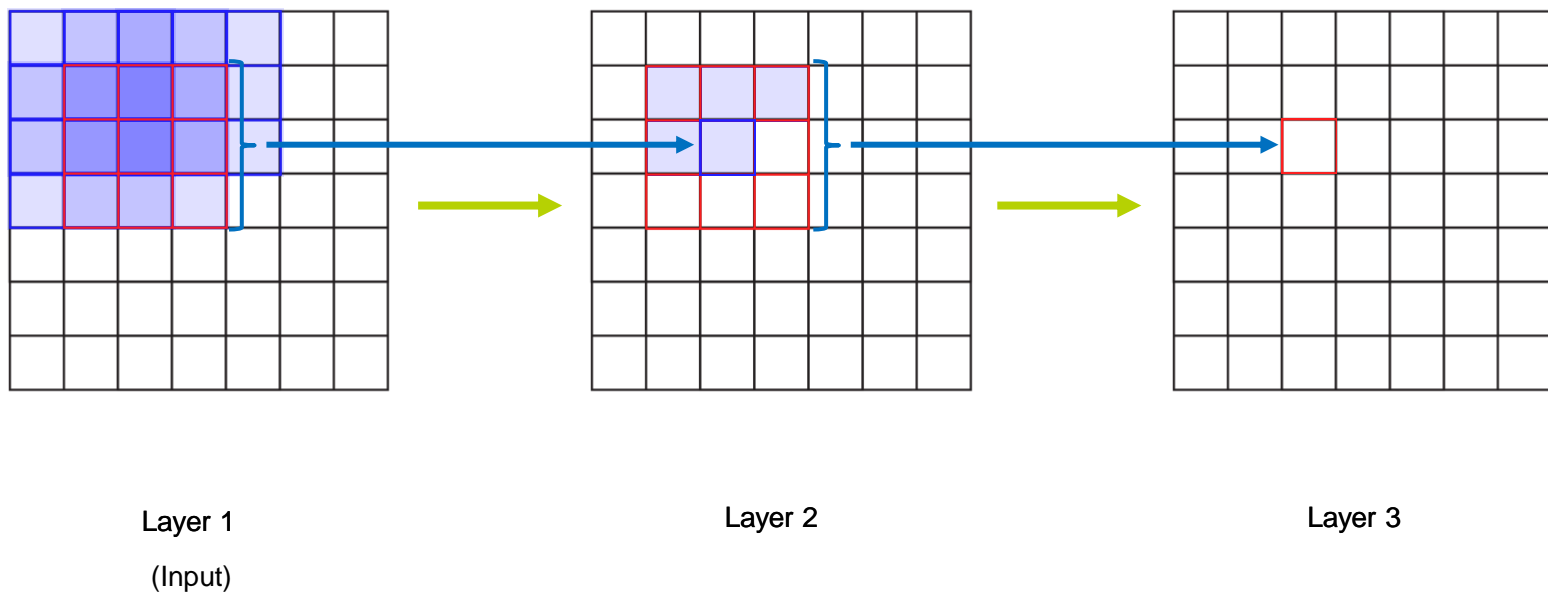
# VGG



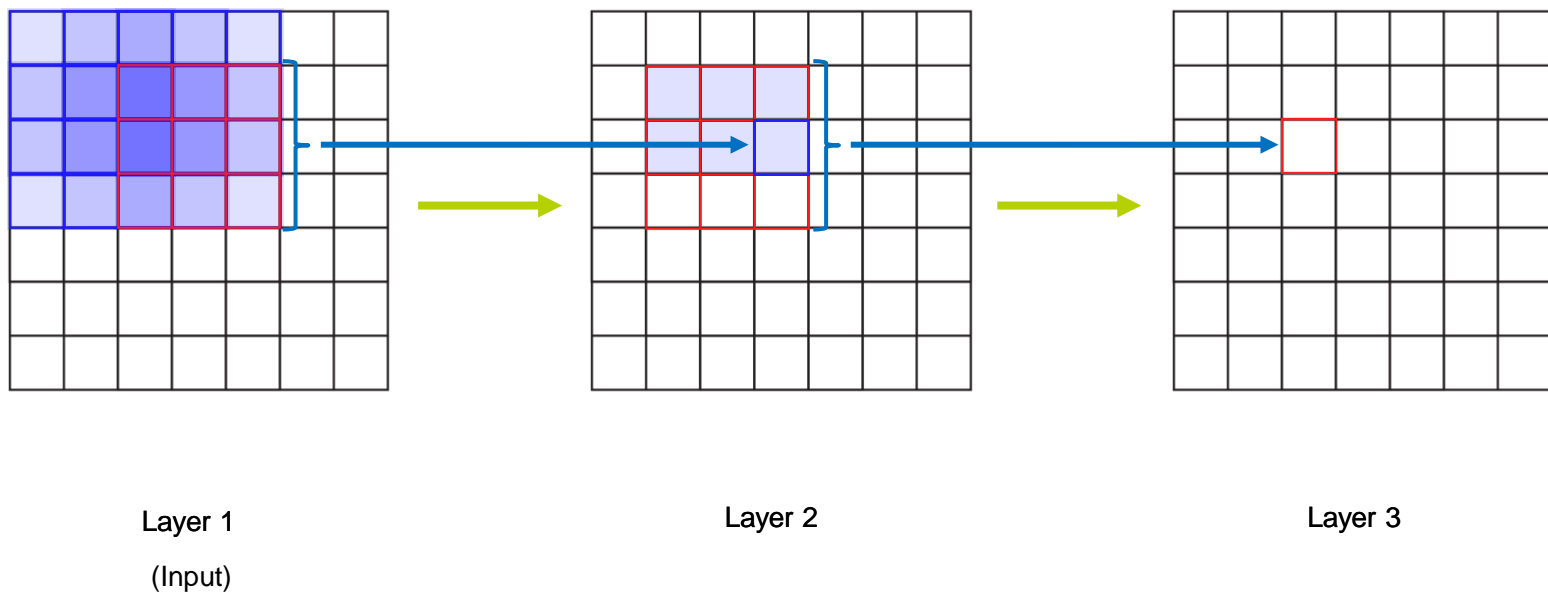
# VGG



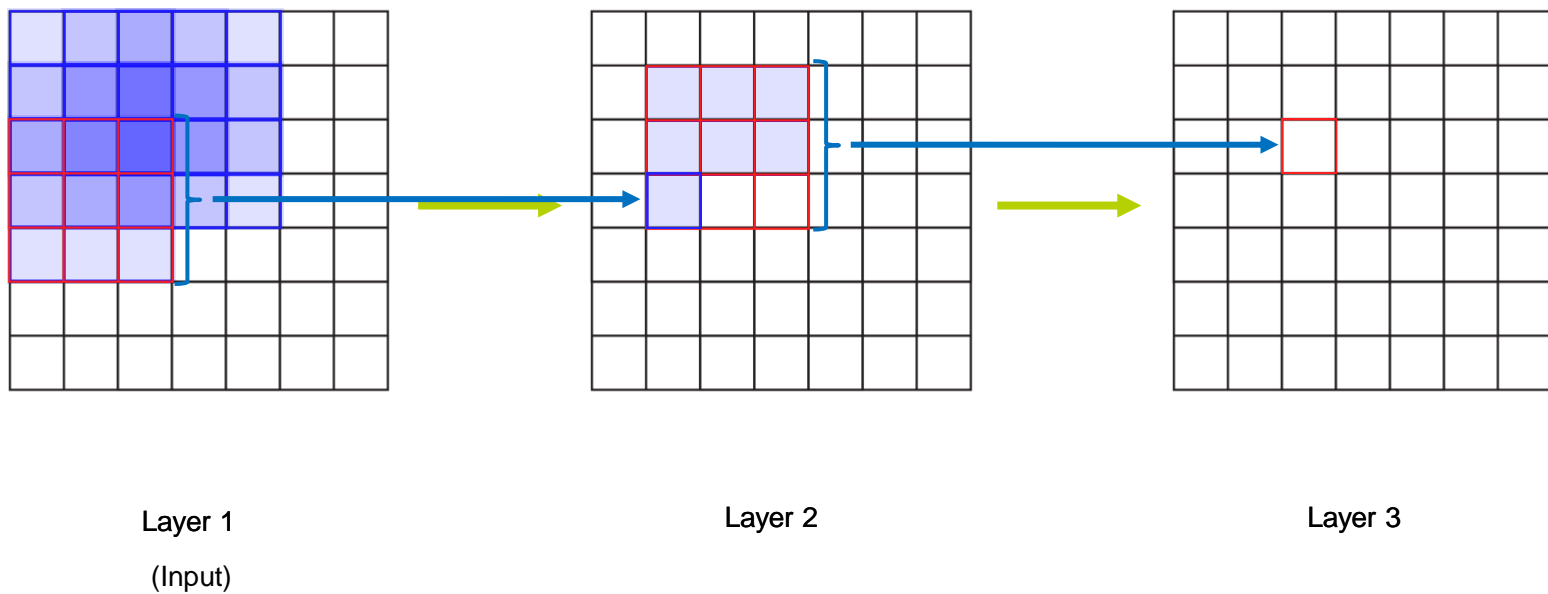
# VGG



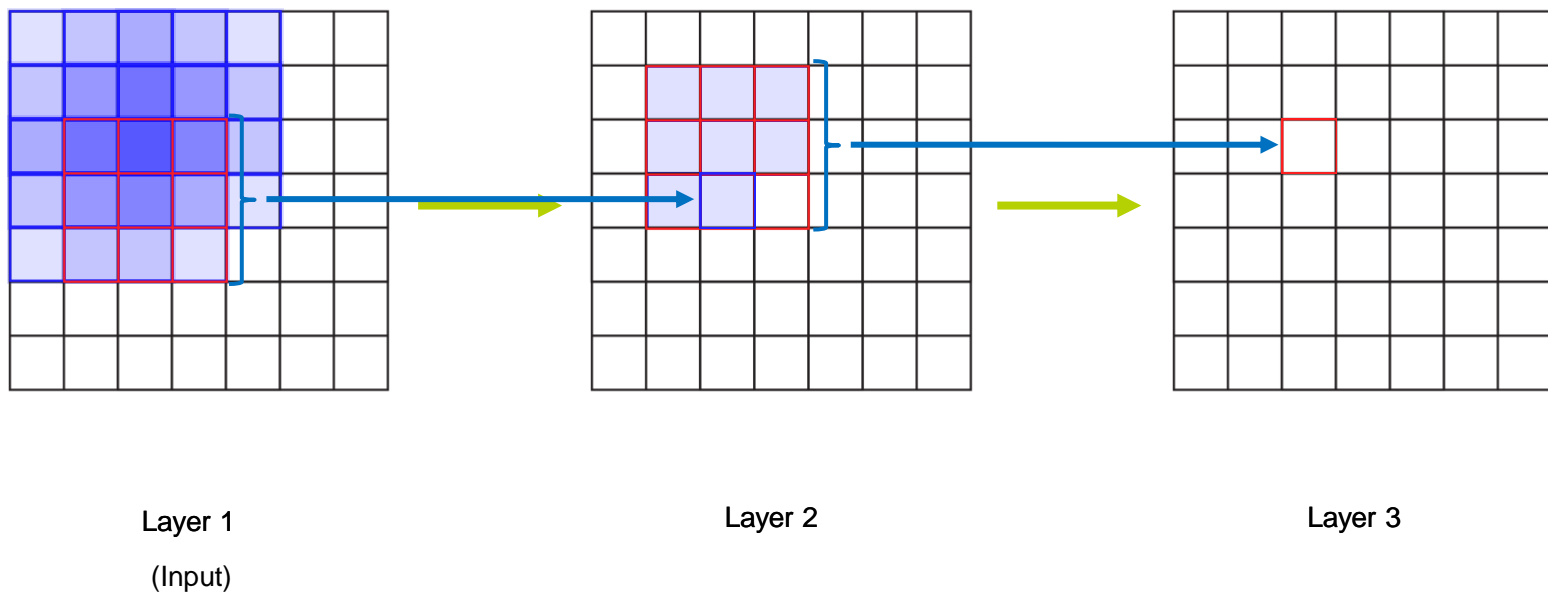
# VGG



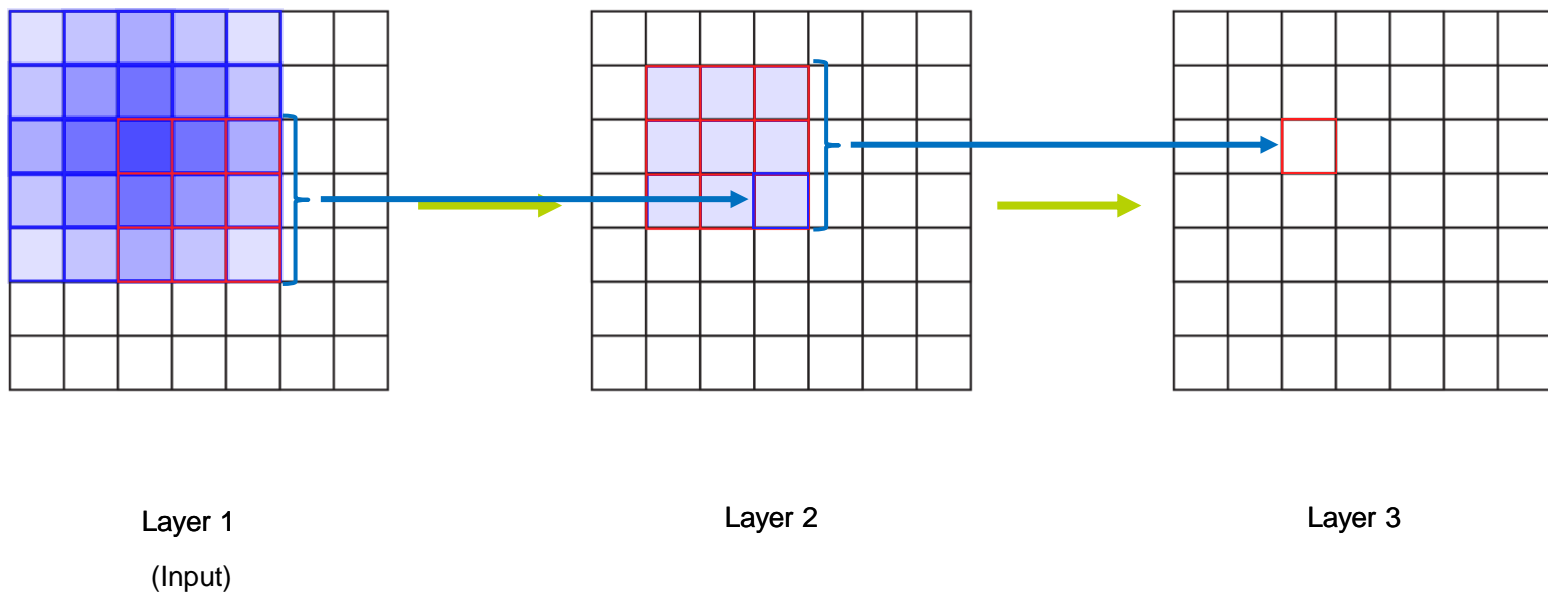
# VGG



# VGG



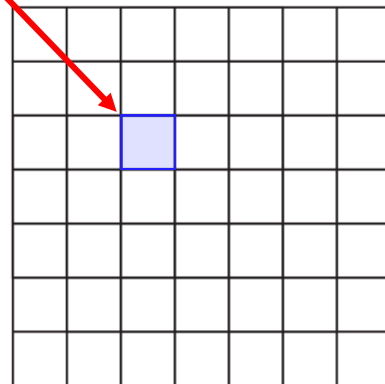
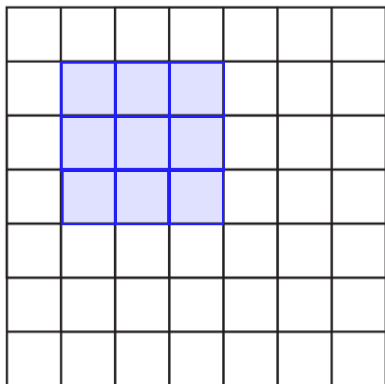
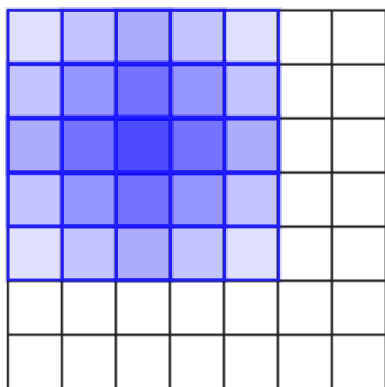
# VGG





# VGG

Each square in Layer 3 “sees” a 5x5 grid from Layer 1



Layer 1  
(Input)

Layer 2

Layer 3

# VGG

Two 3x3, stride 1 convolutions in a row  $\rightarrow$  one 5x5

Three 3x3 convolutions  $\rightarrow$  one 7x7 convolution

Benefit: fewer parameters

One 3x3 layer

$$3 \times 3 \times C = 9C$$

One 7x7 layer

$$7 \times 7 \times C = 49C$$

Three 3x3 layers

$$3 \times (9C) = 27C$$

$$49C \rightarrow 27C \rightarrow \approx 45\% \text{ reduction!}$$

# VGG

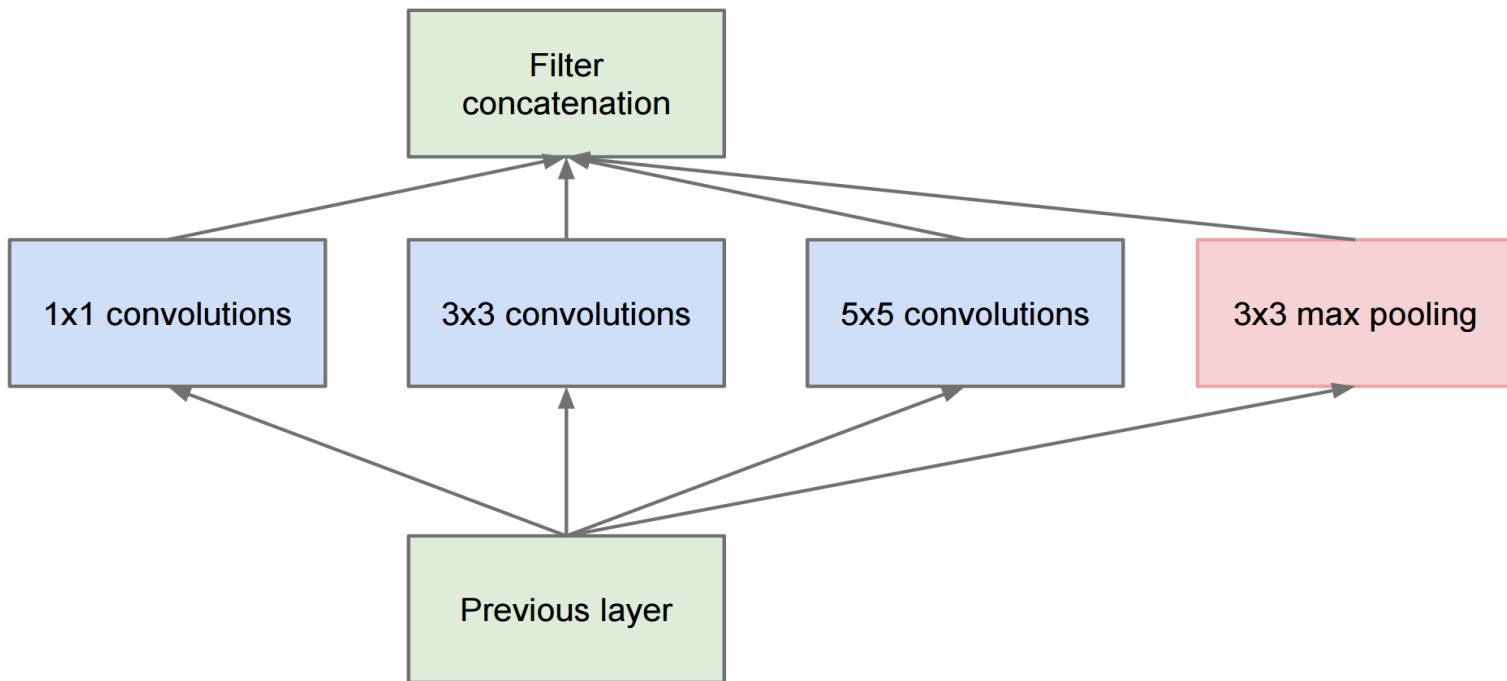
- One of the first architectures to experiment with many layers (More is better!)
- Can use multiple 3x3 convolutions to simulate larger kernels with fewer parameters
- VGGNet has a total of 138 million parameters, which is more than twice that of AlexNet
- Often used as a starting point for transfer learning because of its simple architecture
- Served as "base model" for future works

# Inception

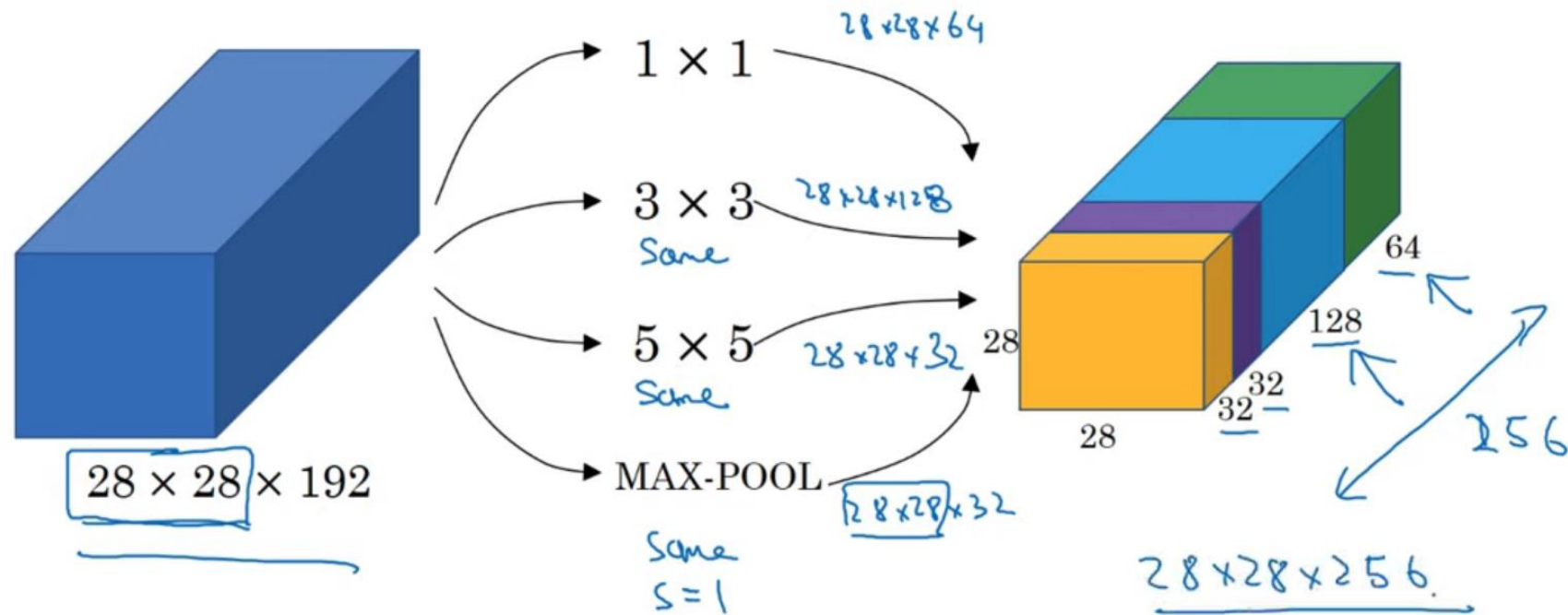
- Szegedy et al 2014
- Inception architecture-based [GoogLeNet](#), was the winner of the 2014 ImageNet challenge, with a 6.7% top five error rate
- Idea: network would want to use different receptive fields
- Want computational efficiency
- Solution: Turn each layer into branches of convolutions
- Each branch handles smaller portion of workload
- Concatenate different branches at the end

# Inception

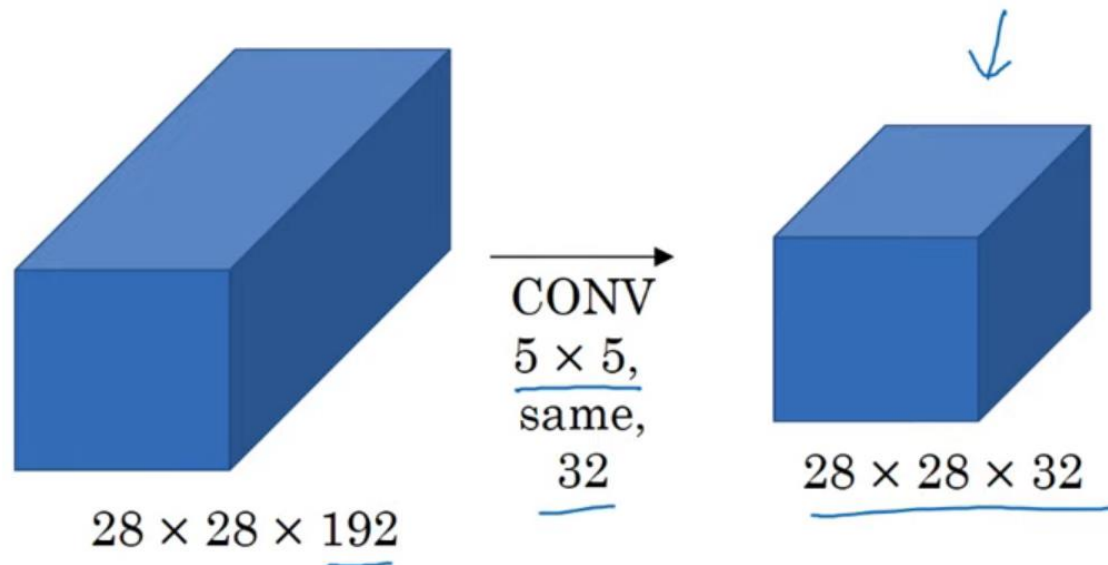
Basic idea: replace single 3x3 convolutions with module



# Motivation for inception network



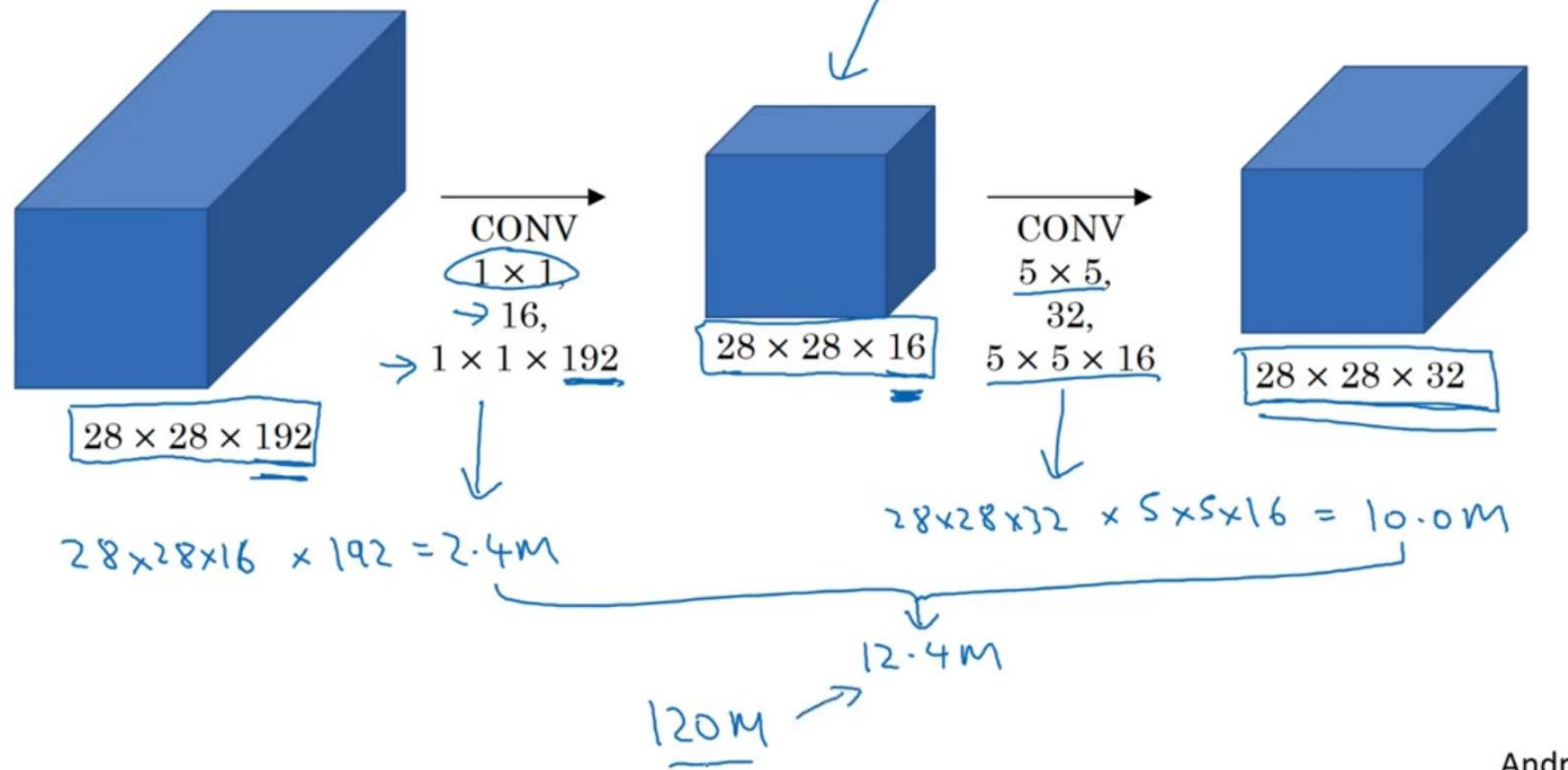
# The problem of computational cost



32 filters. filters are  $5 \times 5 \times 192$  multiplications

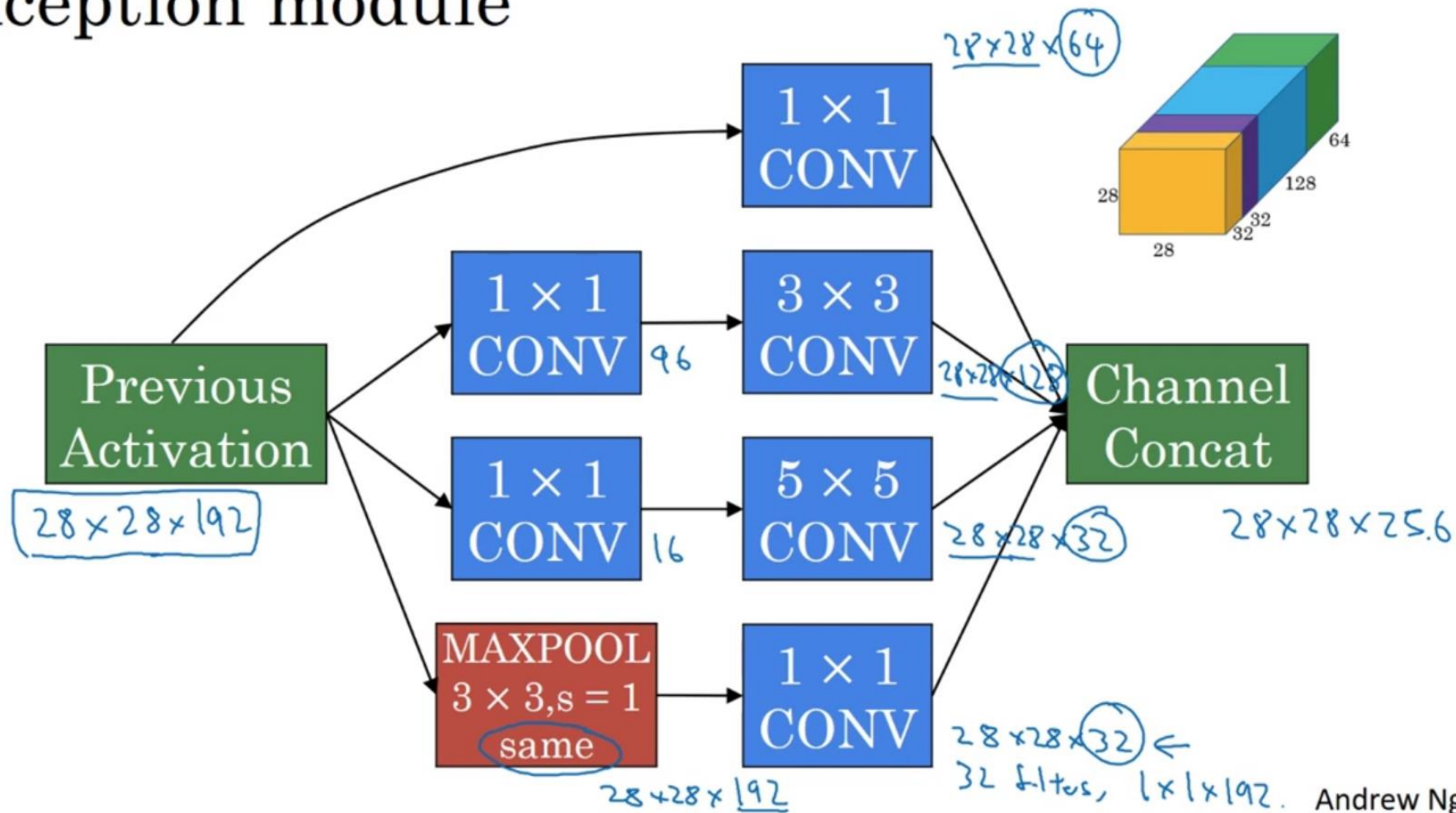
$$\underline{28 \times 28 \times 32} \times \underline{5 \times 5 \times 192} = 120M.$$

# Using 1x1 convolution

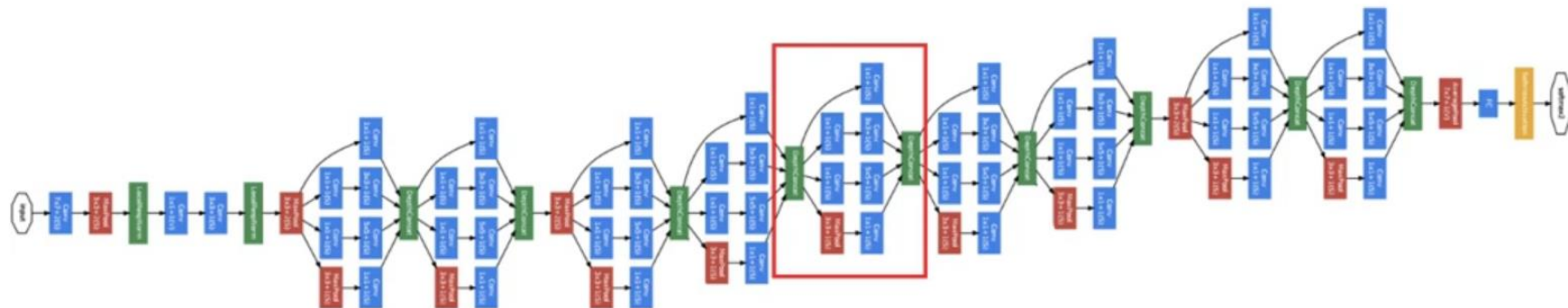


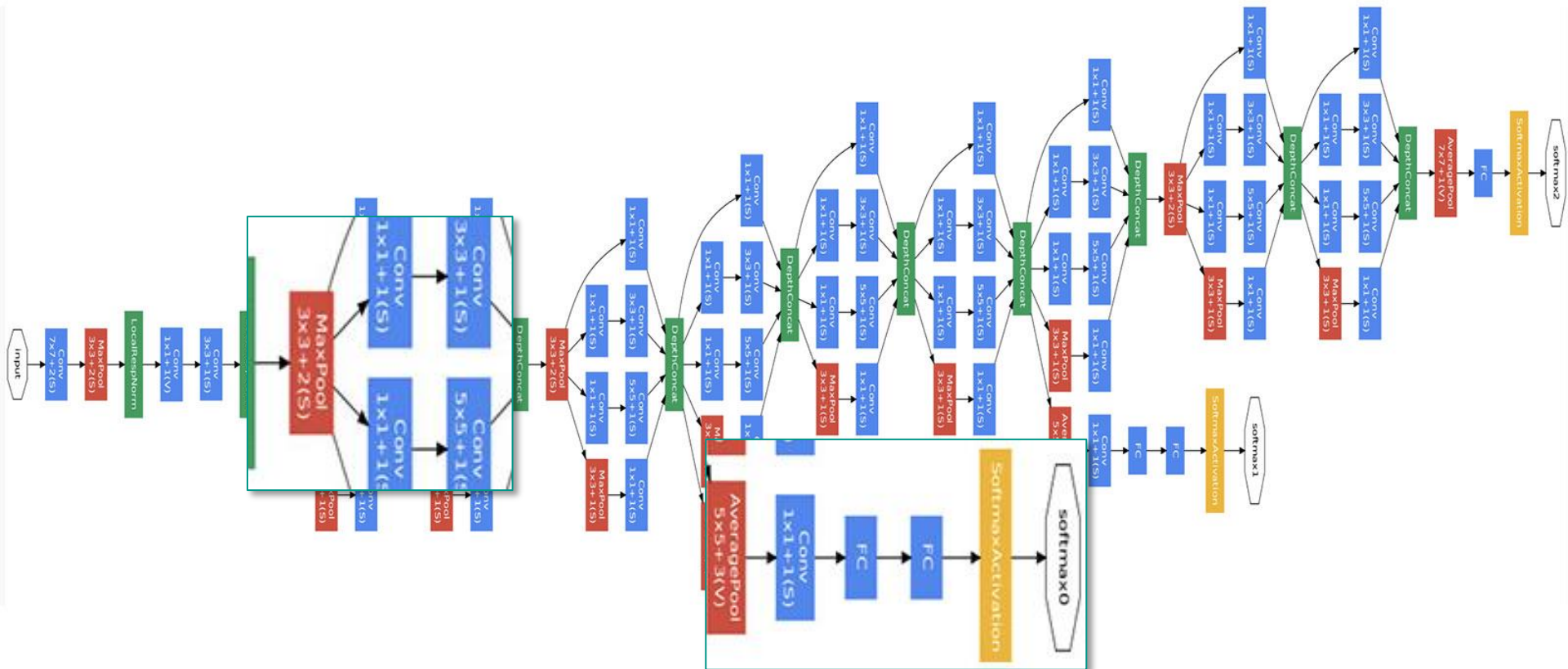


# Inception module



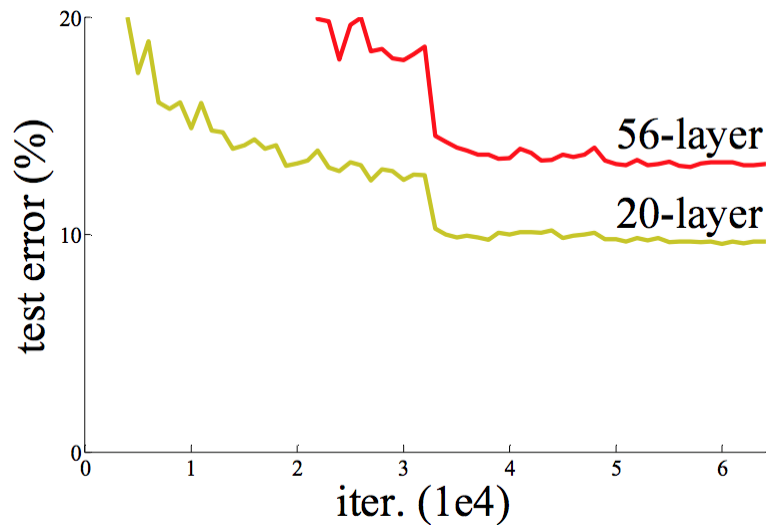
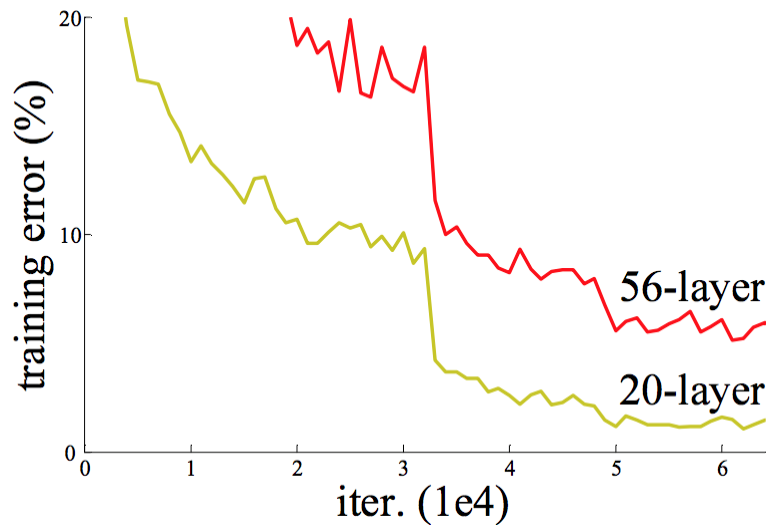
# Inception network





# ResNet - Motivation

Issue: Deeper Networks performing worse on **training** data!  
(as well as test data)

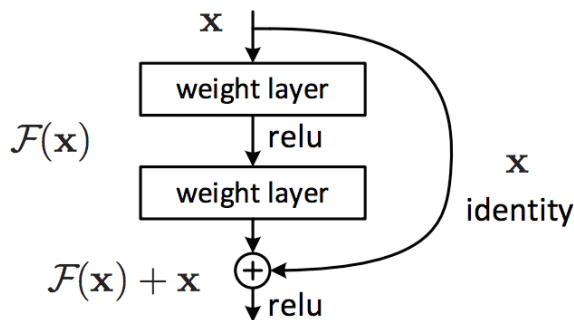


# ResNet

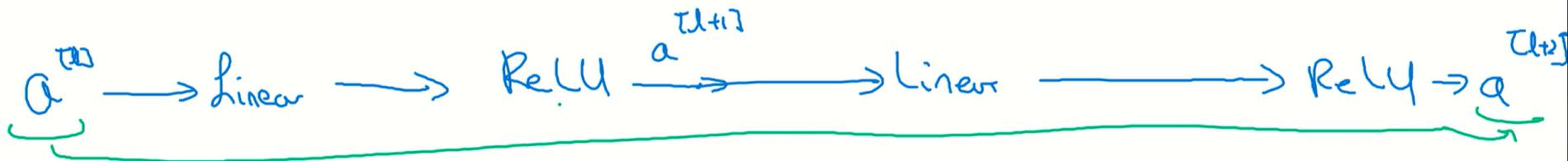
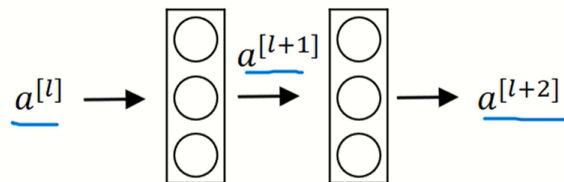
- Surprising because deeper networks should overfit more.
- So what's happening?
- Early layers of Deep Networks are very slow to adjust.
- Analogous to “Vanishing Gradient” issue.
- In theory, should be able to just have an “identity” transformation that makes the deeper network behave like a shallower one

# ResNet

- Assumption: best transformation over multiple layers is close to  $\mathcal{F}(x) + x$
- $x \rightarrow$  input to series of layers
- $\mathcal{F}(x) \rightarrow$  function represented by several layers (such as convs)
- Enforce this by adding “shortcut connections”
- Add the inputs from an earlier layer to the output of current layer



# Residual block

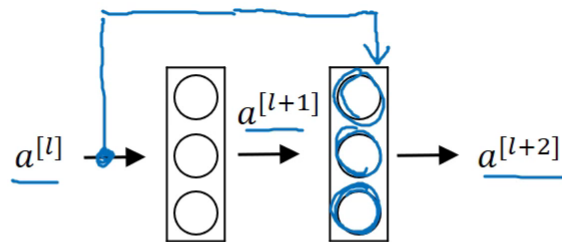


$$\underline{z^{[l+1]} = W^{[l+1]} \underline{a^{[l]} + b^{[l+1]}} \quad \underline{a^{[l+1]} = g(z^{[l+1]})} \quad \underline{z^{[l+2]} = W^{[l+2]} a^{[l+1]} + b^{[l+2]}} \quad \underline{a^{[l+2]} = g(z^{[l+2]})}$$

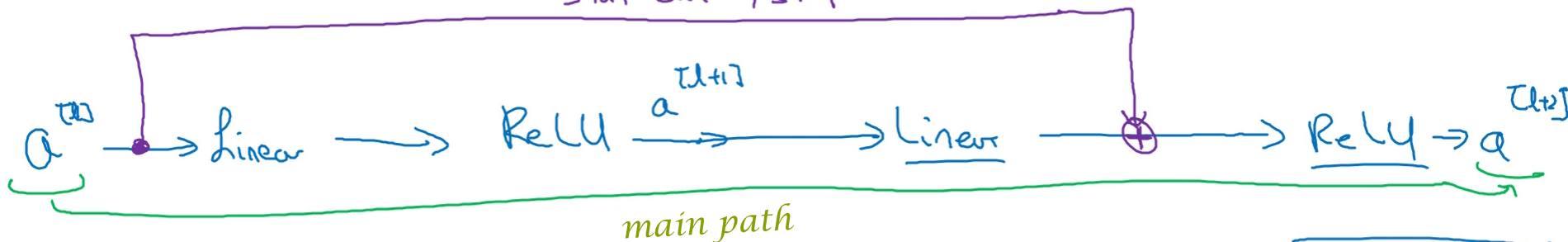
$\uparrow$                        $\uparrow$

---

# Residual block



"short cut" / skip connection



$$z^{[l+1]} = W^{[l+1]} a^{[l]} + b^{[l+1]}$$

$$a^{[l+1]} = g(z^{[l+1]})$$

$$z^{[l+2]} = W^{[l+2]} a^{[l+1]} + b^{[l+2]}$$

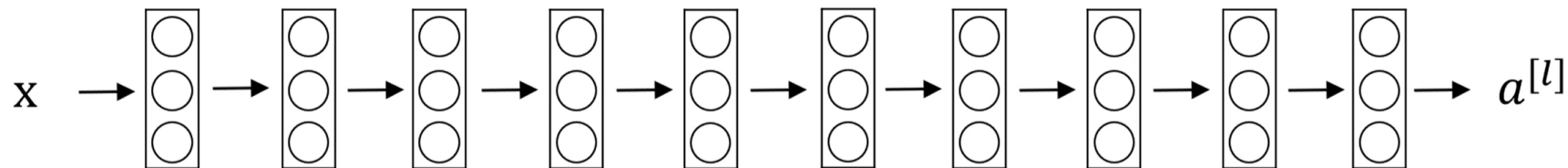
~~$$a^{[l+2]} = g(z^{[l+2]})$$~~

$$a^{[l+1]} = g(z^{[l+2]} + \underbrace{a^{[l]}})$$



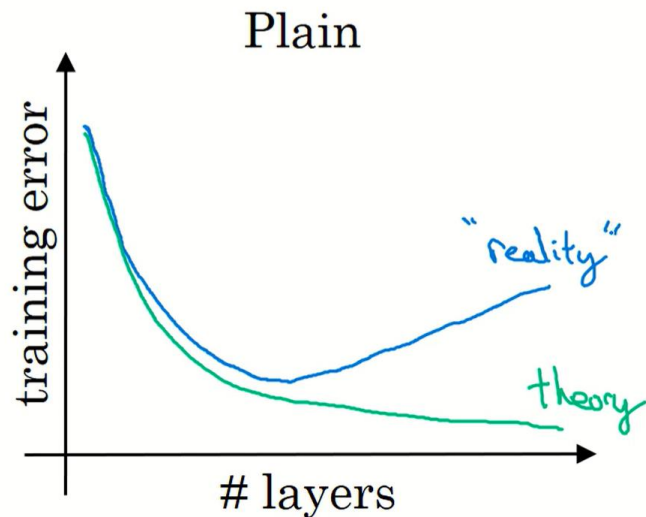
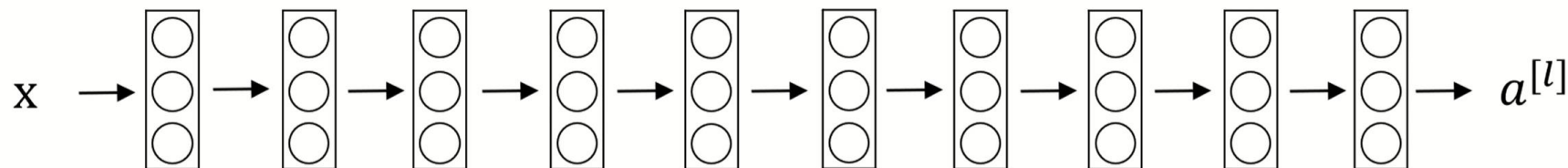
# Residual Network

*plain network*



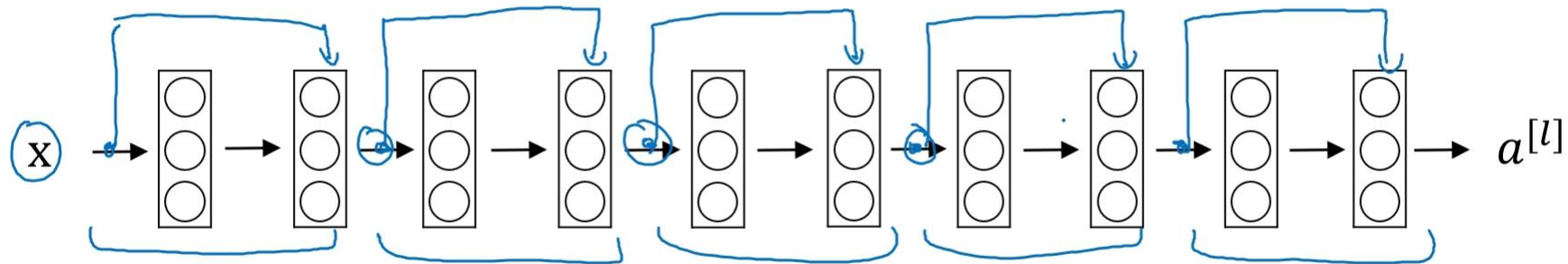
# Residual Network

*plain network*

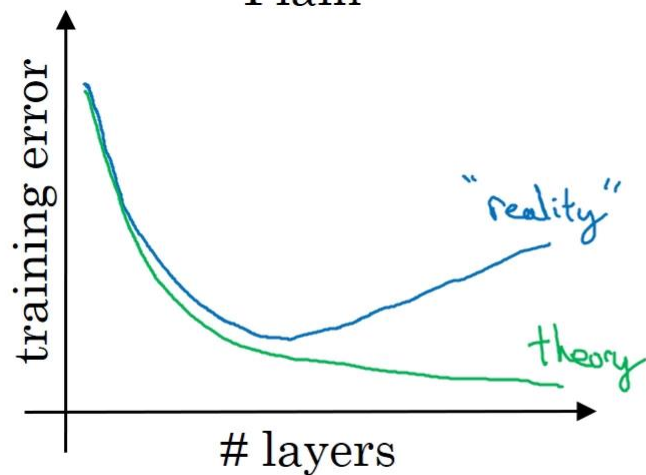


# Residual Network

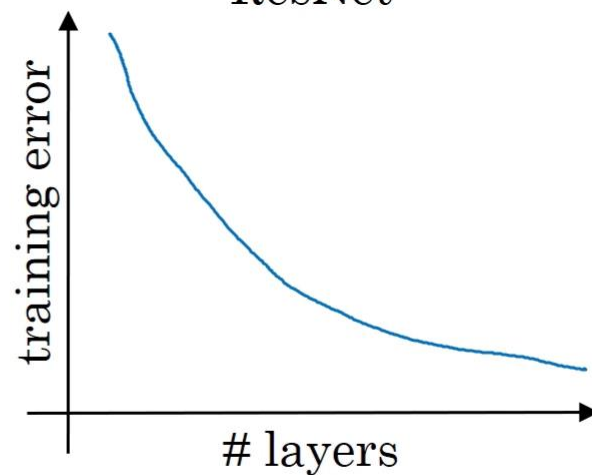
*with residual blocks*



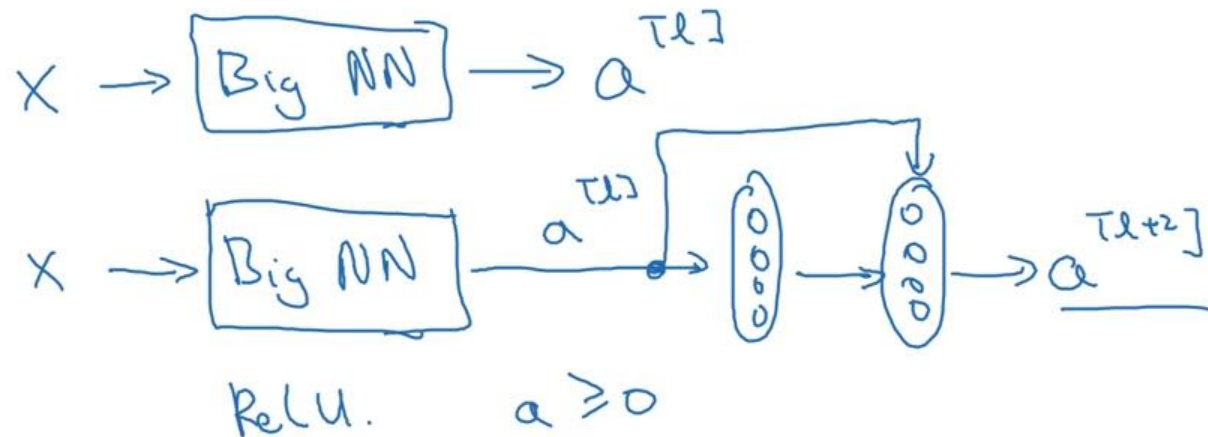
Plain



ResNet

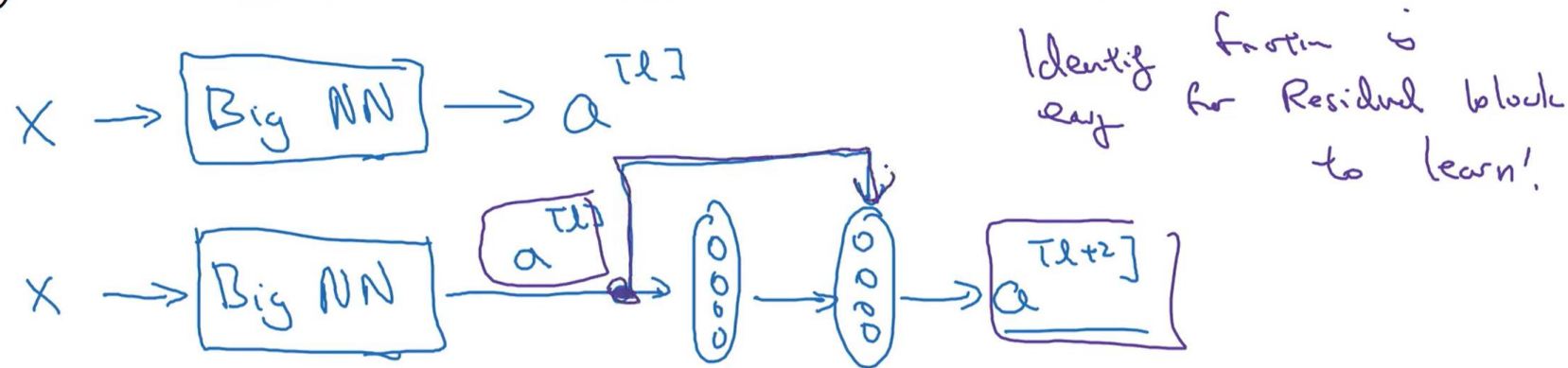


# Why do residual networks work?



$$\begin{aligned} a^{[l+2]} &= g\left(\underline{z^{[l+2]}} + \underline{a^{[l]}}\right) \\ &= g\left(\underline{w^{[l+2]} a^{[l+1]} + b^{[l+2]}}\right) + \underline{a^{[l]}} \end{aligned}$$

# Why do residual networks work?



ReLU.       $a \geq 0$

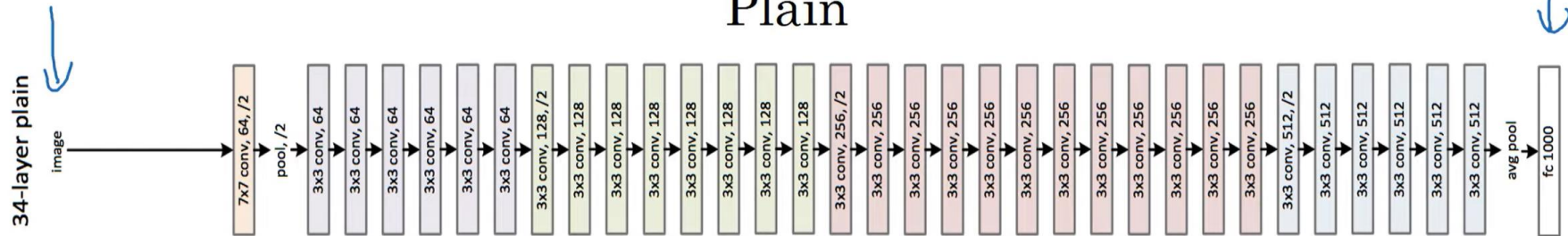
$$a^{[L+2]} = g\left(\underline{z^{[L+2]}} + \underline{a^{[L]}}\right)$$

$$= g\left(\cancel{w^{[L+2]} a^{[L+1]} + b^{[L+2]}} + a^{[L]}\right) = g(a^{[L]}) = \underline{a^{[L]}}$$

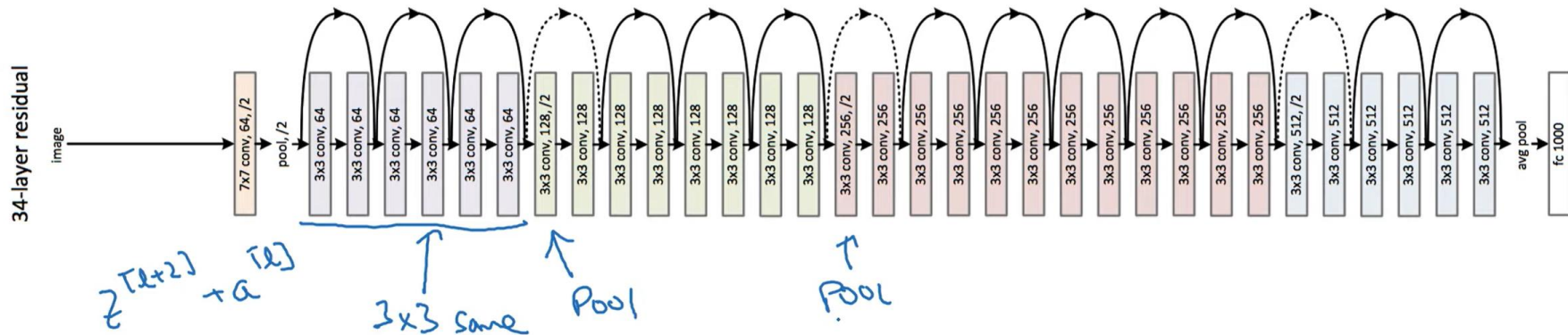
If  $w^{[L+2]} = 0, b^{[L+2]} = 0$

# ResNet

## Plain

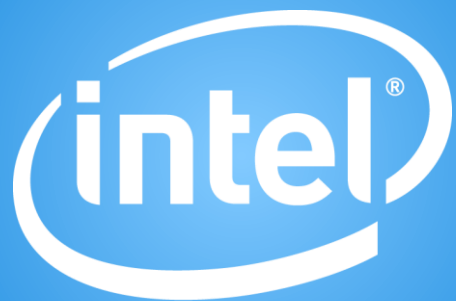


## ResNet



# ResNet

- ResNet uses 152 layers and achieved 3.6% top five error to win the 2015 ImageNet competition.
- This exceeds average human ability.
- ResNet was faster to train than VGGNet, which has only 20 layers.



Software