# Lecture 15: CNNs –II Architectures

Yujiao Shi

SIST, ShanghaiTech

Spring, 2025

# Outline

- **CNN architectures**

  - □ Sequential structure: LeNet/AlexNet/VGGNet

  - □ Parallel branches: GoogLeNet

  - □ Residual structure: ResNet/DenseNet

  - □ Network Architecture Search

*Acknowledgement: Zemel et al's CSC411 and Feifei Li et al's cs231n notes*

ShanghaiTech University

- **Problem Setup**
  - Input: Image
  - Output: Object class

# LeNet-5

- Handwritten digit recognition
- LeCun et al., 1998

# ImageNet (ILSVRC)

# AlexNet

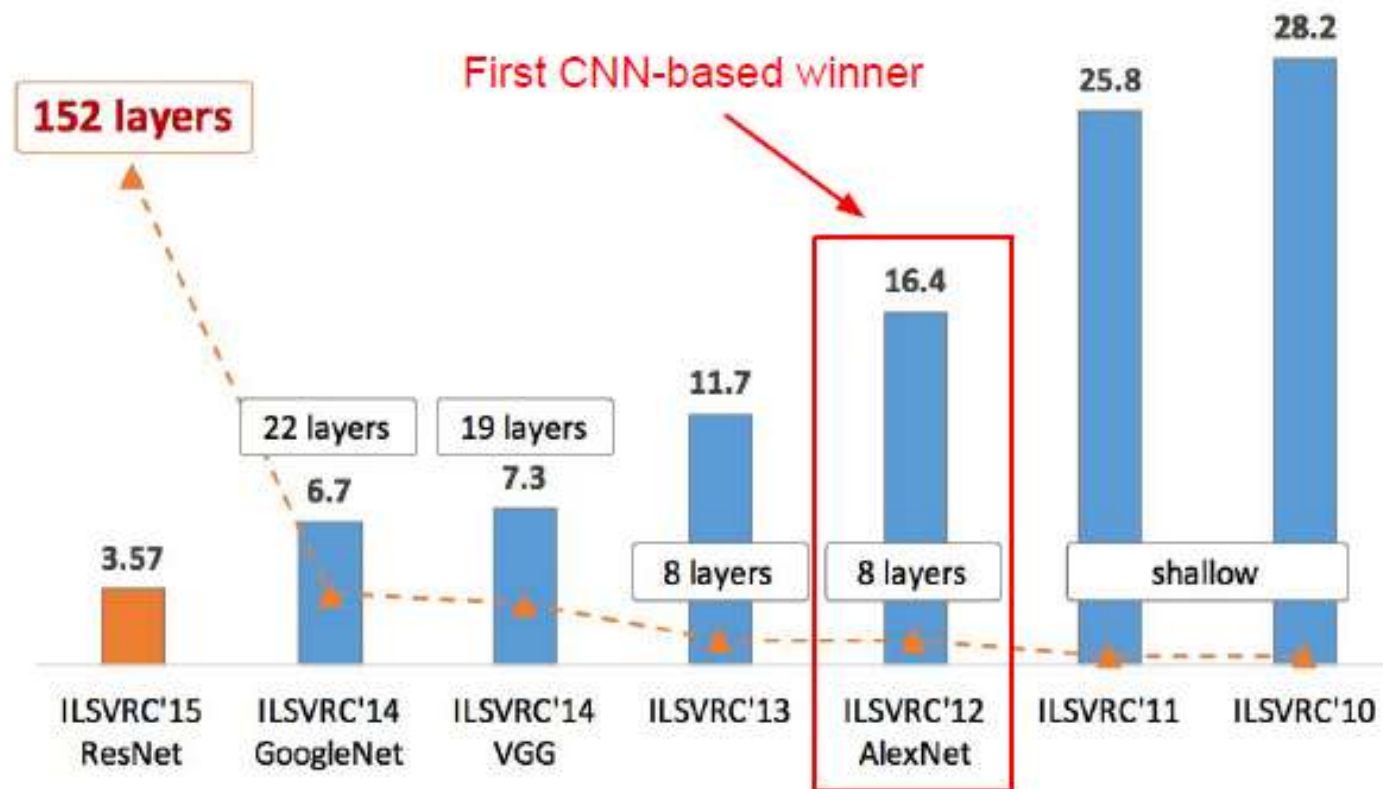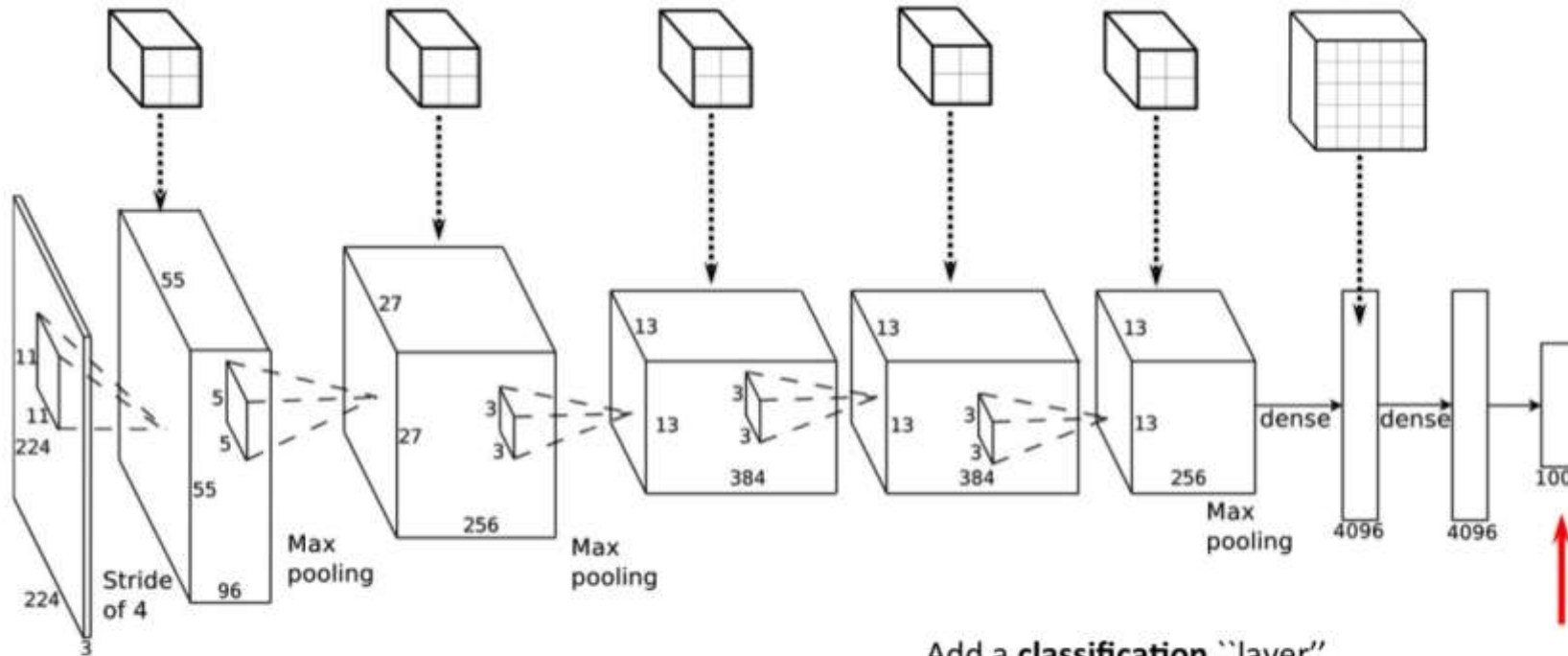- heavy data augmentation
- dropout 0.5
- batch size 128
- SGD Momentum 0.9
- Learning rate 1e-2, reduced by 10 manually when val accuracy plateaus
- L2 weight decay 5e-4
- 7 CNN ensemble: 18.2% -> 15.4%

Add a **classification** ``layer''.

For an input image, the value in a particular dimension of this vector tells you the probability of the corresponding object class.
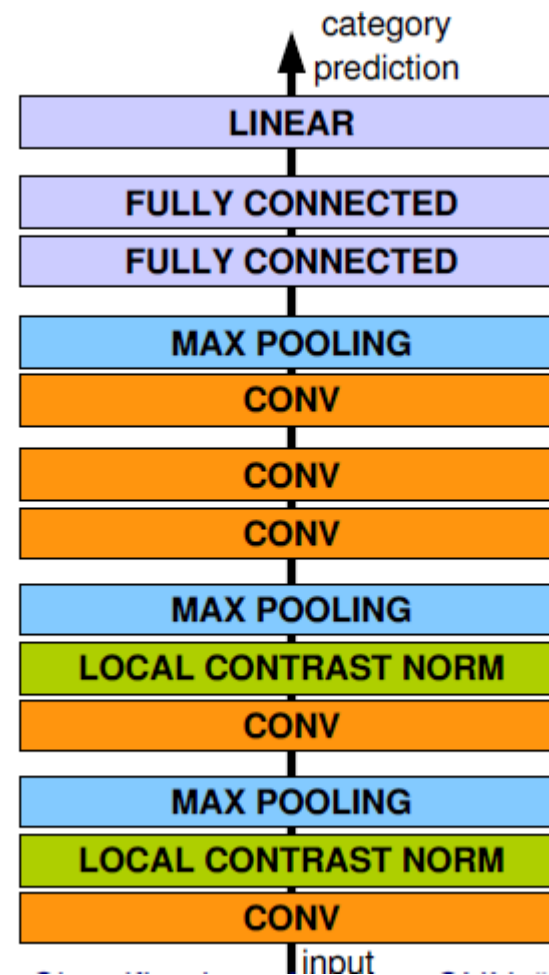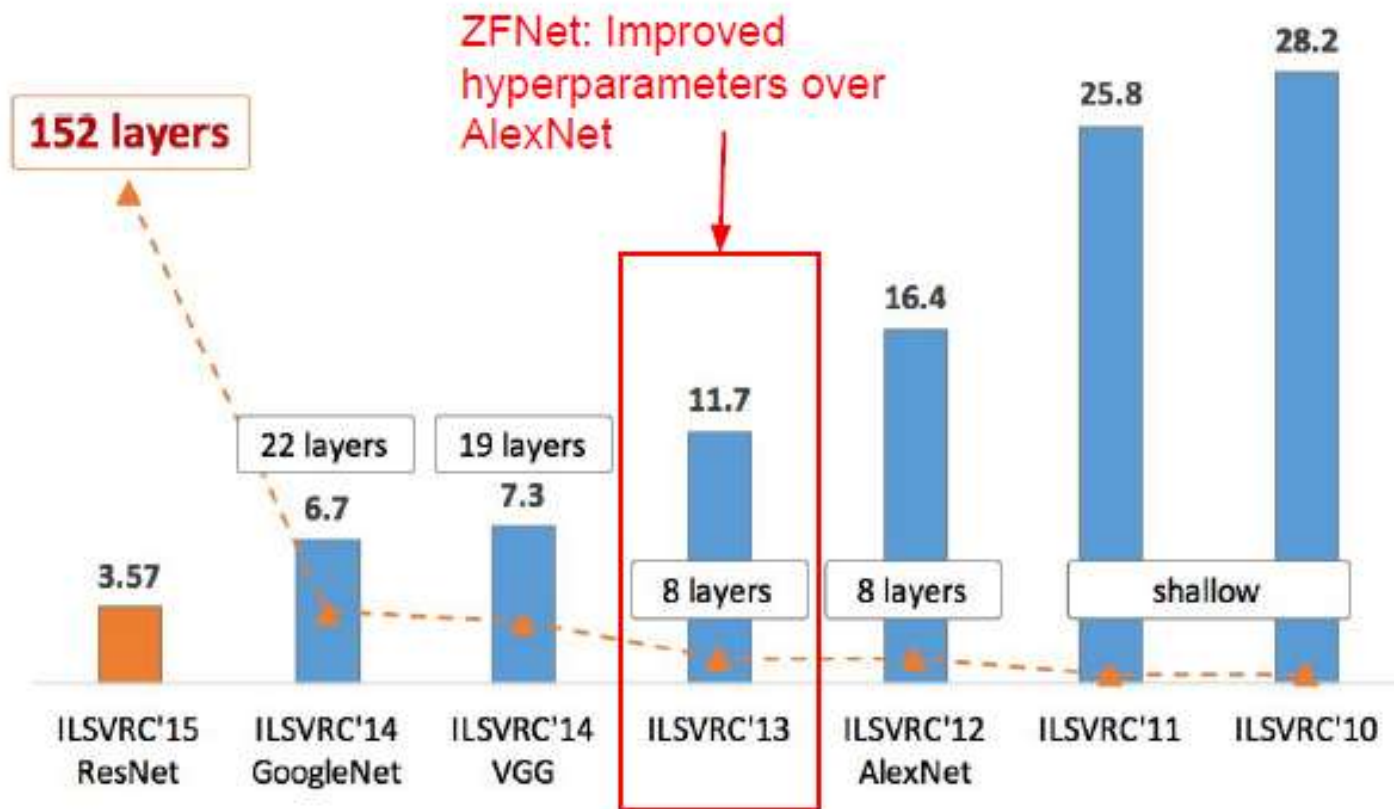
# AlexNet

- Deeper network structure
  - More convolution layers
  - Local contrast normalization
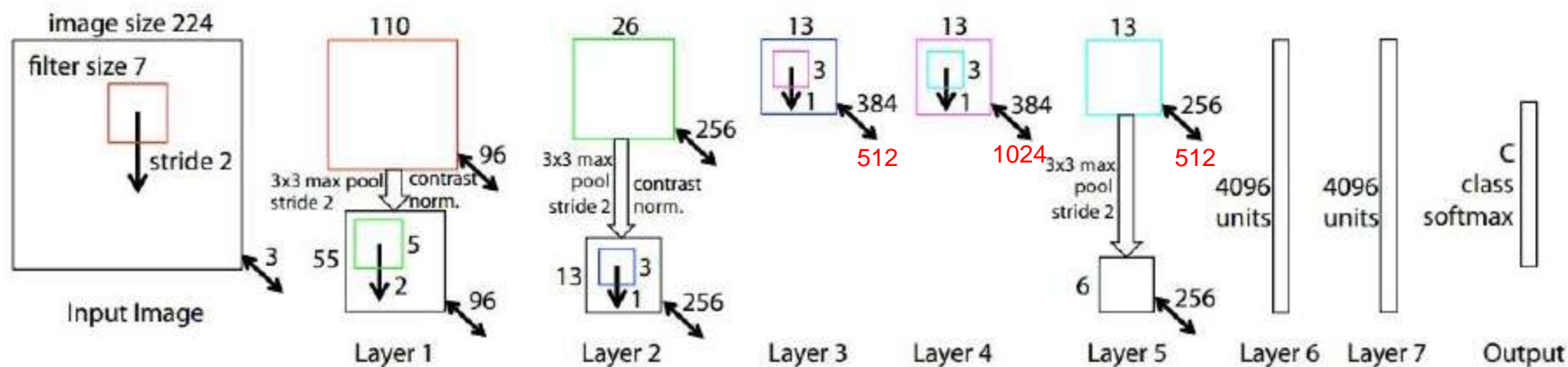  - ReLu instead of Tanh
  - Dropout as regularization

[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)

category prediction

LINEAR

FULLY CONNECTED

FULLY CONNECTED

MAX POOLING

CONV

CONV

CONV

MAX POOLING

LOCAL CONTRAST NORM

CONV

MAX POOLING

LOCAL CONTRAST NORM
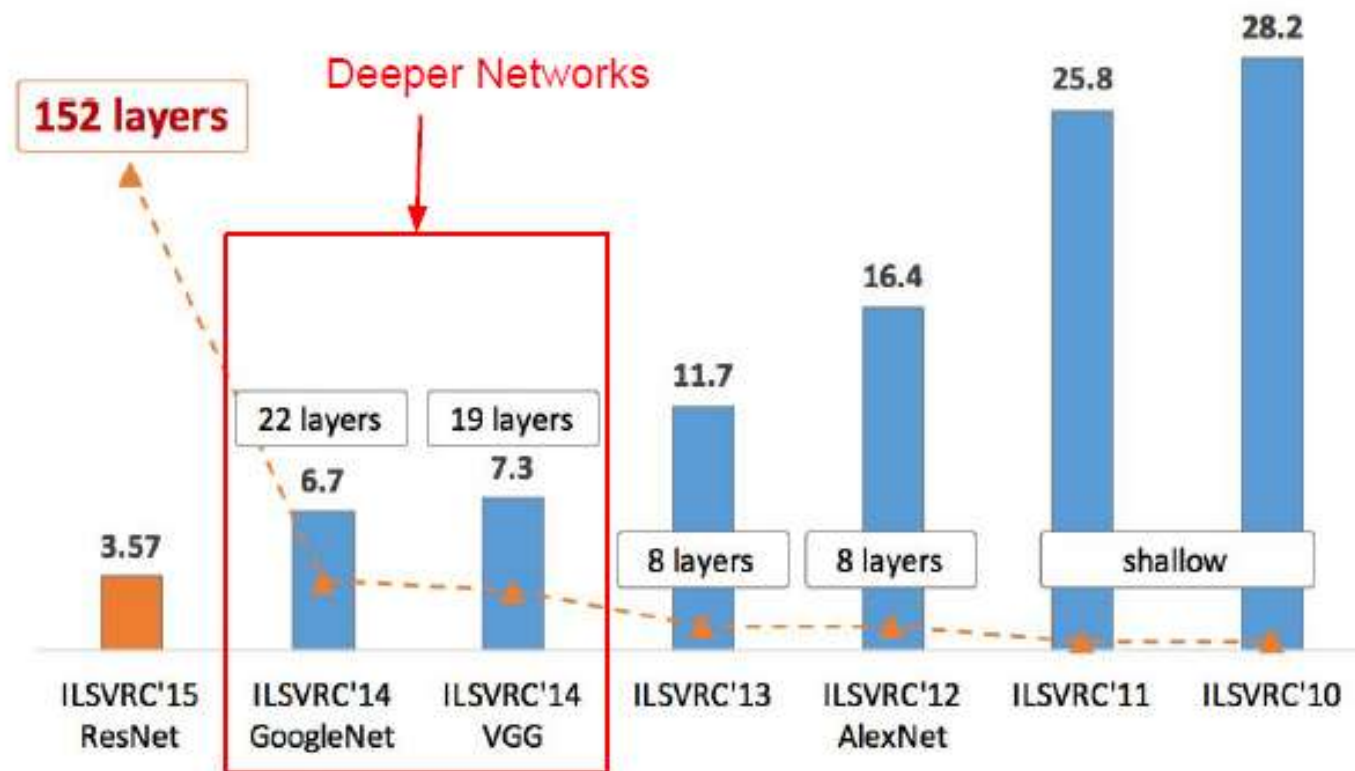
CONV

input

# ImageNet (ILSVRC)

# ZFNet

AlexNet but:
CONV1: change from (11x11 stride 4) to (7x7 stride 2)
CONV3,4,5: instead of 384, 384, 256 filters use 512, 1024, 512

ImageNet top 5 error: 16.4% -> 11.7%

# ImageNet (ILSVRC)

**Deeper Networks**

| | | | | | | |
|---|---|---|---|---|---|---|
| 152 layers | 22 layers | 19 layers | | | | |
| | 6.7 | 7.3 | 11.7 | 16.4 | 25.8 | 28.2 |
| 3.57 | | | 8 layers | 8 layers | shallow | |
| ILSVRC'15 ResNet | ILSVRC'14 GoogleNet | ILSVRC'14 VGG | ILSVRC'13 | ILSVRC'12 AlexNet | ILSVRC'11 | ILSVRC'10 |

上海科技大学
ShanghaiTech University

## Case Study: VGGNet

*[Simonyan and Zisserman, 2014]*
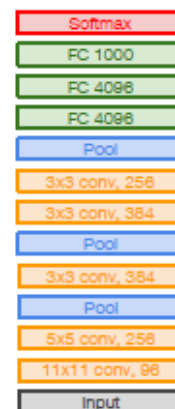
Small filters, Deeper networks

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)
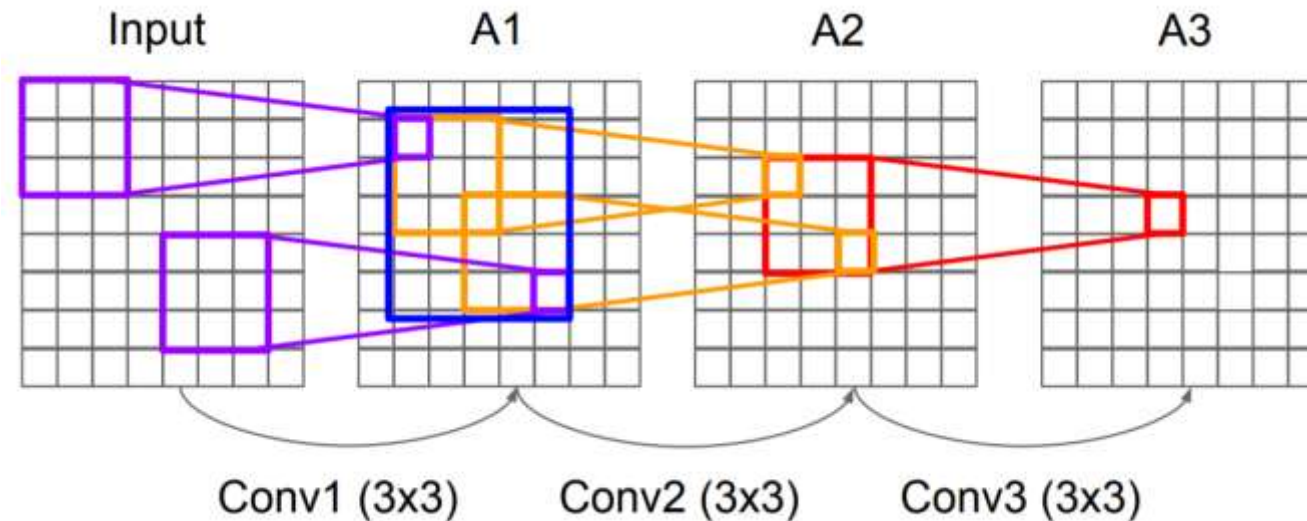-> 7.3% top 5 error in ILSVRC'14



AlexNet        VGG16        VGG19

上海科技大学
**ShanghaiTech University**

## 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

But deeper, more non-linearities

And fewer parameters: $3 * (3^2 C^2)$ vs. $7^2 C^2$ for C channels per layer



Input          A1          A2          A3

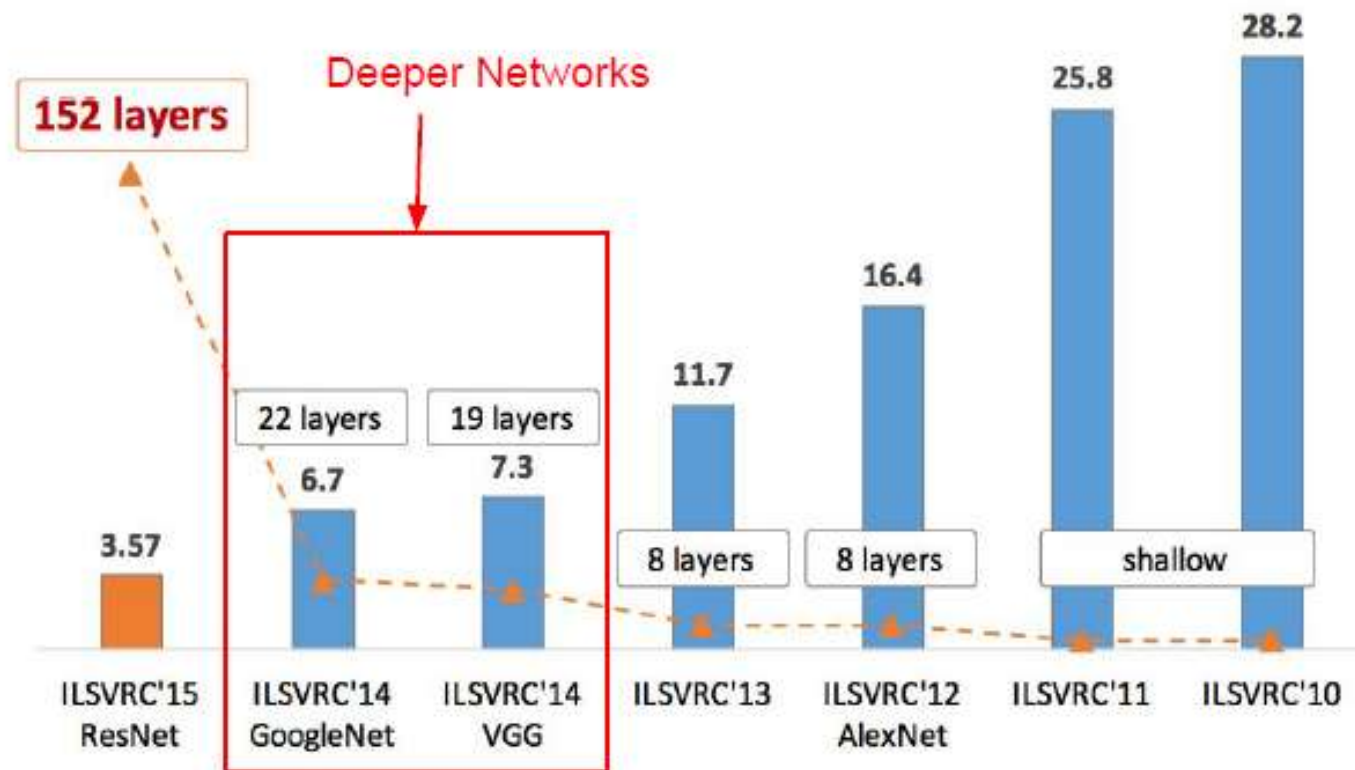Conv1 (3x3)     Conv2 (3x3)     Conv3 (3x3)

# Outline

- **CNN architectures**

  - ☐ Sequential structure: LeNet/AlexNet/VGGNet

  - ☐ Parallel branches: GoogLeNet

  - ☐ Residual structure: ResNet/DenseNet

  - ☐ Network Architecture Search

*Acknowledgement: Zemel et al's CSC411 and Feifei Li et al's cs231n notes*
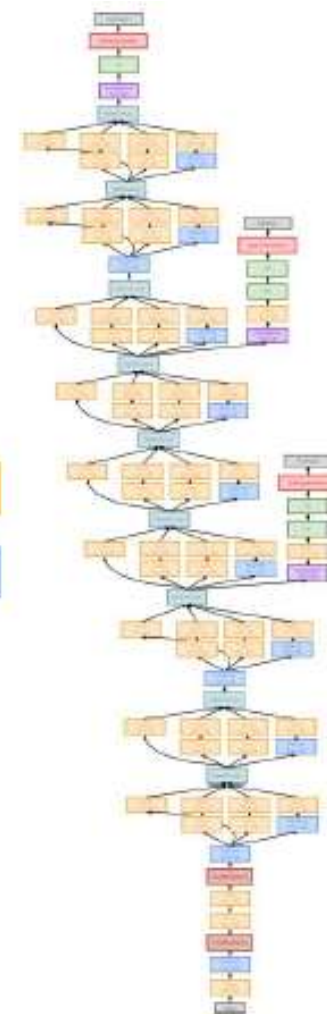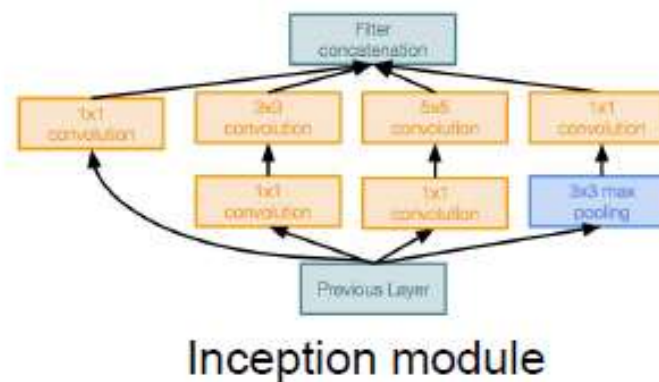
上海科技大学
**ShanghaiTech University**

# Case Study: GoogLeNet

[Szegedy et al., 2014]

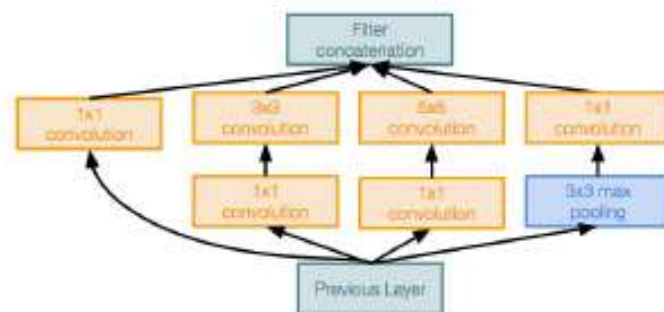Deeper networks, with computational efficiency

- 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)

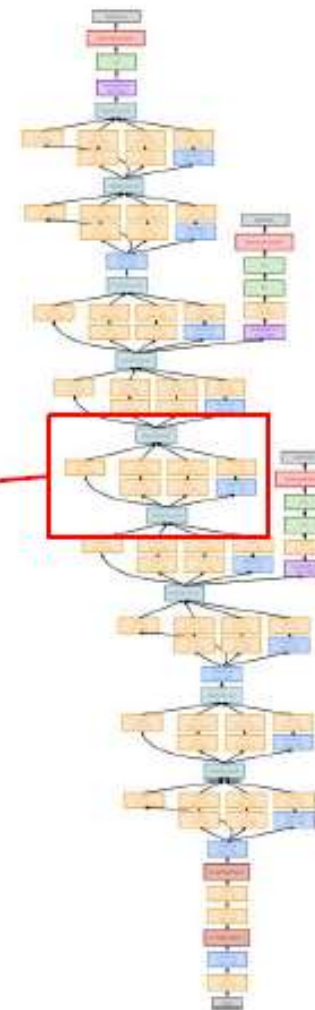Inception module

上海科技大学
ShanghaiTech University

# Case Study: GoogLeNet

*[Szegedy et al., 2014]*

"Inception module": design a good local network topology (network within a network) and then stack these modules on top of each other
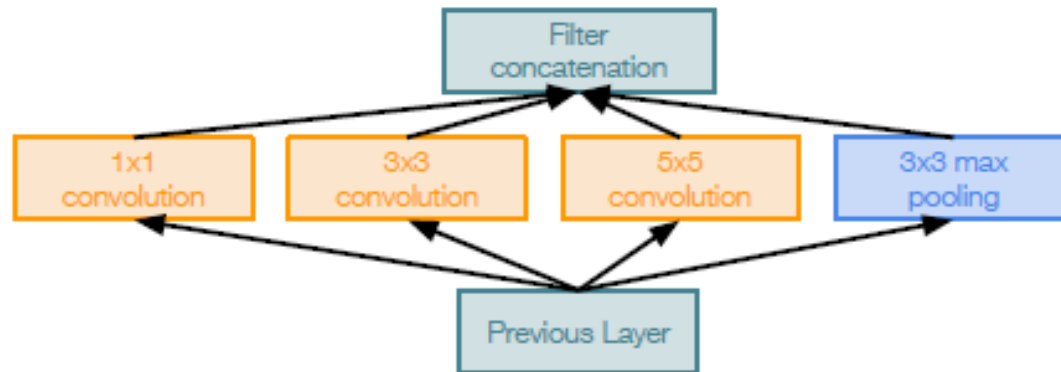
Inception module

# GoogLeNet

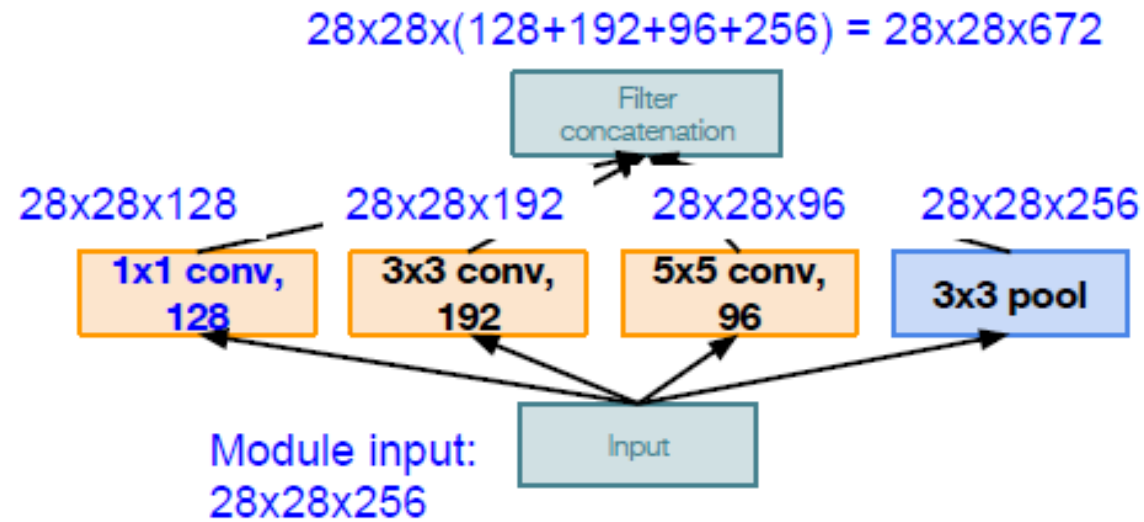- Inception Module



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

# GoogLeNet

- Inception Module

$28 \times 28 \times (128 + 192 + 96 + 256) = 28 \times 28 \times 672$

```
              Filter
           concatenation

28x28x128    28x28x192    28x28x96    28x28x256

1x1 conv,    3x3 conv,    5x5 conv,    3x3 pool
  128          192           96

Module input:            Input
28x28x256
```

Naive Inception module

**Conv Ops:**
[1x1 conv, 128]  28x28x128x1x1x256
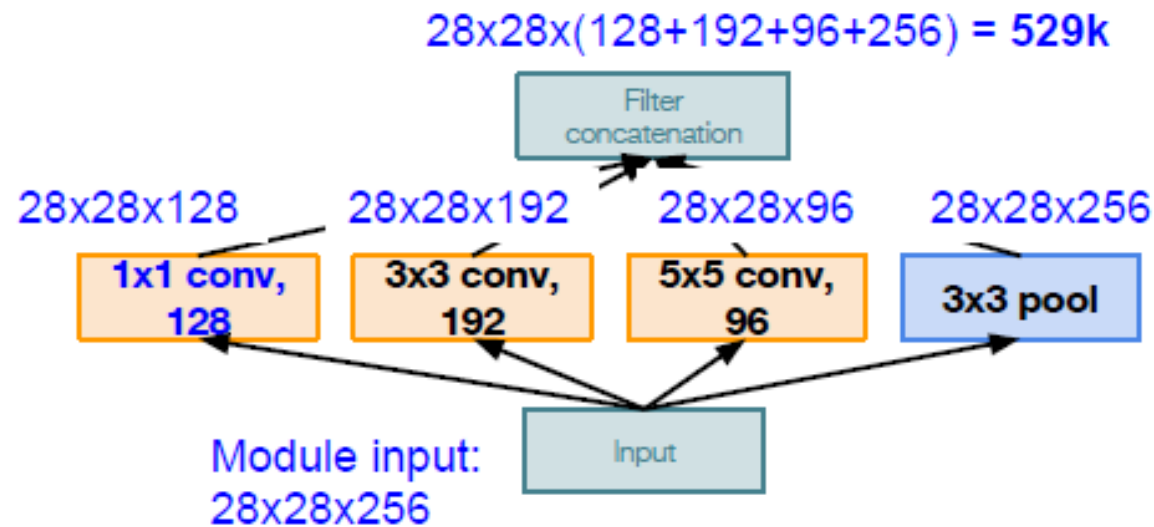[3x3 conv, 192]  28x28x192x3x3x256
[5x5 conv, 96]  28x28x96x5x5x256
**Total: 854M ops**

Very expensive compute

Pooling layer also preserves feature depth, which means total depth after concatenation can only grow at every layer!

# GoogLeNet

- Inception Module

$28 \times 28 \times (128+192+96+256) = 529k$

Filter concatenation

$28 \times 28 \times 128$    $28 \times 28 \times 192$    $28 \times 28 \times 96$    $28 \times 28 \times 256$

| 1x1 conv, 128 | 3x3 conv, 192 | 5x5 conv, 96 | 3x3 pool |

Module input: $28 \times 28 \times 256$

Input

Naive Inception module

Solution: "bottleneck" layers that use 1x1 convolutions to reduce feature depth

# GoogLeNet

- Bottleneck layer



56

**1x1 CONV
with 32 filters**

56

(each filter has size
1x1x64, and performs a
64-dimensional dot
product)

56

56

64

32

preserves spatial
dimensions, reduces depth!

Projects depth to lower
dimension (combination of
feature maps)

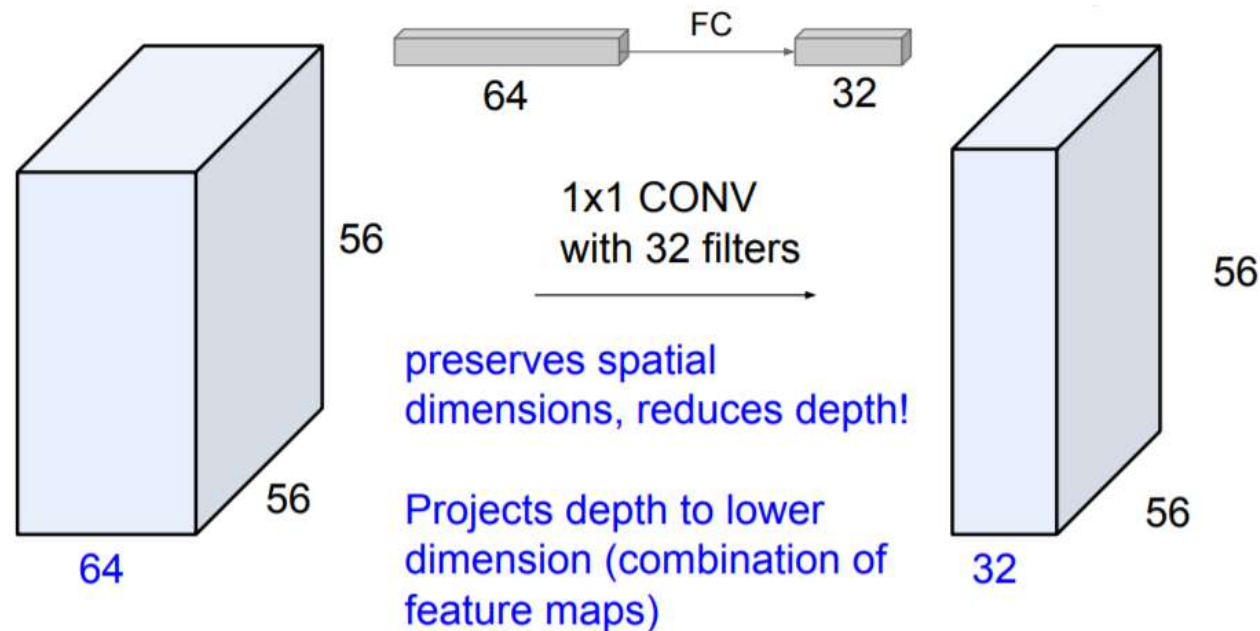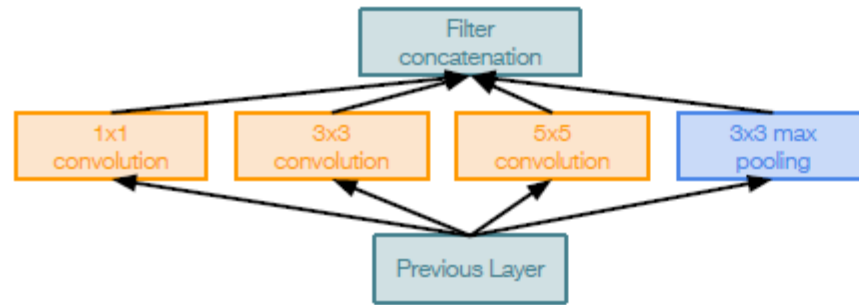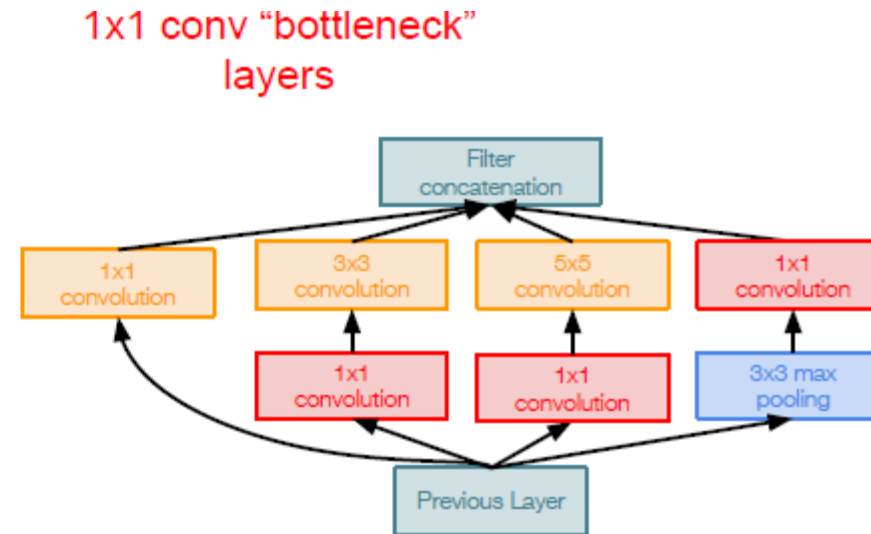- **1x1 Convolutions**
  - ☐ Alternatively, interpret it as applying the same FC layer on each input pixel

- Inception Module

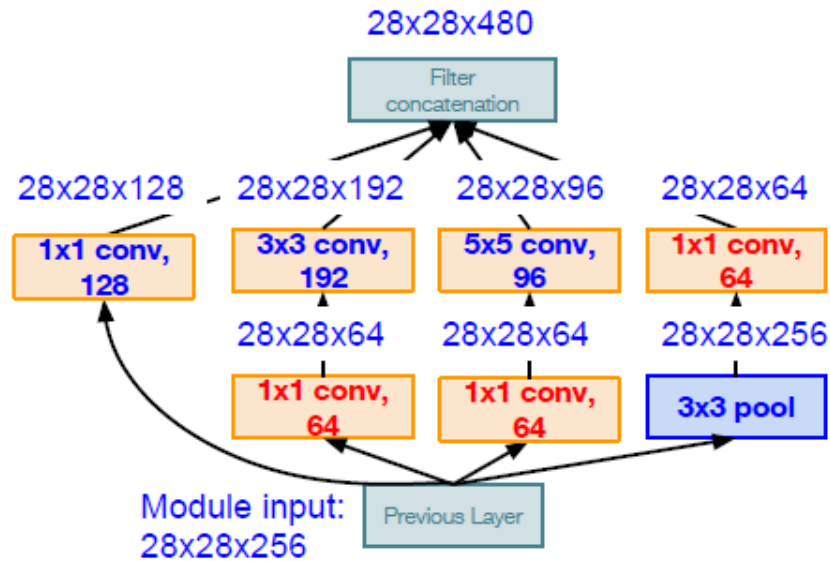

Naive Inception module

Inception module with dimension reduction

# GoogLeNet

- Inception Module



Inception module with dimension reduction

**Conv Ops:**
[1x1 conv, 64] 28x28x64x1x1x256
[1x1 conv, 64] 28x28x64x1x1x256
[1x1 conv, 128] 28x28x128x1x1x256
[3x3 conv, 192] 28x28x192x3x3x64
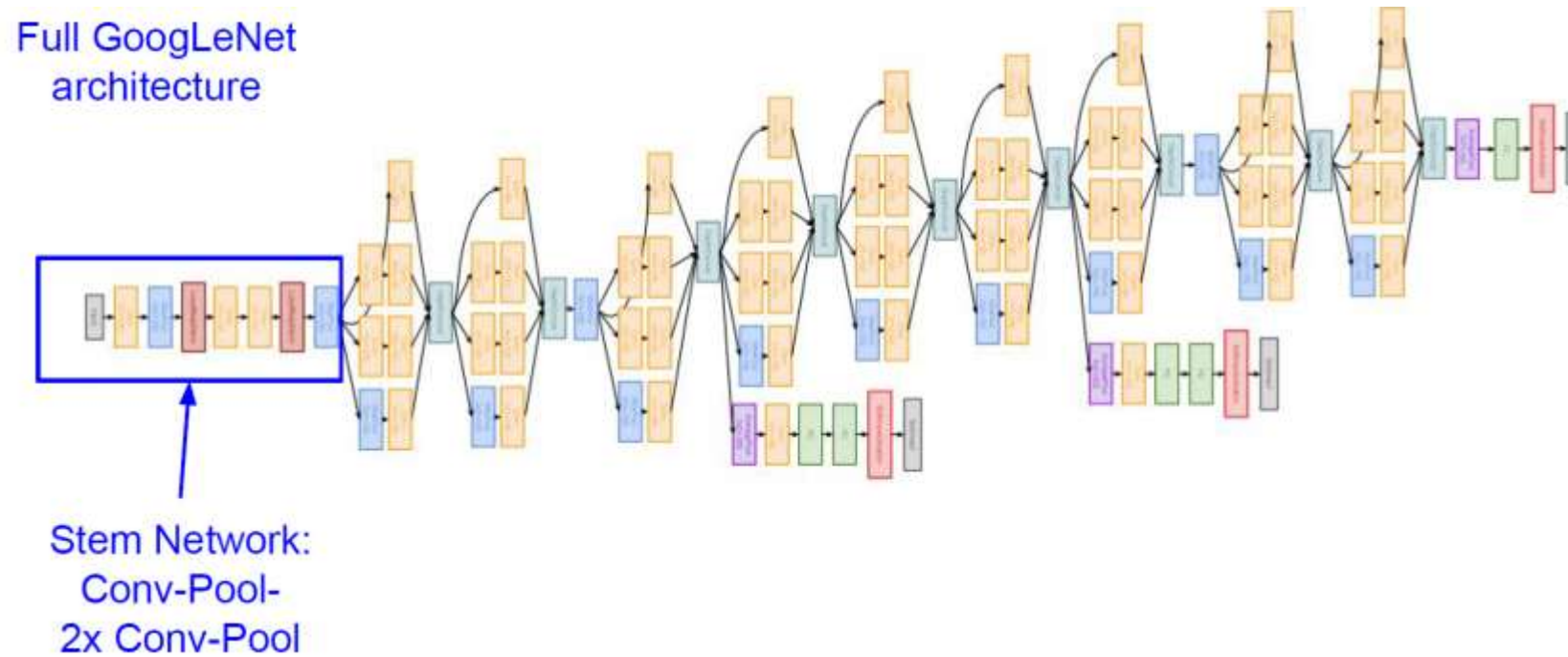[5x5 conv, 96] 28x28x96x5x5x64
[1x1 conv, 64] 28x28x64x1x1x256
**Total: 358M ops**

Compared to 854M ops for naive version
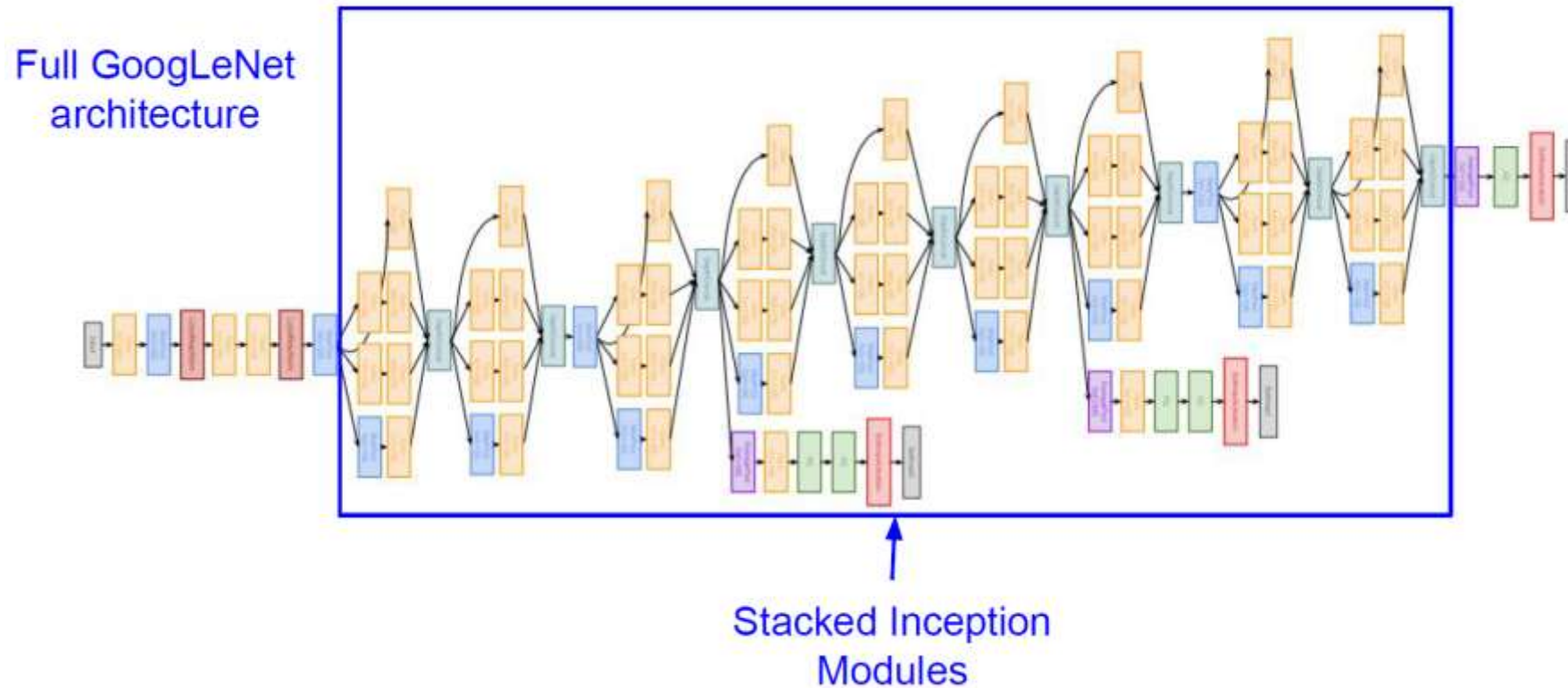Bottleneck can also reduce depth after
pooling layer

# GoogLeNet

- Overall network structure



Full GoogLeNet architecture

Stem Network:
Conv-Pool-
2x Conv-Pool

# GoogLeNet

- Overall network structure



Full GoogLeNet architecture

Stacked Inception Modules
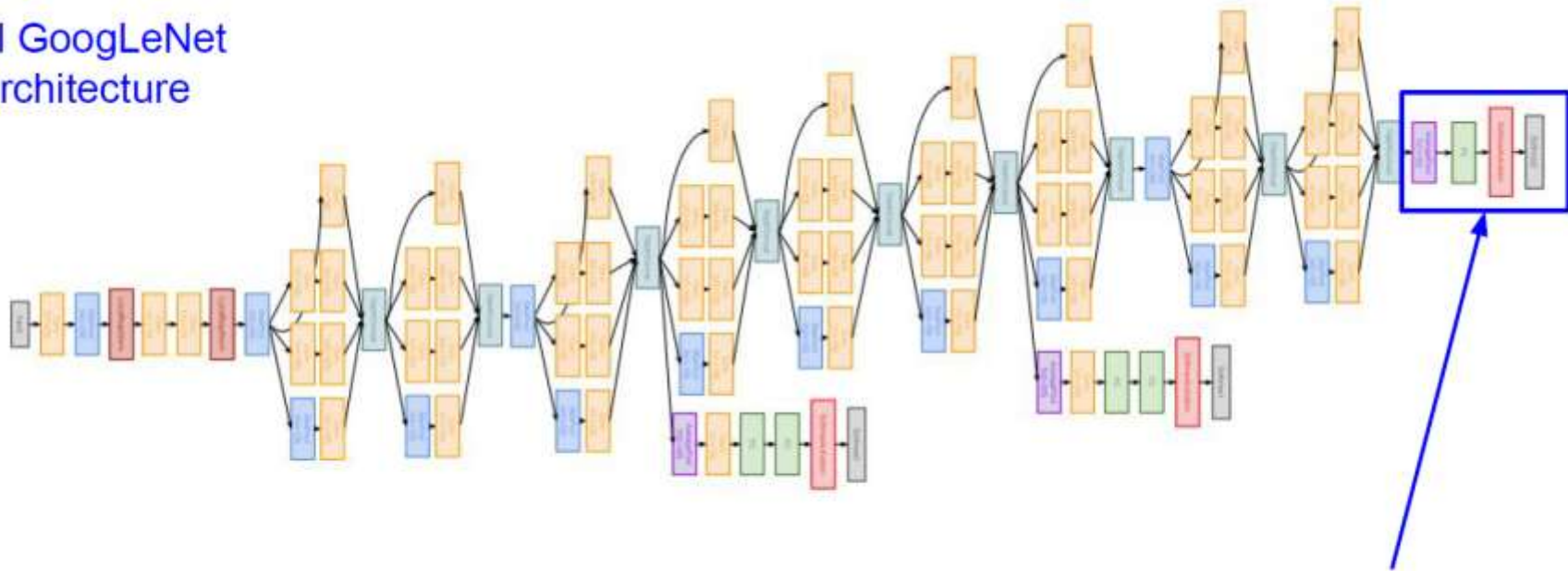
# GoogLeNet

- Overall network structure

Full GoogLeNet
architecture
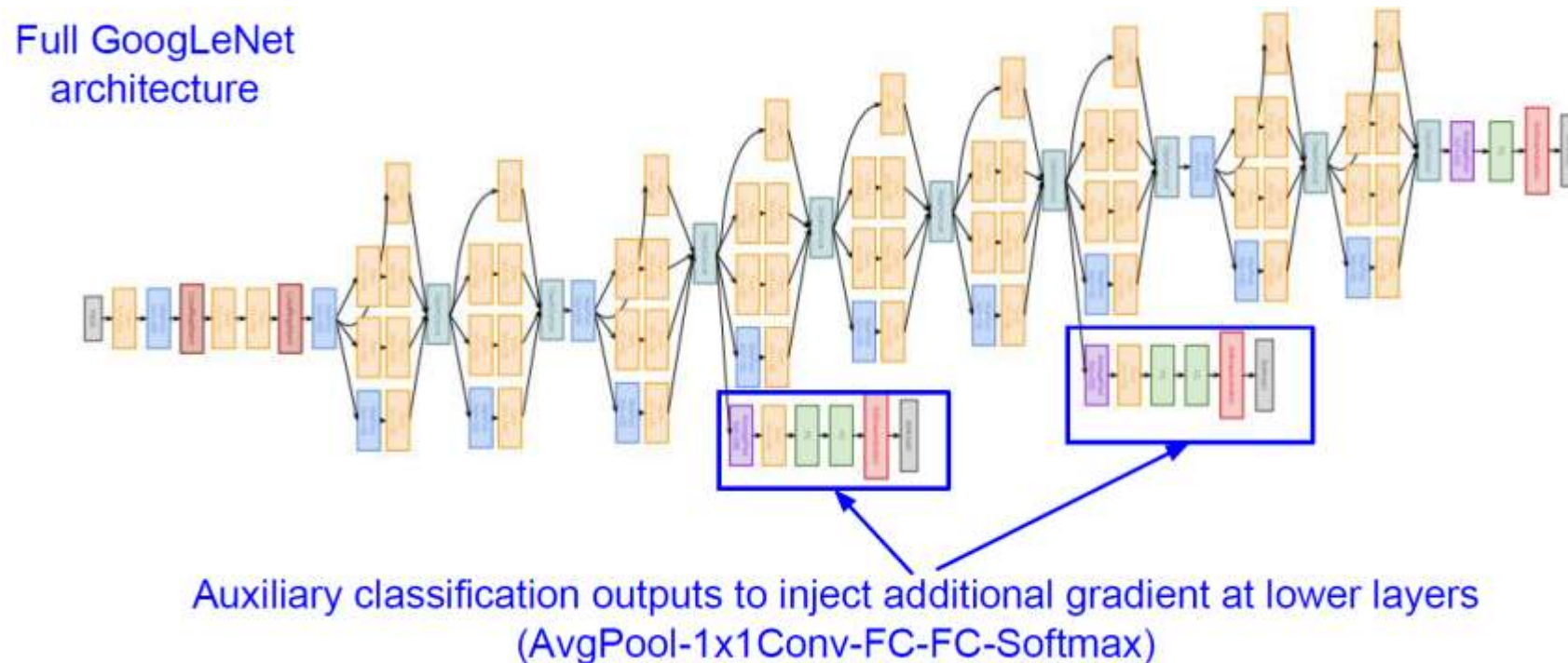
Classifier output
(removed expensive FC layers!)

# GoogLeNet

- Overall network structure



Full GoogLeNet architecture

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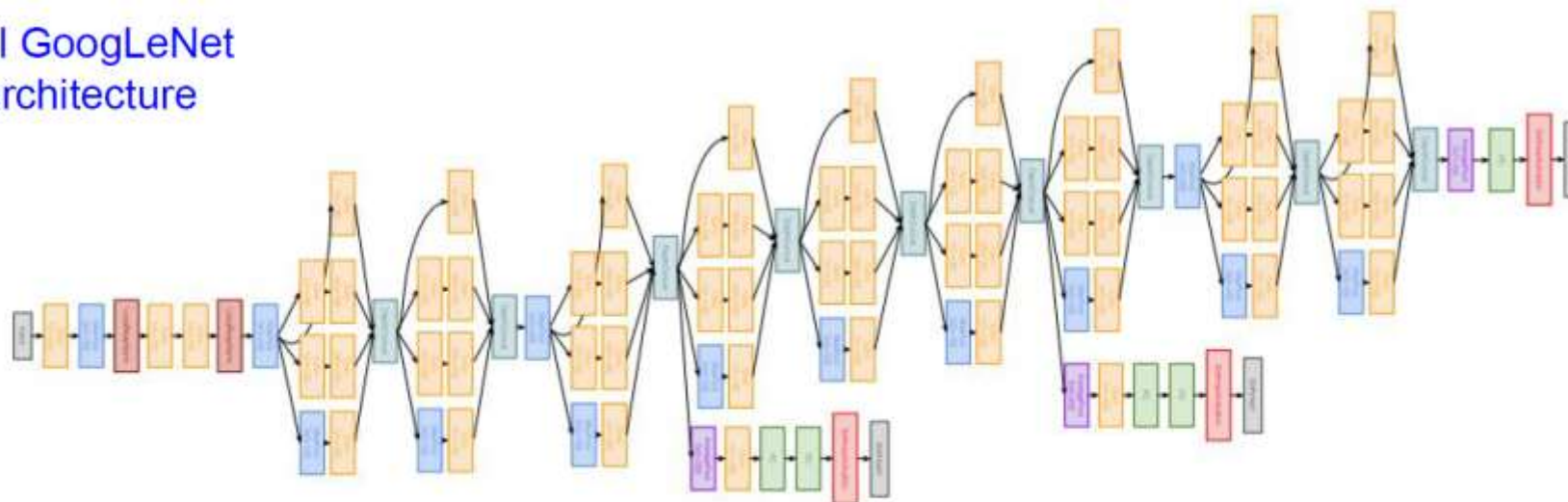
# GoogLeNet

- Overall network structure
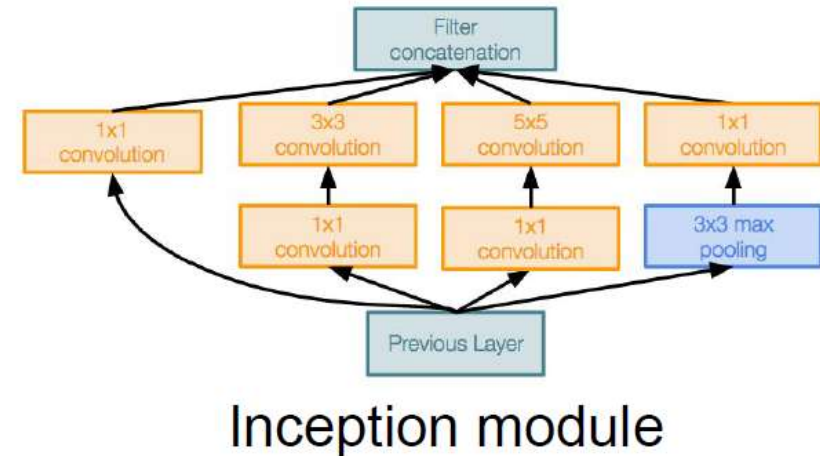
Full GoogLeNet architecture



22 total layers with weights (including each parallel layer in an Inception module)

# GoogLeNet

- Summary

Deeper networks, with computational efficiency

- 22 layers
- Efficient "Inception" module
- No FC layers
- 12x less params than AlexNet
- ILSVRC'14 classification winner (6.7% top 5 error)



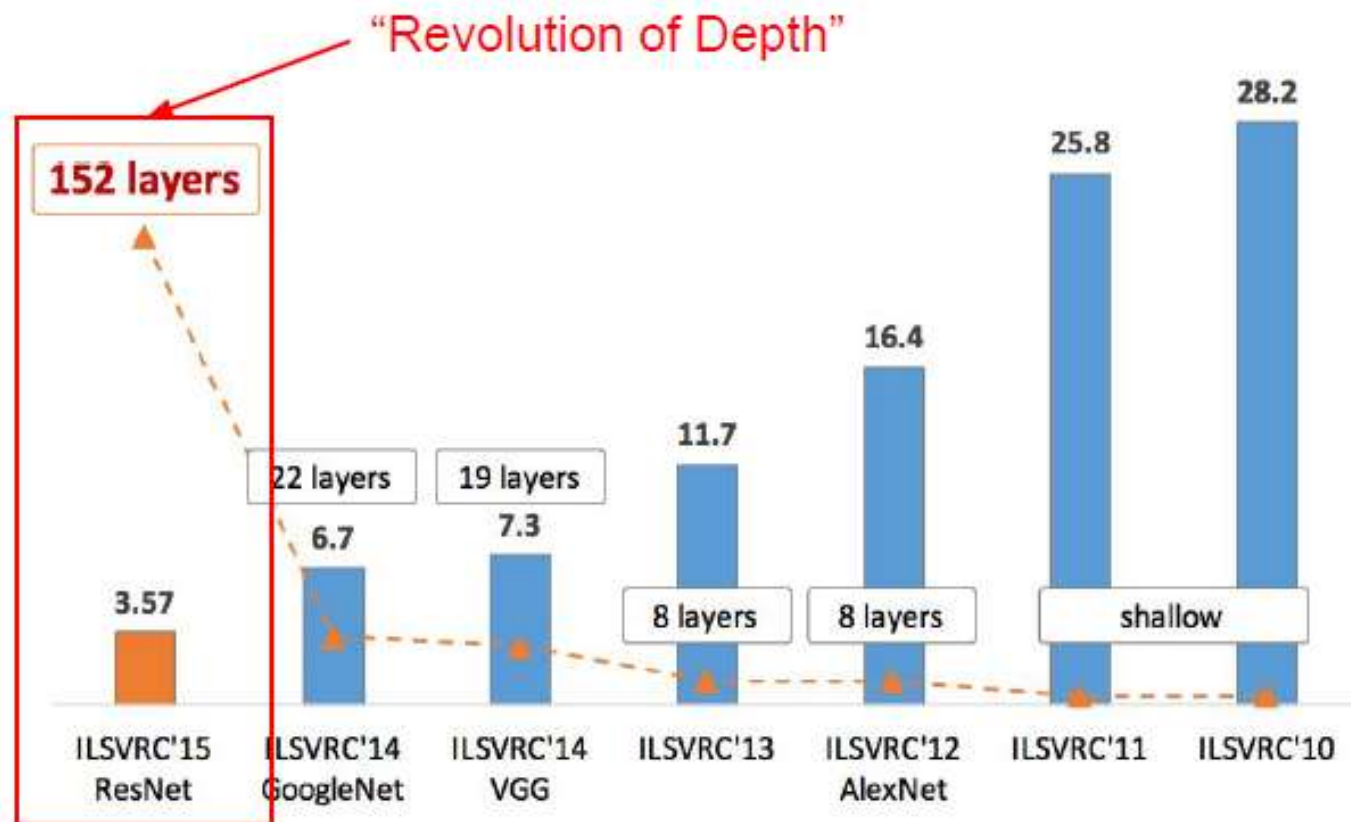Inception module

# Outline

- CNN architectures

  ☐ Sequential structure: LeNet/AlexNet/VGGNet

  ☐ Parallel branches: GoogLeNet

  ☐ Residual structure: ResNet/DenseNet

  ☐ Network Architecture Search

*Acknowledgement: Zemel et al's CSC411 and Feifei Li et al's cs231n notes*
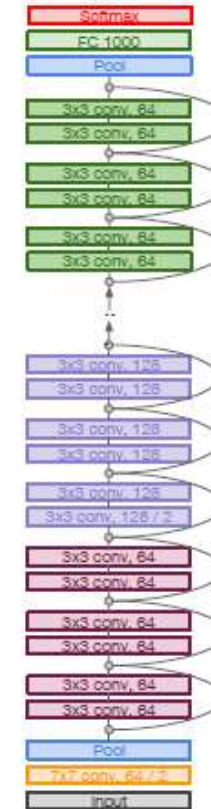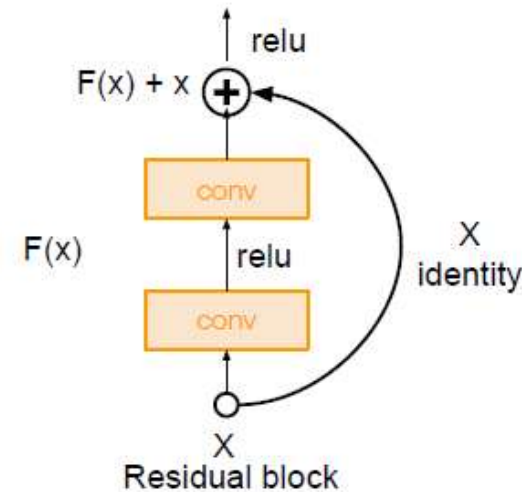
# ResNet

## 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!



Residual block

# ResNet

- What happens when stacking deeper plain conv layers?



56-layer model performs worse on both training and test error
-> The deeper model performs worse, but it's not caused by overfitting!
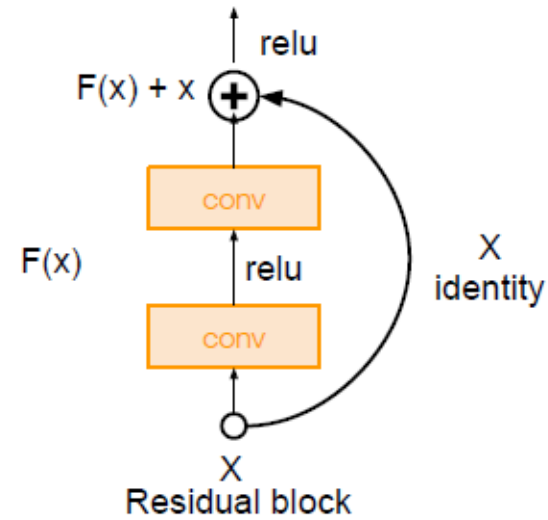
- Hypothesis:
  - The problem is an optimization problem, deeper models are harder to optimize

The deeper model should be able to perform at least as well as the shallower model.

A solution by construction is copying the learned layers from the shallower model and setting additional layers to identity mapping.

# ResNet

- **Solution:**
  - ☐ Use network layers to fit a residual mapping
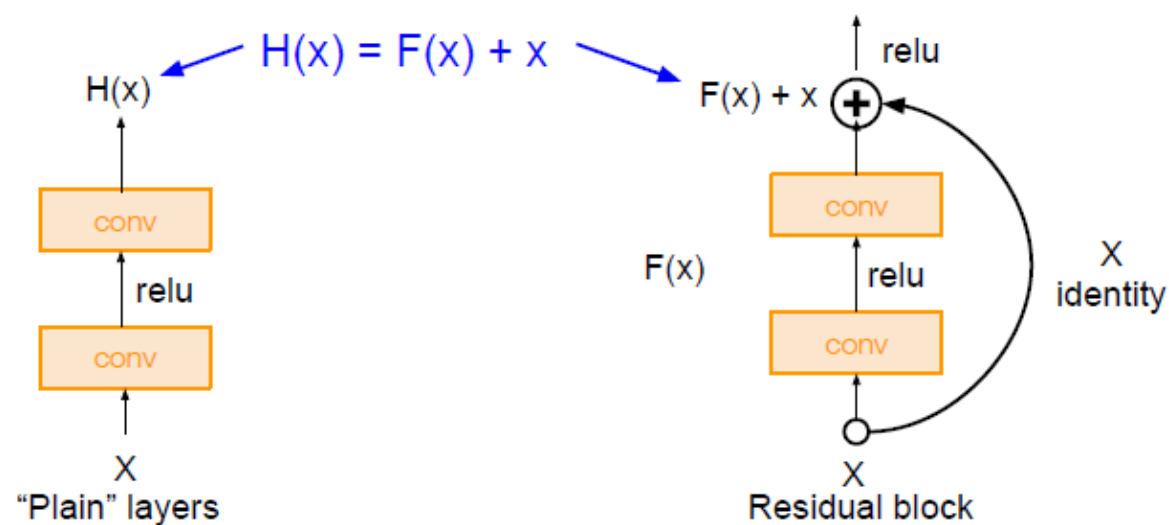


He et al "Deep Residual Learning for Image Recognition", CVPR 2016

# ResNet

- **Solution:**
  - ☐ Use network layers to fit a residual mapping



$H(x) = F(x) + x$

$H(x)$

conv

relu

conv

X
"Plain" layers

relu

$F(x) + x$ ⊕

conv

$F(x)$       relu

conv

X
identity

X
Residual block

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)

relu

$F(x) + x$

3x3 conv

$F(x)$    relu

3x3 conv

X

X identity

Residual block

3x3 conv, 128 filters, /2 spatially with stride 2
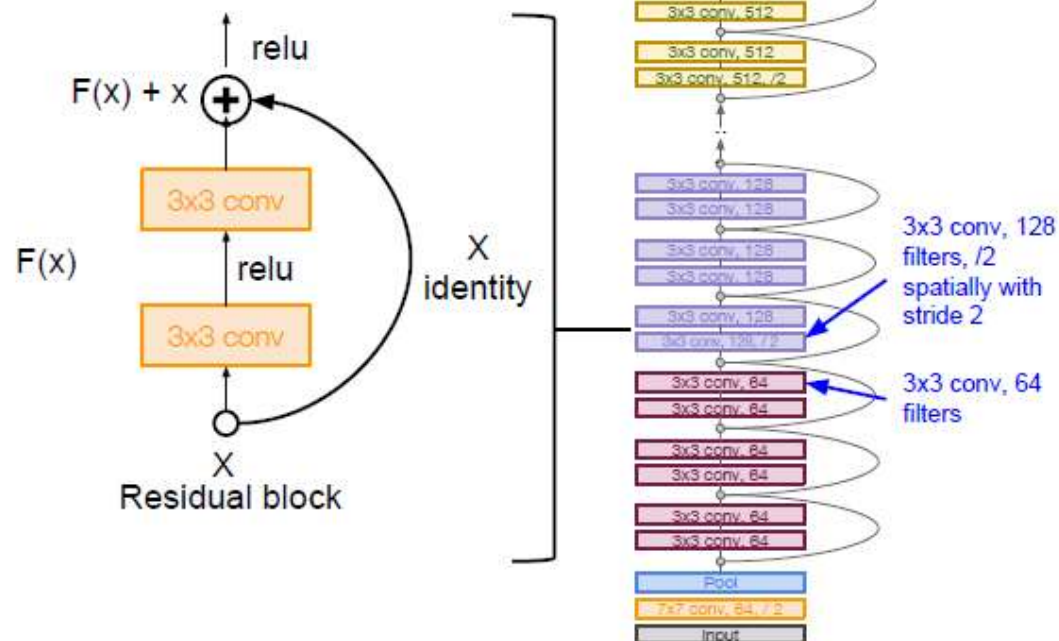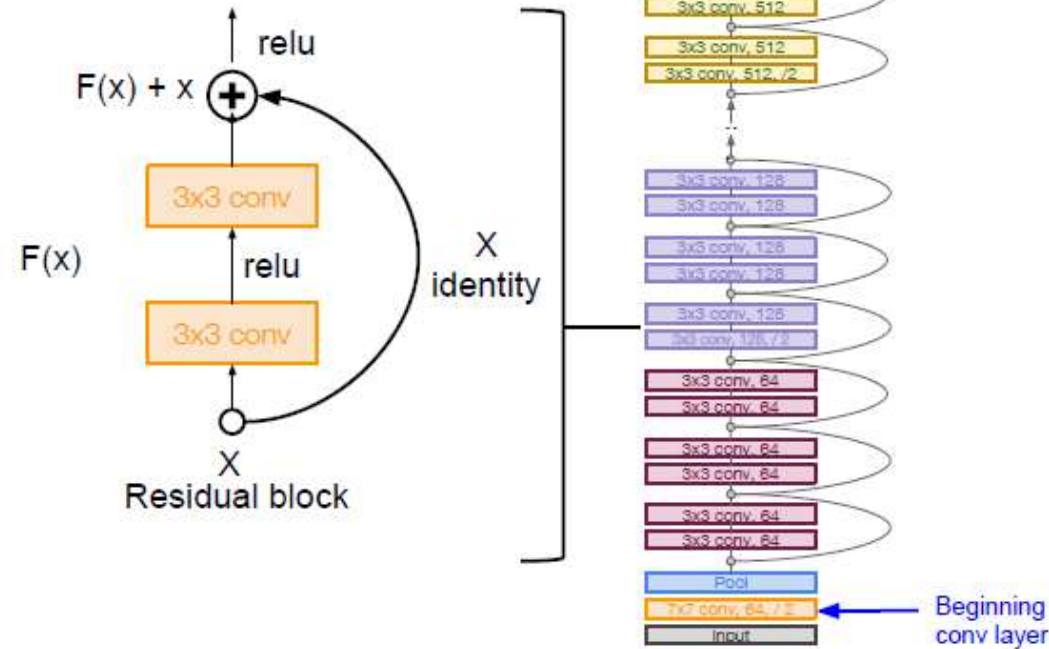
3x3 conv, 64 filters
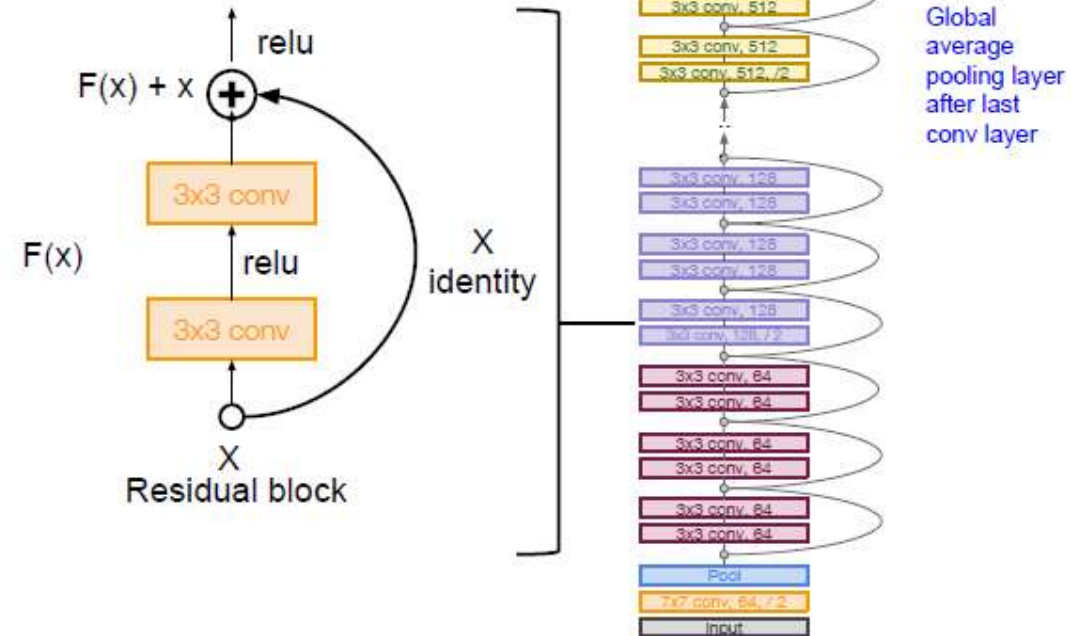
# ResNet



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

$F(x) + x$

relu

$F(x)$

relu

X identity

3x3 conv

3x3 conv

X

Residual block

Beginning conv layer

# ResNet

## 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
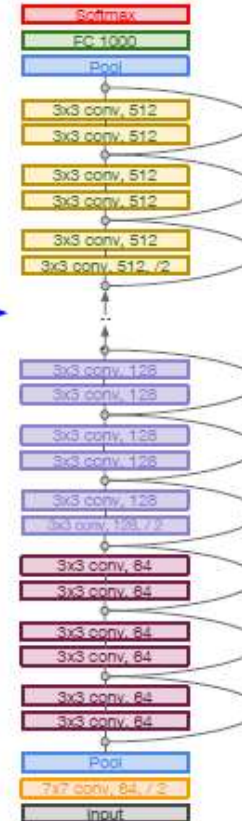- No FC layers at the end (only FC 1000 to output classes)

relu

$F(x) + x$

3x3 conv

$F(x)$     relu

3x3 conv

X identity

X

Residual block

No FC layers besides FC 1000 to output classes

Global average pooling layer after last conv layer

Softmax
FC 1000
Pool
3x3 conv, 512
3x3 conv, 512
3x3 conv, 512
3x3 conv, 512
3x3 conv, 512
3x3 conv, 512, /2
3x3 conv, 128
3x3 conv, 128
3x3 conv, 128
3x3 conv, 128
3x3 conv, 128
3x3 conv, 128, /2
3x3 conv, 64
3x3 conv, 64
3x3 conv, 64
3x3 conv, 64
3x3 conv, 64
3x3 conv, 64
Pool
7x7 conv, 64, / 2
Input

# ResNet



Case Study: ResNet

[He et al., 2015]

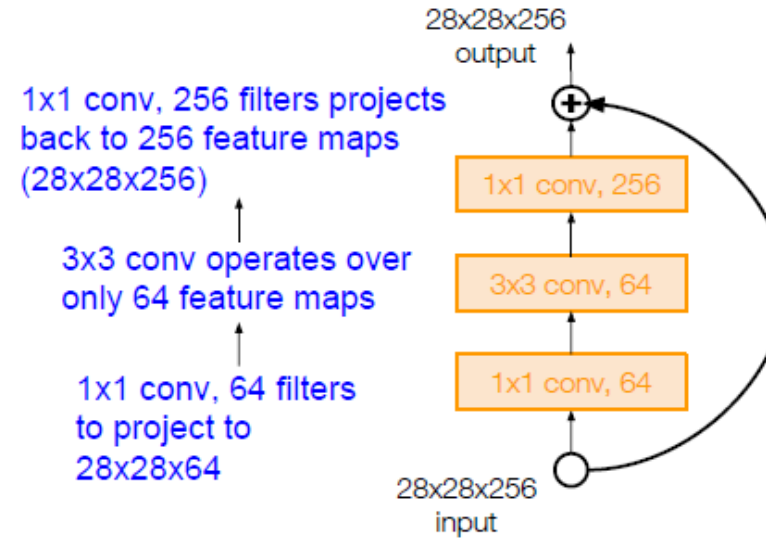Total depths of 34, 50, 101, or 152 layers for ImageNet

# ResNet

## Case Study: ResNet

*[He et al., 2015]*

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

1x1 conv, 256 filters projects
back to 256 feature maps
(28x28x256)

3x3 conv operates over
only 64 feature maps

1x1 conv, 64 filters
to project to
28x28x64

28x28x256
output

1x1 conv, 256

3x3 conv, 64

1x1 conv, 64

28x28x256
input

# ResNet

- **Training details**

  Training ResNet in practice:

  - Batch Normalization after every CONV layer
  - Xavier/2 initialization from He et al.
  - SGD + Momentum (0.9)
  - Learning rate: 0.1, divided by 10 when validation error plateaus
  - Mini-batch size 256
  - Weight decay of 1e-5
  - No dropout used

# ResNet

■ Results

Experimental Results
- Able to train very deep networks without degrading (152 layers on ImageNet, 1202 on Cifar)
- Deeper networks now achieve lowing training error as expected
- Swept 1st place in all ILSVRC and COCO 2015 competitions

MSRA @ ILSVRC & COCO 2015 Competitions

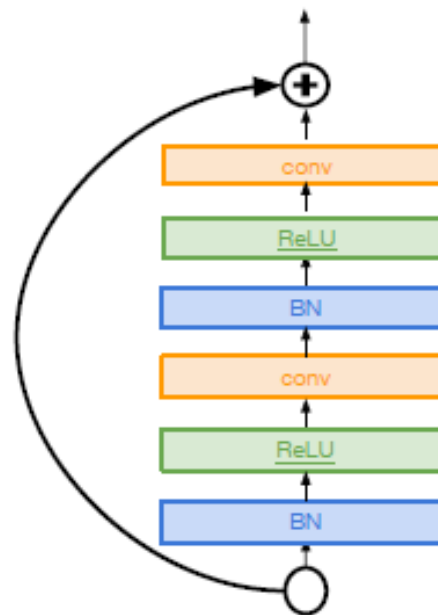- **1st places** in all five main tracks
  - ImageNet Classification: *"Ultra-deep"* (quote Yann) 152-layer nets
  - ImageNet Detection: 16% better than 2nd
  - ImageNet Localization: 27% better than 2nd
  - COCO Detection: 11% better than 2nd
  - COCO Segmentation: 12% better than 2nd

ILSVRC 2015 classification winner (3.6% top 5 error) -- better than "human performance"! (Russakovsky 2014)

上海科技大学
ShanghaiTech University

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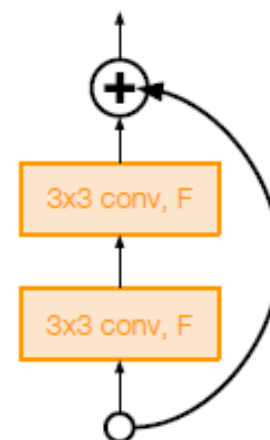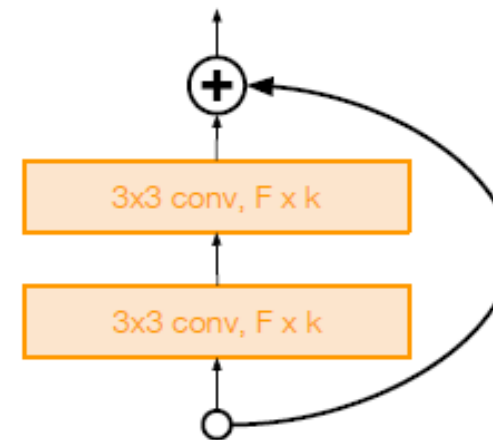- Gives better performance

上海科技大学
ShanghaiTech University

# Wide Residual Networks

*[Zagoruyko et al. 2016]*

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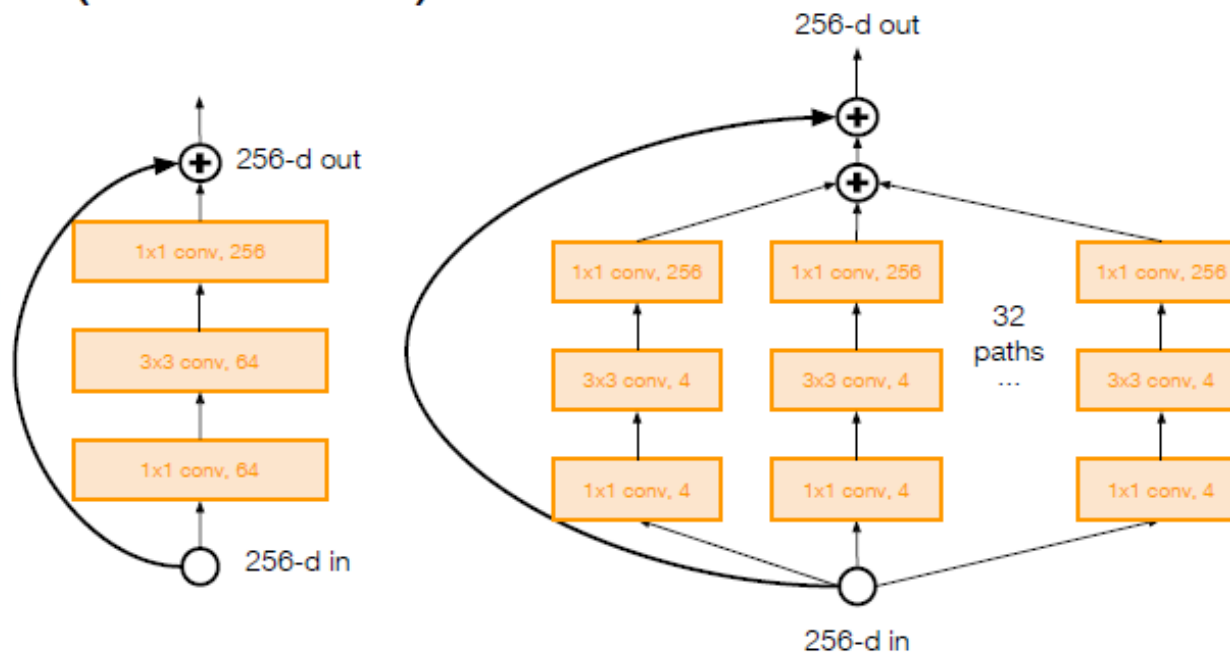- Increasing width instead of depth more computationally efficient (parallelizable)



| 3x3 conv, F |
| 3x3 conv, F |

Basic residual block

| 3x3 conv, F x k |
| 3x3 conv, F x k |

Wide residual block

上海科技大学
ShanghaiTech University

## 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")
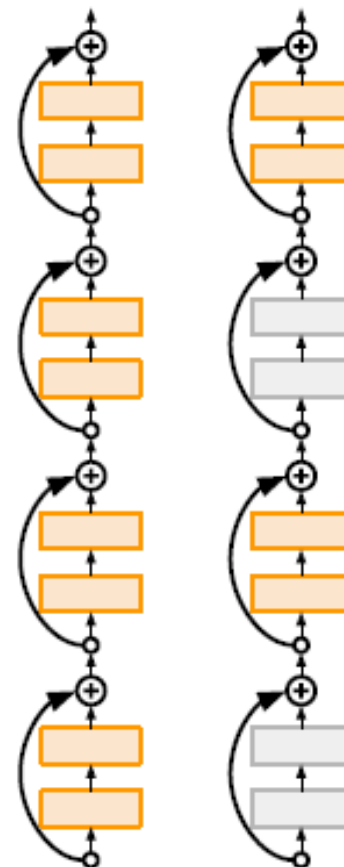- Parallel pathways similar in spirit to Inception module

上 海 科 技 大 学
ShanghaiTech University

# 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
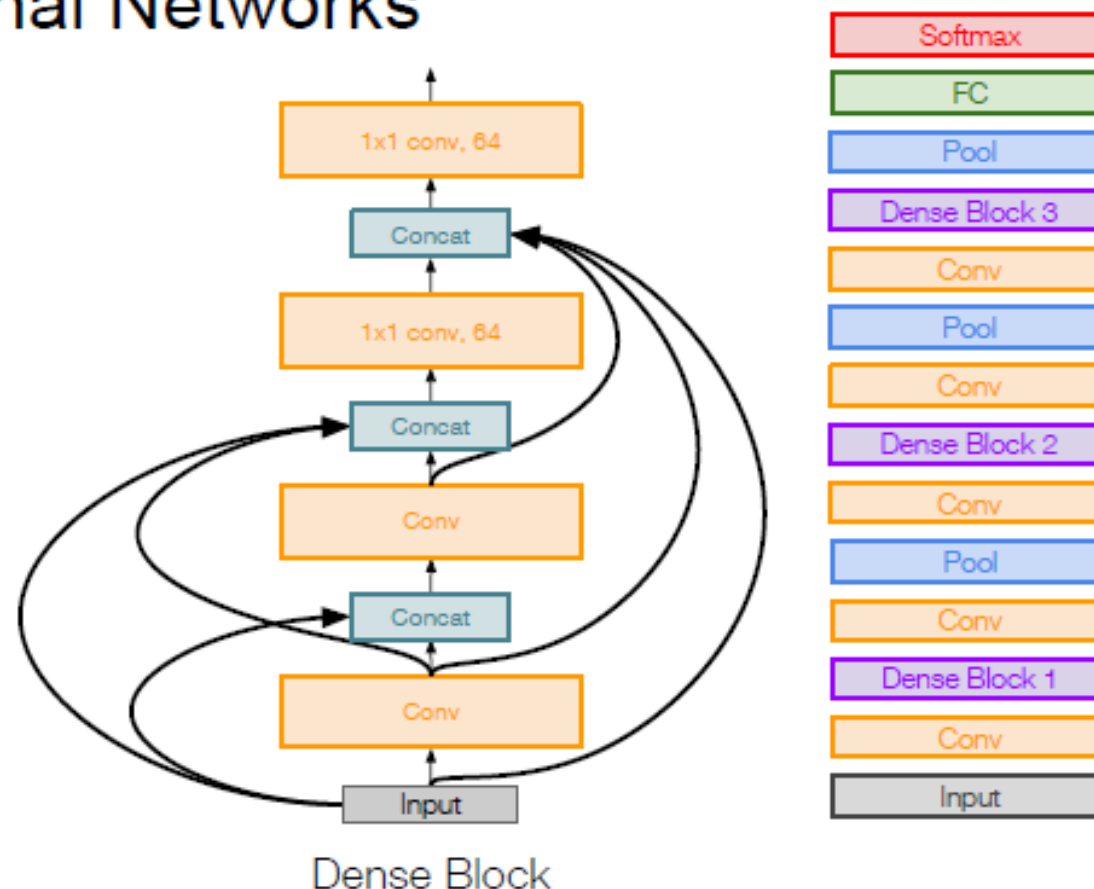- Use full deep network at test time

# DenseNet

## Densely Connected Convolutional Networks

*[Huang et al. 2017]*

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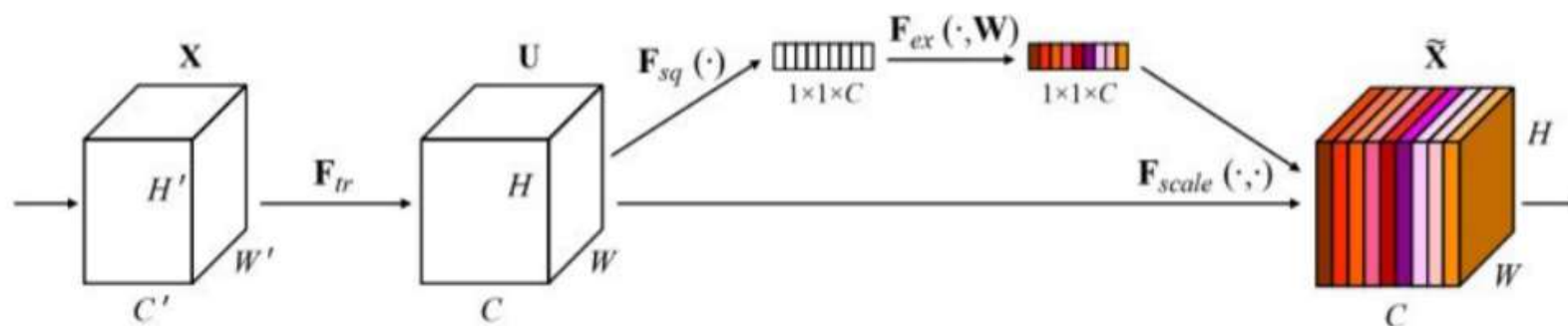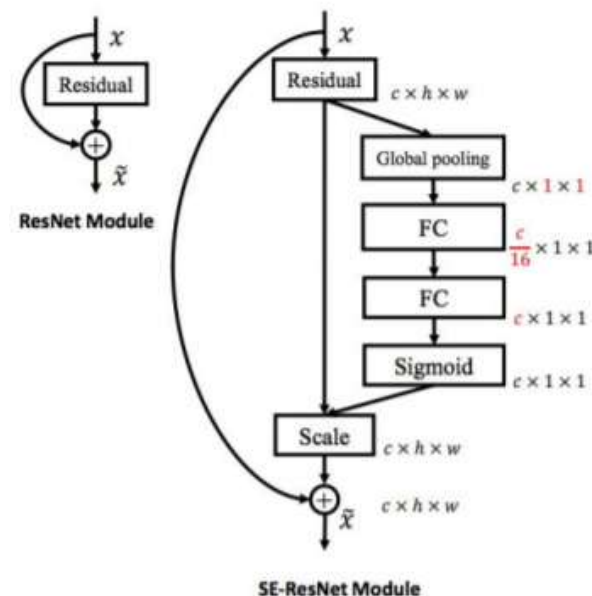- Alleviates vanishing gradient, strengthens feature propagation, encourages feature reuse



Dense Block

**48**

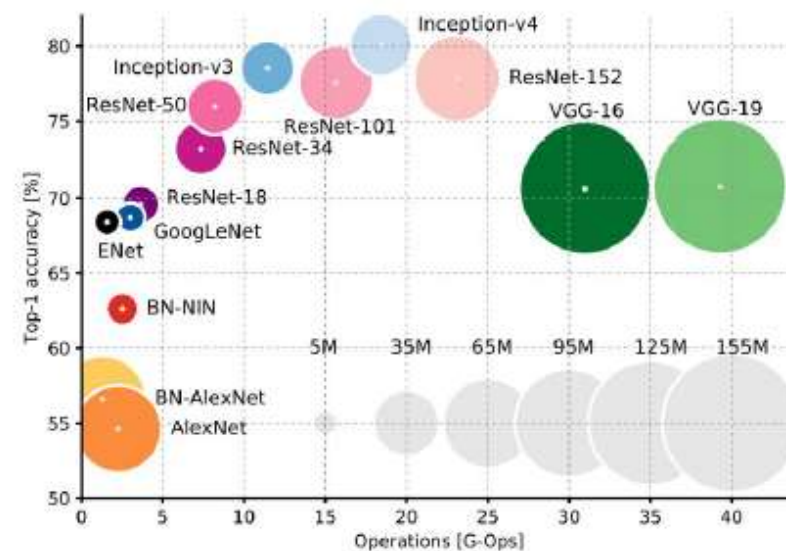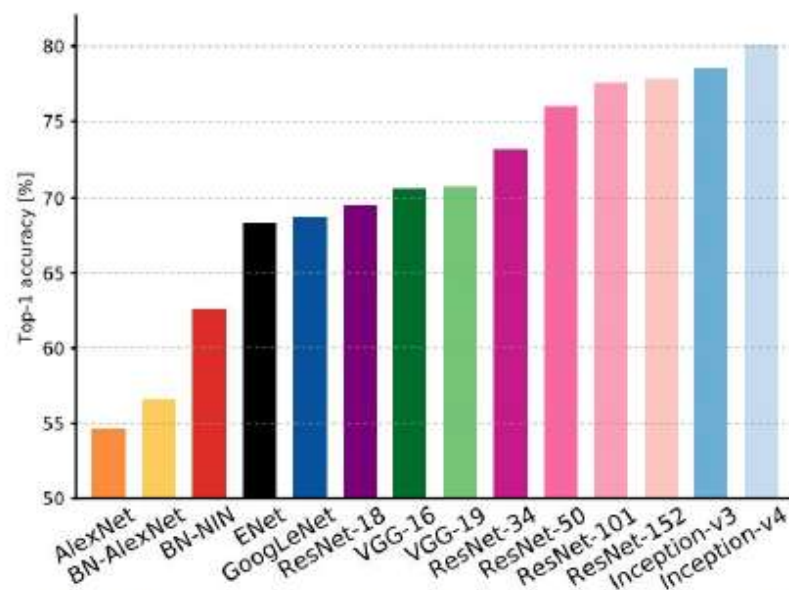# Squeeze-and-Excitation Networks (SENet)

*[Hu et al. 2017]*

- Add a "feature recalibration" module that learns to adaptively reweight feature maps
- Global information (global avg. pooling layer) + 2 FC layers used to determine feature map weights
- ILSVRC'17 classification winner (using ResNeXt-152 as a base architecture)
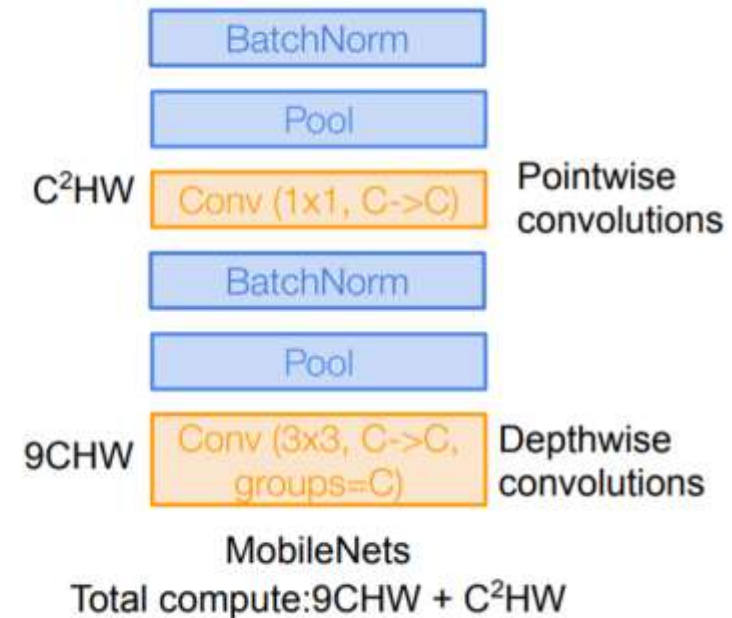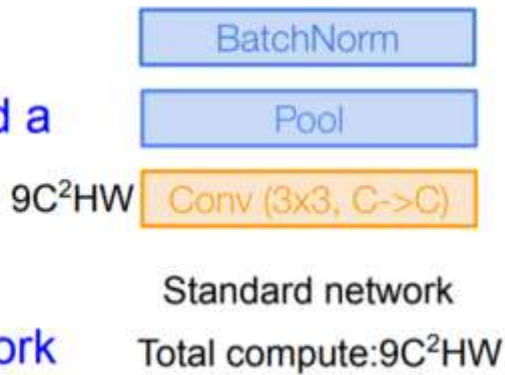
# Model complexity

Comparing complexity...

An Analysis of Deep Neural Network Models for Practical Applications, 2017.

# Efficient networks

- **MobileNets: Efficient Convolutional Neural Networks for Mobile Applications**
  [Howard et al. 2017]

- Depthwise separable convolutions replace standard convolutions by factorizing them into a depthwise convolution and a 1x1 convolution
- Much more efficient, with little loss in accuracy
- Follow-up MobileNetV2 work in 2018 (Sandler et al.)
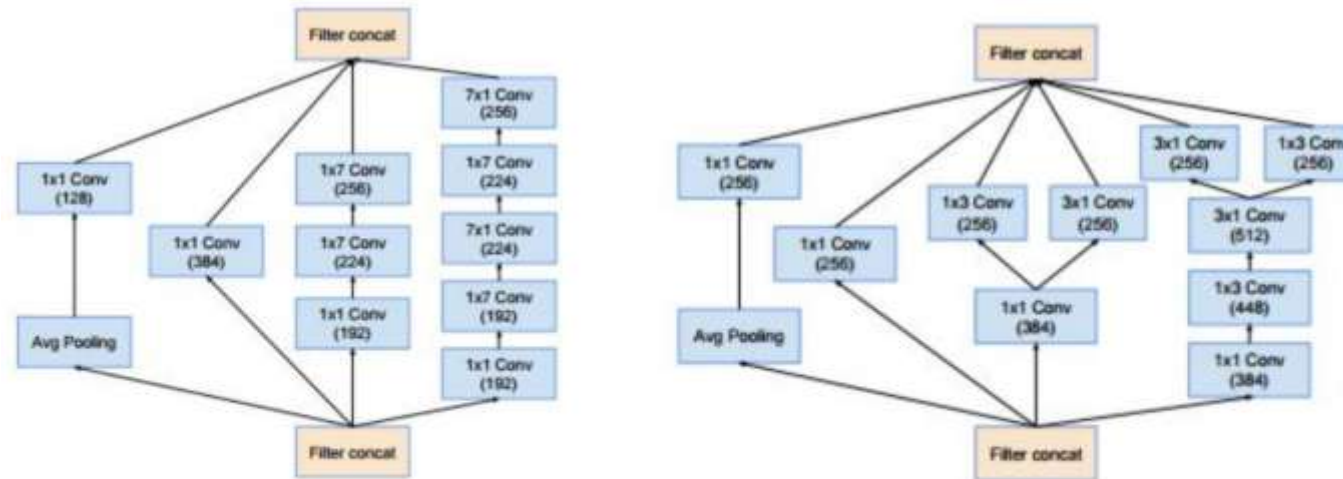- ShuffleNet: Zhang et al, CVPR 2018

| BatchNorm |
|---|
| Pool |
| $9C^2HW$ | Conv (3x3, C->C) |

Standard network
Total compute: $9C^2HW$

| BatchNorm |
|---|
| Pool |
| $C^2HW$ | Conv (1x1, C->C) | Pointwise convolutions |
| BatchNorm |
| Pool |
| $9CHW$ | Conv (3x3, C->C, groups=C) | Depthwise convolutions |

MobileNets
Total compute: $9CHW + C^2HW$

- **CNN architectures**

  - ☐ Sequential structure: LeNet/AlexNet/VGGNet

  - ☐ Parallel branches: GoogLeNet

  - ☐ Residual structure: ResNet/DenseNet

  - ☐ **Network Architecture Search**

*Acknowledgement: Zemel et al's CSC411 and Feifei Li et al's cs231n notes*

- **Problems with network architecture**
  - Designing NA is hard
  - Lots of human efforts go into tuning them
  - Not a lot of intuition into how to design them well
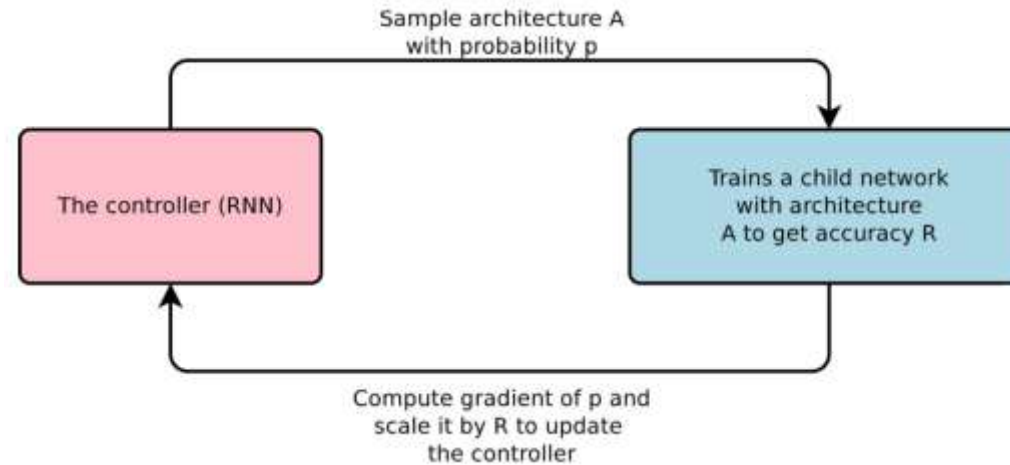  - Can we learn good architectures automatically?



Two layers from the famous Inception V4 computer vision model.
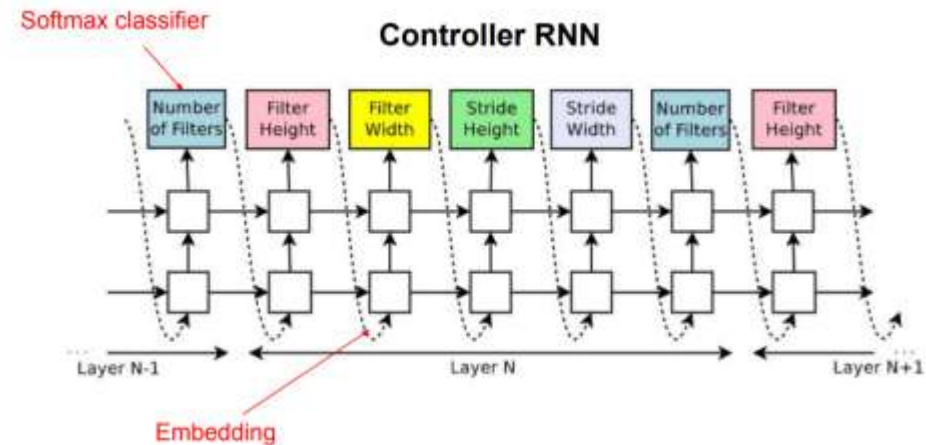Szegedy et al, 2017

# Network Architecture

- Neural architecture search (Zoph and Le, ICLR 2016)

Sample architecture A
with probability p

The controller (RNN)

Trains a child network
with architecture
A to get accuracy R

Compute gradient of p and
scale it by R to update
the controller

- "Controller" network that learns to design a good network architecture (output a string corresponding to network design)
- Iterate:
  1) Sample an architecture from search space
  2) Train the architecture to get a "reward" R corresponding to accuracy
  3) Compute gradient of sample probability, and scale by R to perform controller parameter update (i.e. increase likelihood of good architecture being sampled, decrease likelihood of bad architecture)

Softmax classifier

**Controller RNN**

| Number of Filters | Filter Height | Filter Width | Stride Height | Stride Width | Number of Filters | Filter Height |

Layer N-1          Layer N          Layer N+1

Embedding

# Network structure summary

- AlexNet showed that you can use CNNs to train Computer Vision models.

- ZFNet, VGG shows that bigger networks work better

- GoogLeNet is one of the first to focus on efficiency using 1x1 bottleneck convolutions and global avg pool instead of FC layers

- ResNet showed us how to train extremely deep networks

  - Limited only by GPU & memory!

  - Showed diminishing returns as networks got bigger

- After ResNet: CNNs were *better than the human metric* and focus shifted to Efficient networks:

  - Lots of tiny networks aimed at mobile devices: MobileNet, ShuffleNet

- Neural Architecture Search can now automate architecture design

# More on classification

- Is image classification a solved problem?
  - □ "(Super-)Human level" performance on some benchmarks
    - Face identification
    - ImageNet 1000 classes
- But compared to human vision…
  - □ Limitations in learning
    - We can learn new classes using one or two examples
    - We can also handle label noises
    - We can generalize to unfamiliar scenes
  - □ Limitation in prediction
    - We can also predict the uncertainty
    - We can easily handle adversarial examples
    - We are much more efficient in power consumption

■ What is semantic segmentation?

■ Network architecture for semantic segmentation

    ☐ Main idea for dense prediction

    ☐ Fully convolutional network

    ☐ Upsampling operators

    ☐ Multiscale context modeling

■ Network training losses

*Acknowledgement:  Feifei Li et al's cs231n notes*

- In general, our goal is to learn a mapping from a signal to a 'semantically meaningful' representation.
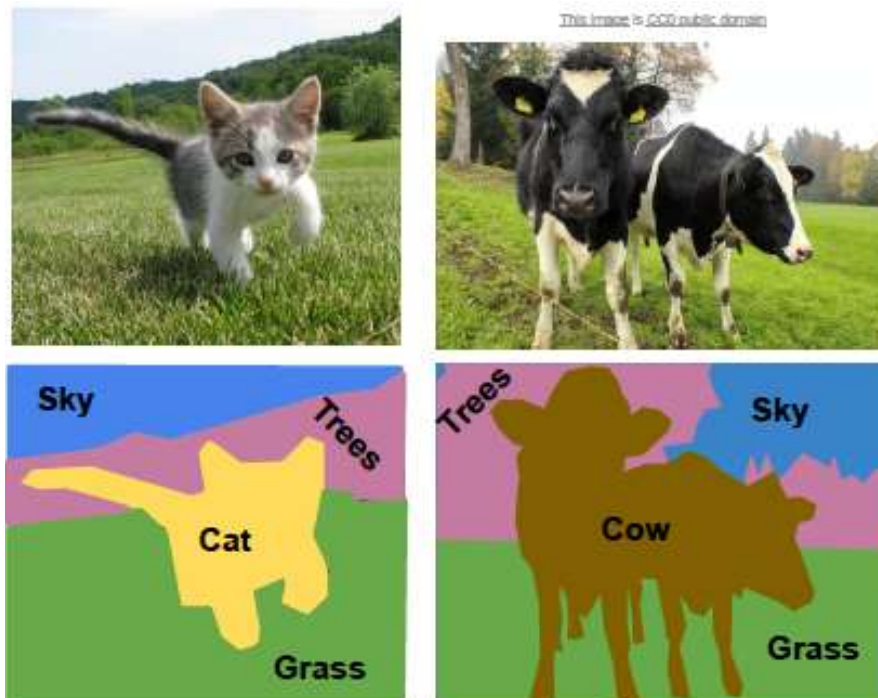    - □ Output can have many different forms:



**red panda** (*Ailurus fulgens*)

# Semantic Segmentation

- Problem setup
  - Label each pixel in the image with an object category label
  - Do not differentiate object instances

# Key to many applications

- Autonomous robots and cars

- Safety and security

- Medical analysis and health

Multi-organ abdominal
CT segmentation

- etc…

Reference standard    NiftyNet segmentation

# Semantic Segmentation

- **Problem formulation**
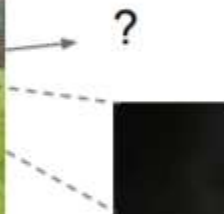  - □ Pixel-wise object classification task



  - □ One-hot encoding



https://www.jeremyjordan.me/semantic-segmentation/

61

- A naïve approach

Full image

?

Impossible to classify without context

Q: how do we include context?

- A naïve approach

- A naïve approach



Problem: Very inefficient! Not reusing shared features between overlapping patches

Farabet et al, "Learning Hierarchical Features for Scene Labeling," TPAMI 2013
Pinheiro and Collobert, "Recurrent Convolutional Neural Networks for Scene Labeling", ICML 2014

■ **Main idea for dense prediction**
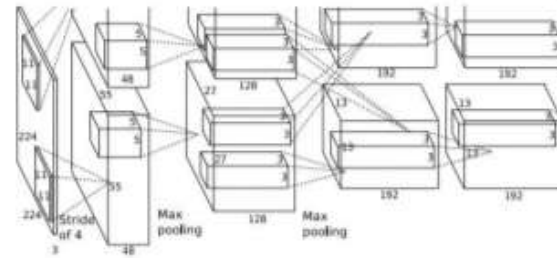
■ Fully convolutional network

■ Upsampling operators

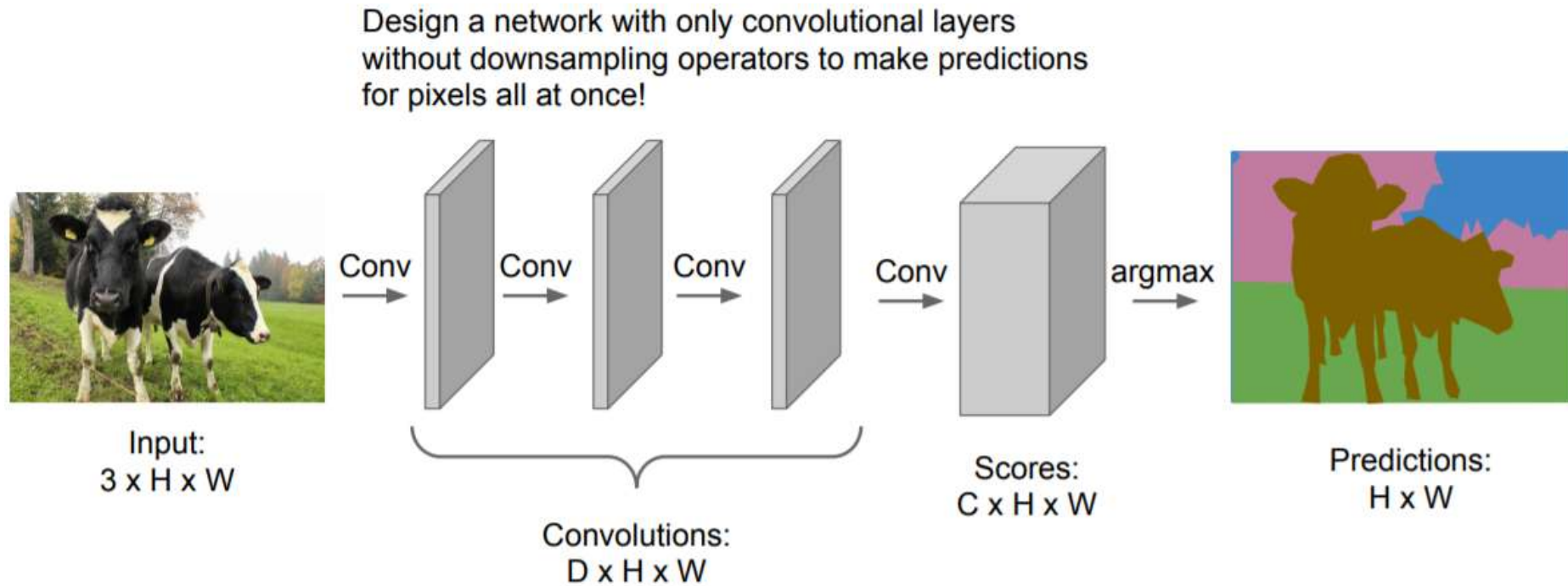■ Multiscale context modeling

- First idea



Full image

An intuitive idea: encode the entire image with conv net, and do semantic segmentation on top.

Problem: classification architectures often reduce feature spatial sizes to go deeper, but semantic segmentation requires the output size to be the same as input size.
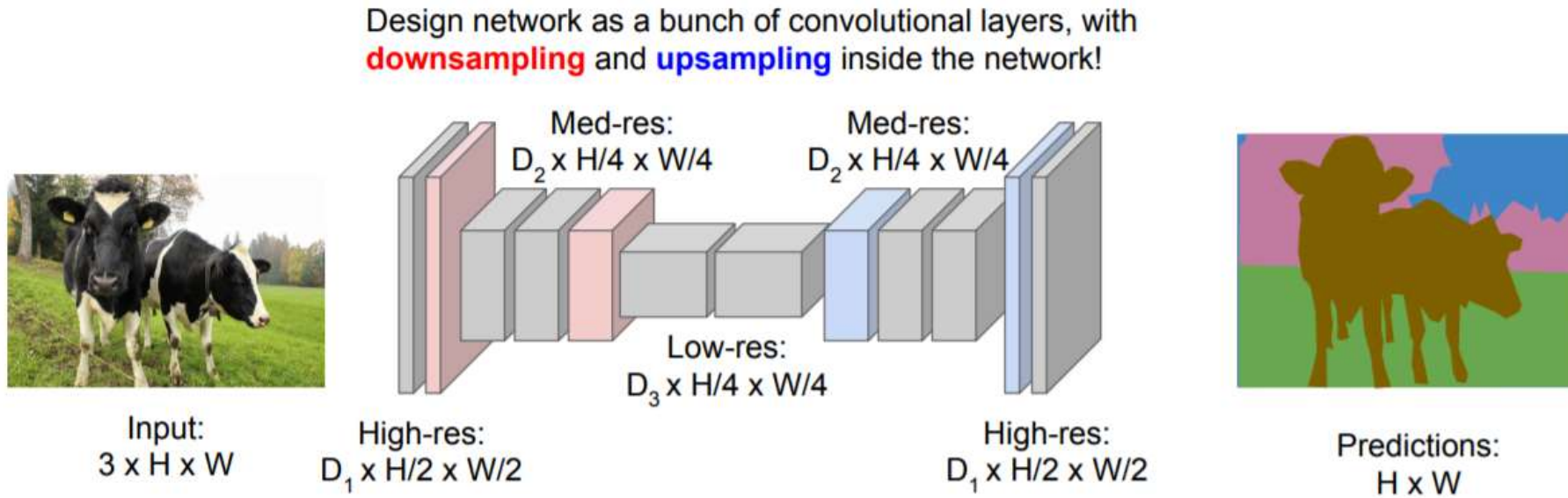
- Second idea



Design a network with only convolutional layers without downsampling operators to make predictions for pixels all at once!

Input: 3 x H x W

Conv → Conv → Conv → Conv → argmax

Convolutions: D x H x W

Scores: C x H x W

Predictions: H x W

■ Second idea improved



Design network as a bunch of convolutional layers, with **downsampling** and **upsampling** inside the network!

Input: $3 \times H \times W$

High-res: $D_1 \times H/2 \times W/2$

Med-res: $D_2 \times H/4 \times W/4$

Low-res: $D_3 \times H/4 \times W/4$

Med-res: $D_2 \times H/4 \times W/4$

High-res: $D_1 \times H/2 \times W/2$

Predictions: $H \times W$

上海科技大学
ShanghaiTech University

- Third idea



Input:
3 x H x W

High-res:
$D_1$ x H/2 x W/2

Med-res:
$D_2$ x H/4 x W/4

Med-res:
$D_2$ x H/4 x W/4

Low-res:
$D_3$ x H/4 x W/4

High-res:
$D_1$ x H/2 x W/2

Predictions:
H x W

上海科技大学
ShanghaiTech University

- Main idea for dense predictionedicti1o

- **Fully convolutional network**

- Upsampling operators

- Multiscale context modeling
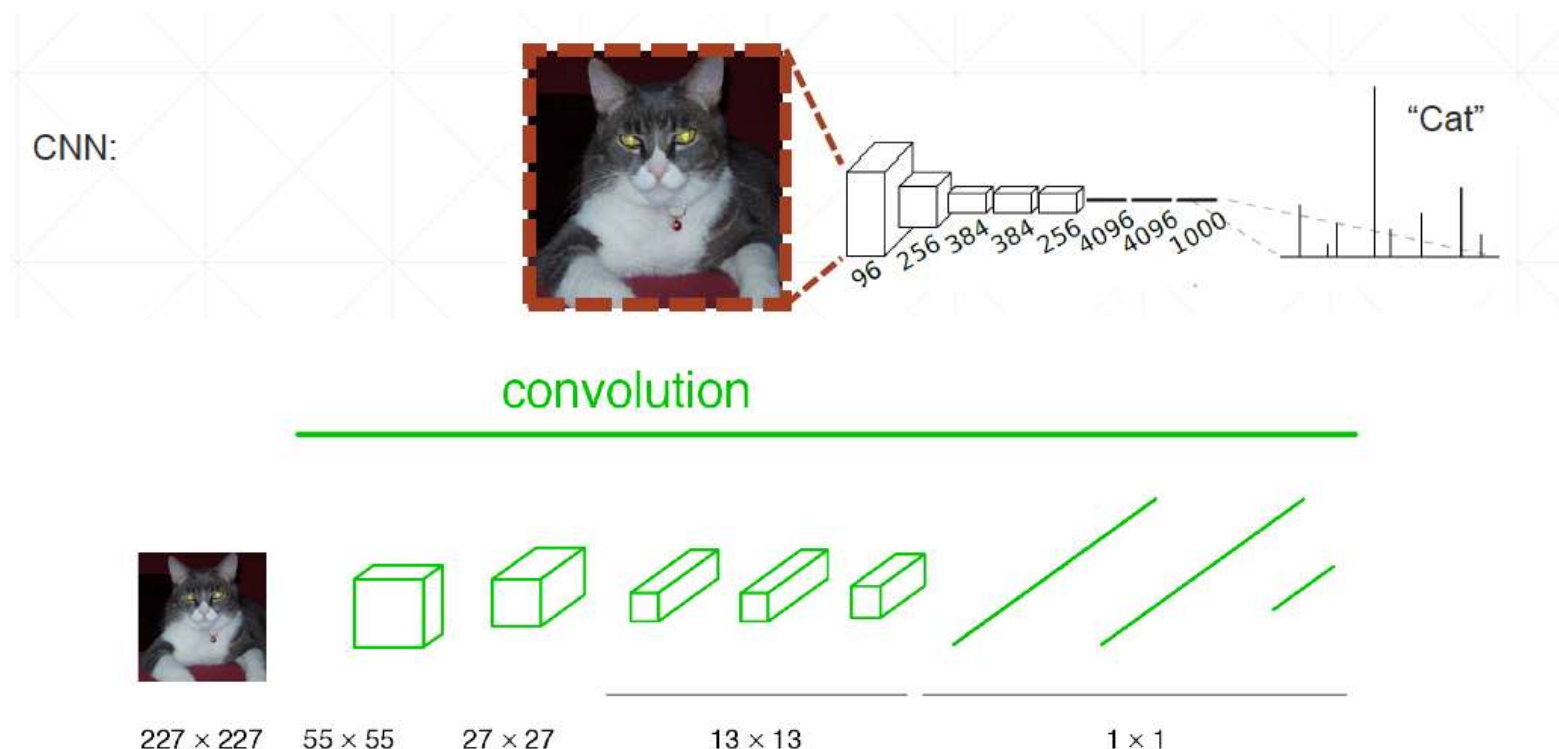
*Acknowledgement:  Feifei Li et al's cs231n notes*

- Starting from a classification network



CNN:

"Cat"

96 256 384 384 256 4096 4096 1000

convolution                    fully connected

"tabby cat"

227 × 227    55 × 55    27 × 27         13 × 13

上海科技大学
ShanghaiTech University

- Interpreting fully connected layers as 1x1 convolution (after reshaping)

- Extending to a complete image

- Extending to a complete image



- Keep kernel sizes and strides
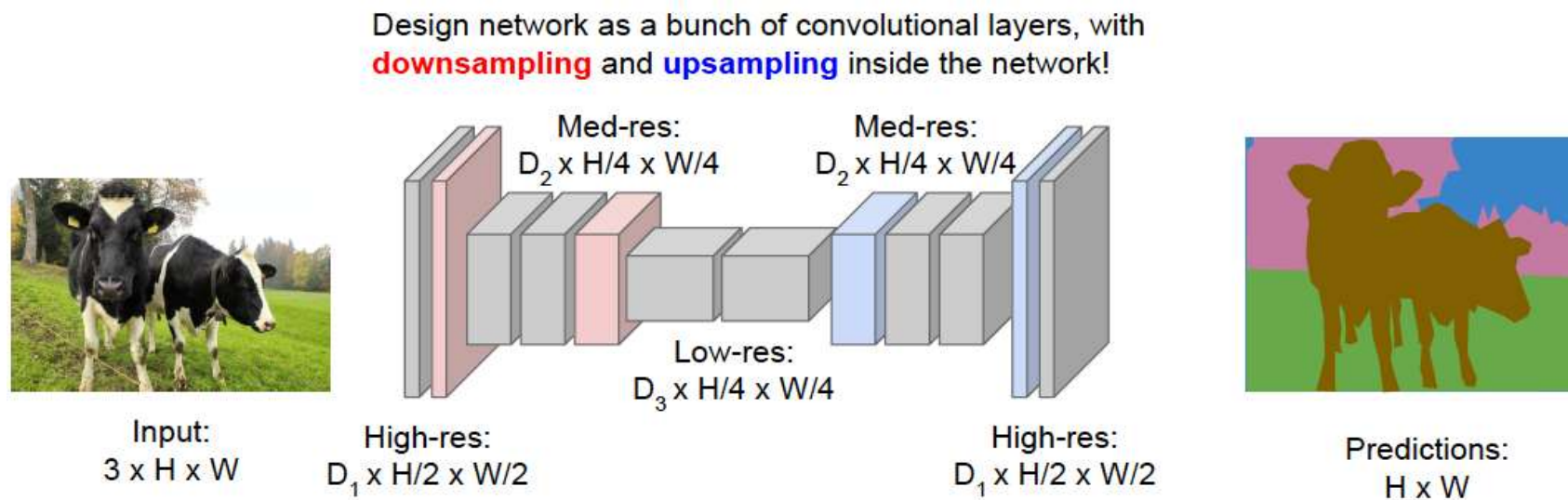- Replace dense layer with convolution

上 海 科 技 大 学
ShanghaiTech University

- Main idea for dense predictionedicti1o

- Fully convolutional network

- Upsampling operators

- Multiscale context modeling

*Acknowledgement:  Feifei Li et al's cs231n notes*

上 海 科 技 大 学
ShanghaiTech University

- General encoder-decoder architecture



Design network as a bunch of convolutional layers, with
**downsampling** and **upsampling** inside the network!

Med-res: $D_2 \times H/4 \times W/4$

Med-res: $D_2 \times H/4 \times W/4$

Low-res: $D_3 \times H/4 \times W/4$

Input: $3 \times H \times W$

High-res: $D_1 \times H/2 \times W/2$

High-res: $D_1 \times H/2 \times W/2$

Predictions: $H \times W$

**76**

# In-Network upsampling

- Unpooling

**Nearest Neighbor**

| 1 | 2 |
|---|---|
| 3 | 4 |

→

| 1 | 1 | 2 | 2 |
|---|---|---|---|
| 1 | 1 | 2 | 2 |
| 3 | 3 | 4 | 4 |
| 3 | 3 | 4 | 4 |

Input: 2 x 2      Output: 4 x 4

**"Bed of Nails"**

| 1 | 2 |
|---|---|
| 3 | 4 |

→

| 1 | 0 | 2 | 0 |
|---|---|---|---|
| 0 | 0 | 0 | 0 |
| 3 | 0 | 4 | 0 |
| 0 | 0 | 0 | 0 |

Input: 2 x 2      Output: 4 x 4

# In-Network upsampling

- Max Unpooling



**Max Pooling**
Remember which element was max!

Input: 4 x 4    Output: 2 x 2

Rest of the network

**Max Unpooling**
Use positions from pooling layer

Input: 2 x 2    Output: 4 x 4

Corresponding pairs of downsampling and upsampling layers

■ Learnable Upsampling: Transpose convolution

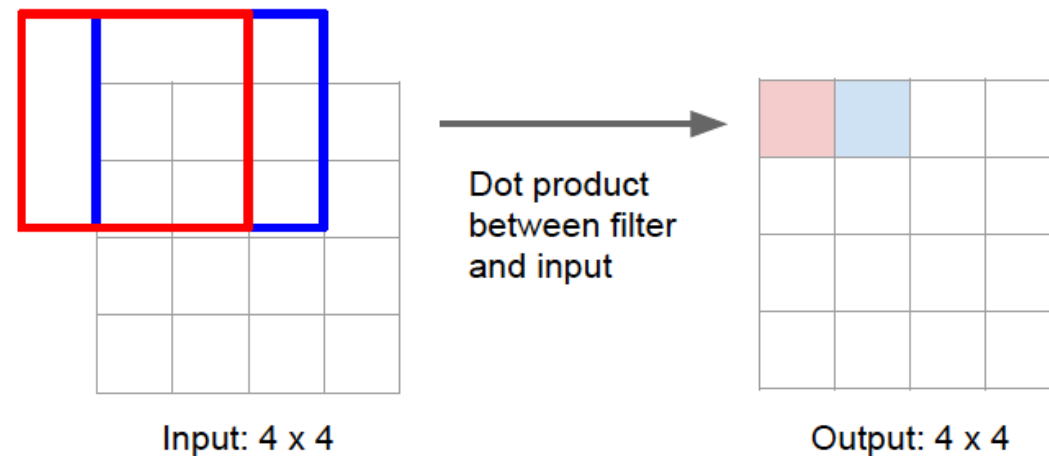**Recall:**Typical 3 x 3 convolution, stride 1 pad 1

Input: 4 x 4

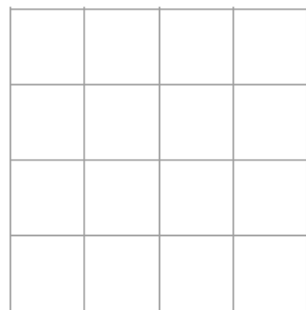Output: 4 x 4

■ Learnable Upsampling: Transpose convolution



**Recall:** Normal 3 x 3 convolution, stride 1 pad 1

Dot product between filter and input

Input: 4 x 4

Output: 4 x 4

# In-Network upsampling

- **Learnable Upsampling: Transpose convolution**

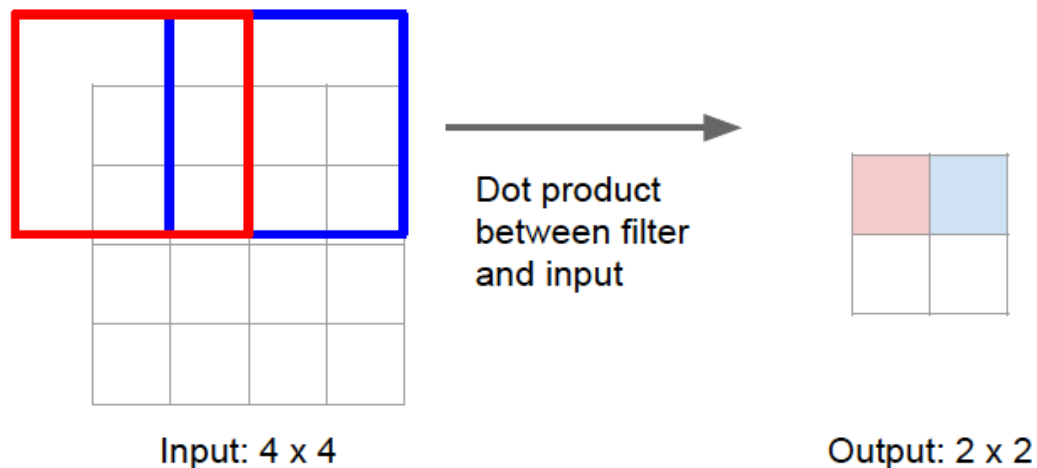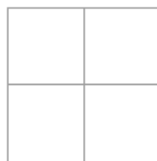**Recall:** Normal 3 x 3 convolution, <u>stride 2</u> pad 1

Input: 4 x 4                    Output: 2 x 2

# In-Network upsampling

- Learnable Upsampling: Transpose convolution

**Recall:** Normal 3 x 3 convolution, <u>stride 2</u> pad 1

Dot product between filter and input

Input: 4 x 4

Output: 2 x 2

Filter moves 2 pixels in the input for every one pixel in the output

Stride gives ratio between movement in input and output

# In-Network upsampling

- Learnable Upsampling: Transpose convolution

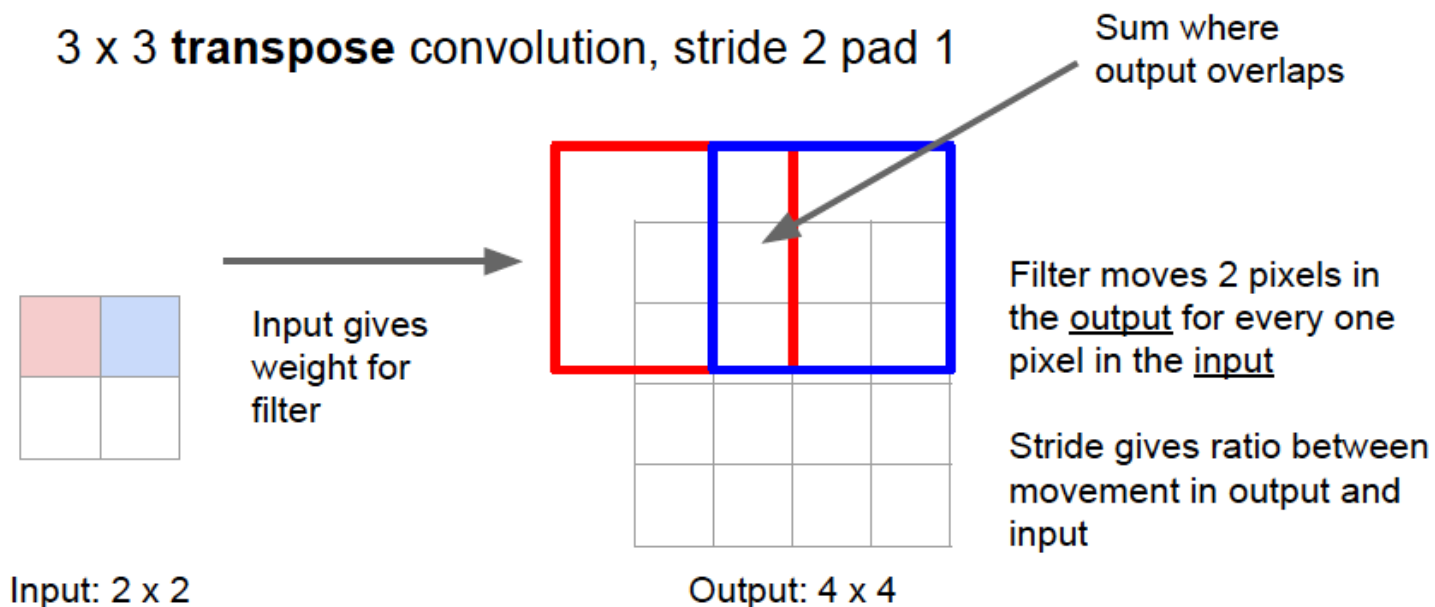3 x 3 **transpose** convolution, stride 2 pad 1

Input: 2 x 2

Output: 4 x 4

# In-Network upsampling
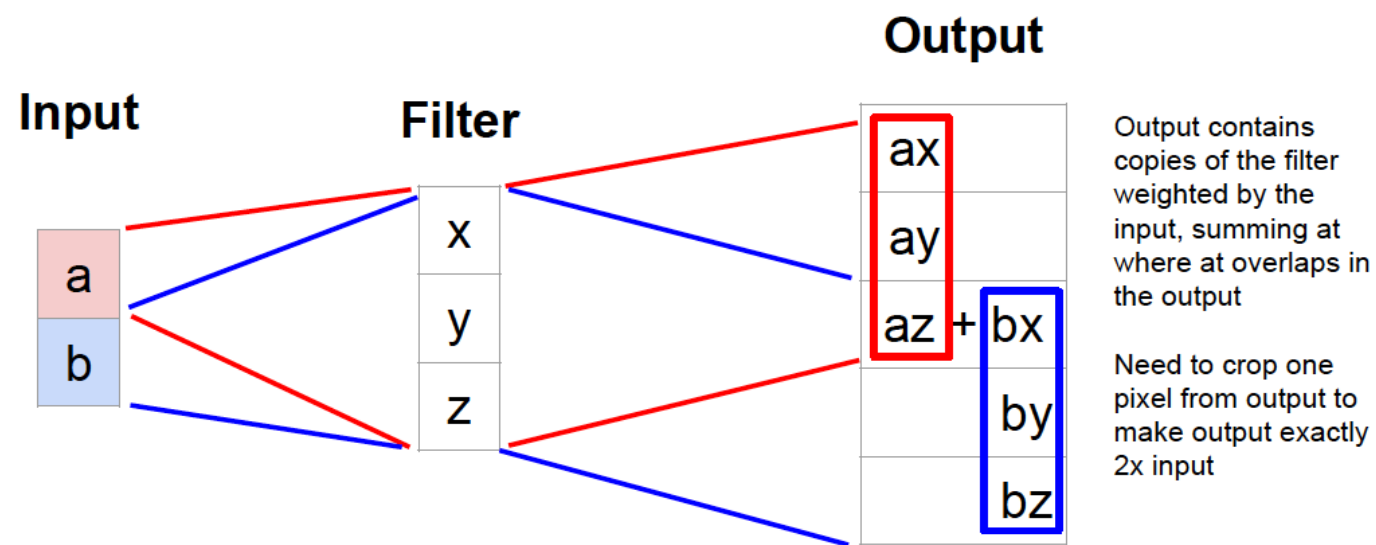
- Learnable Upsampling: Transpose convolution

3 x 3 **transpose** convolution, stride 2 pad 1

Other names:
-Deconvolution (bad)
-Upconvolution
-Fractionally strided convolution
-Backward strided convolution

Input gives weight for filter

Sum where output overlaps

Filter moves 2 pixels in the output for every one pixel in the input

Stride gives ratio between movement in output and input

Input: 2 x 2

Output: 4 x 4

- Learnable Upsampling: Transpose convolution
  - 1D example



Output contains copies of the filter weighted by the input, summing at where at overlaps in the output

Need to crop one pixel from output to make output exactly 2x input

# In-Network upsampling

- **Learnable Upsampling: Transpose convolution**
  - □ 1D example

We can express convolution in terms of a matrix multiplication

$$\vec{x} * \vec{a} = X\vec{a}$$

$$\begin{bmatrix} x & y & z & 0 & 0 & 0 \\ 0 & 0 & x & y & z & 0 \end{bmatrix} \begin{bmatrix} 0 \\ a \\ b \\ c \\ d \\ 0 \end{bmatrix} = \begin{bmatrix} ay + bz \\ bx + cy + dz \end{bmatrix}$$

Example: 1D conv, kernel size=3, stride=2, padding=1

Convolution transpose multiplies by the transpose of the same matrix:

$$\vec{x} *^T \vec{a} = X^T \vec{a}$$

$$\begin{bmatrix} x & 0 \\ y & 0 \\ z & x \\ 0 & y \\ 0 & z \\ 0 & 0 \end{bmatrix} \begin{bmatrix} a \\ b \end{bmatrix} = \begin{bmatrix} ax \\ ay \\ az + bx \\ by \\ bz \\ 0 \end{bmatrix}$$

When stride>1, convolution transpose is no longer a normal convolution!

**上海科技大学**
**ShanghaiTech University**

- **Fully Convolutional Network** [Long et al, CVPR 2015]
  - ☐ Upsampling: low-resolution, lack spatial details
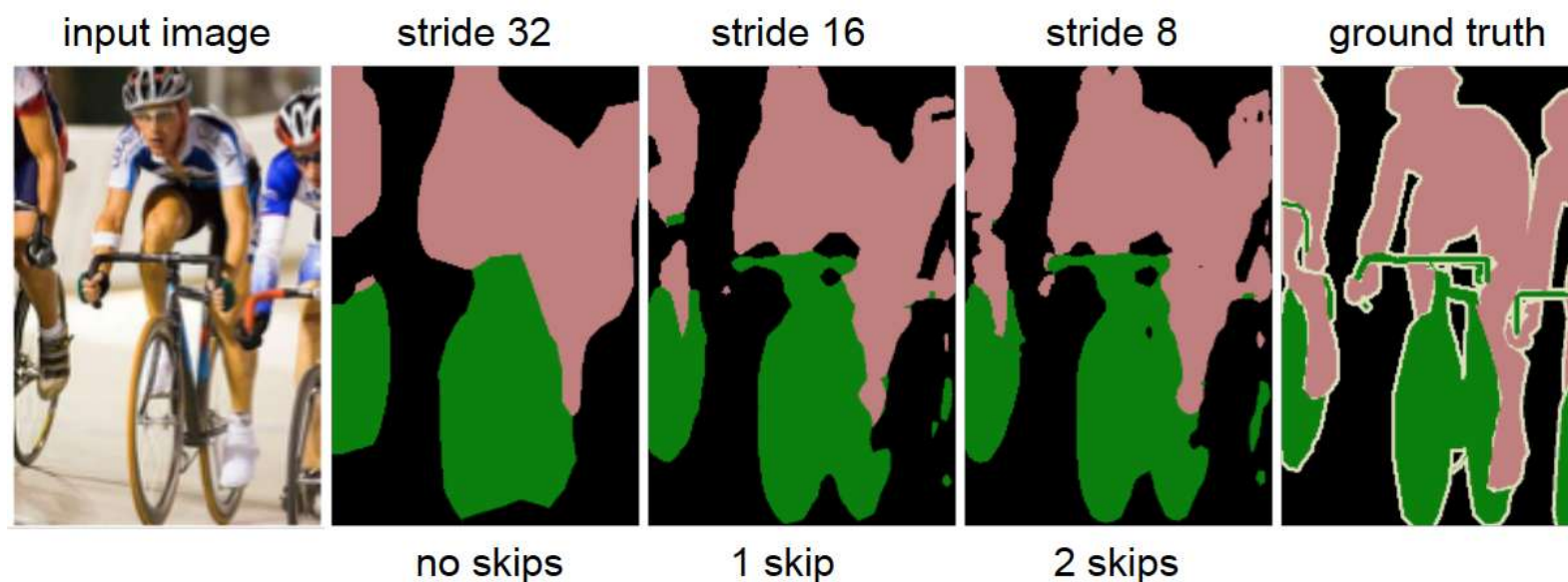  - ☐ Combining *where (local, shallow)* with *what (global, deep)*



image                    intermediate layers

fuse features into **deep jet**

上海科技大学
ShanghaiTech University

- Fully Convolutional Network [Long et al, CVPR 2015]
  - Upsampling: low-resolution, lack spatial details
  - Introducing **skip layers**

上海科技大学
ShanghaiTech University

- **Fully Convolutional Network** [Long et al, CVPR 2015]
  - ☐ Upsampling: low-resolution, lack spatial details
  - ☐ Skip layer refinement

| input image | stride 32 | stride 16 | stride 8 | ground truth |
| --- | --- | --- | --- | --- |
| | no skips | 1 skip | 2 skips | |

- U-Net [Ronneberger et al, MICCAI 2015]
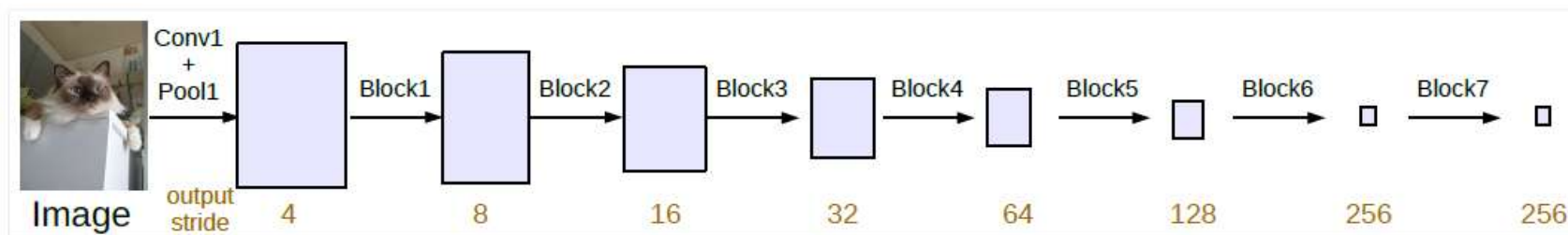
上 海 科 技 大 学
ShanghaiTech University

- **Dilated Convolutional Network** [Yu and Koltun, ICLR 2016]
    - ☐ Dense feature map without upsampling
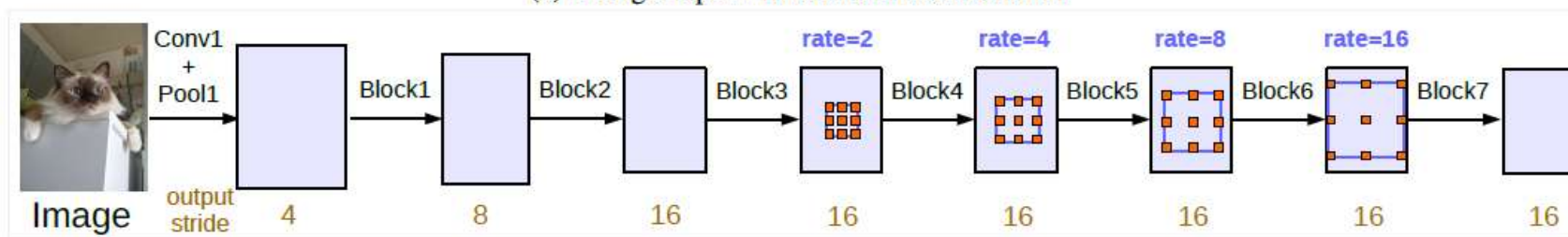    - ☐ ***Dilated (or Atrous) convolution***



$$y[i] = \sum_{k=1}^{K} x[i + r \cdot k] w[k].$$

上海科技大学
ShanghaiTech University

- **Dilated Convolutional Network** [Yu and Koltun, ICLR 2016]
  - □ Dense feature map without upsampling
  - □ ***Dilated (or Atrous) convolution***

$$y[i] = \sum_{k=1}^{K} x[i + r \cdot k]w[k].$$

# Network Design: Spatial resolution

- Dilated Convolutional Network [Yu and Koltun, ICLR 2016]
  - □ Dense feature map without upsampling
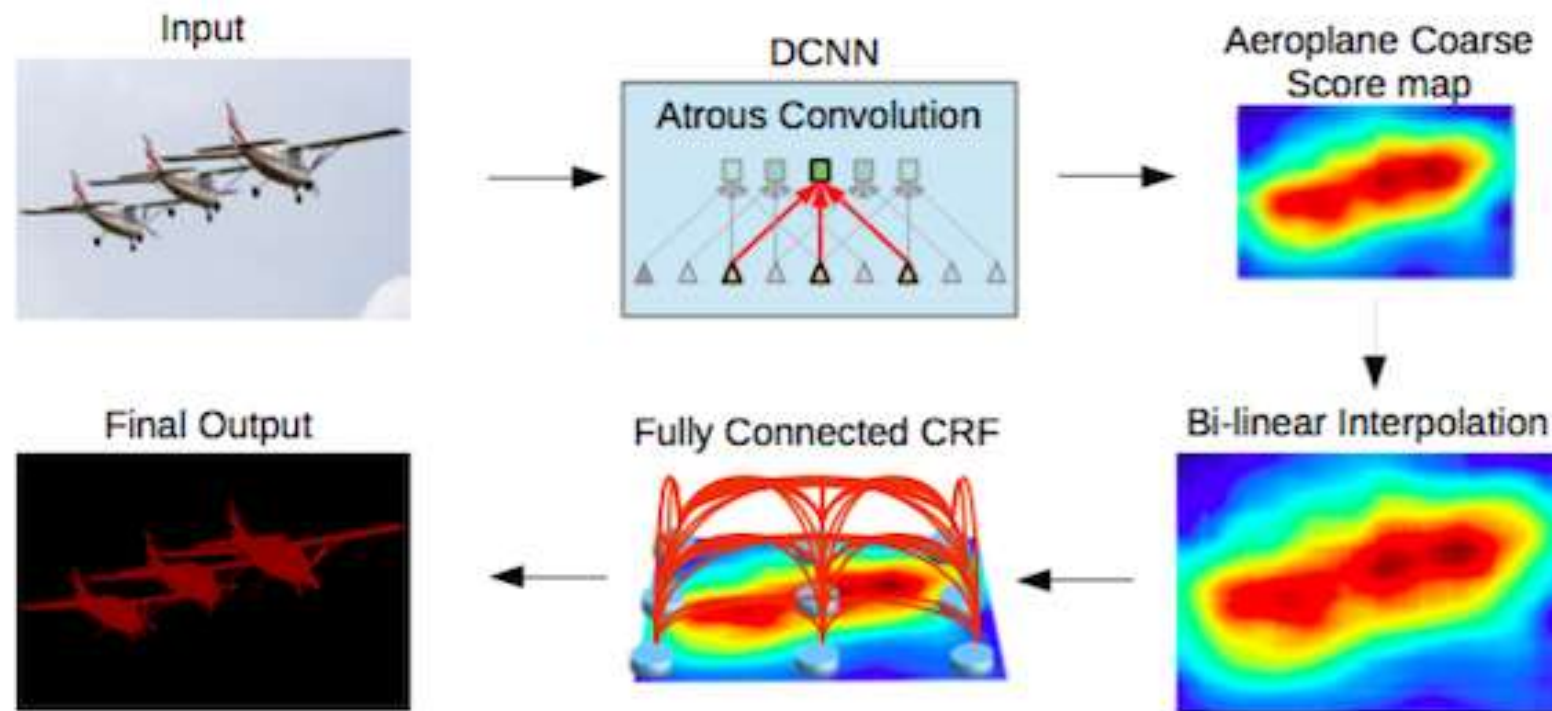  - □ ***Dilated (or Atrous) convolution***

■ Dilated Convolutional Network [Yu and Koltun, ICLR 2016]

　□ Dense feature map without upsampling

　□ ***Dilated (or Atrous) convolution***



(a) Going deeper without atrous convolution.

(b) Going deeper with atrous convolution. Atrous convolution with $rate > 1$ is applied after block3 when $output\_stride = 16$.

上海科技大学
ShanghaiTech University

- **DeepLab v1&v2**
  - ☐ Post-processing with dense CRFs.

上 海 科 技 大 学
ShanghaiTech University

- ## PSPNet [Zhao et al CVPR 2017]
  - □ A pyramid parsing module that carries both local and global context information



(a) Input Image     (b) Feature Map     (c) Pyramid Pooling Module     (d) Final Prediction

- DeepLab v3



- Deeplab v3+



https://ai.googleblog.com/2018/03/semantic-image-segmentation-with.html

■ Main idea: pixel-wise classification
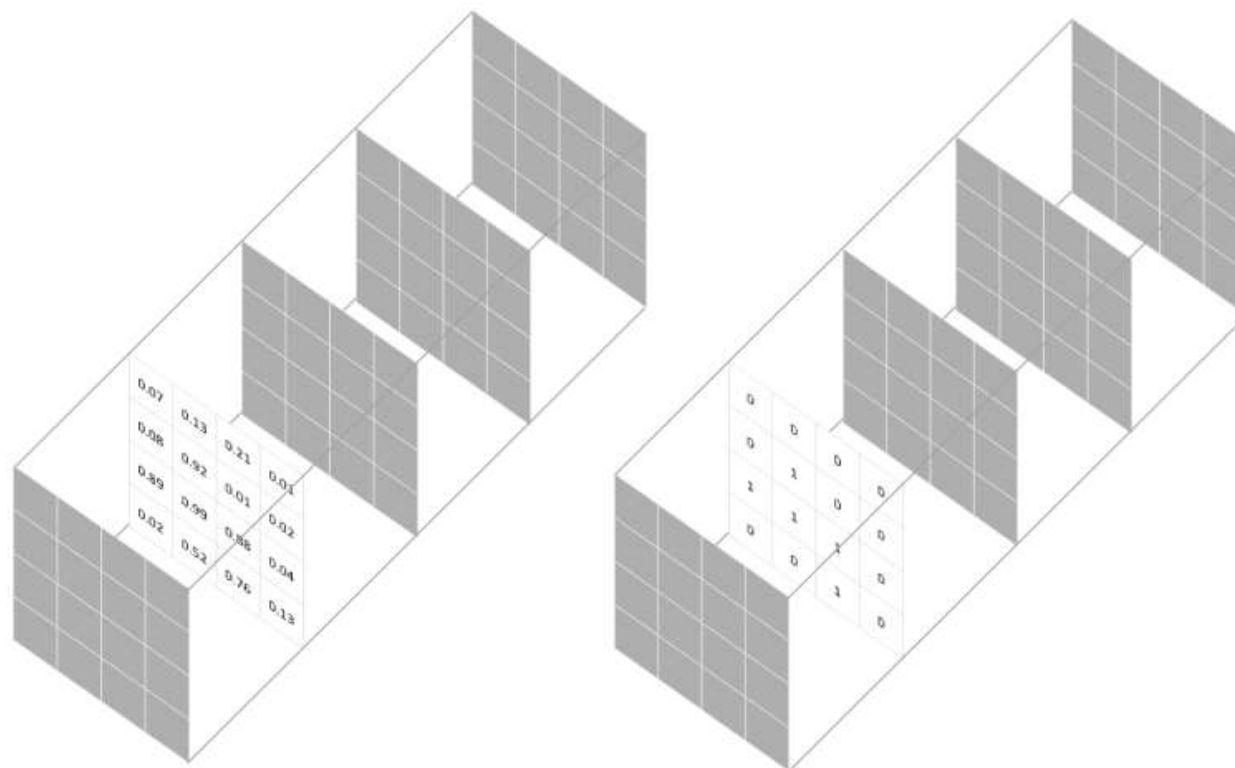
**■ Pixel-wise loss**



Pixel-wise loss is calculated as the log loss, summed over all possible classes

$$-\sum_{classes} y_{true} \log\left(y_{pred}\right)$$

This scoring is repeated over all **pixels** and averaged

Prediction for a selected pixel

Target for the corresponding pixel

# Region-based loss



$$Dice = \frac{2|A \cap B|}{|A| + |B|}$$

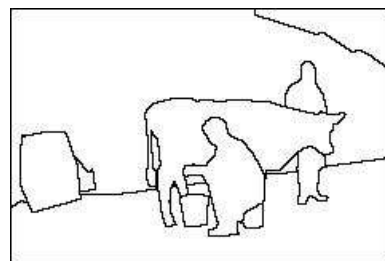Soft Dice coefficient is calculated for each class mask

$$1 - \frac{2\sum\limits_{pixels} y_{true}y_{pred}}{\sum\limits_{pixels} y_{true}^2 + \sum\limits_{pixels} y_{pred}^2}$$

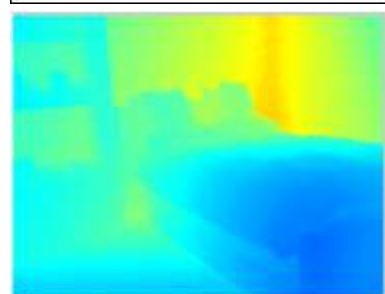This scoring is repeated over all **classes** and averaged

Prediction for a selected class

Target for the corresponding class

# Semantic Segmentation: Summary

上 海 科 技 大 学
ShanghaiTech University

- **Pixel-wise annotation of images**
  - An instance of scene understanding



Boundary

Depth

- **Other research topics (not discussed)**
  - *Low-level vision: superresolution, deblurring, inpainting, depth*
  - *Video: optical flow, action and activity recognition and detection*
  - *Volumetric/Multimodality: RGB-D images, medical imaging, etc.*