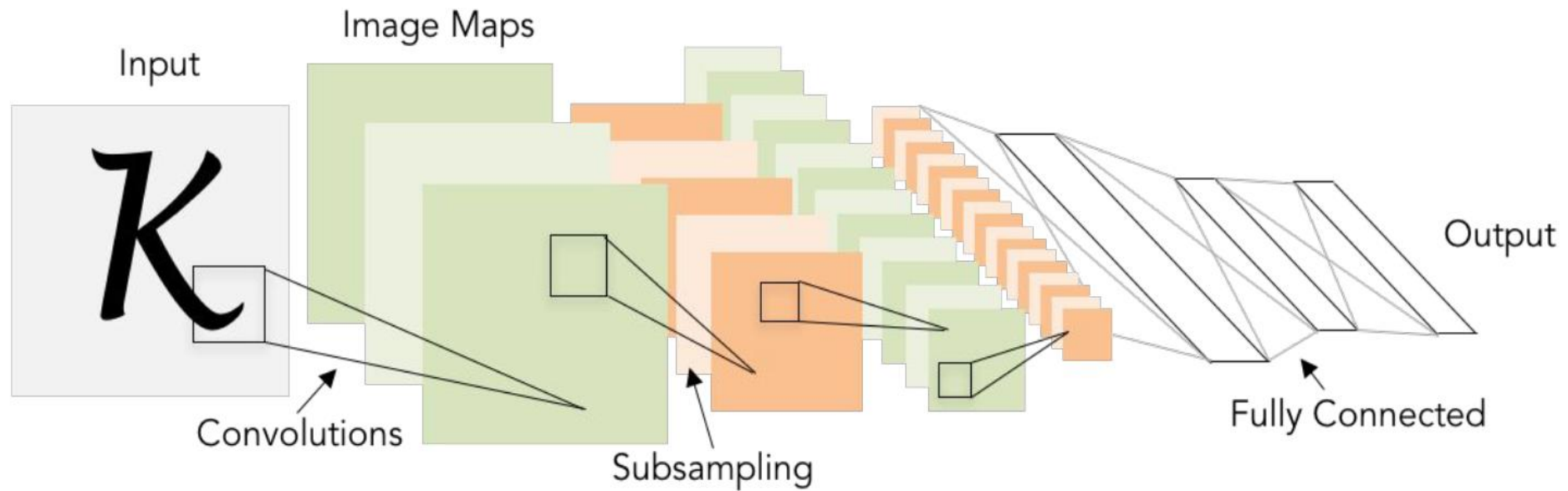




# CS231n Lecture 2

- ☑ BOAZ 10기 박성현
- ☑ BOAZ 11기 김태희
- ☑ BOAZ 11기 홍지민

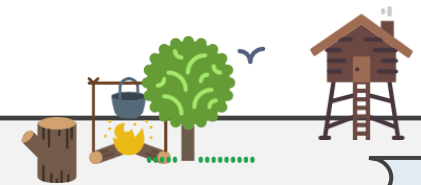
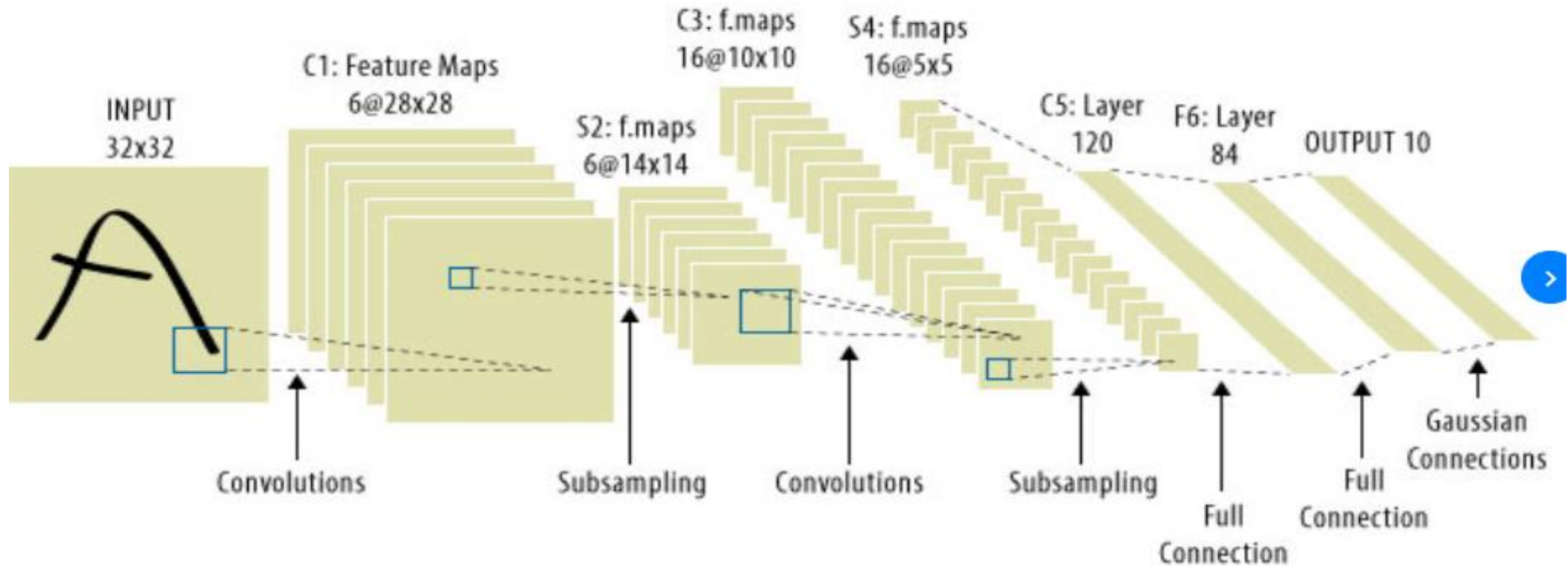
[LeCun et al., 1998]



Conv filters were 5x5, applied at stride 1  
Subsampling (Pooling) layers were 2x2 applied at stride 2  
i.e. architecture is [CONV-POOL-CONV-POOL-FC-FC]



# More details in LeNet





MAX POOL1

## NORM1

CONV2

MAX POOL2

## NORM2

CONV3

CONV4

CONV5

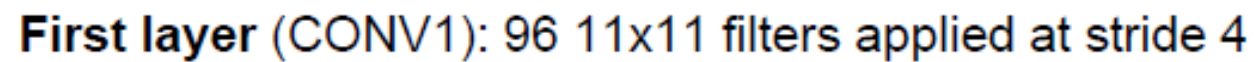
Max POOL3

FC6

FC7

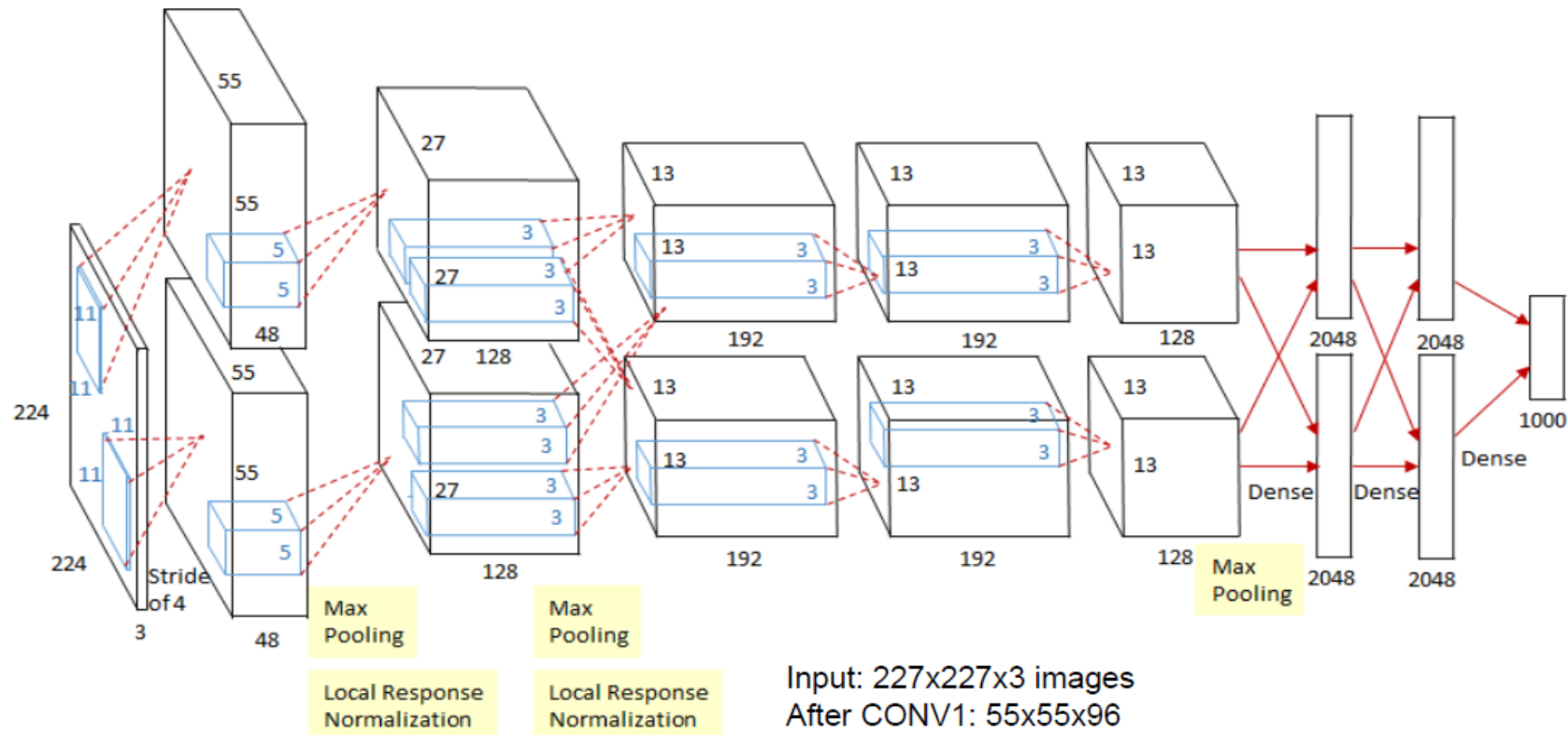
FC8





Q: what is the output volume size?

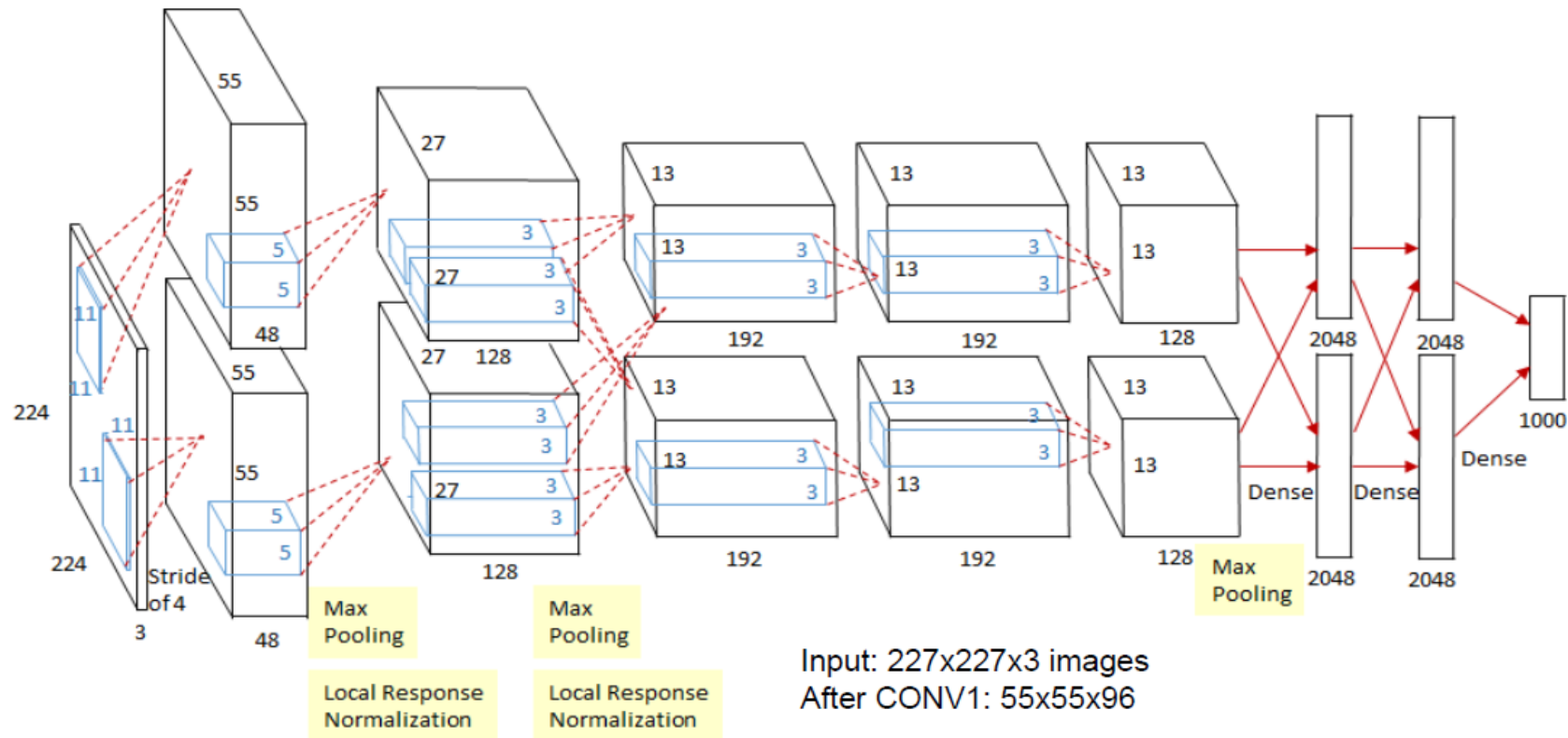




**Second layer (POOL1):** 3x3 filters applied at stride 2

Q: what is the output volume size? Hint:  $(55-3)/2+1 = 27$





Q) Difference btw convolution and pooling?





# Case Study: AlexNet

[Krizhevsky et al. 2012]

Full (simplified) AlexNet architecture:

[227x227x3] INPUT

[55x55x96] CONV1: 96 11x11 filters at stride 4, pad 0

[27x27x96] MAX POOL1: 3x3 filters at stride 2

[27x27x96] NORM1: Normalization layer

[27x27x256] CONV2: 256 5x5 filters at stride 1, pad 2

[13x13x256] MAX POOL2: 3x3 filters at stride 2

[13x13x256] NORM2: Normalization layer

[13x13x384] CONV3: 384 3x3 filters at stride 1, pad 1

[13x13x384] CONV4: 384 3x3 filters at stride 1, pad 1

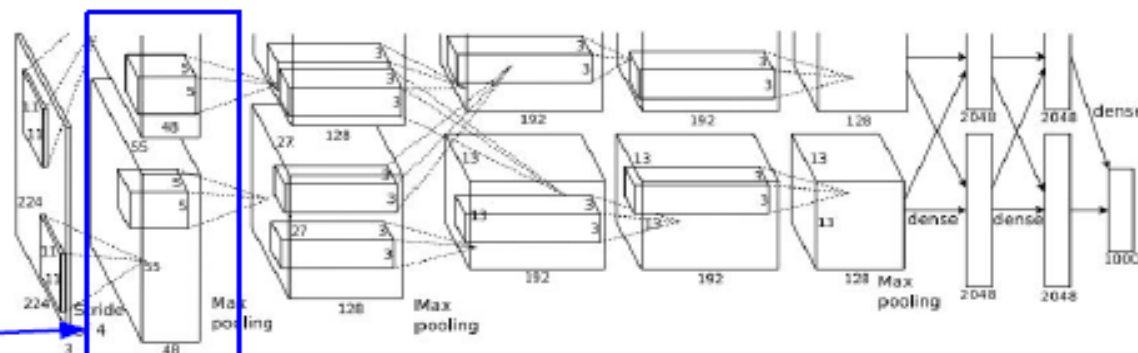
[13x13x256] CONV5: 256 3x3 filters at stride 1, pad 1

[6x6x256] MAX POOL3: 3x3 filters at stride 2

[4096] FC6: 4096 neurons

[4096] FC7: 4096 neurons

[1000] FC8: 1000 neurons (class scores)



[55x55x48] x 2

Historical note: Trained on GTX 580 GPU with only 3 GB of memory.

Network spread across 2 GPUs, half the neurons (feature maps) on each GPU.

Figure copyright Alex Krizhevsky, Ilya Sutskever, and Geoffrey Hinton, 2012. Reproduced with permission.





## Case Study: VGGNet

[Simonyan and Zisserman, 2014]

Small filters, Deeper networks

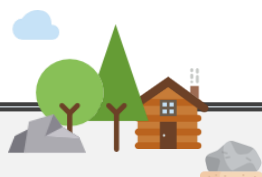
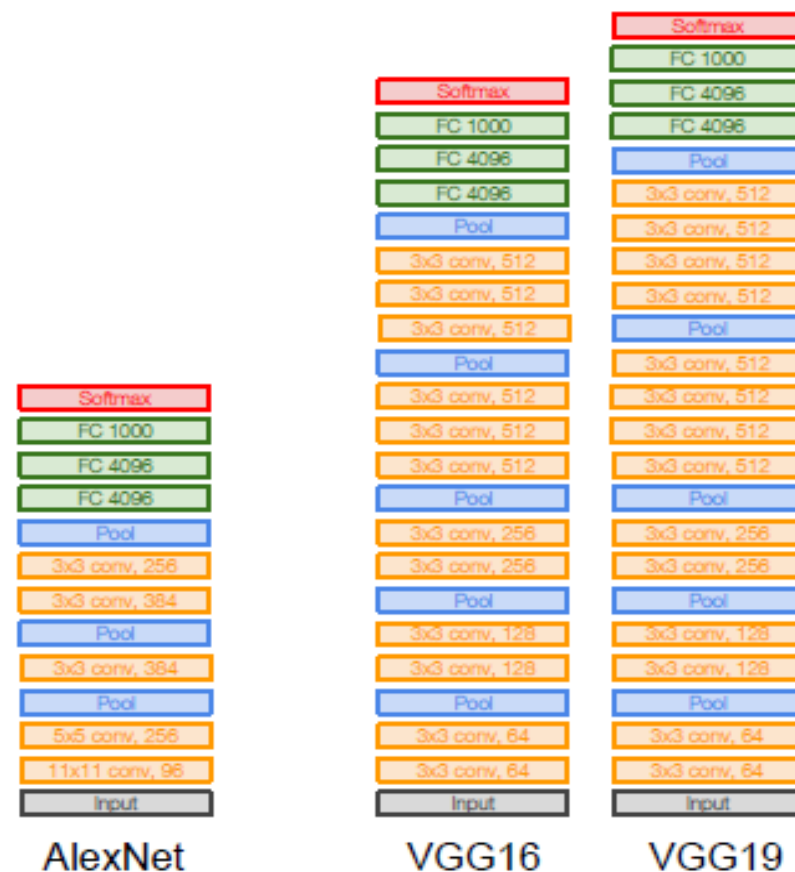
8 layers (AlexNet)

-> 16 - 19 layers (VGG16Net)

Only 3x3 CONV stride 1, pad 1  
and 2x2 MAX POOL stride 2

11.7% top 5 error in ILSVRC'13  
(ZFNet)

-> 7.3% top 5 error in ILSVRC'14



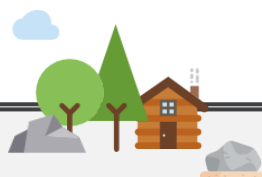
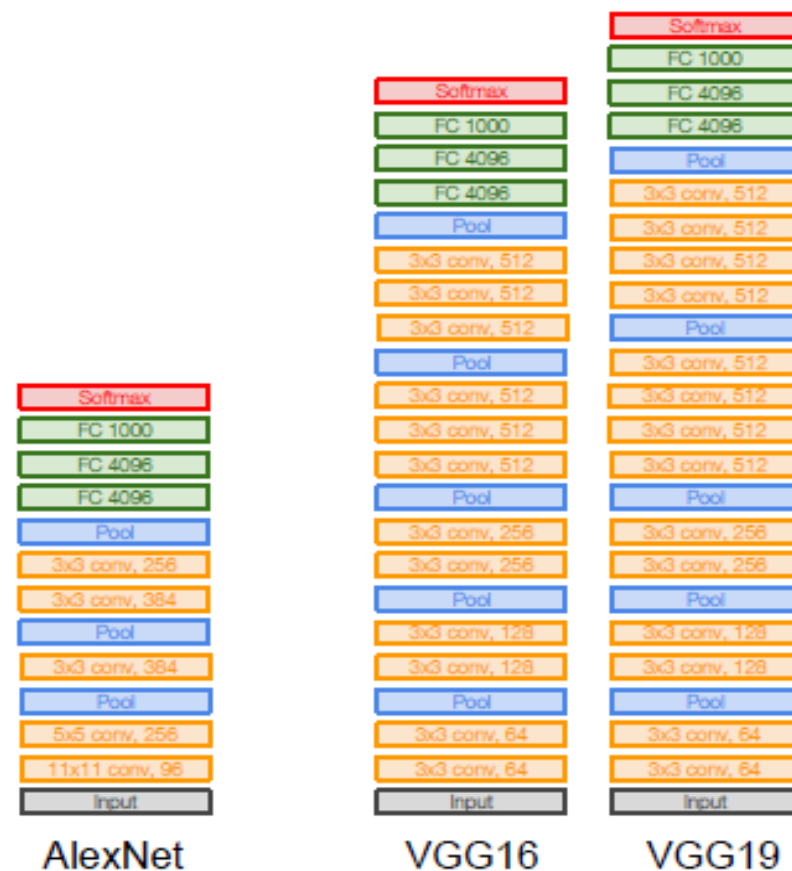
# Case Study: VGGNet

[Simonyan and Zisserman, 2014]

Q: Why use smaller filters? (3x3 conv)

Stack of three 3x3 conv (stride 1) layers has same **effective receptive field** as one 7x7 conv layer

Q: What is the effective receptive field of three 3x3 conv (stride 1) layers?



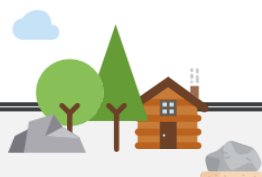
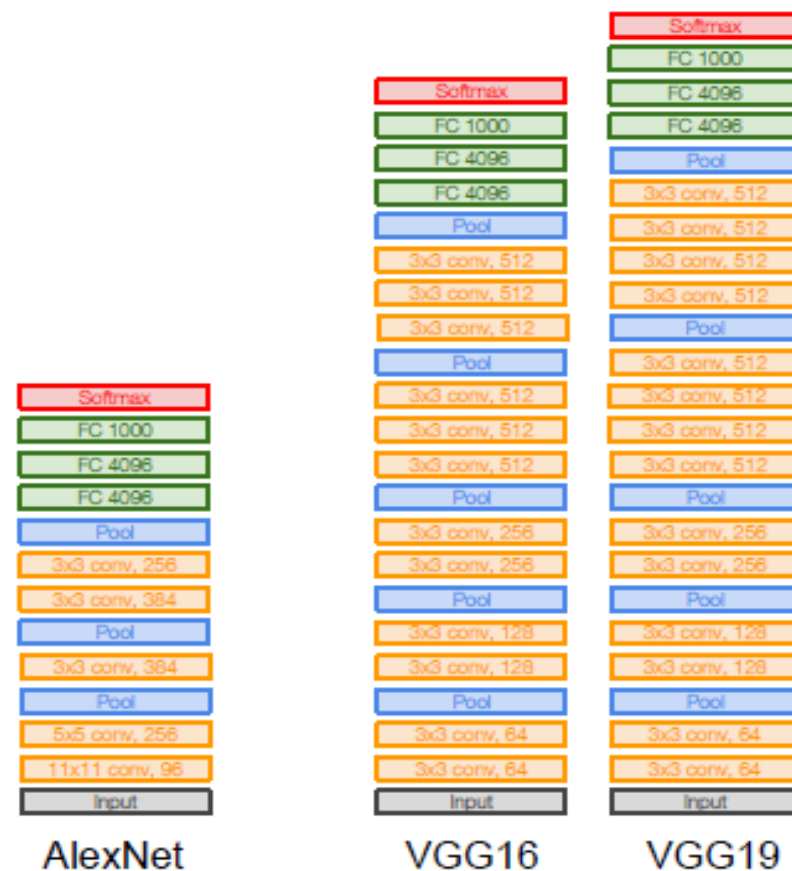
## Case Study: VGGNet

[Simonyan and Zisserman, 2014]

Q: Why use smaller filters? (3x3 conv)

Stack of three 3x3 conv (stride 1) layers has same **effective receptive field** as one 7x7 conv layer

Q: What is the effective receptive field of three 3x3 conv (stride 1) layers?



```

graph TD
    Input[Input] --> C1_1[3x3 conv, 64]
    C1_1 --> C1_2[3x3 conv, 64]
    C1_2 --> P1[Pool]
    P1 --> C2_1[3x3 conv, 128]
    C2_1 --> C2_2[3x3 conv, 128]
    C2_2 --> P2[Pool]
    P2 --> C3_1[3x3 conv, 256]
    C3_1 --> C3_2[3x3 conv, 256]
    C3_2 --> P3[Pool]
    P3 --> C4_1[3x3 conv, 512]
    C4_1 --> C4_2[3x3 conv, 512]
    C4_2 --> P4[Pool]
    P4 --> C5_1[3x3 conv, 512]
    C5_1 --> C5_2[3x3 conv, 512]
    C5_2 --> FC1[FC 4096]
    FC1 --> FC2[FC 4096]
    FC2 --> FC3[FC 4096]
    FC3 --> FC4[FC 1000]
    FC4 --> Softmax[Softmax]
  
```

## VGG16

TOTAL params: 138M parameters

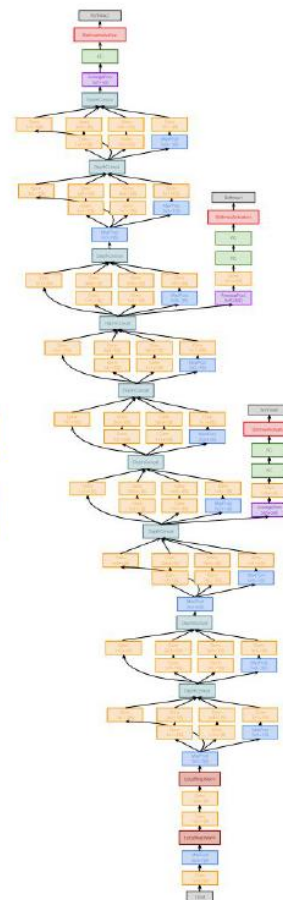
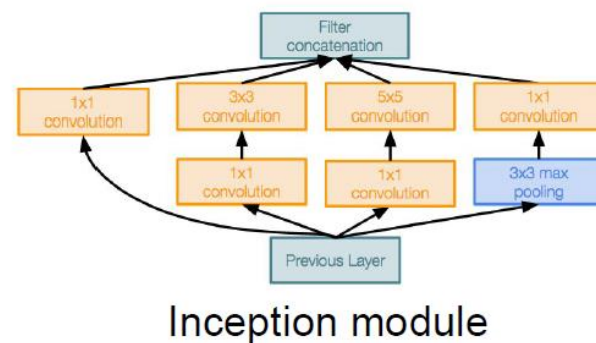


## Case Study: GoogLeNet

[Szegedy et al., 2014]

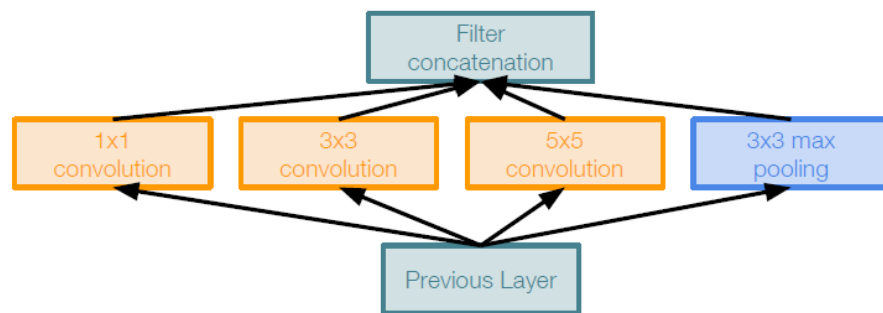
Deeper networks, with computational efficiency

- 22 layers
- Efficient “Inception” module
- No FC layers
- Only 5 million parameters!  
12x less than AlexNet
- ILSVRC’14 classification winner  
(6.7% top 5 error)



## Case Study: GoogLeNet

[Szegedy et al., 2014]



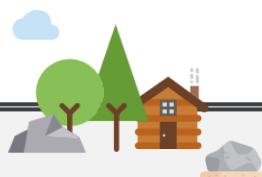
Naive Inception module

Apply parallel filter operations on the input from previous layer:

- Multiple receptive field sizes for convolution (1x1, 3x3, 5x5)
- Pooling operation (3x3)

Concatenate all filter outputs together depth-wise

Q: What is the problem with this?  
[Hint: Computational complexity]

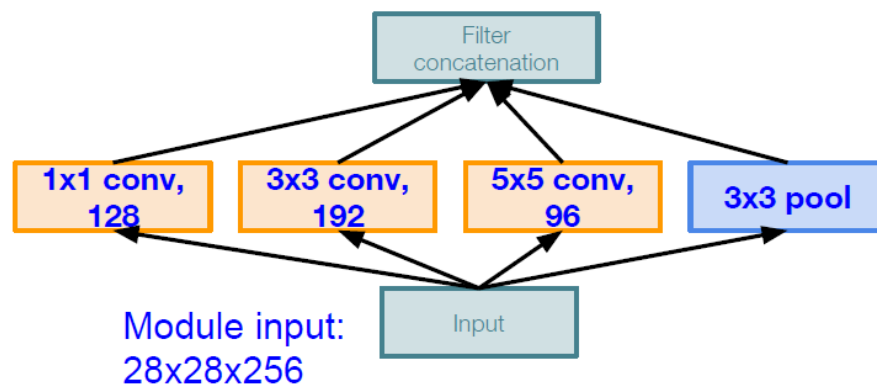




## Case Study: GoogLeNet

[Szegedy et al., 2014]

Example:



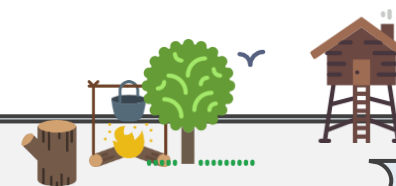
Naive Inception module

Q: What is the problem with this?  
[Hint: Computational complexity]

Q1: What is the output size of the  
1x1 conv, with 128 filters?

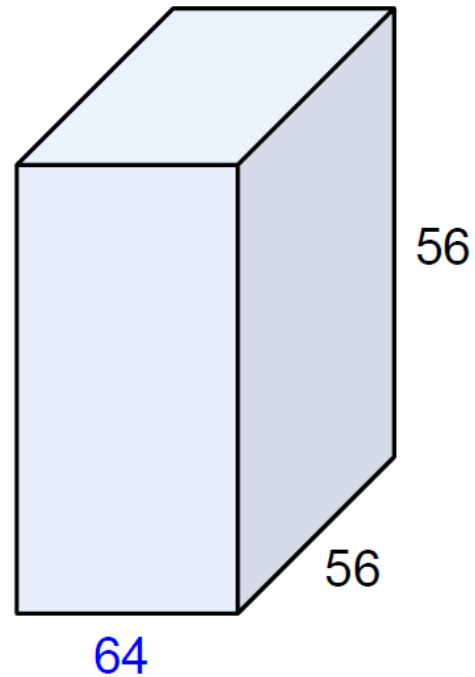
Q2: What are the output sizes of  
all different filter operations?

Q3: What is output size after  
filter concatenation?





## Reminder: 1x1 convolutions

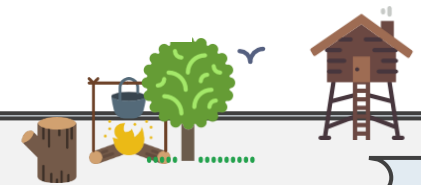
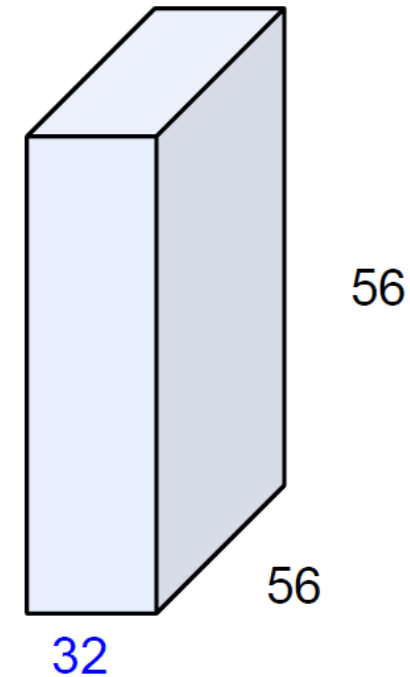


1x1 CONV  
with 32 filters



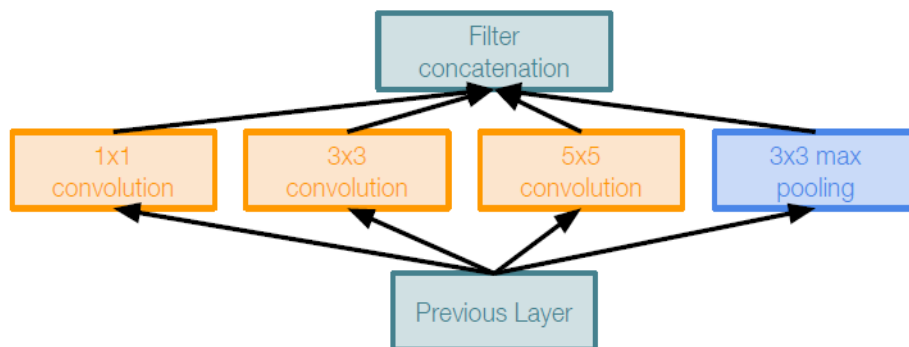
preserves spatial  
dimensions, reduces depth!

Projects depth to lower  
dimension (combination of  
feature maps)



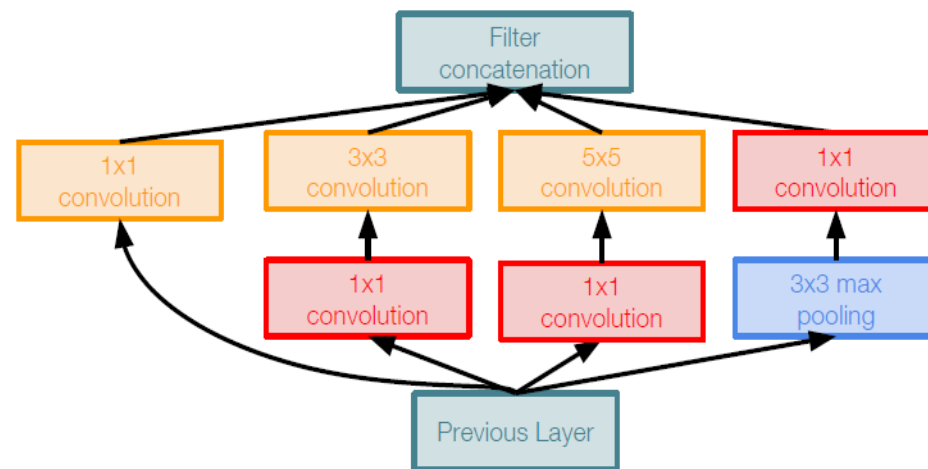
# Case Study: GoogLeNet

[Szegedy et al., 2014]



Naive Inception module

1x1 conv “bottleneck”  
layers



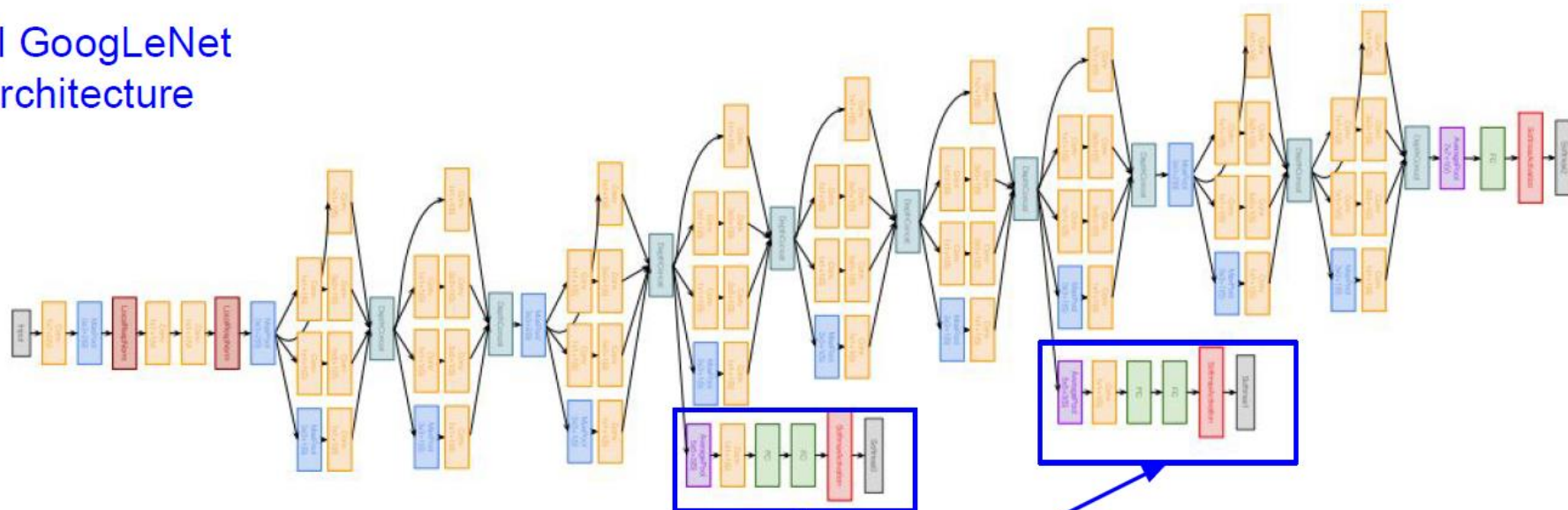
Inception module with dimension reduction



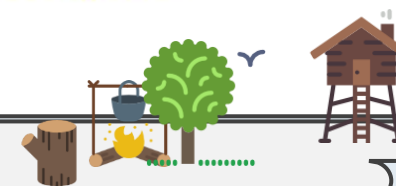
# Case Study: GoogLeNet

[Szegedy et al., 2014]

Full GoogLeNet  
architecture



Auxiliary classification outputs to inject additional gradient at lower layers  
(AvgPool-1x1Conv-FC-FC-Softmax)

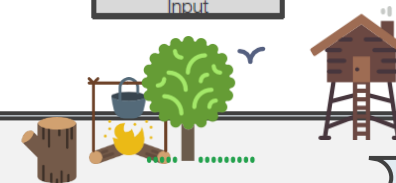
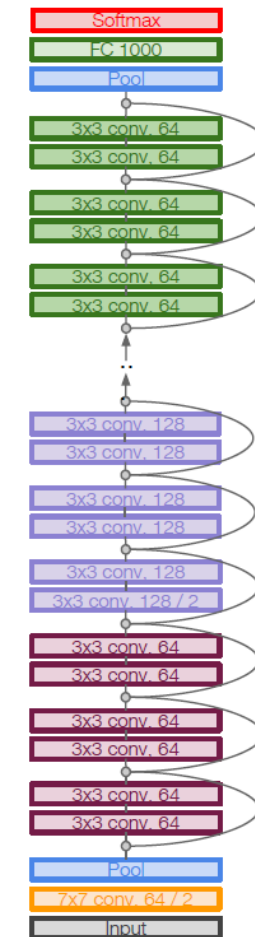
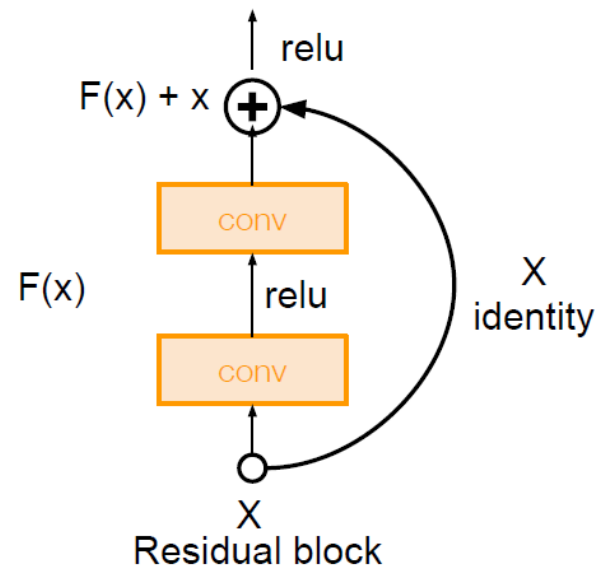


## Case Study: ResNet

[He et al., 2015]

Very deep networks using residual connections

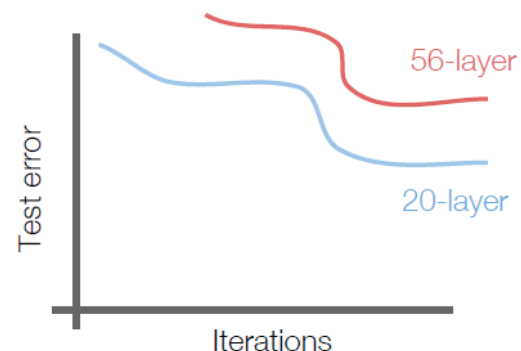
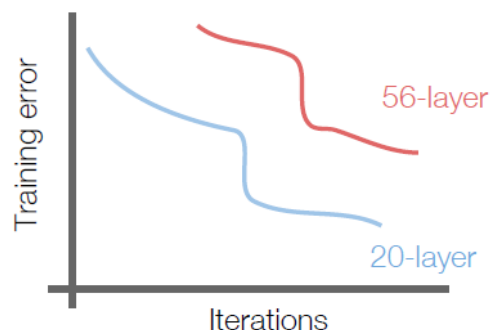
- 152-layer model for ImageNet
- ILSVRC'15 classification winner (3.57% top 5 error)
- Swept all classification and detection competitions in ILSVRC'15 and COCO'15!



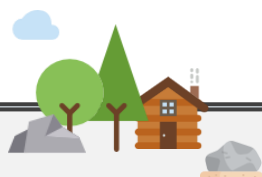
## Case Study: ResNet

[He et al., 2015]

What happens when we continue stacking deeper layers on a “plain” convolutional neural network?



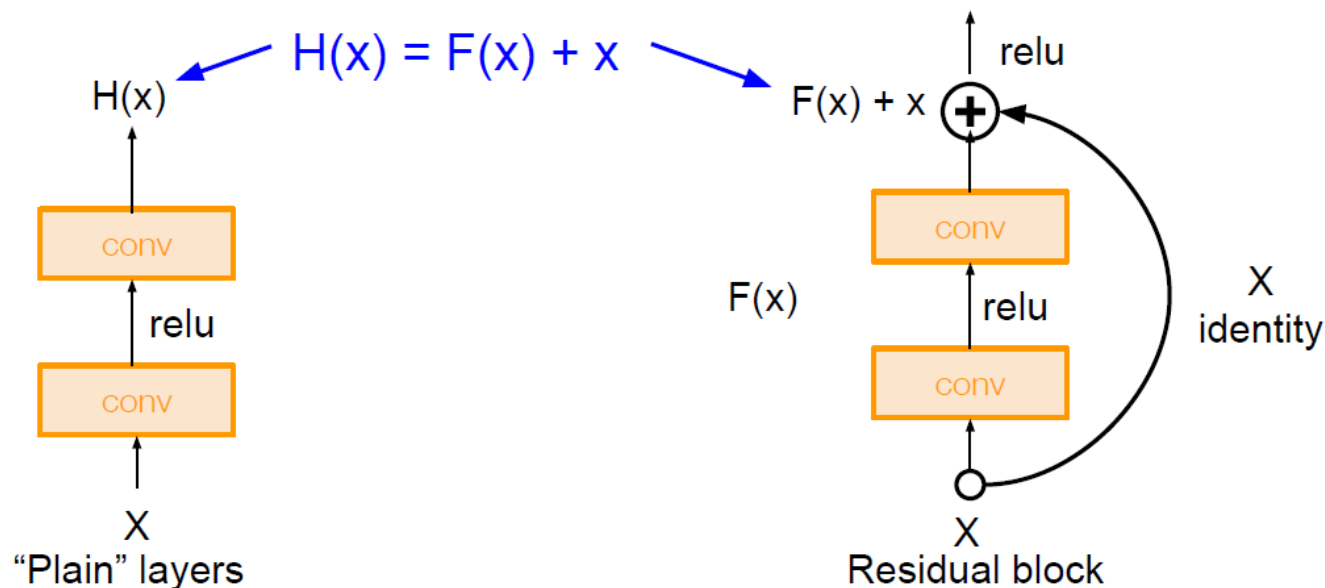
56-layer model performs worse on both training and test error  
-> The deeper model performs worse, but it's not caused by overfitting!



# Case Study: ResNet

[He et al., 2015]

Solution: Use network layers to fit a residual mapping instead of directly trying to fit a desired underlying mapping



Use layers to fit residual  $F(x) = H(x) - x$  instead of  $H(x)$  directly

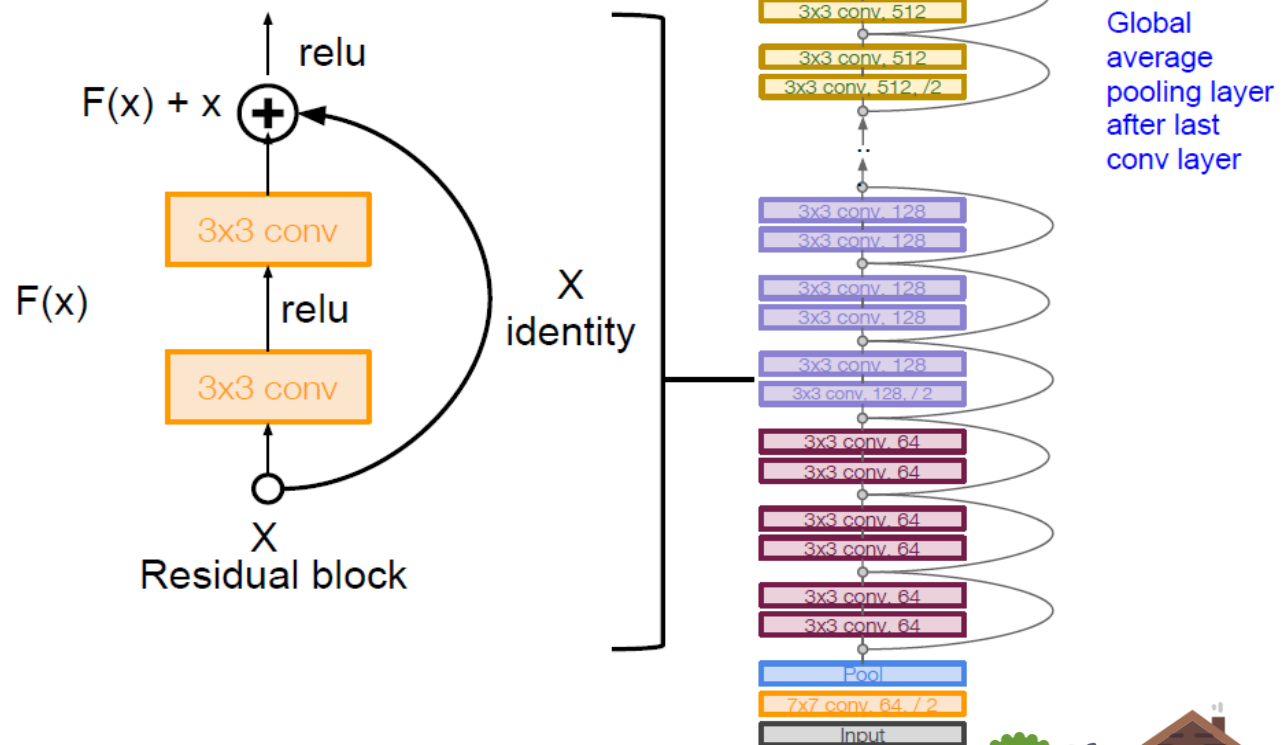


# Case Study: ResNet

[He et al., 2015]

Full ResNet architecture:

- Stack residual blocks
- Every residual block has two 3x3 conv layers
- Periodically, double # of filters and downsample spatially using stride 2 (/2 in each dimension)
- Additional conv layer at the beginning
- No FC layers at the end (only FC 1000 to output classes)

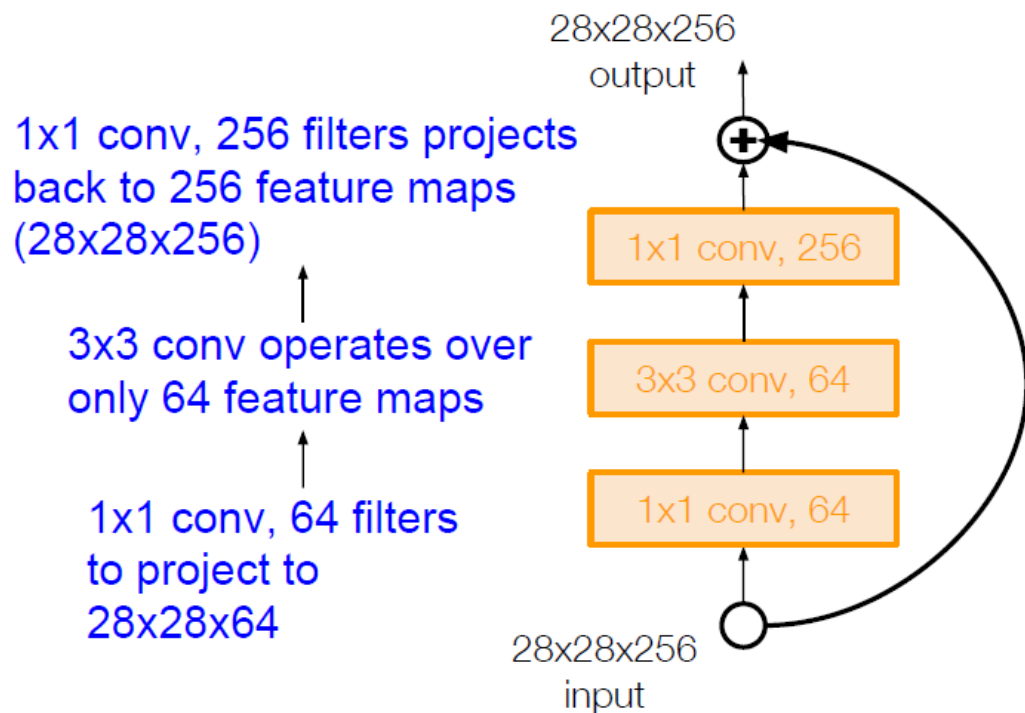




# Case Study: ResNet

[He et al., 2015]

For deeper networks  
(ResNet-50+), use “bottleneck”  
layer to improve efficiency  
(similar to GoogLeNet)

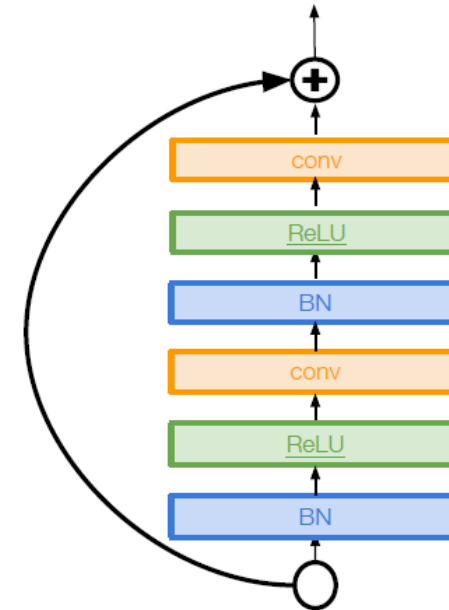


## Improving ResNets...

# Identity Mappings in Deep Residual Networks

[He et al. 2016]

- Improved ResNet block design from creators of ResNet
- Creates a more direct path for propagating information throughout network (moves activation to residual mapping pathway)
- Gives better performance

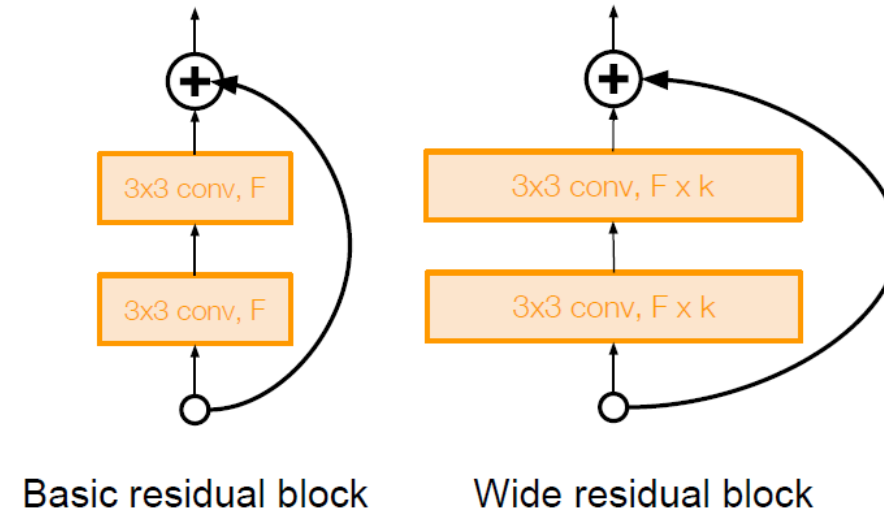


Improving ResNets...

## Wide Residual Networks

[Zagoruyko et al. 2016]

- Argues that residuals are the important factor, not depth
- User wider residual blocks ( $F \times k$  filters instead of  $F$  filters in each layer)
- 50-layer wide ResNet outperforms 152-layer original ResNet
- Increasing width instead of depth more computationally efficient (parallelizable)

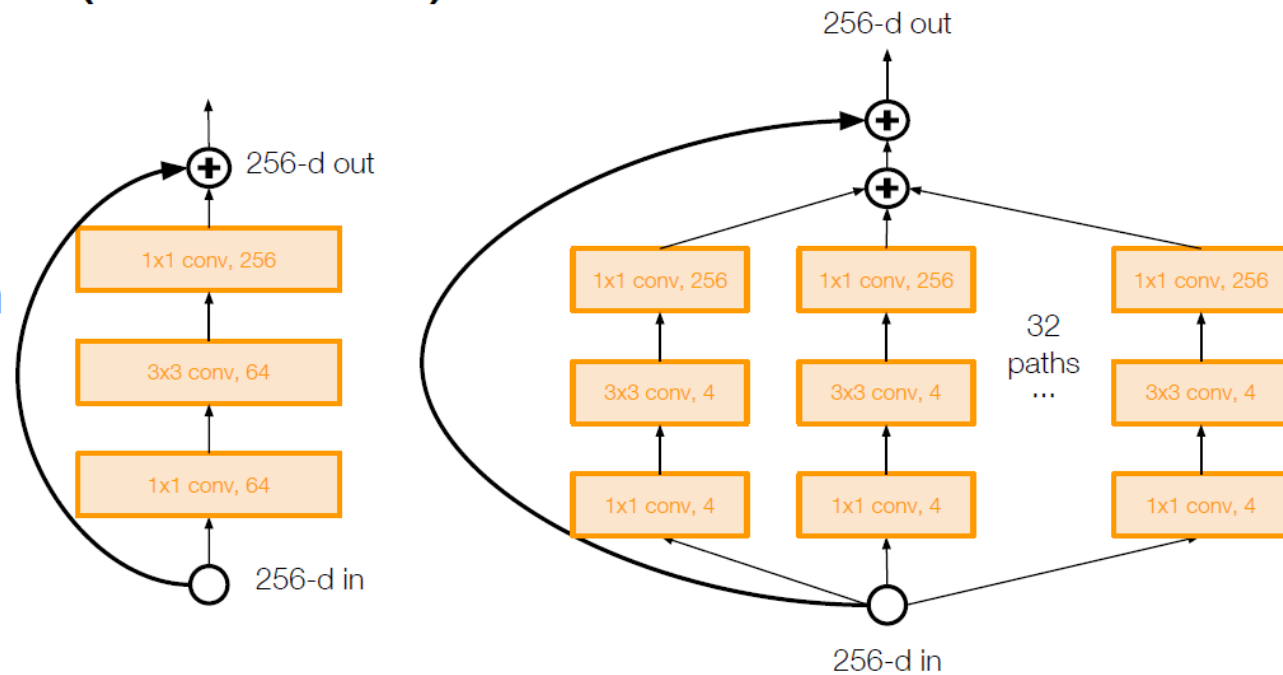


Improving ResNets...

## Aggregated Residual Transformations for Deep Neural Networks (ResNeXt)

[Xie et al. 2016]

- Also from creators of ResNet
- Increases width of residual block through multiple parallel pathways (“cardinality”)
- Parallel pathways similar in spirit to Inception module

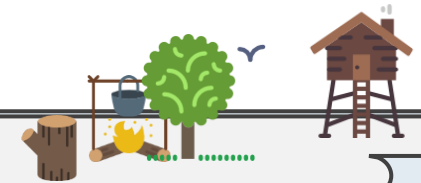
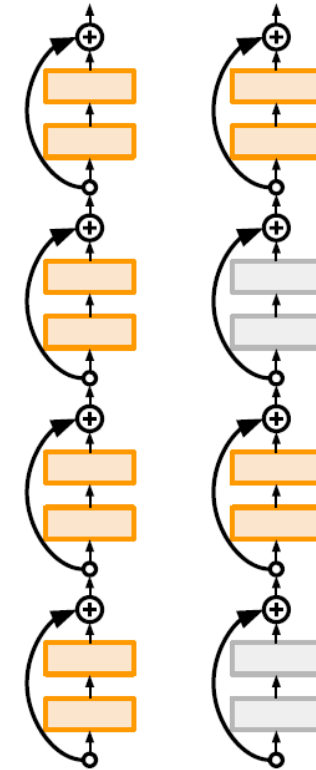


Improving ResNets...

## Deep Networks with Stochastic Depth

[Huang et al. 2016]

- Motivation: reduce vanishing gradients and training time through short networks during training
- Randomly drop a subset of layers during each training pass
- Bypass with identity function
- Use full deep network at test time

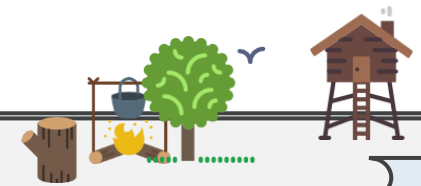
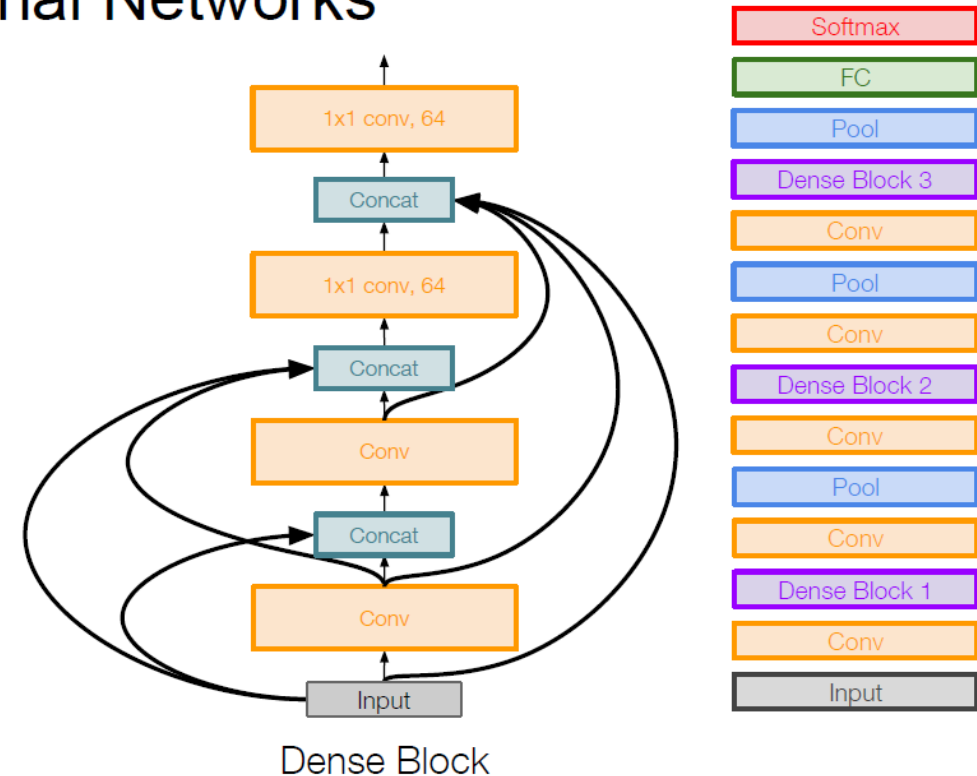


## Beyond ResNets...

### Densely Connected Convolutional Networks

[Huang et al. 2017]

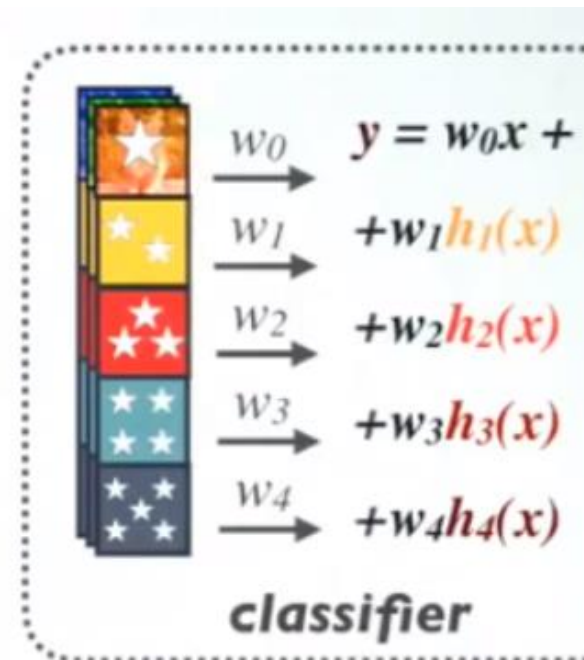
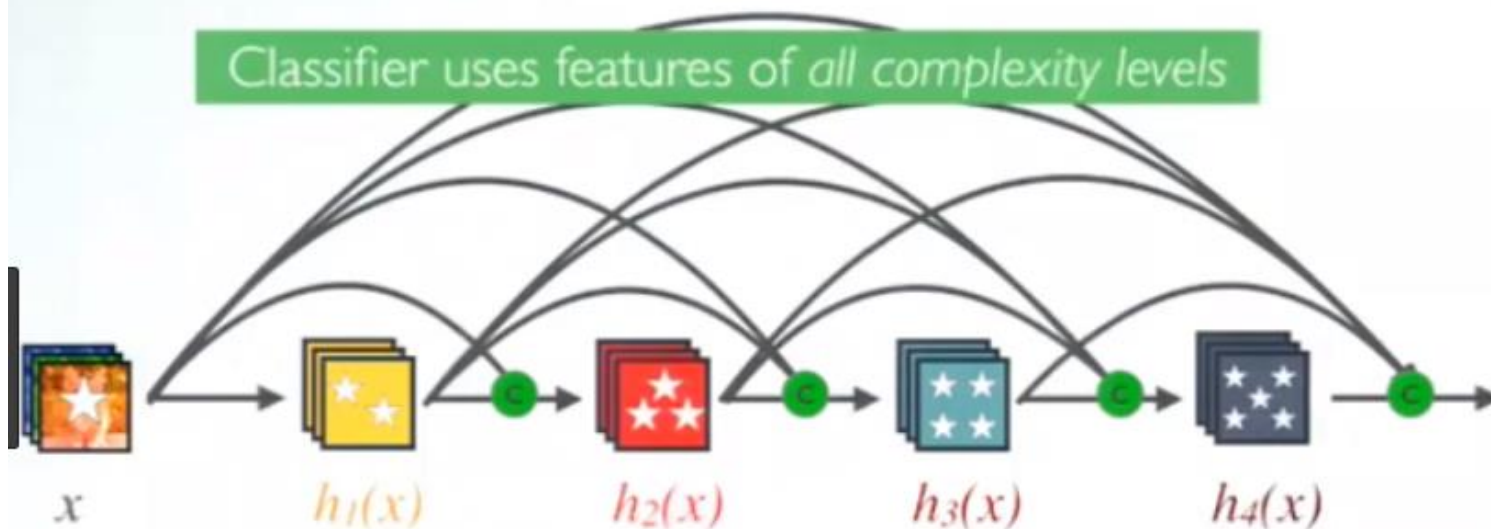
- Dense blocks where each layer is connected to every other layer in feedforward fashion
- Alleviates vanishing gradient, strengthens feature propagation, encourages feature reuse



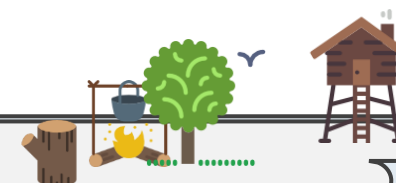
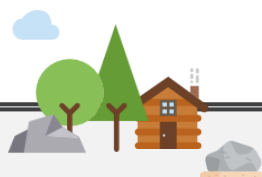
# Beyond ResNet – DenseNet

## Dense Connectivity:

Classifier uses features of *all* complexity levels

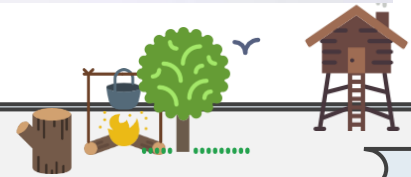
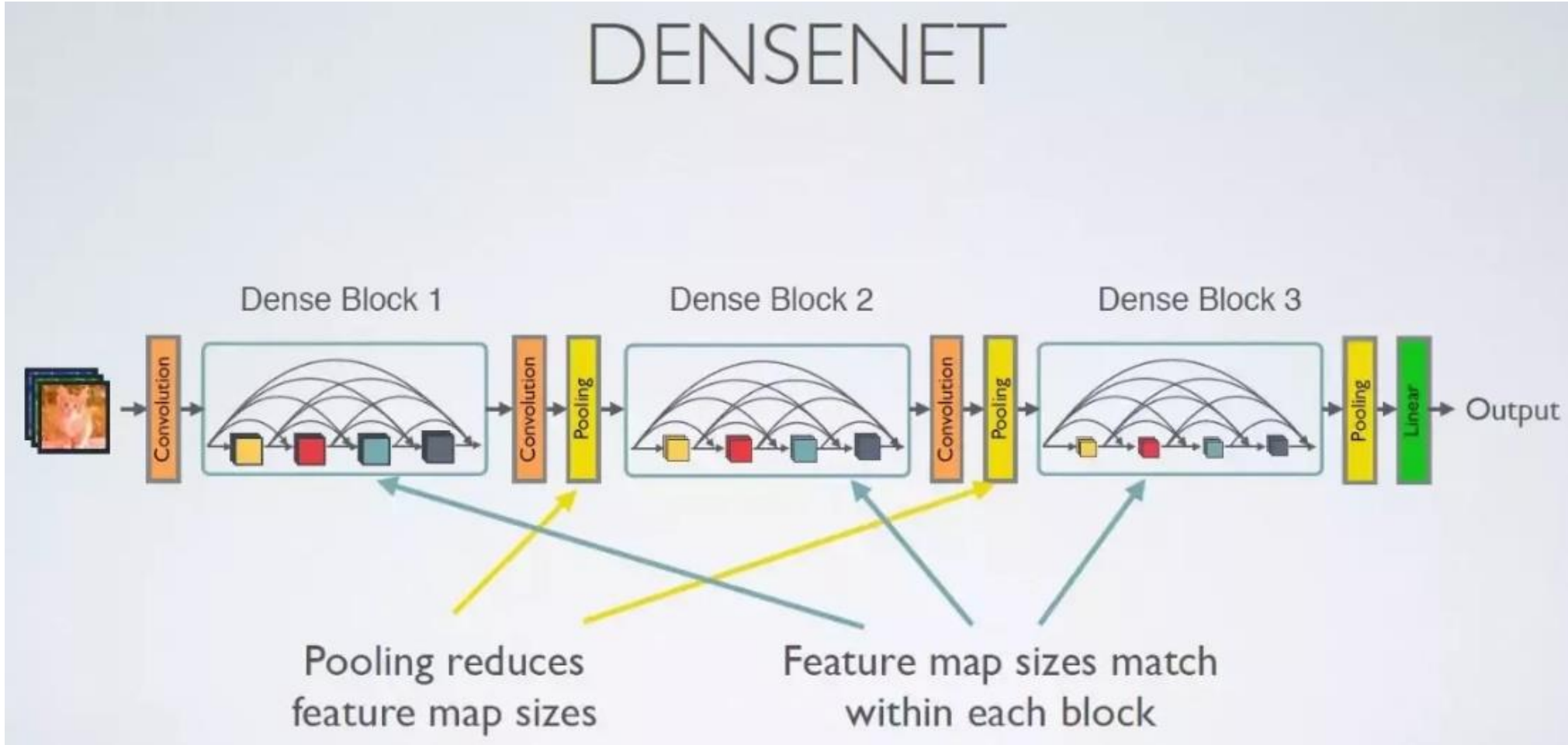


★ Increasingly complex features





# Beyond ResNet - DenseNet



CS231n : <http://cs231n.stanford.edu/syllabus.html>

