# Topics in AI (CPSC 532S):
# Multimodal Learning with Vision, Language and Sound

**Lecture 6: Convolutional Neural Networks (Part 3)**

# Logistics:

**Assignment 2** is due on **Monday**
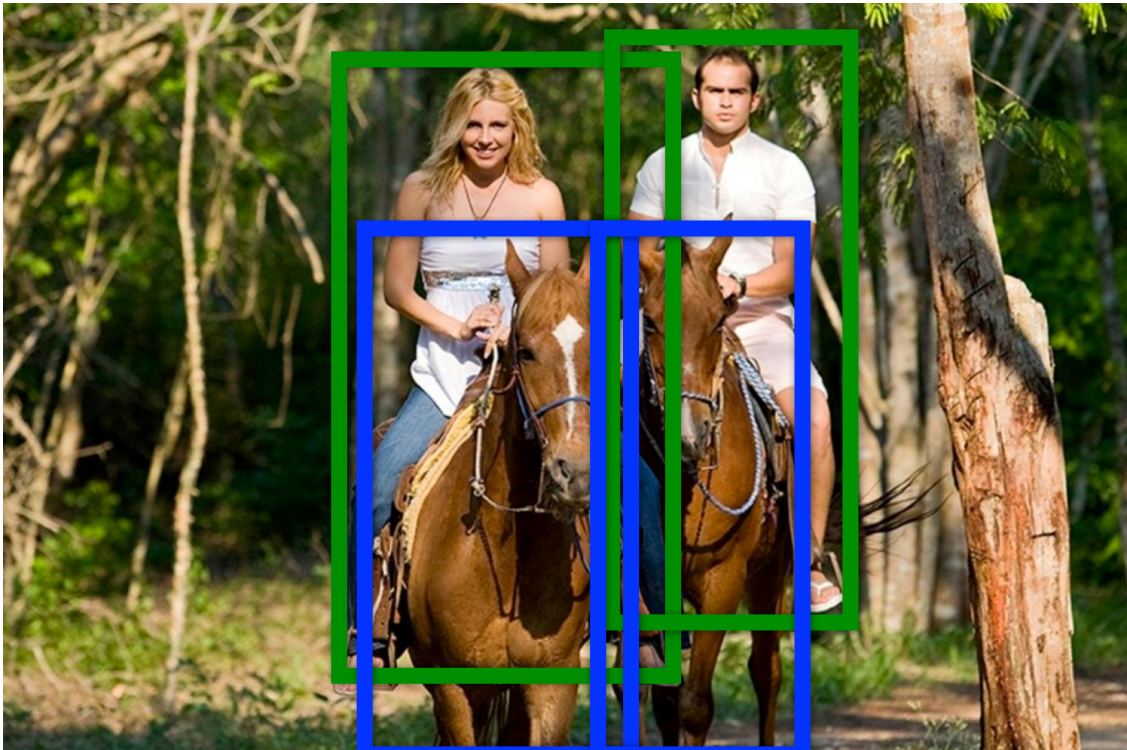
# Computer **Vision Problems** (no language for now)

## Categorization



Multi-**class:**  Horse
Church
Toothbrush
**Person**

IMAGENET

Multi-**label**:  **Horse**
Church
Toothbrush
**Person**

## Detection



Horse (x, y, w, h)
Horse (x, y, w, h)
Person (x, y, w, h)
Person (x, y, w, h)

COCO
Common Objects in Context

## Segmentation



Horse
Person

COCO
Common Objects in Context

## Instance Segmentation



Horse1
Horse2
Person1
Person2

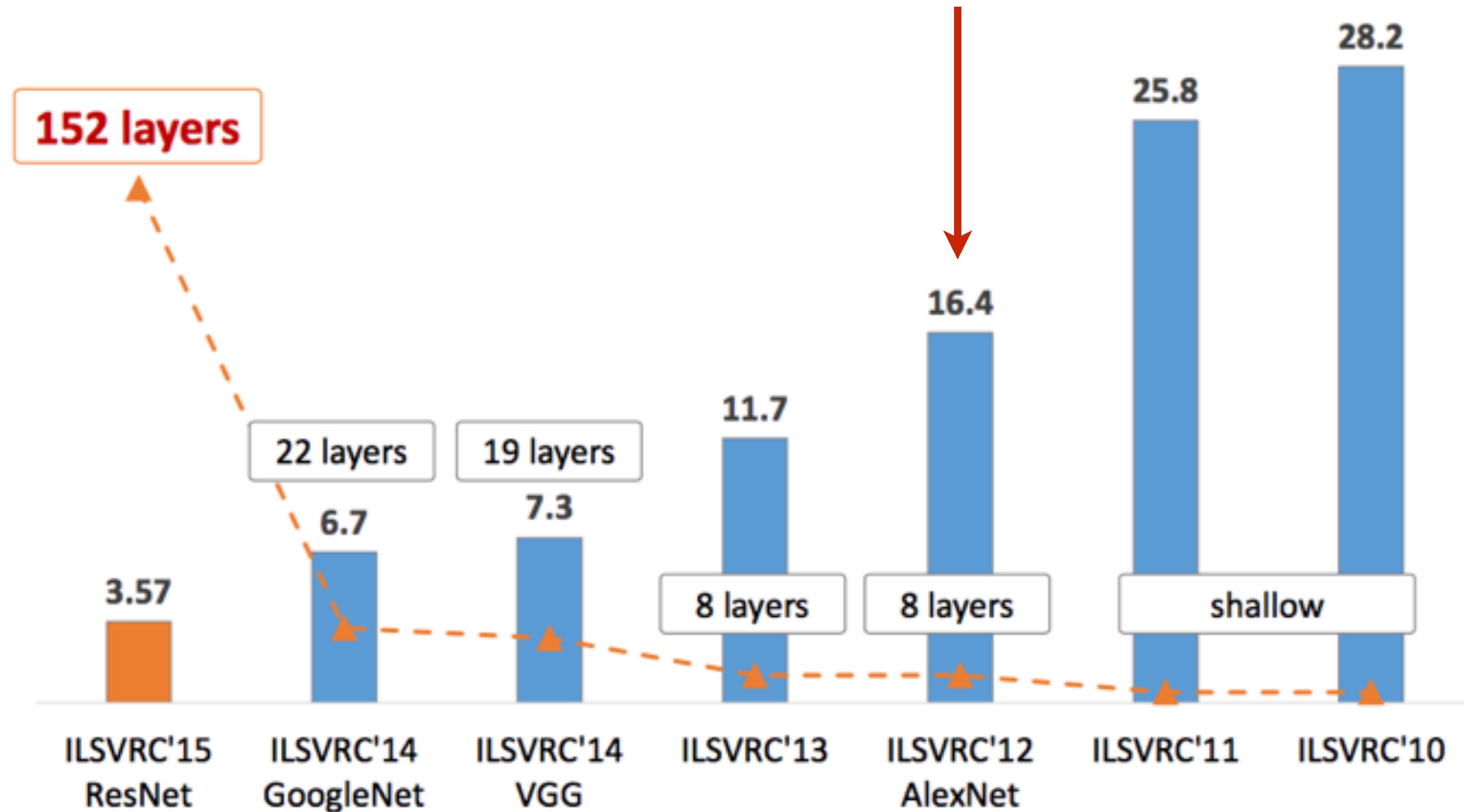# Computer **Vision Problems** (no language for now)

Categorization



Multi-**class:**  Horse
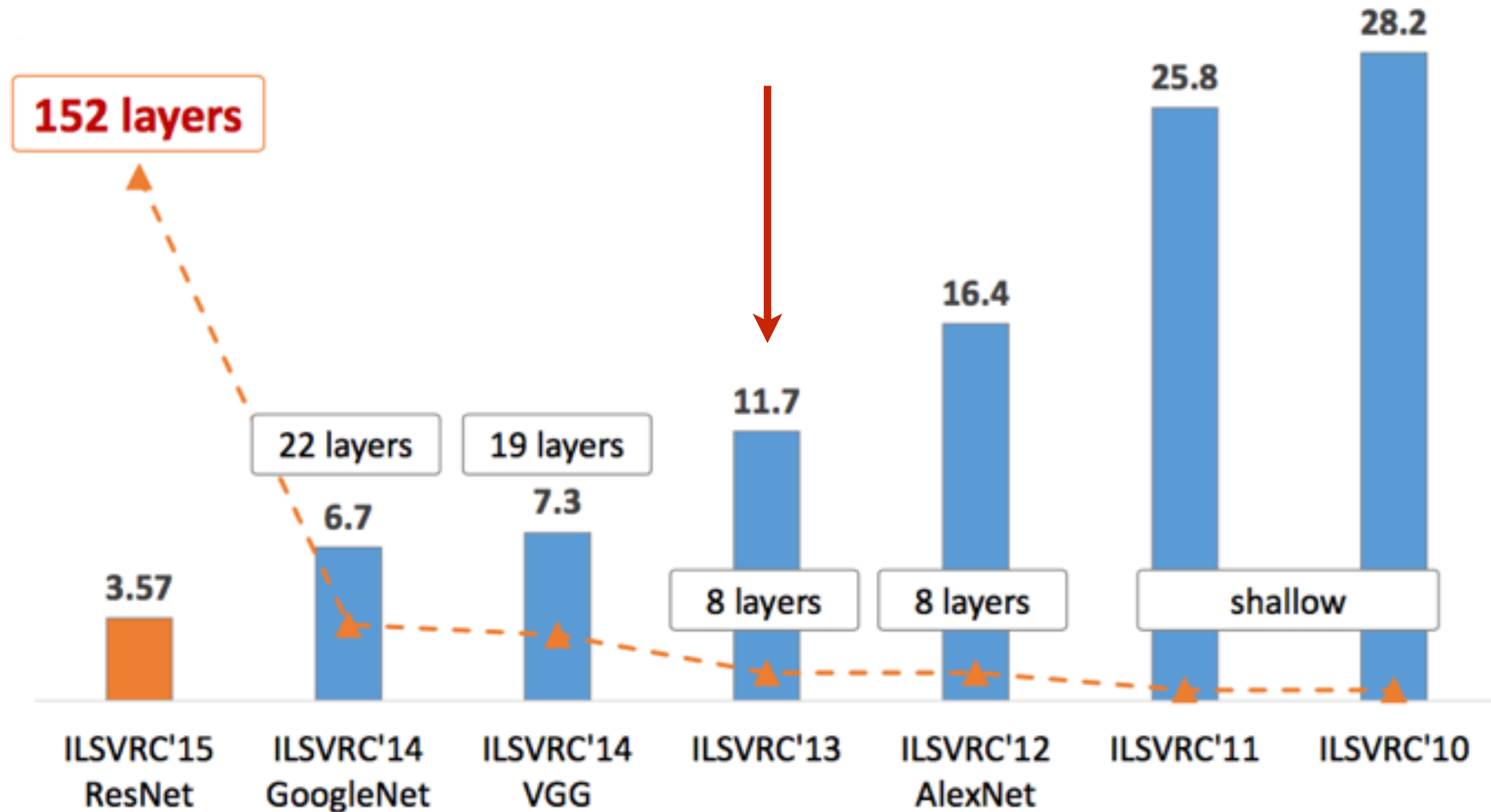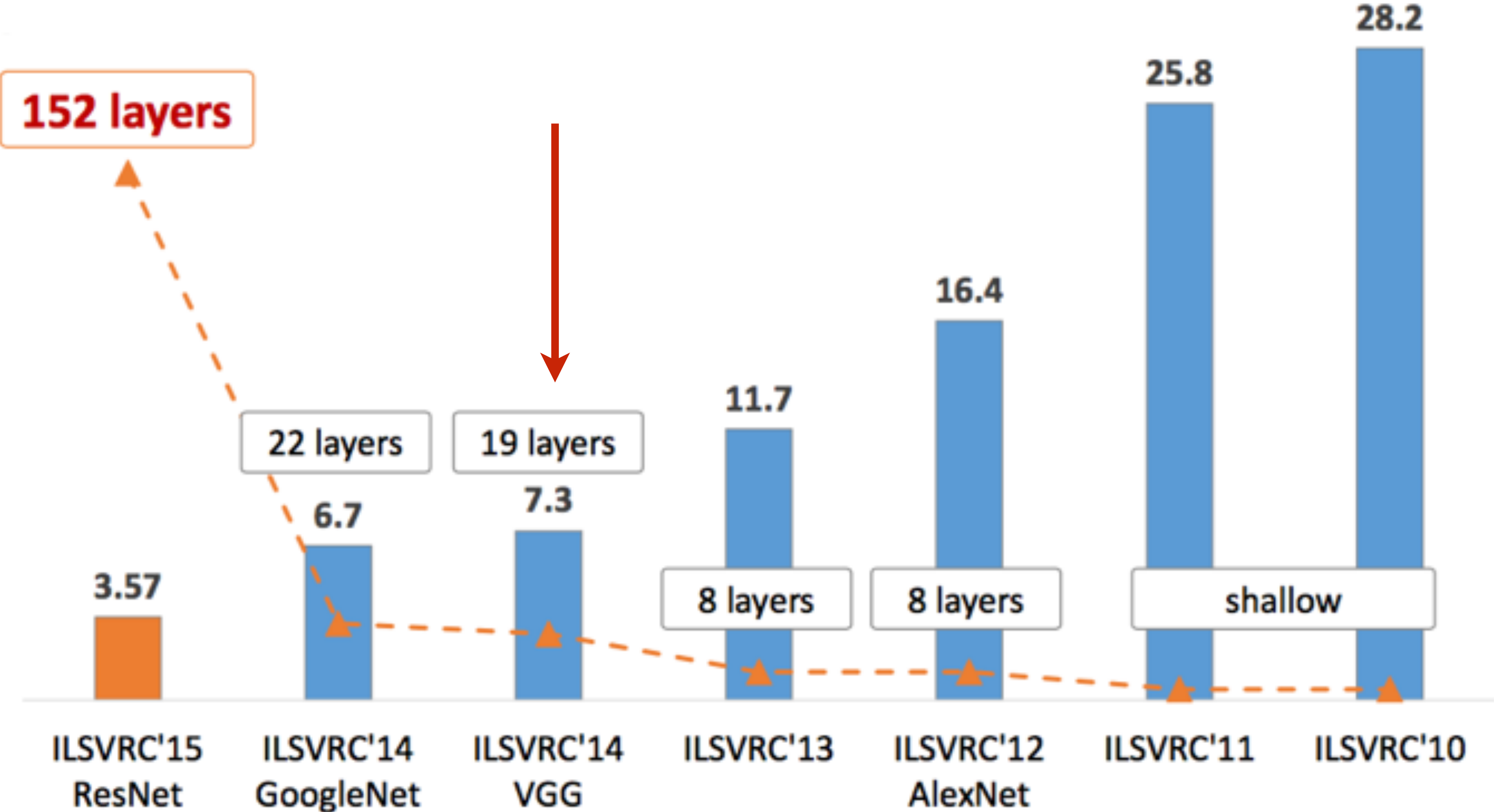Church
Toothbrush
**Person**

IM🔲GENET

# ILSVRC winner 2012

# **ILSVRC** winner 2012

# ILSVRC winner 2012
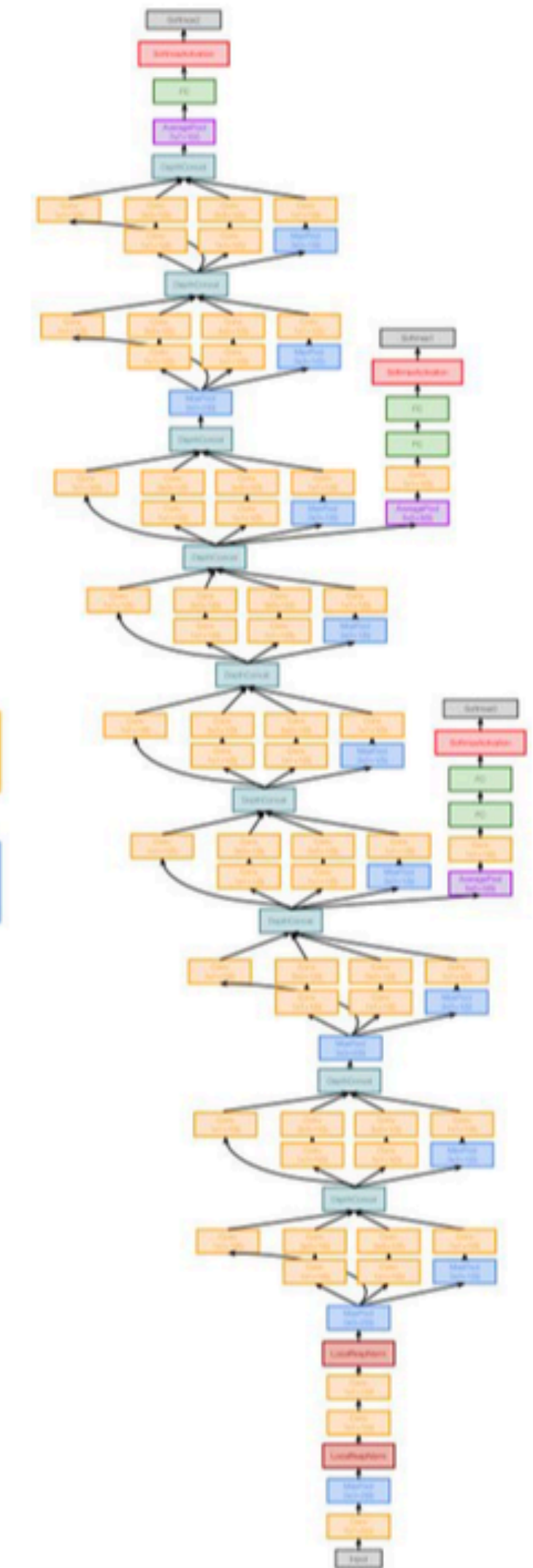
# **Google**LeNet

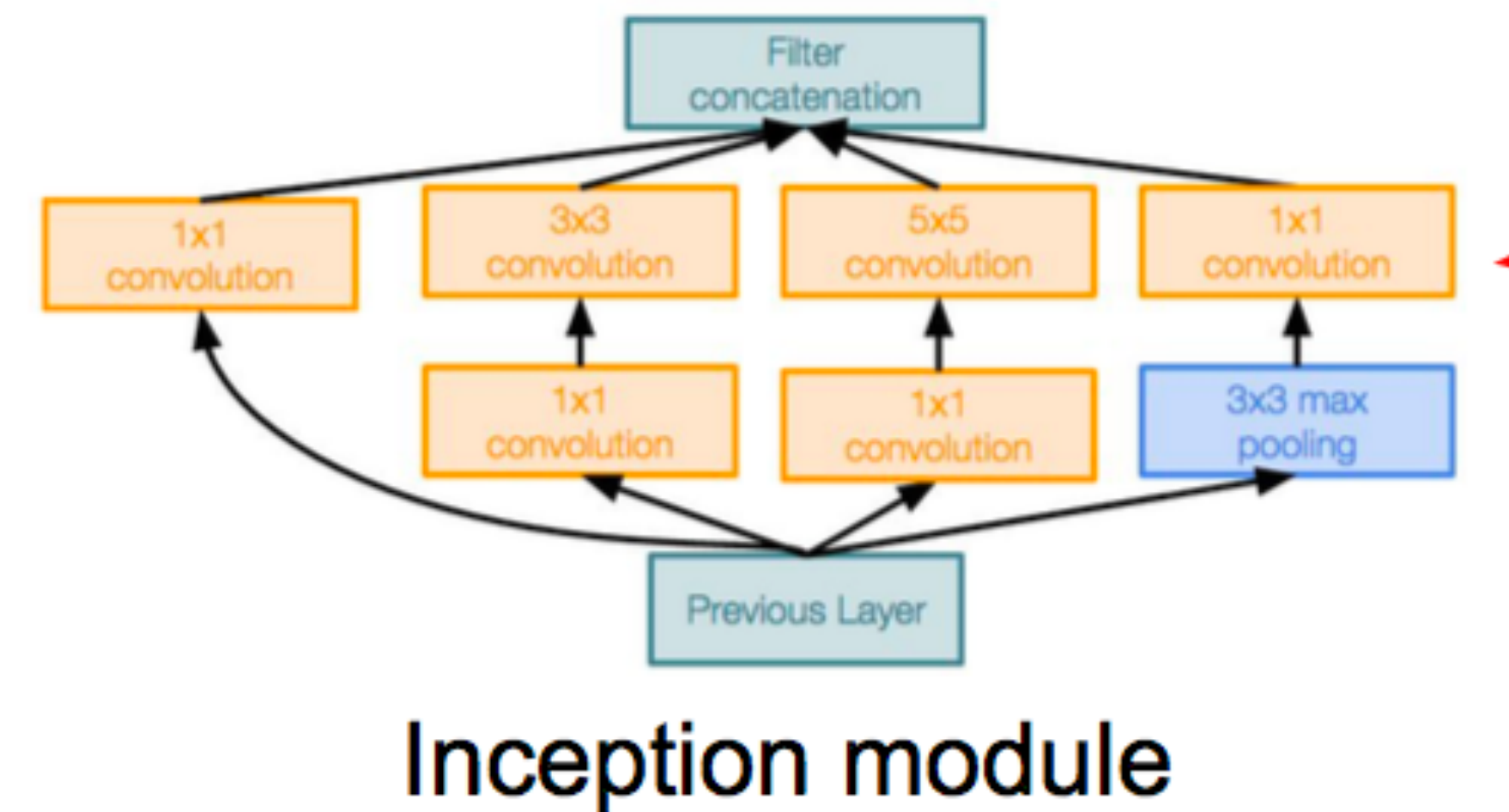even deeper network with **computational efficiency**

— 22 layers

— Efficient "Inception" module

— No FC layers

— Only 5 million parameters!

(12x less than AlexNet!)

— Better performance (@6.7 top 5 error)

**Inception module**

# **Google**LeNet: Inception Module

**Idea:** design good local topology ("network within network")  and then stack these modules



**Inception module**

# **Google**LeNet: Inception Module

**Idea:** design good local topology ("network within network")  and then stack these modules

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 at output depth-wise
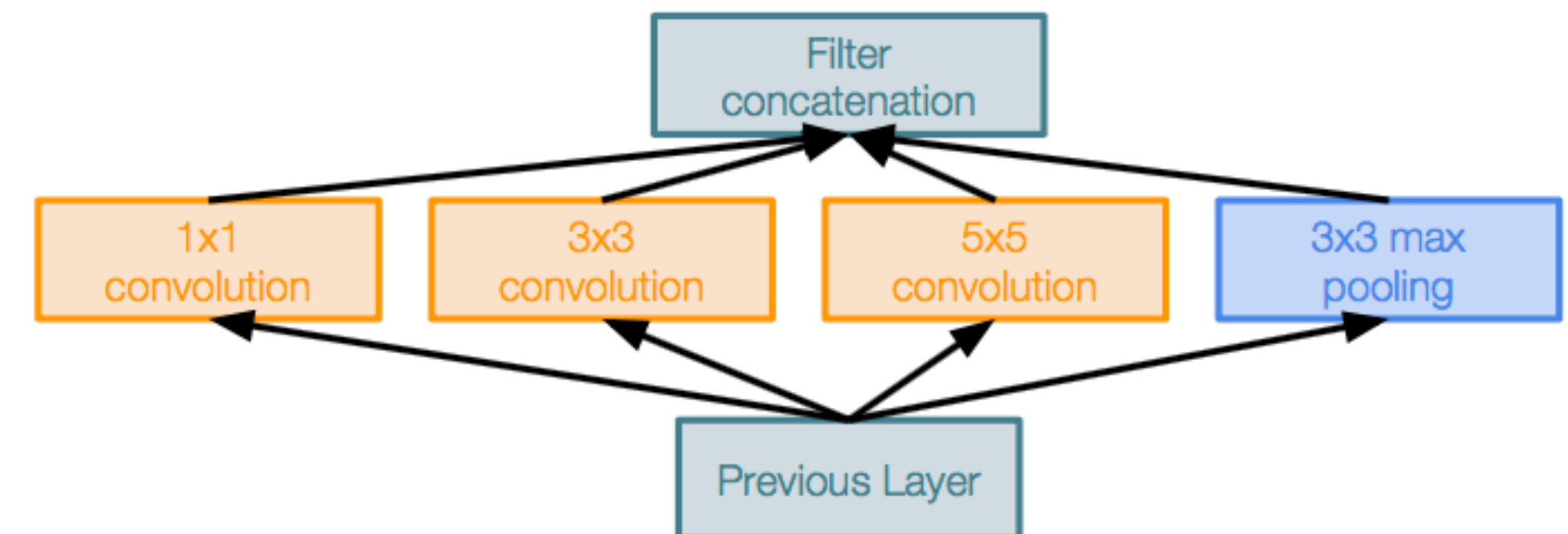


Naive Inception module

# **Google**LeNet: Inception Module

**Idea:** design good local topology ("network within network")  and then stack these modules

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 at output depth-wise

What's the problem?



**Naive Inception module**
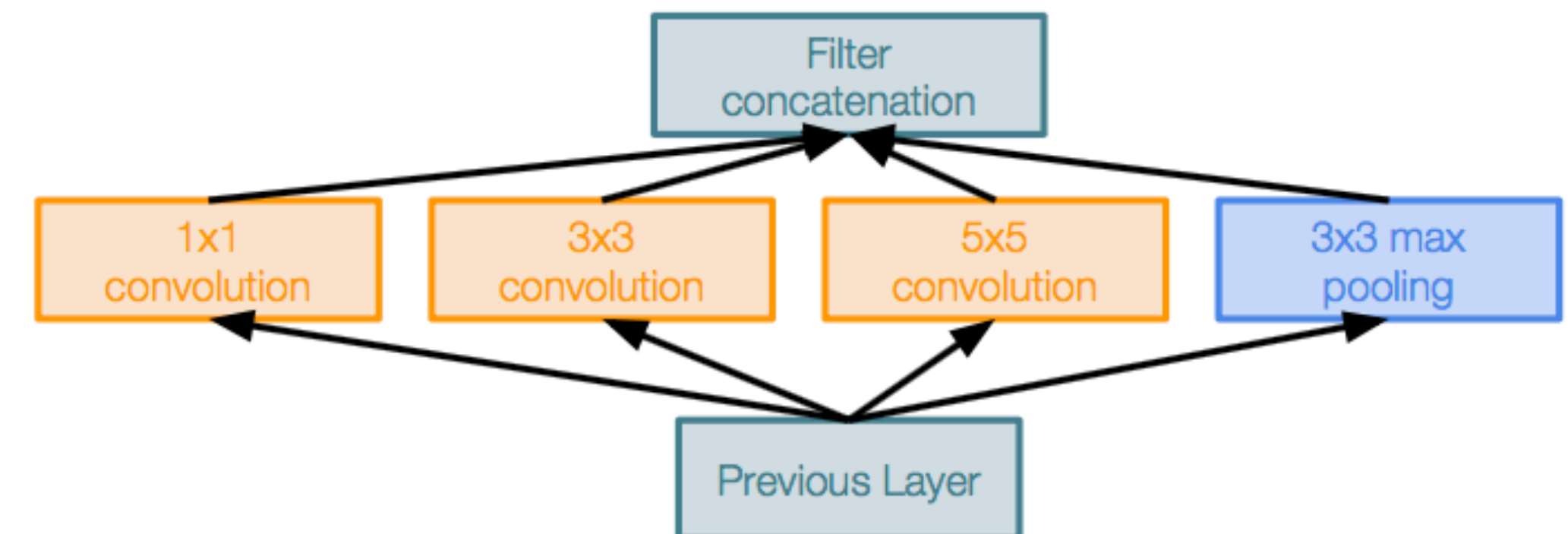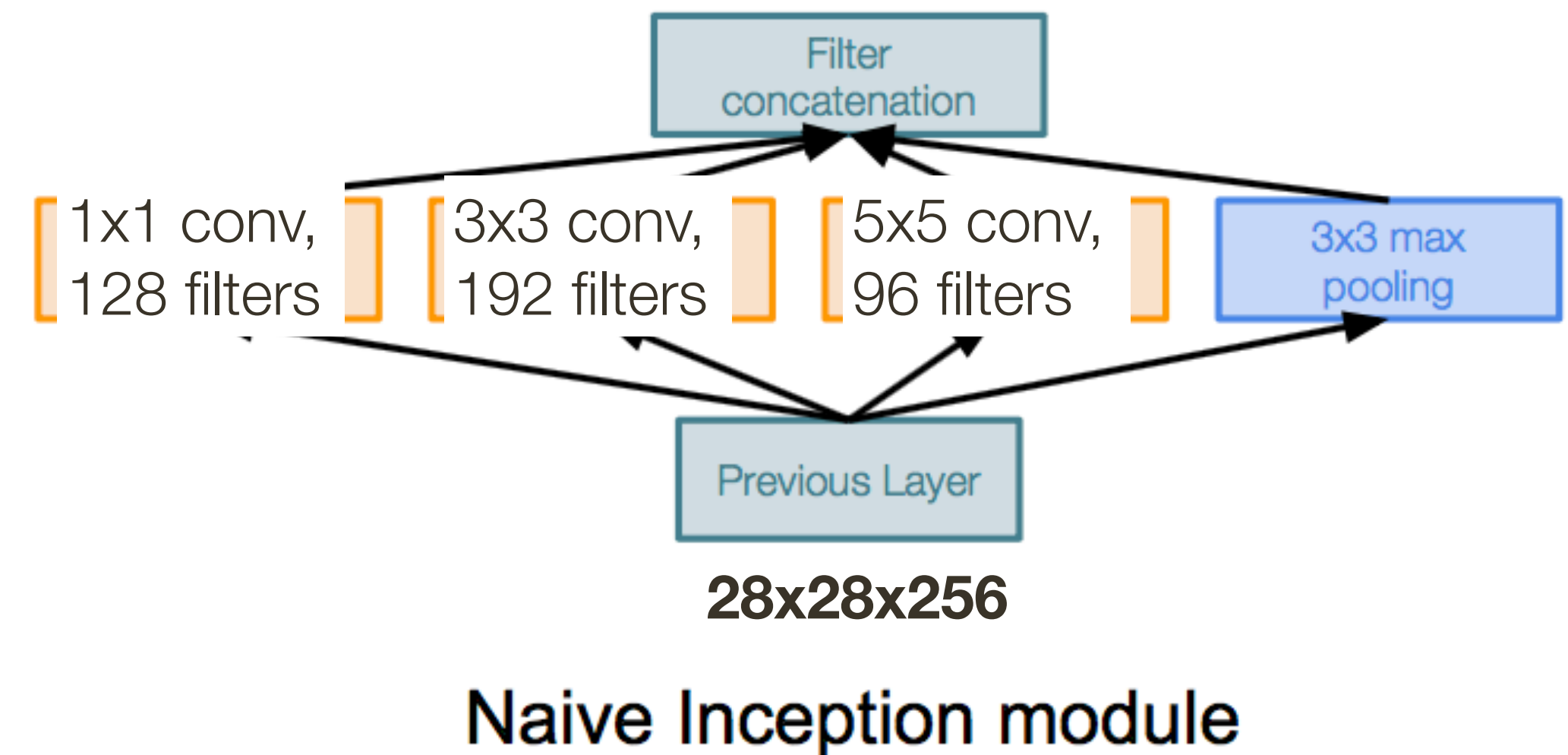
# **Google**LeNet: Inception Module

**Idea:** design good local topology ("network within network")  and then stack these modules

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 at output depth-wise



Filter concatenation

1x1 conv, 128 filters

3x3 conv, 192 filters

5x5 conv, 96 filters

3x3 max pooling

Previous Layer

**28x28x256**

## Naive Inception module

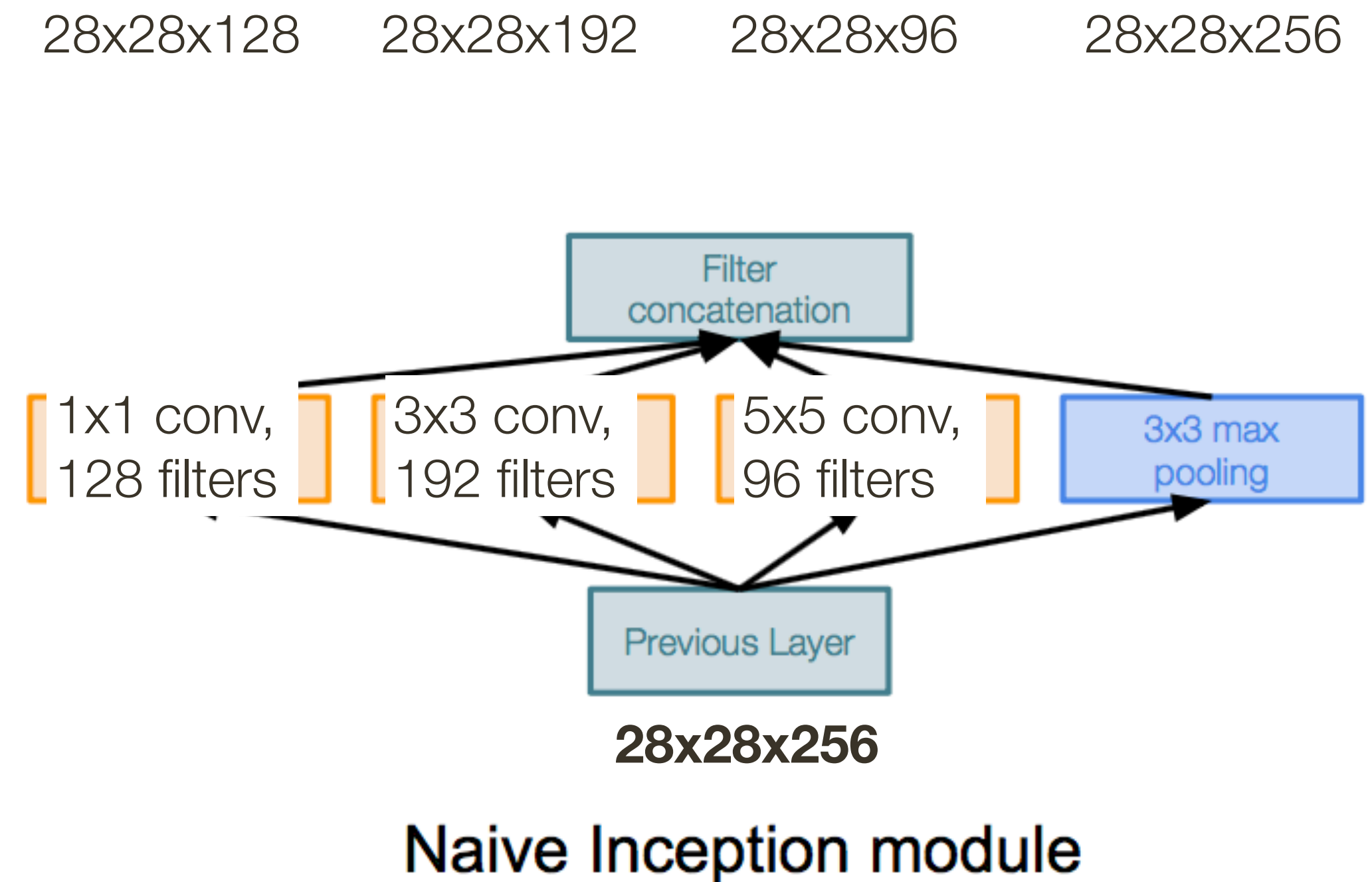* slide from Fei-Dei Li, Justin Johnson, Serena Yeung, **cs231n Stanford**

# **Google**LeNet: Inception Module

**Idea:** design good local topology ("network within network")  and then stack these modules

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 at output depth-wise

28x28x128    28x28x192    28x28x96    28x28x256



Filter concatenation

1x1 conv, 128 filters    3x3 conv, 192 filters    5x5 conv, 96 filters    3x3 max pooling

Previous Layer

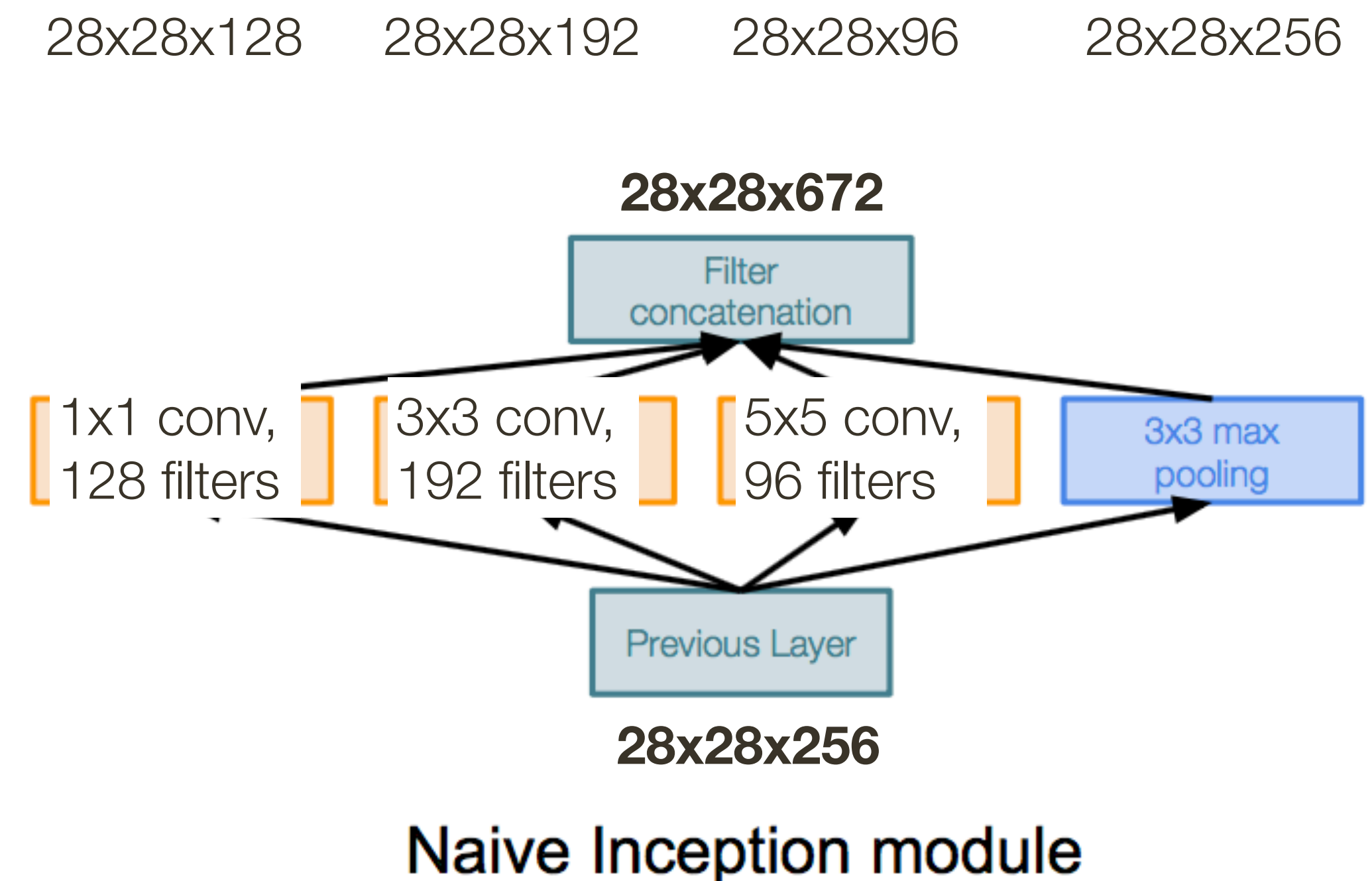**28x28x256**

## Naive Inception module

# **Google**LeNet: Inception Module

**Idea:** design good local topology ("network within network") and then stack these modules

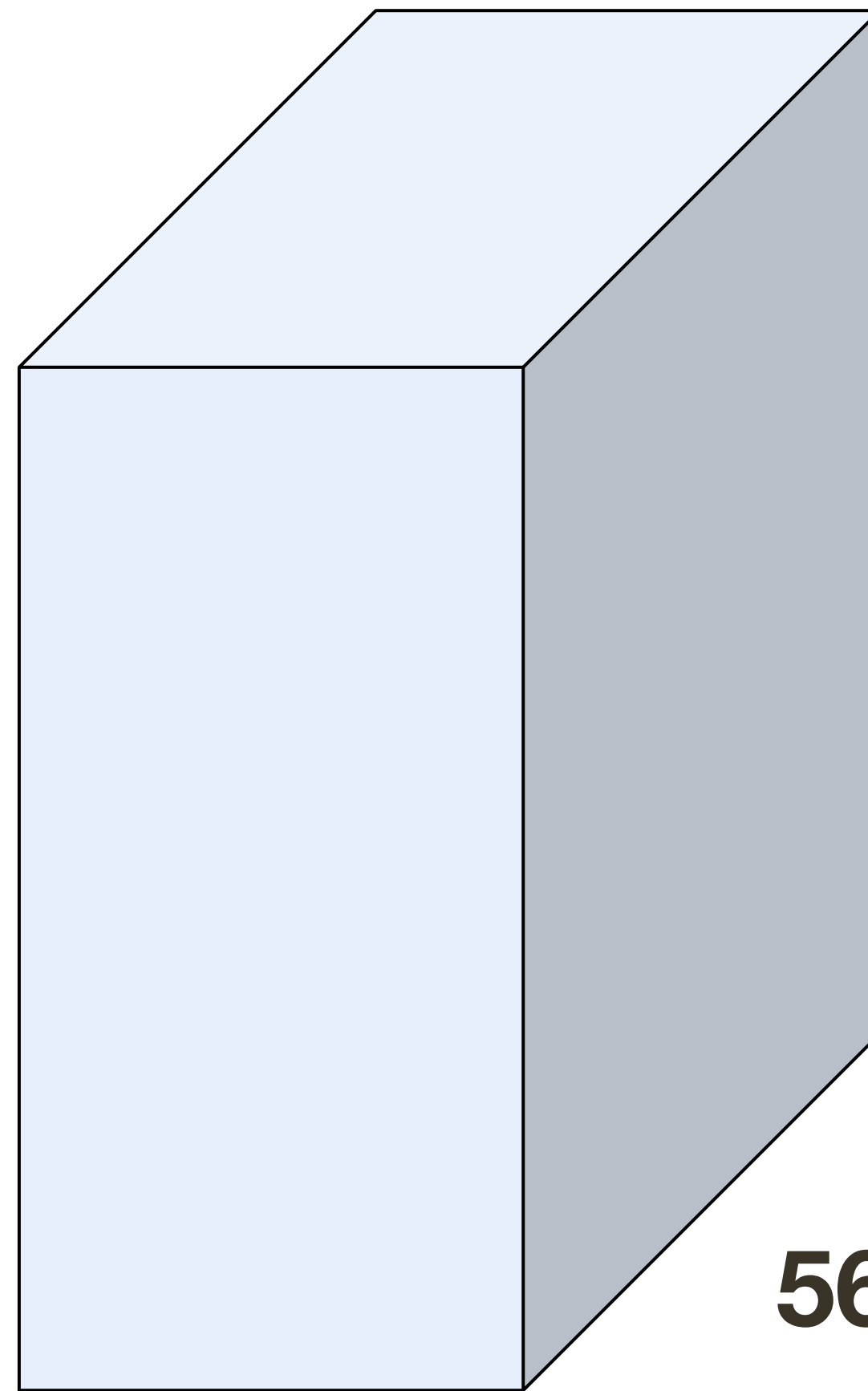Apply **parallel filter operations** on the input from previous layer

— Multiple receptive field sizes for convolution (1x1, 3x3, 5x5)

— Pooling operation (3x3)

**Concatenate** all filter outputs together at output depth-wise

28x28x128    28x28x192    28x28x96    28x28x256

**28x28x672**

Filter concatenation

1x1 conv, 128 filters    3x3 conv, 192 filters    5x5 conv, 96 filters    3x3 max pooling

Previous Layer

**28x28x256**

## Naive Inception module

# Convolutional Layer: **1x1 convolutions**

56 x 56 x 64 **image**

56 x 56 x 32 **image**



**56** height

**56** width

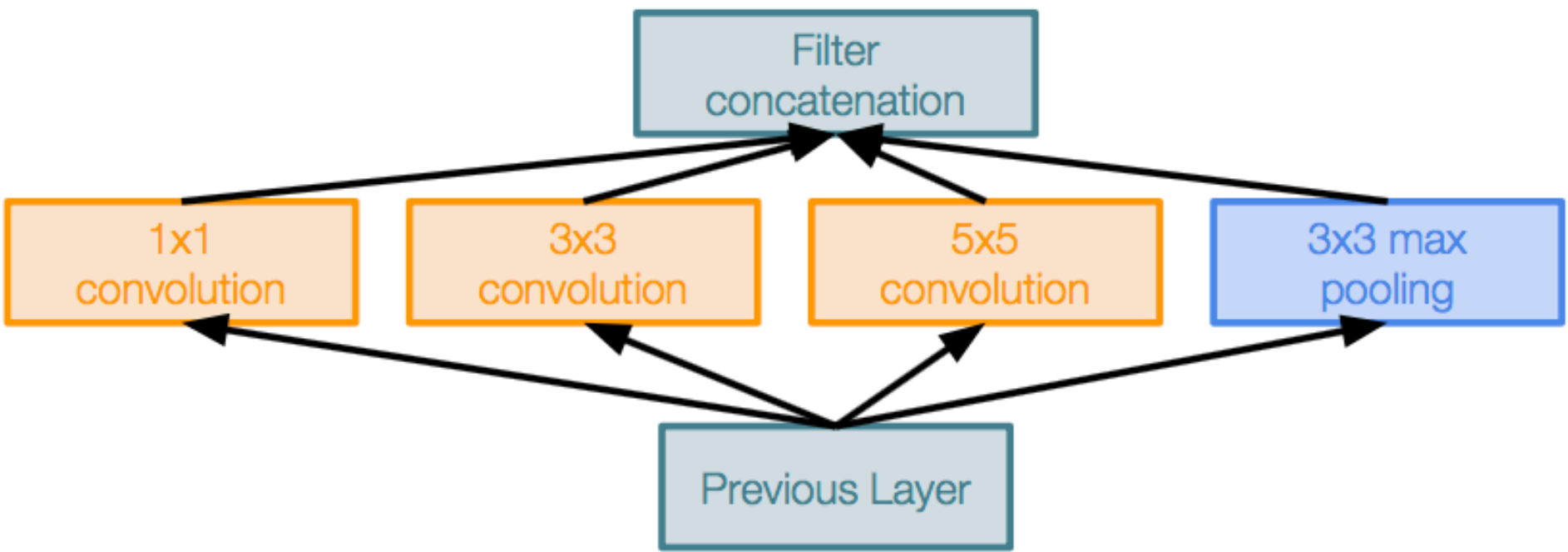**64** depth

32 **filters** of size, 1 x 1 x 64

**56** height

**56** width

**32** depth

# **Google**LeNet: Inception Module

**Idea:** design good local topology ("network within network")  and then stack these modules

1x1 "bottleneck" layers



Naive Inception module

Inception module with dimension reduction

saves approximately 60% of computations

# **Google**LeNet

even deeper network with **computational efficiency**

— 22 layers

— Efficient "Inception" module

— No FC layers

— Only 5 million parameters!

(12x less than AlexNet!)

— Better performance (@6.7 top 5 error)



**Inception module**

# **Google**LeNet

even deeper network with **computational efficiency**

— 22 layers

— Efficient "Inception" module

— No FC layers

— Only 5 million parameters!

(12x less than AlexNet!)

— Better performance (@6.7 top 5 error)



**Inception module**

# **Google**LeNet

even deeper network with **computational efficiency**

— 22 layers

— Efficient "Inception" module

— No FC layers

— Only 5 million parameters!

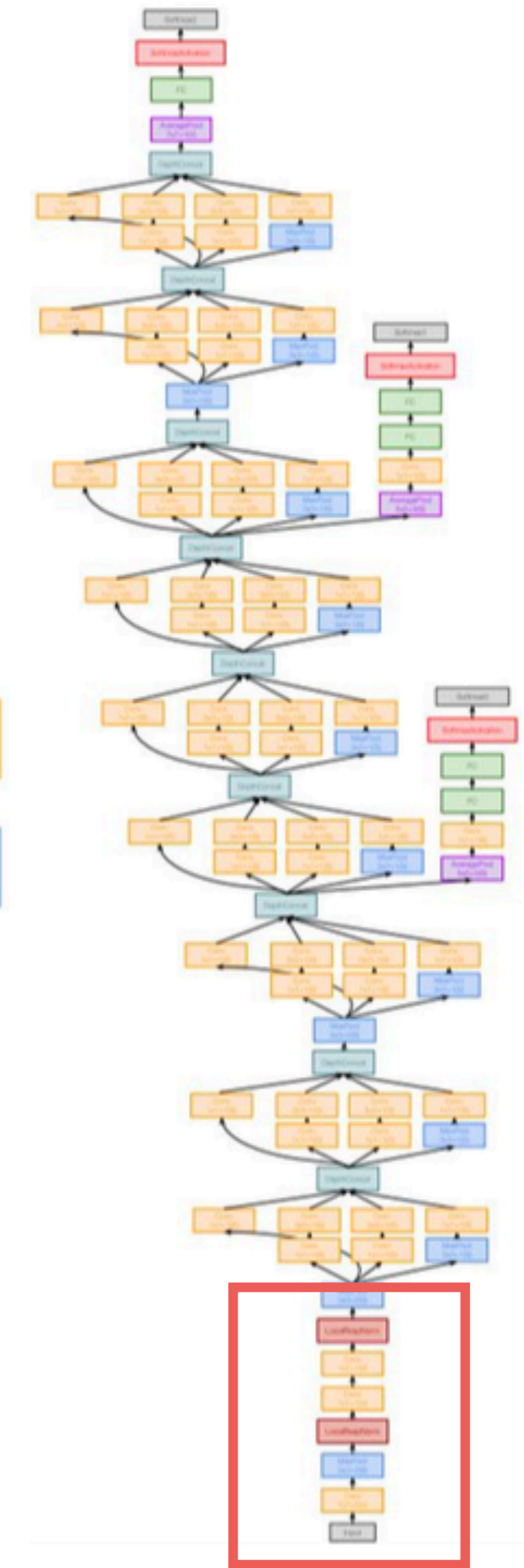(12x less than AlexNet!)

— Better performance (@6.7 top 5 error)



**Inception module**

# **Google**LeNet

even deeper network with **computational efficiency**

— 22 layers

— Efficient "Inception" module

— No FC layers

— Only 5 million parameters!

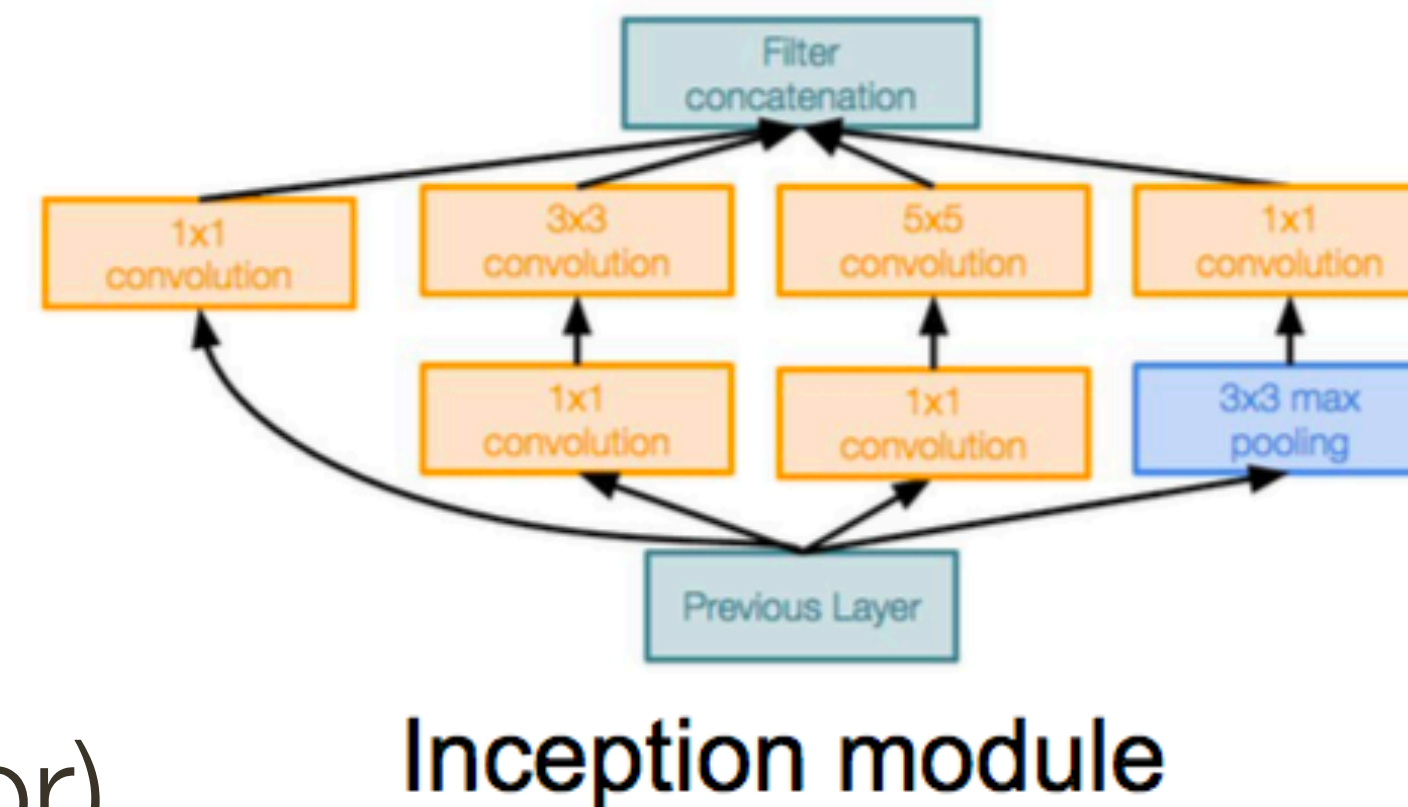(12x less than AlexNet!)
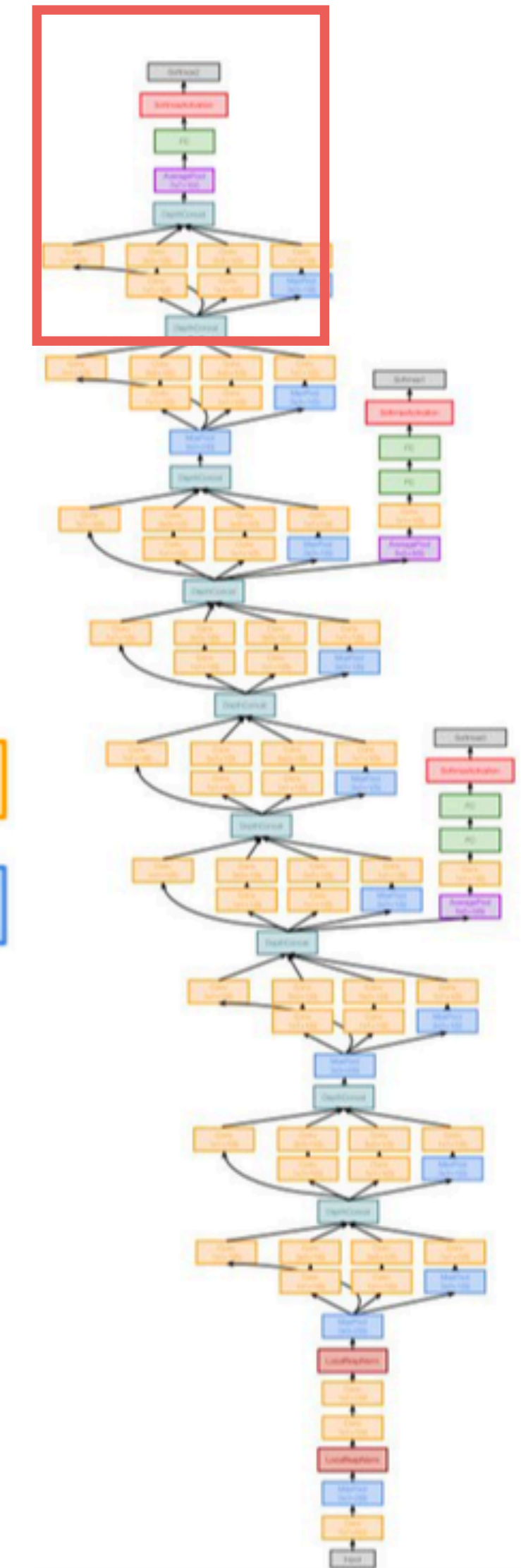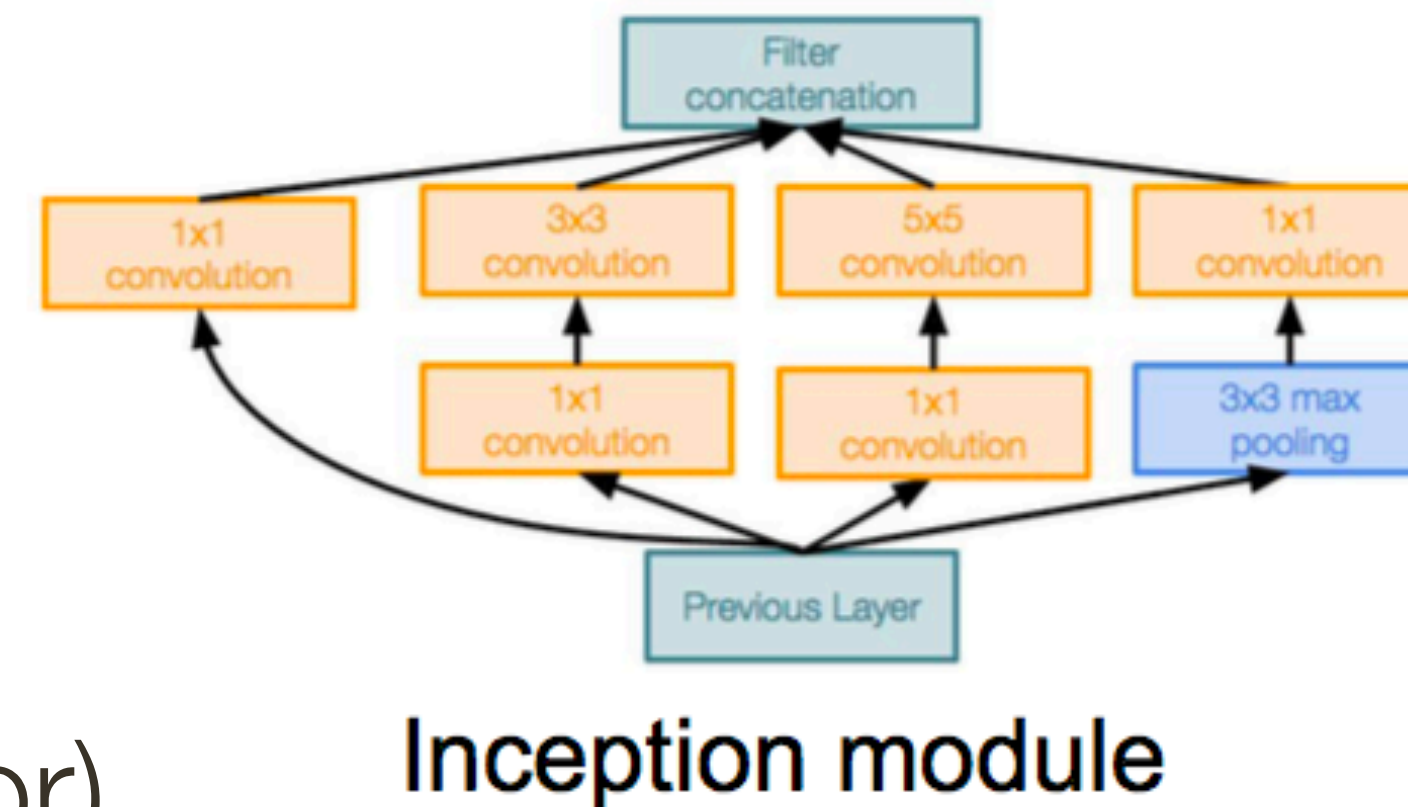
— Better performance (@6.7 top 5 error)



**Inception module**

# Optimizing **Deep** Neural Networks

Consider multi-layer neural network with sigmoid activations and loss C

# Optimizing **Deep** Neural Networks

$$\frac{\partial C}{\partial b_1} = \sigma'(z_1) \times w_2 \times \sigma'(z_2) \times w_3 \times \sigma'(z_3) \times w_4 \times \sigma'(z_4) \times \frac{\partial C}{\partial a_4}$$

# Optimizing **Deep** Neural Networks

$$\frac{\partial C}{\partial b_1} = \sigma'(z_1) \times w_2 \times \sigma'(z_2) \times w_3 \times \sigma'(z_3) \times w_4 \times \sigma'(z_4) \times \frac{\partial C}{\partial a_4}$$



Expression for **gradient** of bias in **Layer 1**: $\frac{\partial C}{\partial b_1} = \sigma'(z_1)\, w_2 \sigma'(z_2)\, w_3 \sigma'(z_3)\, w_4 \sigma'(z_4) \frac{\partial C}{\partial a_4}$

Expression for **gradient** of bias in **Layer 3**: $\frac{\partial C}{\partial b_3} = \sigma'(z_3)\, w_4 \sigma'(z_4) \frac{\partial C}{\partial a_4}$

# Optimizing **Deep** Neural Networks

$$\frac{\partial C}{\partial b_1} = \sigma'(z_1) \times w_2 \times \sigma'(z_2) \times w_3 \times \sigma'(z_3) \times w_4 \times \sigma'(z_4) \times \frac{\partial C}{\partial a_4}$$



Expression for **gradient** of bias in **Layer 1**:  $\frac{\partial C}{\partial b_1} = \sigma'(z_1)\, w_2 \sigma'(z_2)\, w_3 \sigma'(z_3)\, w_4 \sigma'(z_4) \frac{\partial C}{\partial a_4}$

common terms

Expression for **gradient** of bias in **Layer 3**:  $\frac{\partial C}{\partial b_3} = \sigma'(z_3)\, w_4 \sigma'(z_4) \frac{\partial C}{\partial a_4}$

# Optimizing **Deep** Neural Networks

$$\frac{\partial C}{\partial b_1} = \sigma'(z_1) \times w_2 \times \sigma'(z_2) \times w_3 \times \sigma'(z_3) \times w_4 \times \sigma'(z_4) \times \frac{\partial C}{\partial a_4}$$
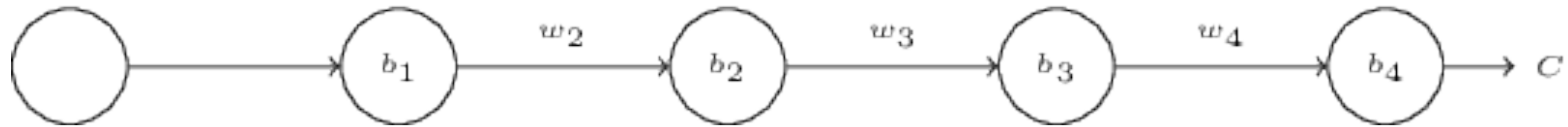


$$\frac{\partial C}{\partial b_1} = \sigma'(z_1)\, w_2 \sigma'(z_2)\, w_3 \sigma'(z_3)\, \underbrace{w_4 \sigma'(z_4) \frac{\partial C}{\partial a_4}}$$

common terms

$$\frac{\partial C}{\partial b_3} = \sigma'(z_3)\, w_4 \sigma'(z_4) \frac{\partial C}{\partial a_4}$$

**Observations**:
 |weight| < 1 (due to initialization)
 max of derivative of sigmoid = 1/4 @ 0

# Optimizing **Deep** Neural Networks

$$\frac{\partial C}{\partial b_1} = \sigma'(z_1) \times w_2 \times \sigma'(z_2) \times w_3 \times \sigma'(z_3) \times w_4 \times \sigma'(z_4) \times \frac{\partial C}{\partial a_4}$$



$$\frac{\partial C}{\partial b_1} = \sigma'(z_1) \overbrace{w_2 \sigma'(z_2)}^{< \frac{1}{4}} \overbrace{w_3 \sigma'(z_3)}^{< \frac{1}{4}} \underbrace{w_4 \sigma'(z_4) \frac{\partial C}{\partial a_4}}$$
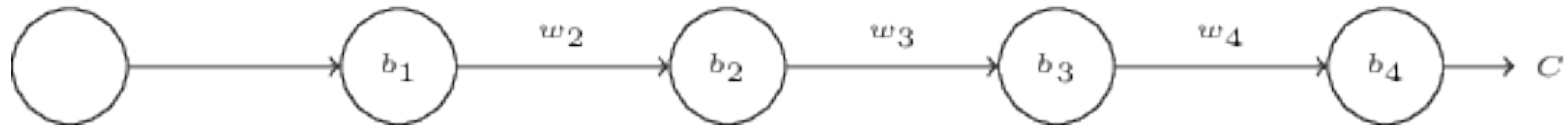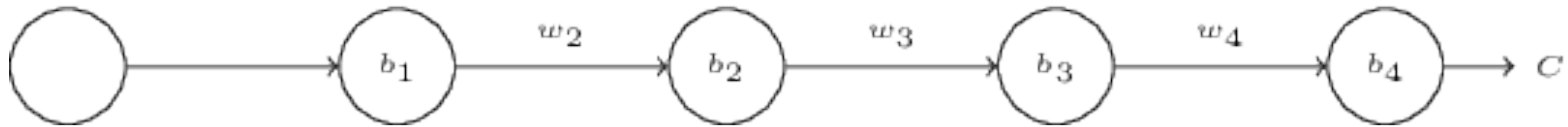
**Observations**:
   |weight| < 1 (due to initialization)
   max of derivative of sigmoid = 1/4 @ 0

common terms

$$\frac{\partial C}{\partial b_3} = \sigma'(z_3) \overbrace{w_4 \sigma'(z_4) \frac{\partial C}{\partial a_4}}$$

# Optimizing **Deep** Neural Networks

$$\frac{\partial C}{\partial b_1} = \sigma'(z_1) \times w_2 \times \sigma'(z_2) \times w_3 \times \sigma'(z_3) \times w_4 \times \sigma'(z_4) \times \frac{\partial C}{\partial a_4}$$



$$\frac{\partial C}{\partial b_1} = \sigma'(z_1) \overbrace{w_2 \sigma'(z_2)}^{<\frac{1}{4}} \overbrace{w_3 \sigma'(z_3)}^{<\frac{1}{4}} \underbrace{w_4 \sigma'(z_4) \frac{\partial C}{\partial a_4}}$$

common terms

This is called **vanishing gradient** problem

— makes deep networks hard to train
— later layers learn faster than earlier ones

$$\frac{\partial C}{\partial b_3} = \sigma'(z_3) \overbrace{w_4 \sigma'(z_4) \frac{\partial C}{\partial a_4}}$$

# Optimizing **Deep** Neural Networks

$$\frac{\partial C}{\partial b_1} = \sigma'(z_1) \times w_2 \times \sigma'(z_2) \times w_3 \times \sigma'(z_3) \times w_4 \times \sigma'(z_4) \times \frac{\partial C}{\partial a_4}$$



$$\frac{\partial C}{\partial b_1} = \sigma'(z_1) \overbrace{w_2 \sigma'(z_2)}^{>1} \overbrace{w_3 \sigma'(z_3)}^{>1} \underbrace{w_4 \sigma'(z_4) \frac{\partial C}{\partial a_4}}$$

common terms

**Exploding gradient** problem

— makes weights large (e.g., 100)
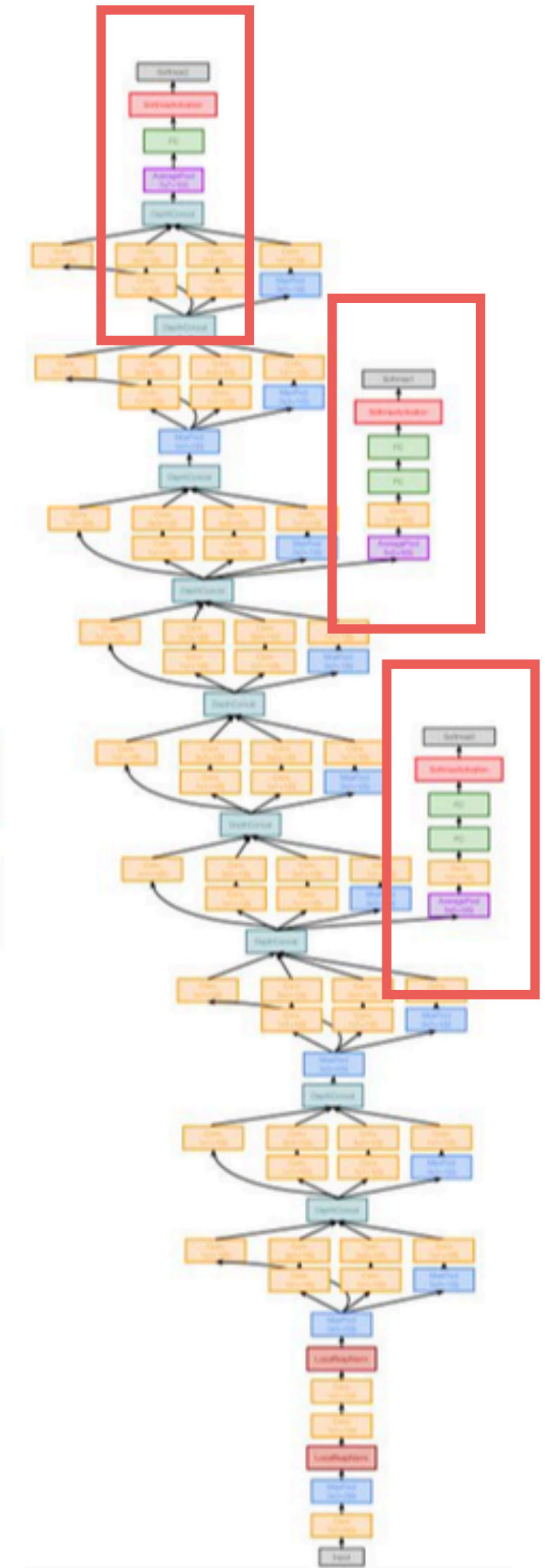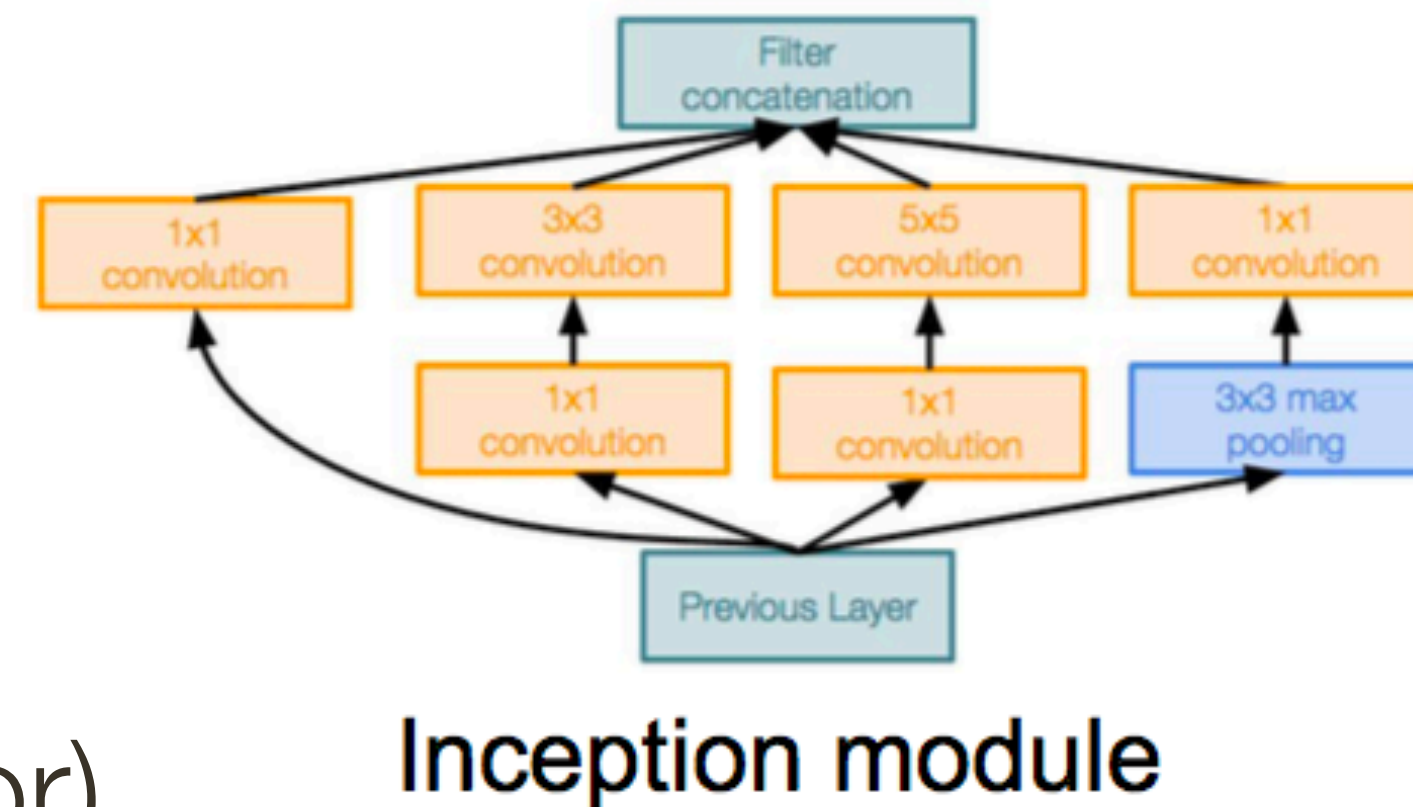— make bias such that pre-activation = 0

$$\frac{\partial C}{\partial b_3} = \sigma'(z_3) \overbrace{w_4 \sigma'(z_4) \frac{\partial C}{\partial a_4}}$$

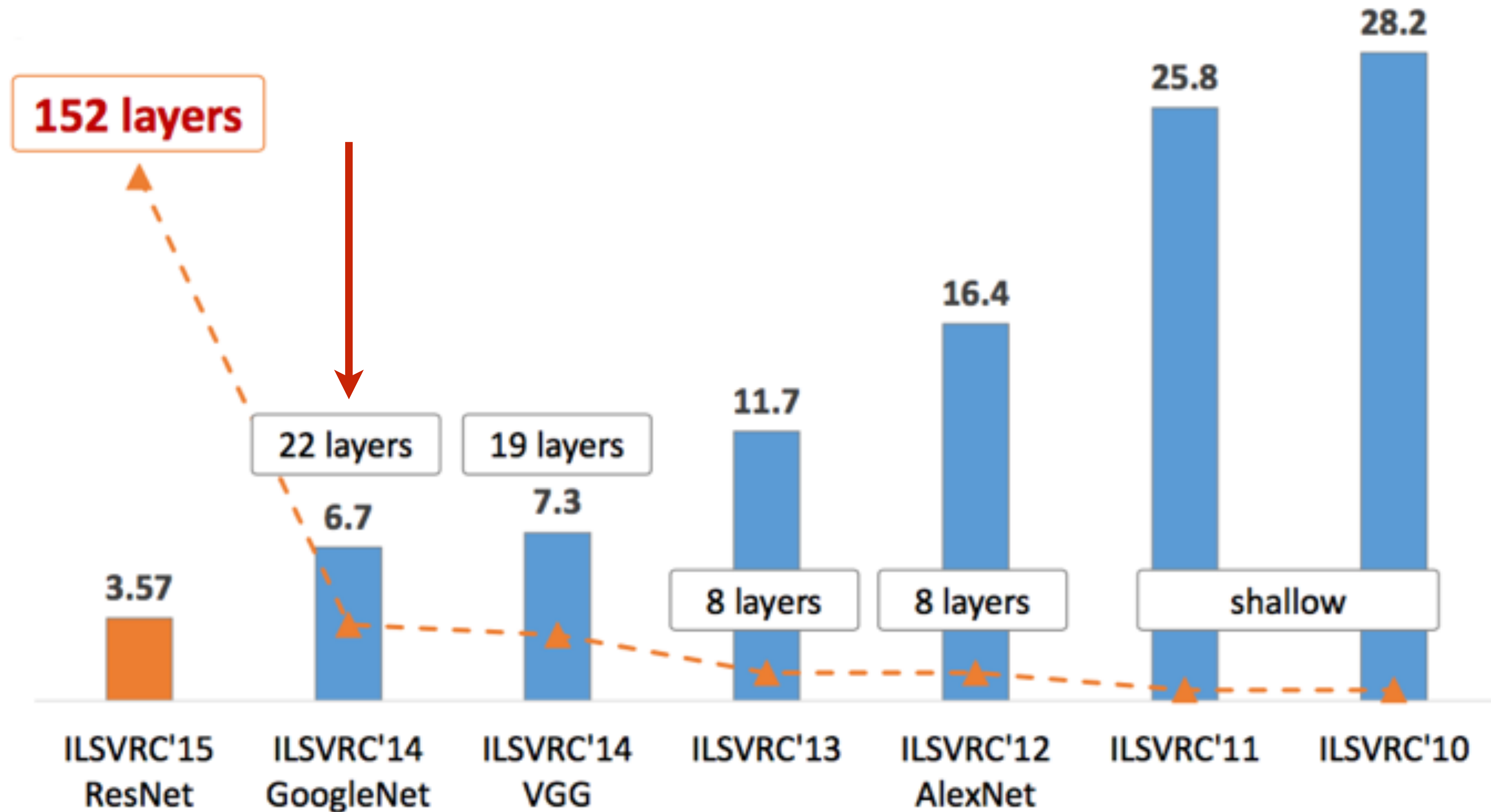**Source**: http://neuralnetworksanddeeplearning.com/chap5.html

# **Google**LeNet

even deeper network with **computational efficiency**

— 22 layers

— Efficient "Inception" module

— No FC layers

— Only 5 million parameters!

(12x less than AlexNet!)
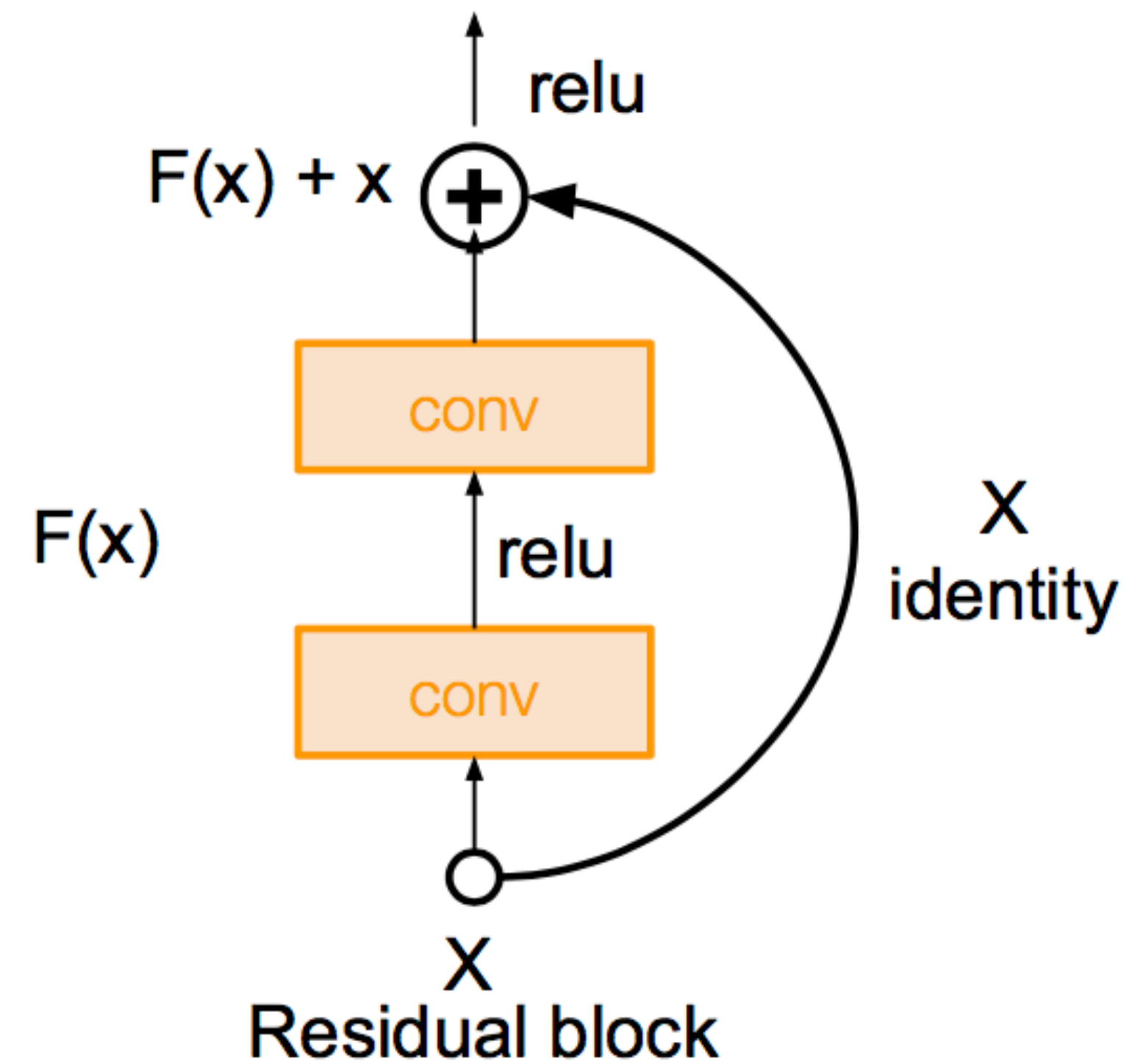
— Better performance (@6.7 top 5 error)



Inception module

# **ILSVRC** winner 2012

# ResNet

even deeper — **152 layers**!

using residual connections



F(x) + x

relu

conv

F(x)                relu

conv

X
identity

X
**Residual block**

Softmax
FC 1000
Pool
3x3 conv, 64
3x3 conv, 64
3x3 conv, 64
3x3 conv, 64
3x3 conv, 64
3x3 conv, 64
...
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: Motivation

What happens when we continue to stacking deeper layers on a "plain" CNN
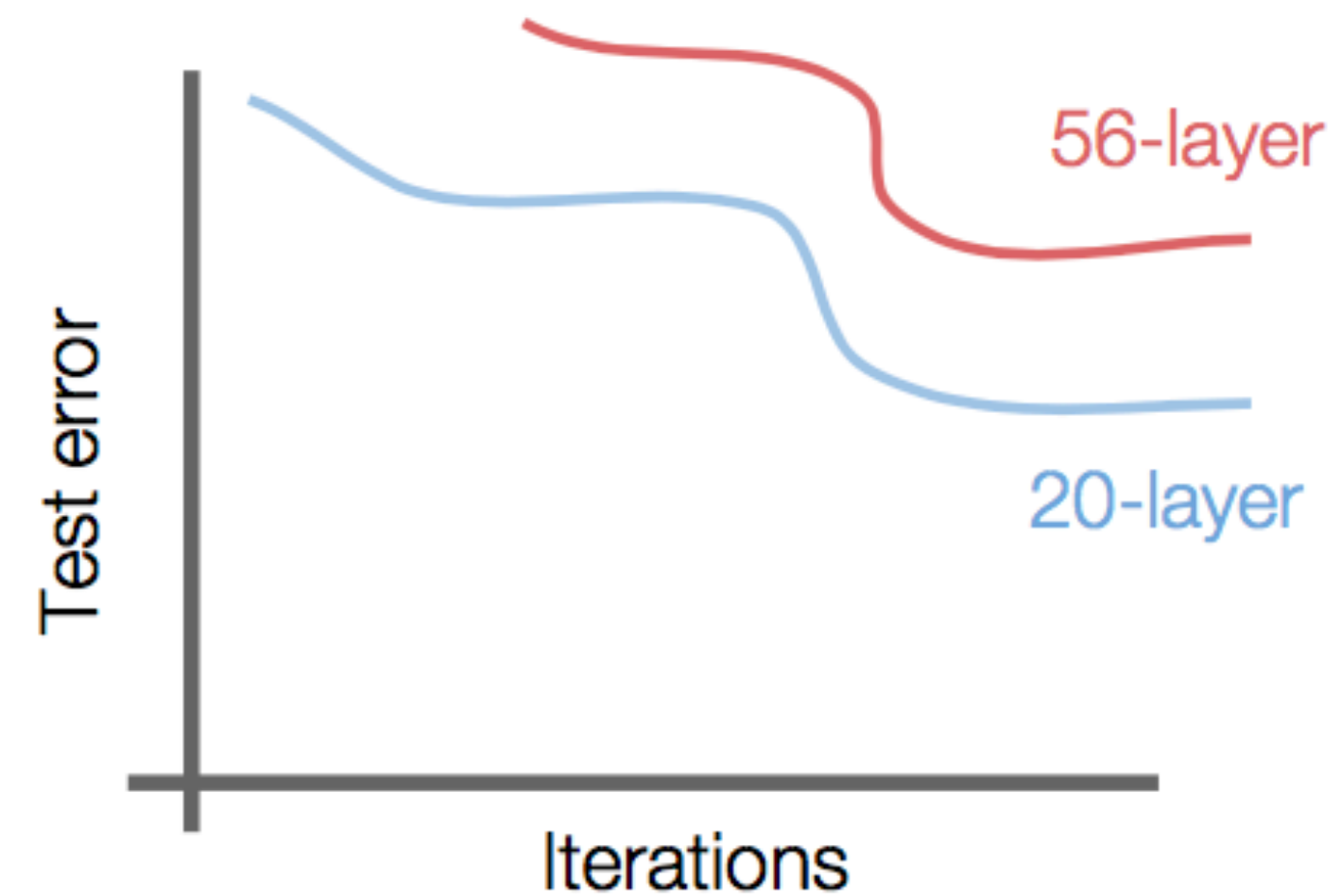
# ResNet: Motivation

What happens when we continue to stacking deeper layers on a "plain" CNN
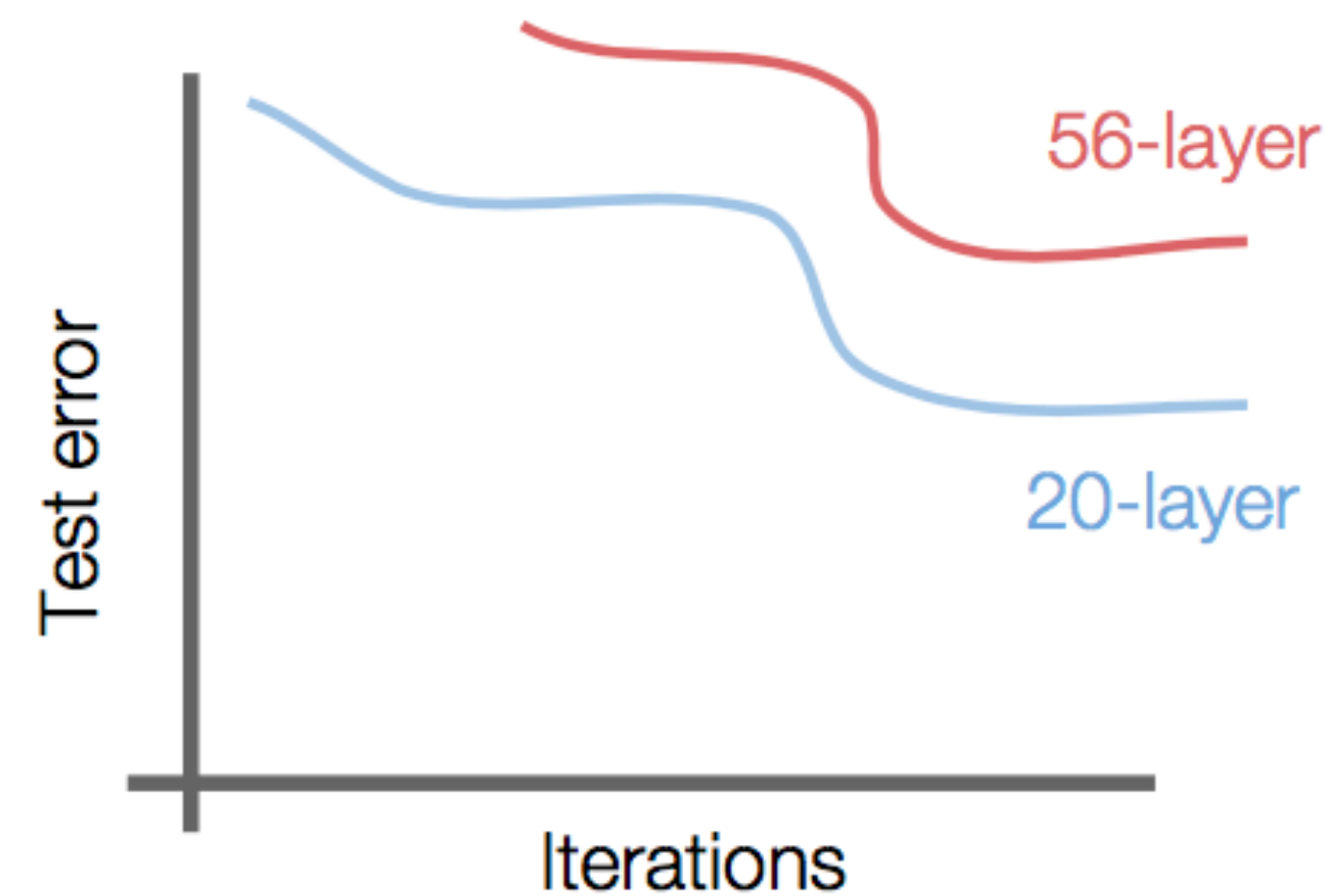


Whats the **problem**?

# ResNet: Motivation

What happens when we continue to stacking deeper layers on a "plain" CNN



Whats the **problem**?

# ResNet: Motivation

**Hypothesis:** deeper models are harder to optimize (optimization problem)

# ResNet: Motivation

**Hypothesis:** deeper models are harder to optimize (optimization problem)

**Observation:** the deeper model should (conceptually) perform just as well (e.g., take shallower model and use identity for all remaining layers)

# **ResNet:** Motivation

**Hypothesis:** deeper models are harder to optimize (optimization problem)

**Observation:** the deeper model should (conceptually) perform just as well (e.g., take shallower model and use identity for all remaining layers)

How do we implement this idea in practice

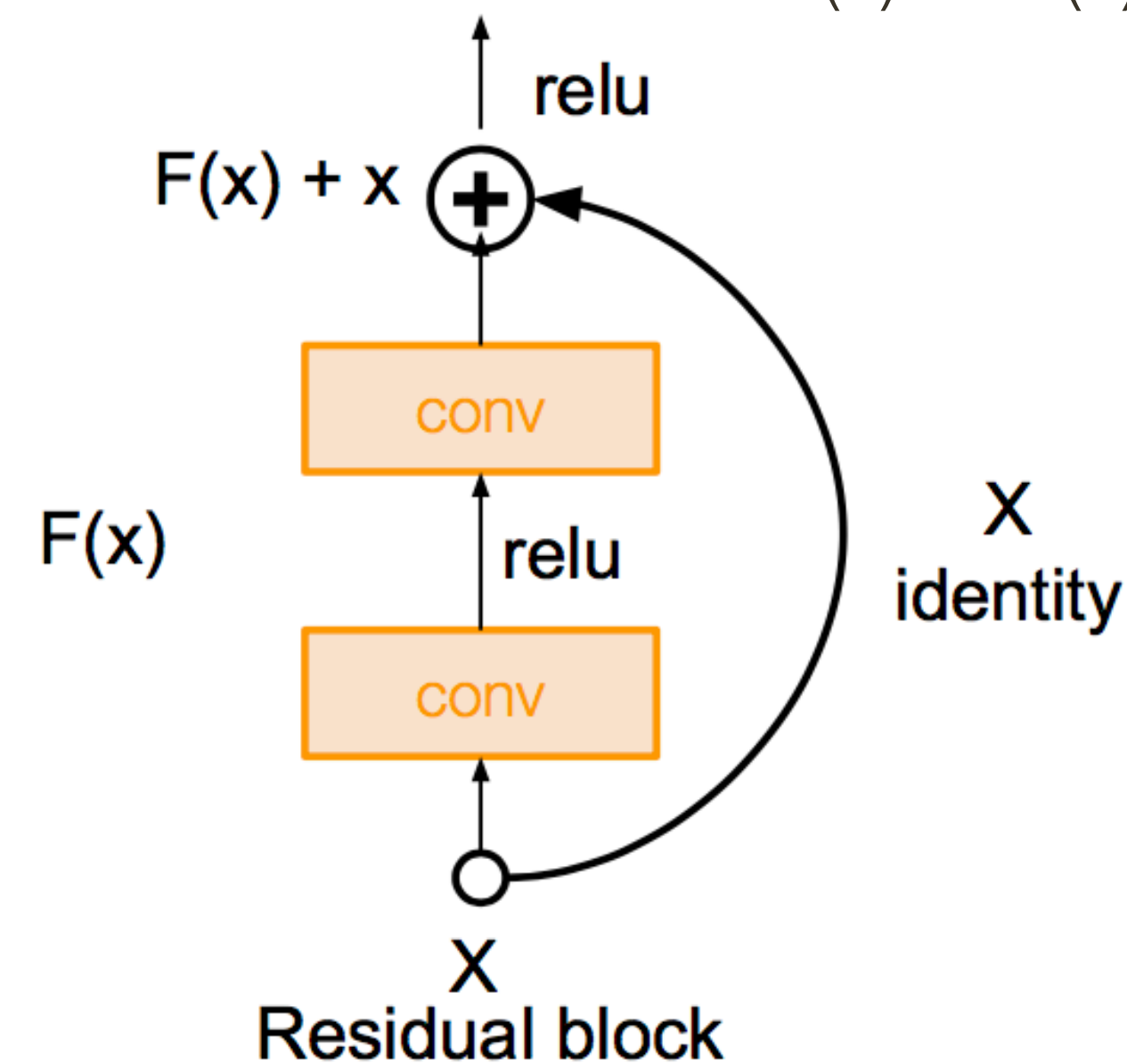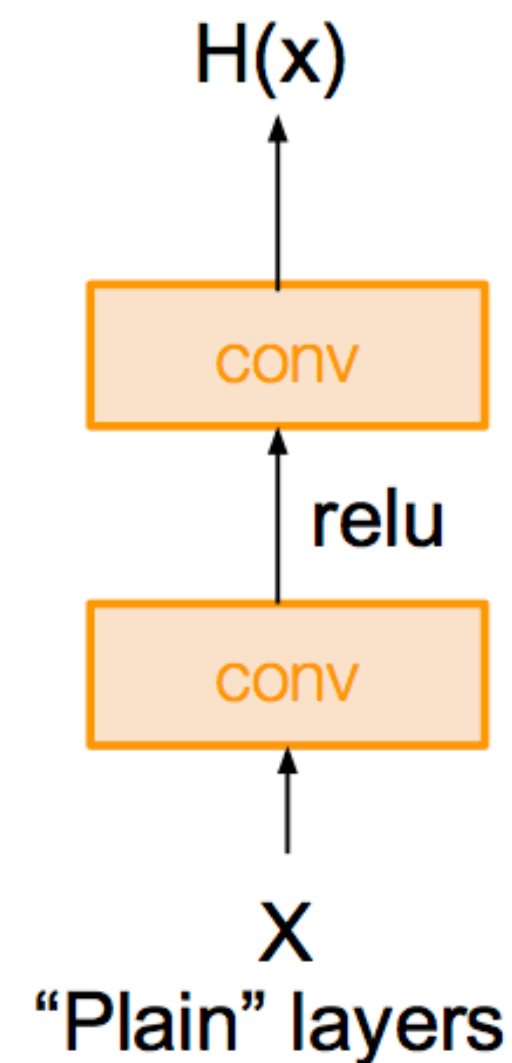# ResNet

**Solution:** use network to fit residual mapping instead of directly trying to fit a desired underlying mapping

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

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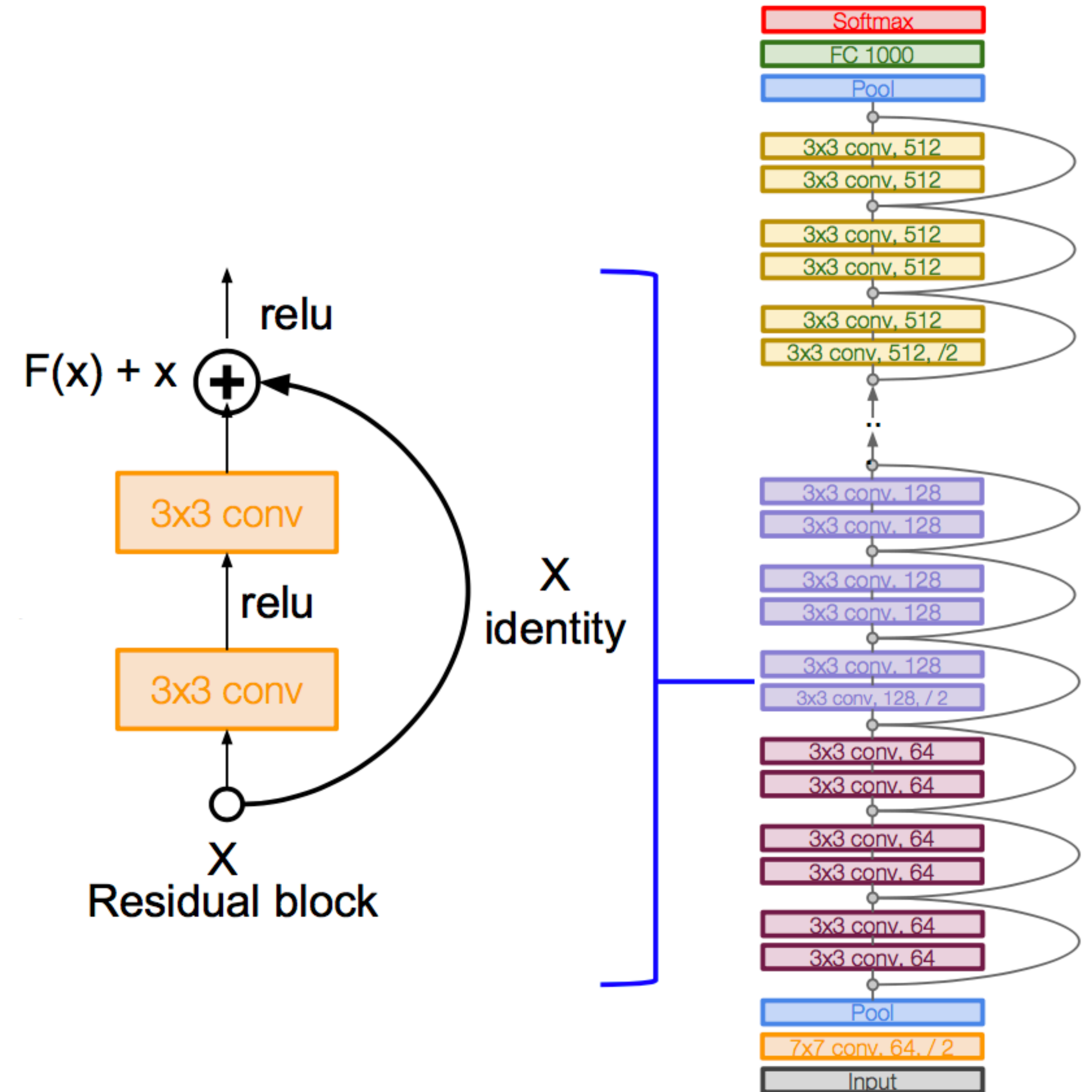$F(x) = H(x) - X$ instead of $H(x)$ directly



H(x)

conv

relu

conv

X
"Plain" layers

relu
F(x) + x ⊕
F(x)      relu

conv

conv

X
identity

X
Residual block

# ResNet

## Full details

— Stacked **residual blocks**

— Every residual block consists of **two 3x3 filters**

— Periodically double # of filters and downsample spatially using stride of 2

— Additional convolutional layer in the beginning

— **No FC layers** at the end (only FC to output 1000 classes)



relu

$F(x) + x \quad \bigoplus$

3x3 conv

relu

3x3 conv

X
identity

X
**Residual block**

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

# ILSVRC winner 2012

# ILSVRC winner 2012



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

152 layers

22 layers    19 layers

8 layers     8 layers     shallow

3.57    6.7    7.3    11.7    25.8    28.2

ILSVRC'15    ILSVRC'14    ILSVRC'14    ILSVRC'13    ILSVRC'12    ILSVRC'11    ILSVRC'10
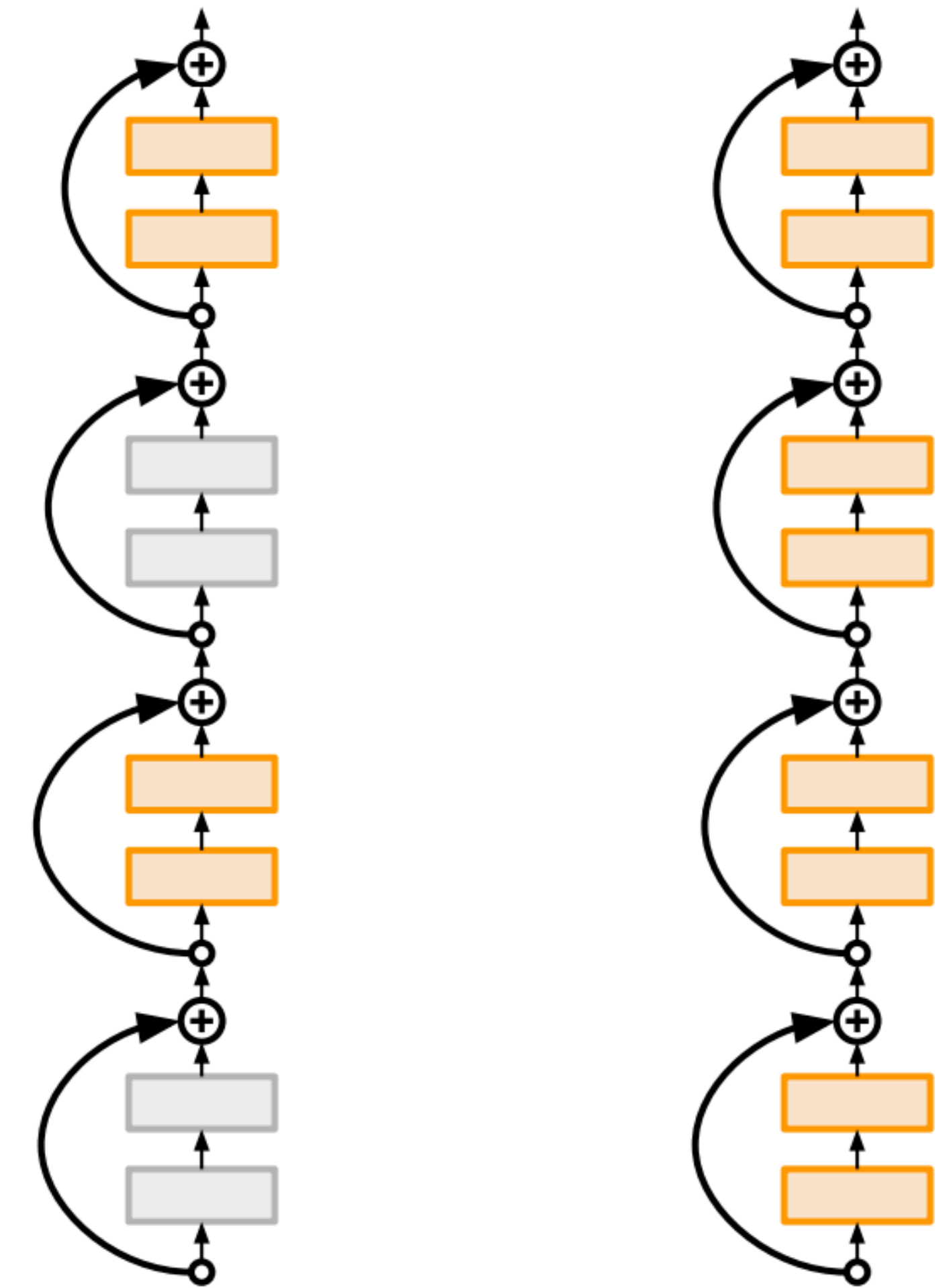ResNet       GoogleNet    VGG                       AlexNet

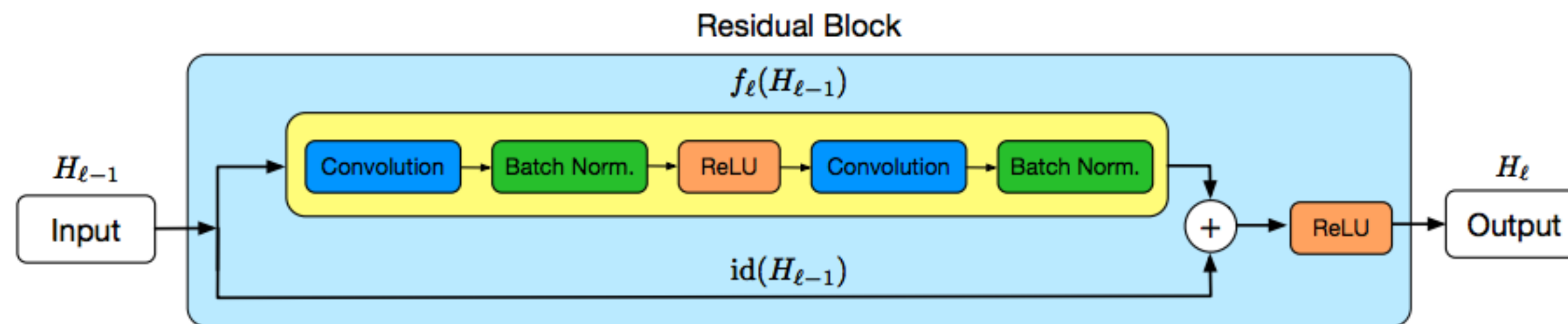# Regularization: Stochastic Depth

Effectively "dropout" but for layers
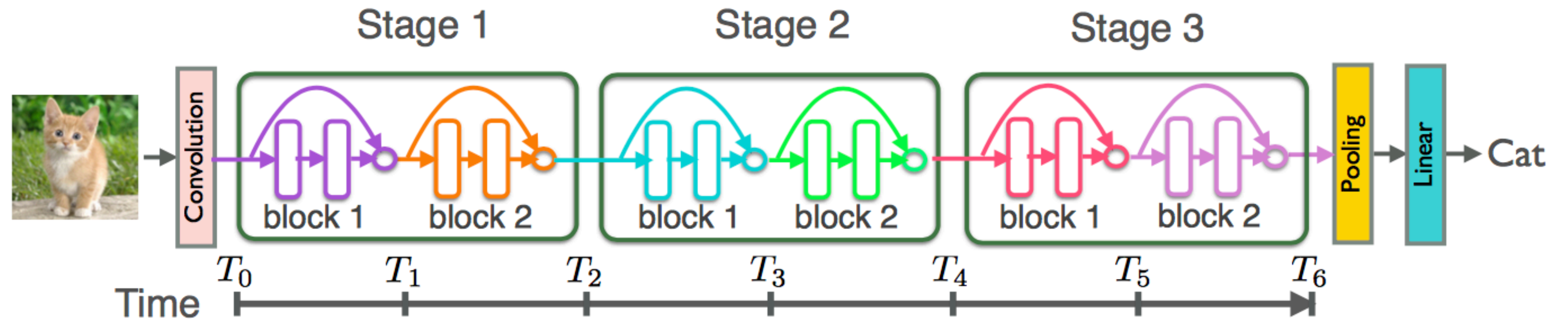
Stochastically with some probability **turn off some layer** (for each batch)

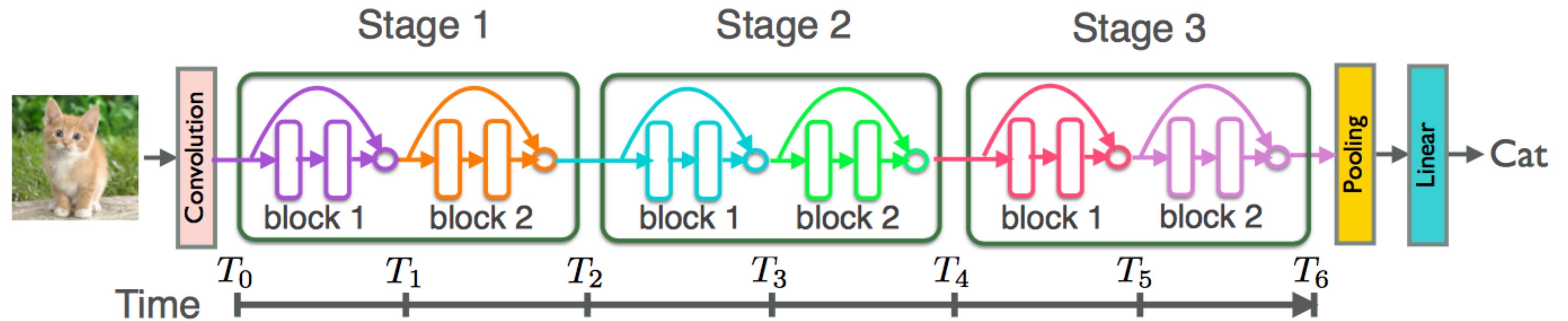Effectively trains a collection of neural networks



**Residual Block**

$f_\ell(H_{\ell-1})$

$H_{\ell-1}$ → Input → Convolution → Batch Norm. → ReLU → Convolution → Batch Norm. → + → ReLU → Output $H_\ell$

$id(H_{\ell-1})$

# ResNet: A little theory

One can view a sequence of outputs from residual layers as a **Dynamical System**
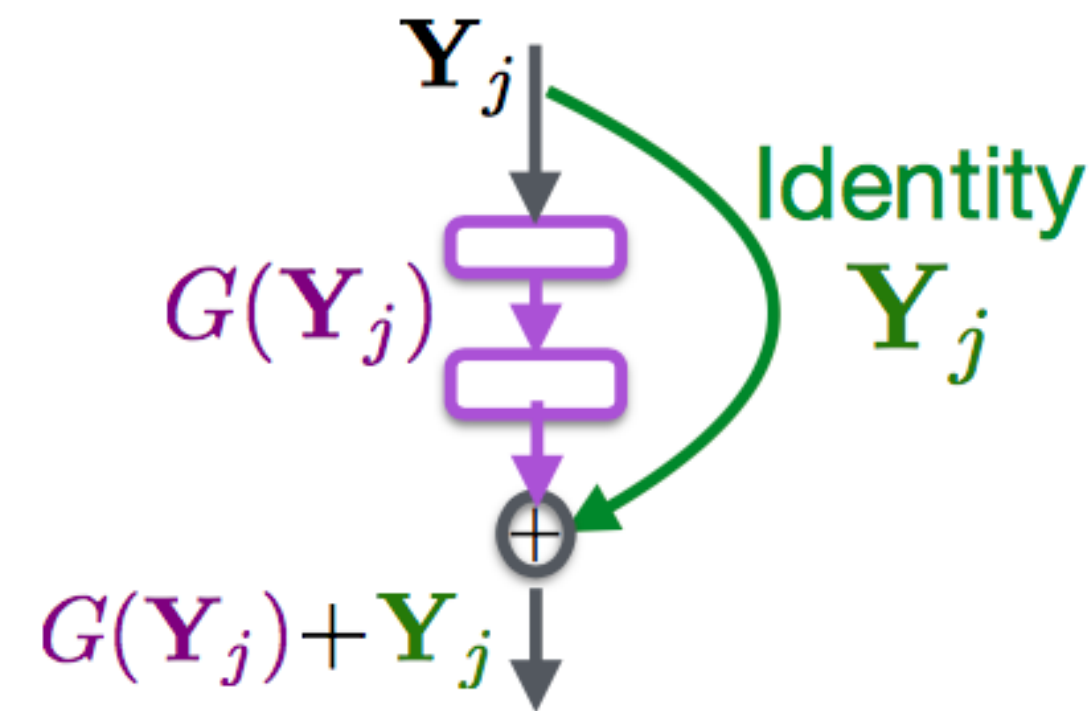


[ Cheng et al., ICLR 2018 ]

# **ResNet**: A little theory

One can view a sequence of outputs from residual layers as a **Dynamical System**



$$\mathbf{Y}_{j+1} = \mathbf{Y}_j + G(\mathbf{Y}_j, \boldsymbol{\theta}_j)$$

[ Cheng et al., ICLR 2018 ]

# **ResNet**: A little theory

One can view a sequence of outputs from residual layers as a **Dynamical System**



What happens if you take more layers and take smaller steps?

[ Chen et al., NIPS 2018 **best paper** ]

# **ResNet**: A little theory

One can view a sequence of outputs from residual layers as a **Dynamical System**



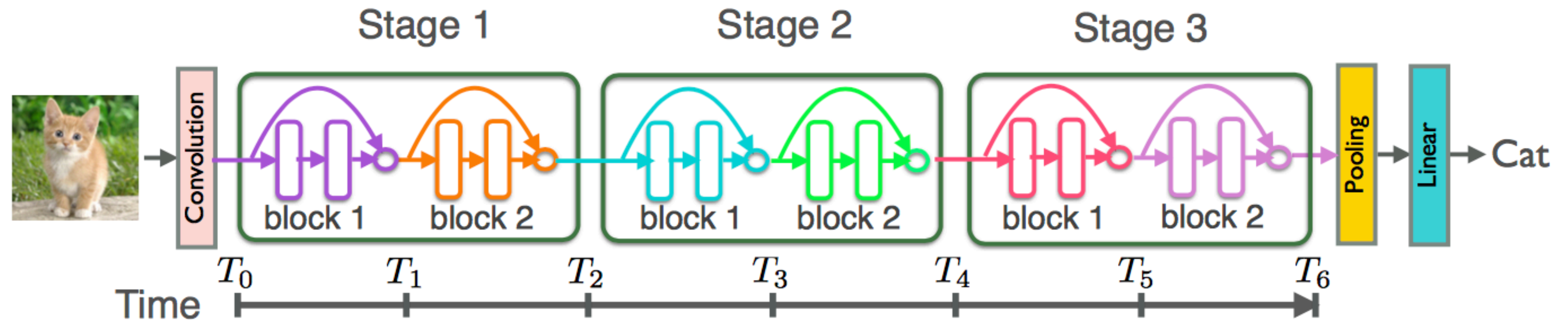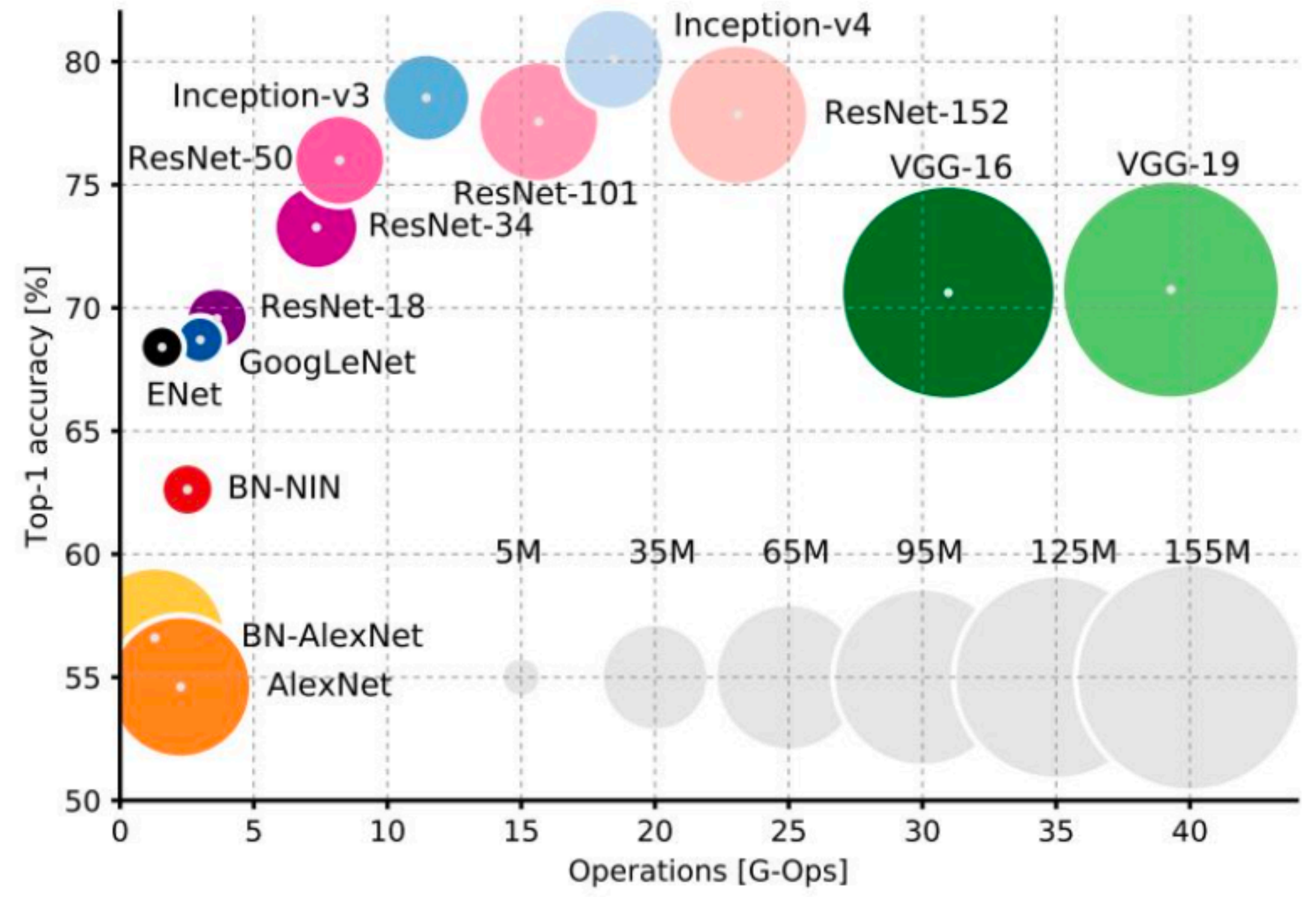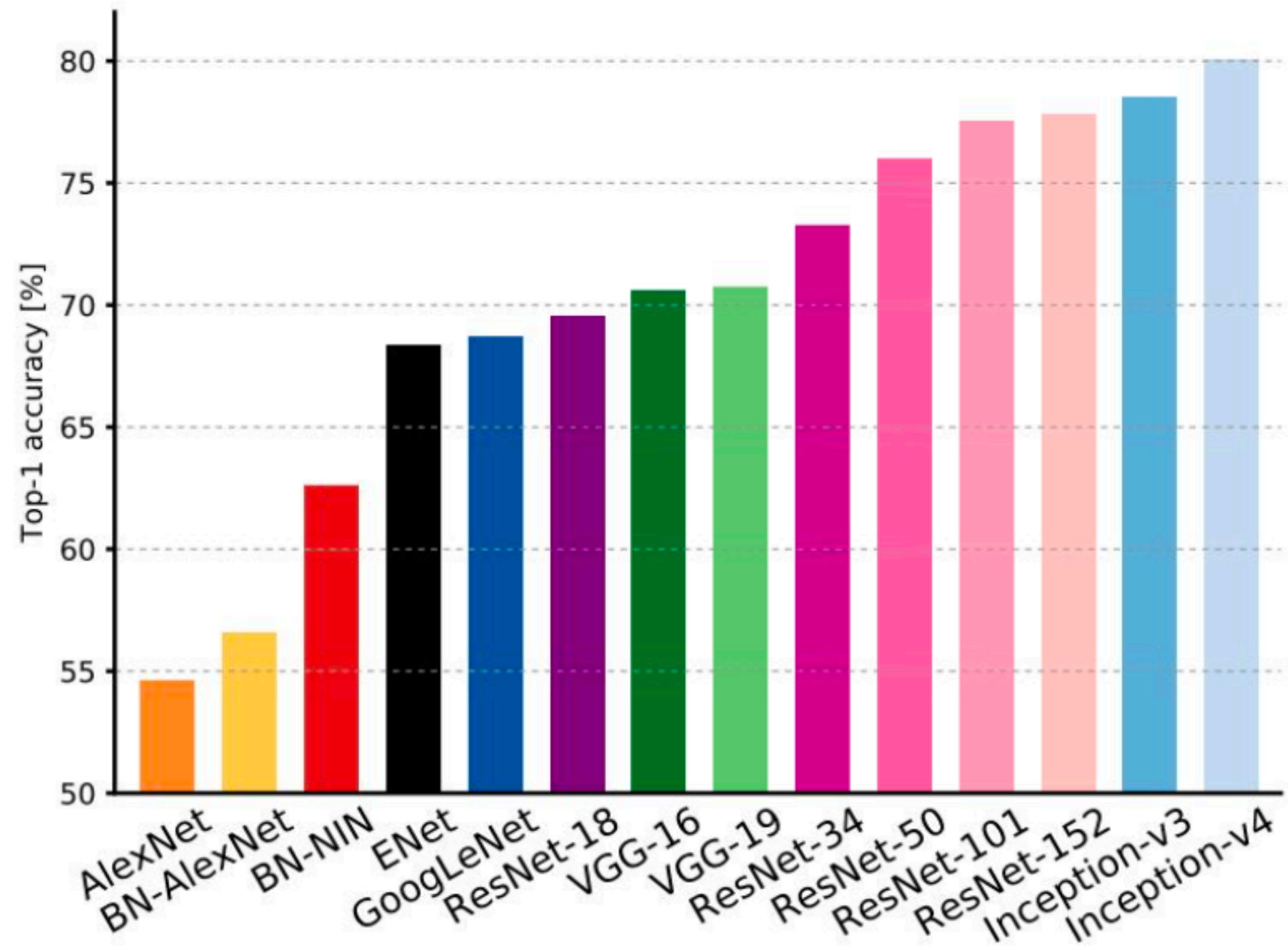What happens if you take more layers and take smaller steps?

You can actually treat a neural network as an **ODE**: $\dfrac{d\mathbf{h}(t)}{dt} = f(\mathbf{h}(t), t, \boldsymbol{\theta})$

[ Chen et al., NIPS 2018 **best paper** ]

# Comparing **Complexity**



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

# Computer **Vision Problems** (no language for now)

## Categorization



Multi-**class:**   Horse
Church
Toothbrush
**Person**

IM**A**GENET

Multi-**label:**   **Horse**
Church
Toothbrush
**Person**

## Detection



Horse (x, y, w, h)
Horse (x, y, w, h)
Person (x, y, w, h)
Person (x, y, w, h)

COCO
Common Objects in Context

## Segmentation



Horse
Person

COCO
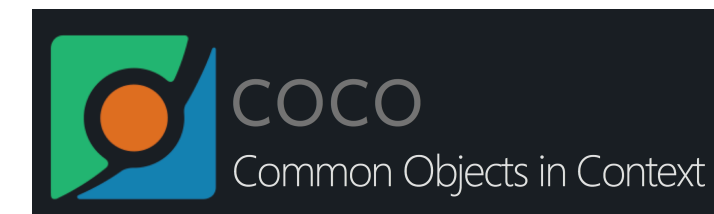Common Objects in Context

## Instance Segmentation



Horse1
Horse2
Person1
Person2

# Computer **Vision Problems** (no language for now)
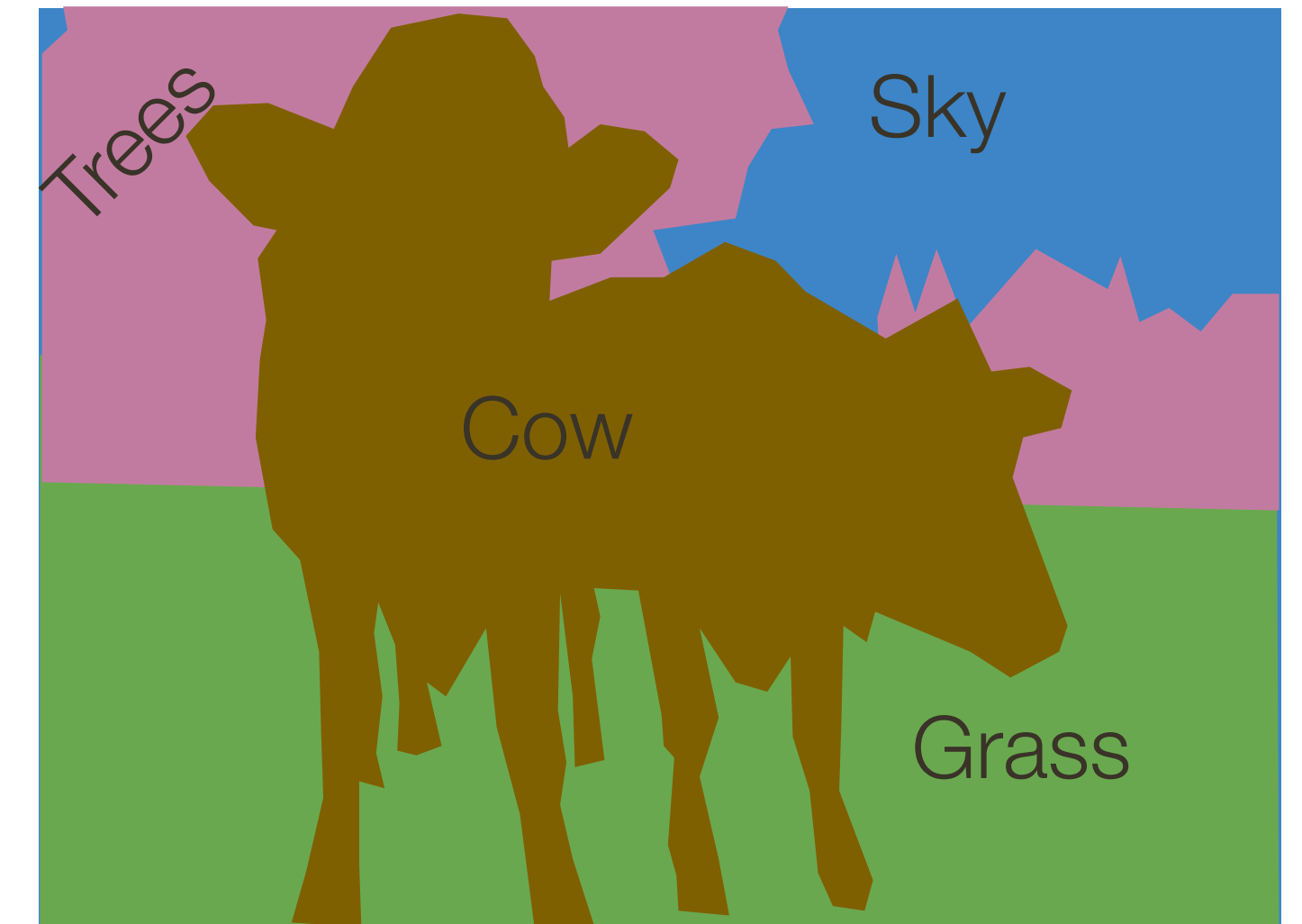
## Segmentation



Horse
Person

# Semantic **Segmentation**

Label **every pixel** with a category label (without differentiating instances)

# Semantic **Segmentation:** Sliding Window

Extract **patches**    **Classify center pixel** with CNN

Cow

Cow

Grass

# Semantic **Segmentation:** Sliding Window

Extract **patches**        **Classify center pixel** with CNN



Cow

Cow

Grass

**Problem:** VERY inefficient, no reuse of computations for overlapping patches

# Semantic **Segmentation:** Fully Convolutional CNNs

Design a network as a number of convolutional layers to make predictions for all pixels at once!



Input **Image**

3 x H x W

CONV, ReLU

CONV, ReLU

CONV, ReLU

**Convolutions**

D x H x W

Argmax

Class **Scores**

C x H x W

**Predicted** Labels

H x W

# Semantic **Segmentation:** Fully Convolutional CNNs

Design a network as a number of convolutional layers to make predictions for all pixels at once!
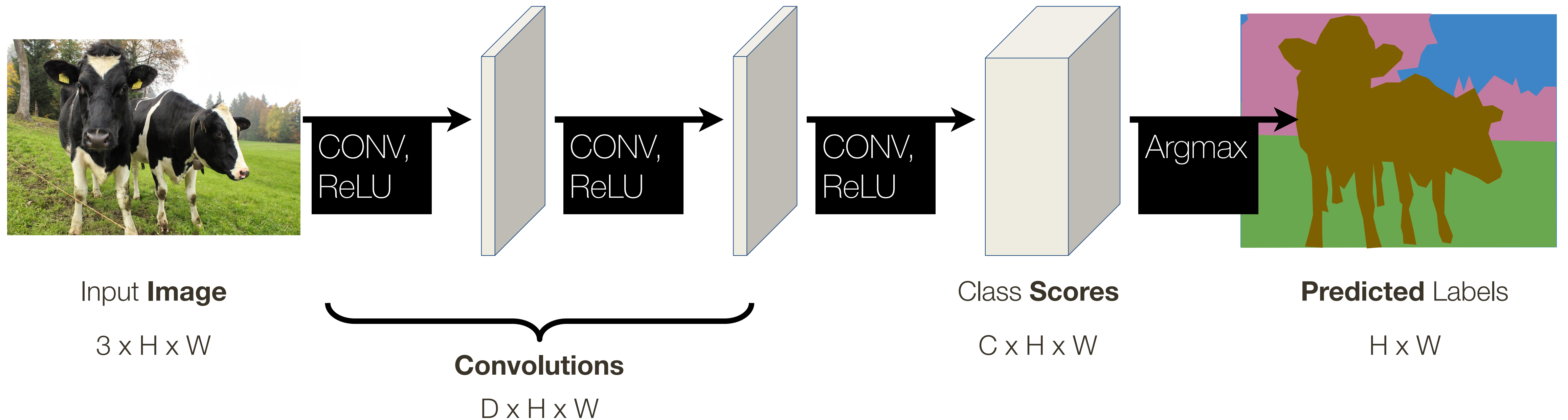


Input **Image**

3 x H x W

**Convolutions**

D x H x W

Class **Scores**

C x H x W

**Predicted** Labels

H x W

**Problem:** Convolutions at the original image scale will be very expensive

# Semantic **Segmentation:** Fully Convolutional CNNs

Design a network as a number 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 Image**

$3 \times H \times W$

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

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

**Predicted** Labels

$H \times W$

[ Long et al, CVPR 2015 ]
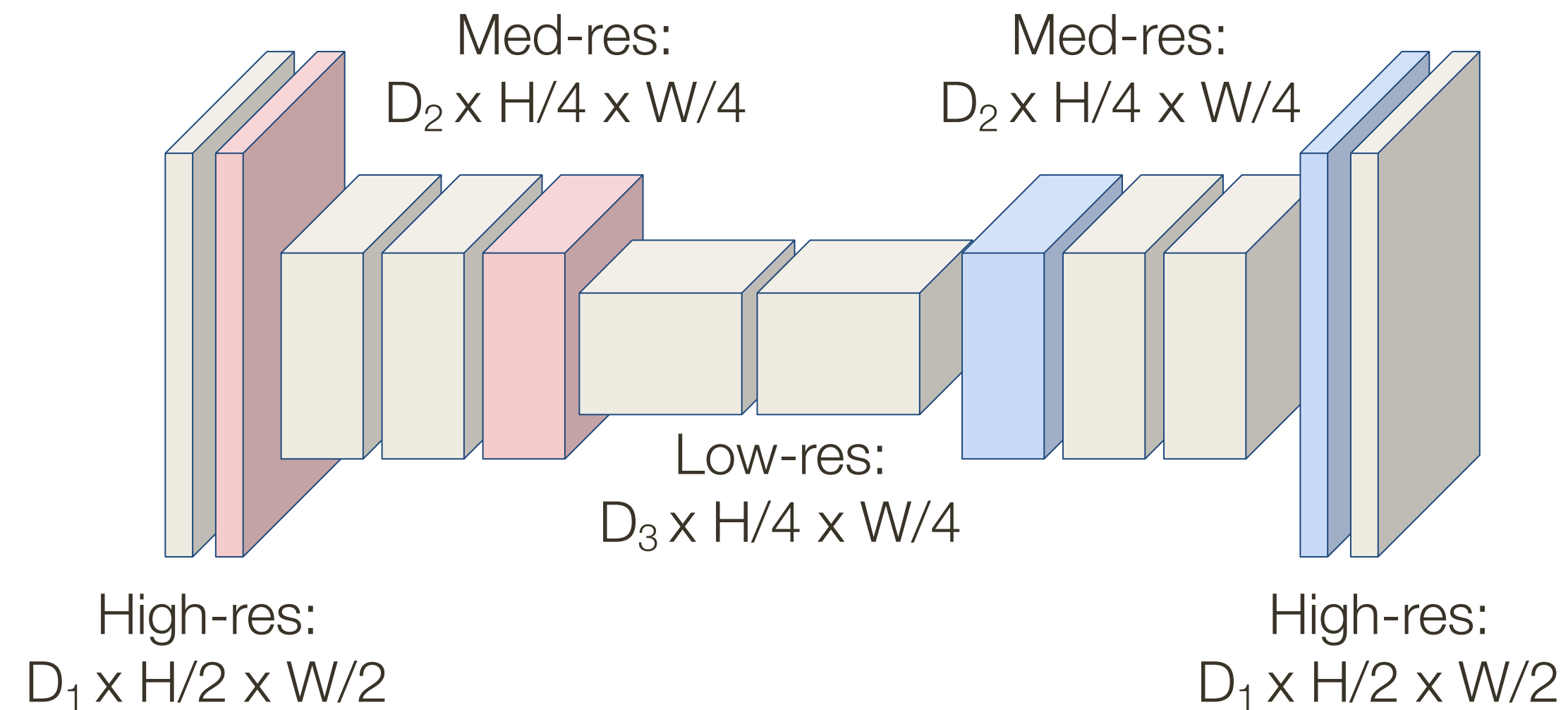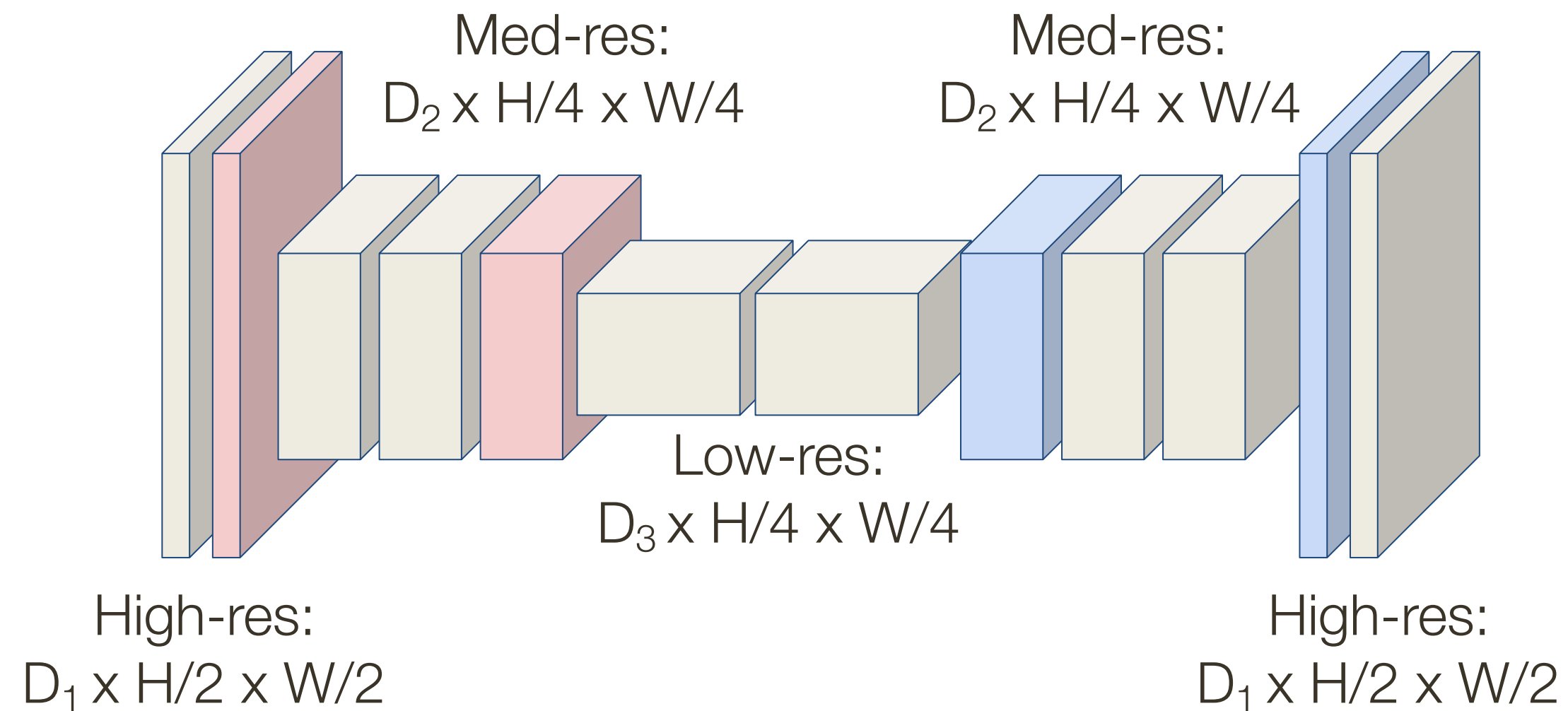[ Noh et al, ICCV 2015 ]

# Semantic **Segmentation:** Fully Convolutional CNNs

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



Input **Image**

$3 \times H \times W$

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

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

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

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

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

**Predicted** Labels

$H \times W$

**Downsampling** = Pooling

**Upsampling** = ???

[ Long et al, CVPR 2015 ]
[ Noh et al, ICCV 2015 ]

* slide from Fei-Dei Li, Justin Johnson, Serena Yeung, **cs231n Stanford**

# In-network **Up Sampling** (a.k.a "Unpooling")

Nearest Neighbor

$$\begin{array}{cc} 1 & 2 \\ 3 & 4 \end{array} \longrightarrow \begin{array}{cc|cc} 1 & 1 & 2 & 2 \\ 1 & 1 & 2 & 2 \\ \hline 3 & 3 & 4 & 4 \\ 3 & 3 & 4 & 4 \end{array}$$

**Input:** 2 x 2          **Output:** 4 x 4

# In-network **Up Sampling** (a.k.a "Unpooling")

**Nearest Neighbor**

$$\begin{matrix} 1 & 2 \\ 3 & 4 \end{matrix} \longrightarrow \begin{matrix} 1 & 1 & 2 & 2 \\ 1 & 1 & 2 & 2 \\ 3 & 3 & 4 & 4 \\ 3 & 3 & 4 & 4 \end{matrix}$$

**Input:** 2 x 2          **Output:** 4 x 4

**"Bed of Nails"**

$$\begin{matrix} 1 & 2 \\ 3 & 4 \end{matrix} \longrightarrow \begin{matrix} 1 & 0 & 2 & 0 \\ 0 & 0 & 0 & 0 \\ 3 & 0 & 4 & 0 \\ 0 & 0 & 0 & 0 \end{matrix}$$

**Input:** 2 x 2          **Output:** 4 x 4

# In-network **Up Sampling:** Max Unpooling

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

**Max Unpooling**
Use positions from pooling layer



**Input:** 4 x 4

**Output:** 2 x 2

Rest of the network

**Input:** 2 x 2

**Output:** 4 x 4

Corresponding pairs of downsampling and upsampling layers

# In-network **Up Sampling:** 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 **Up Sampling:** 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 **Up Sampling:** Transpose Convolution

**Recall:** Normal 3 x 3 convolution, stride 2 pad 1



Dot product between filter and input

**Input:** 4 x 4

**Output:** 2 x 2

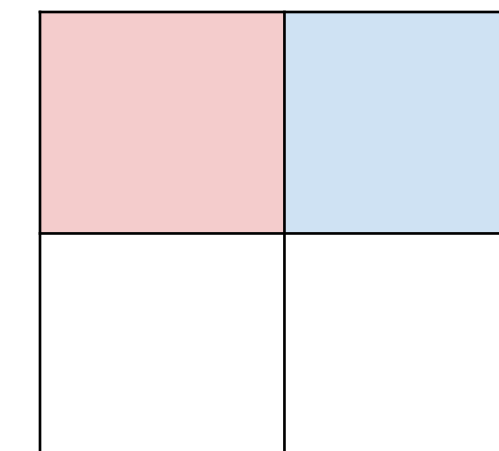# In-network **Up Sampling:** Transpose Convolution

**Recall:** Normal 3 x 3 convolution, stride 2 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 in movement in input vs output
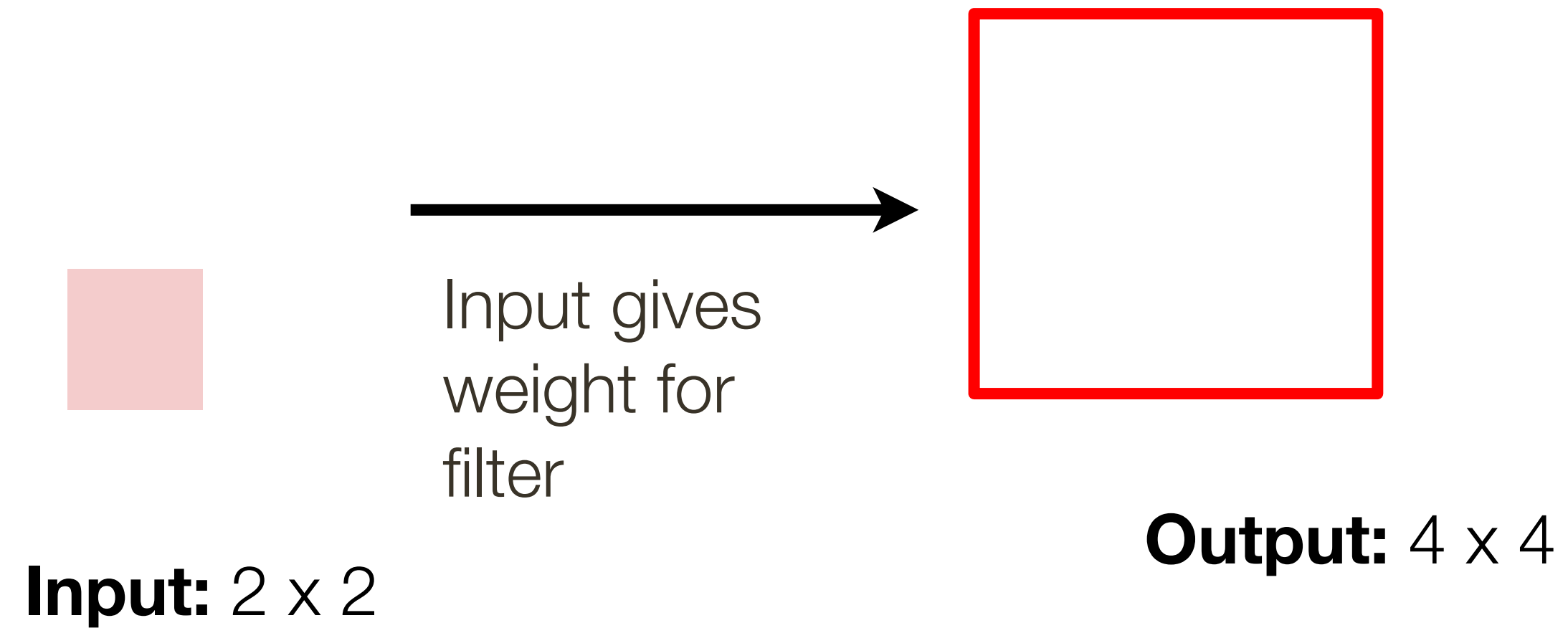
# In-network **Up Sampling:** Transpose Convolution

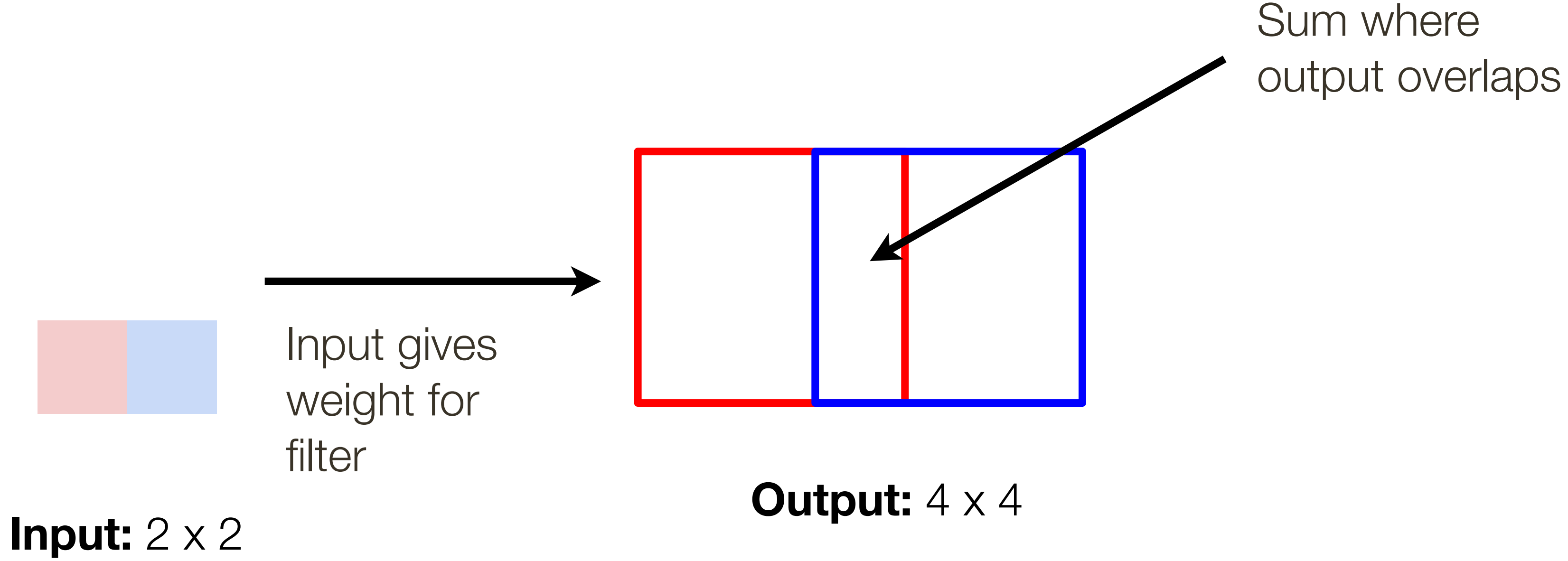3 x 3 **transpose** convolution, stride 2 pad 1

**Input:** 2 x 2

**Output:** 4 x 4

# In-network **Up Sampling:** Transpose Convolution

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

Input gives
weight for
filter

**Output:** 4 x 4

**Input:** 2 x 2

# In-network **Up Sampling:** Transpose Convolution

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

Sum where
output overlaps

Input gives
weight for
filter

**Input:** 2 x 2

**Output:** 4 x 4
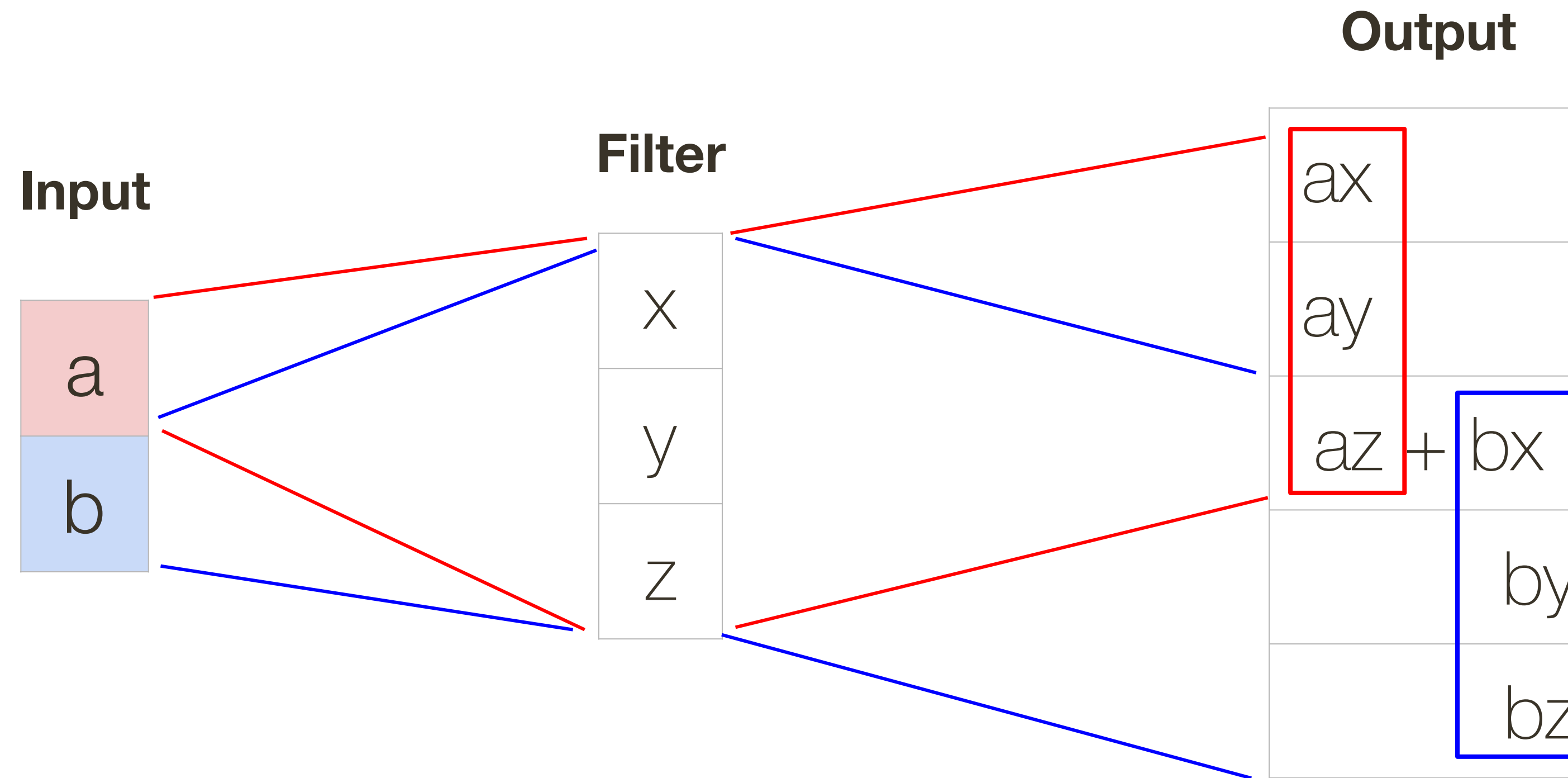
Filter moves 2 pixels in the **output** for every one
pixel in the **input**

Stride gives ratio in movement in output vs input
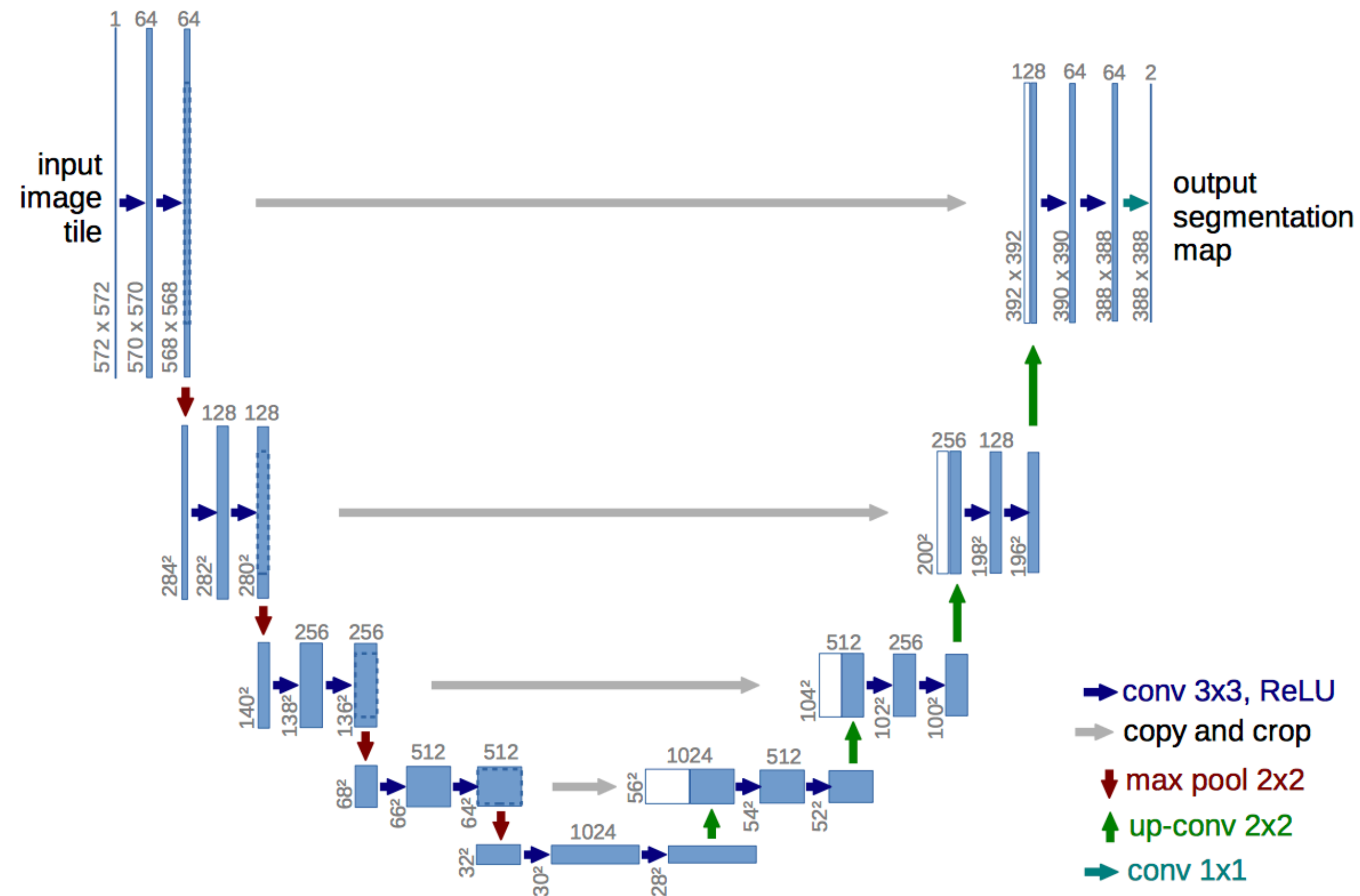
# **Transpose Convolution**: 1-D Example

**Output**

**Filter**

**Input**

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

# **U-Net** Architecture

ResNet-like Fully convolutional CNN

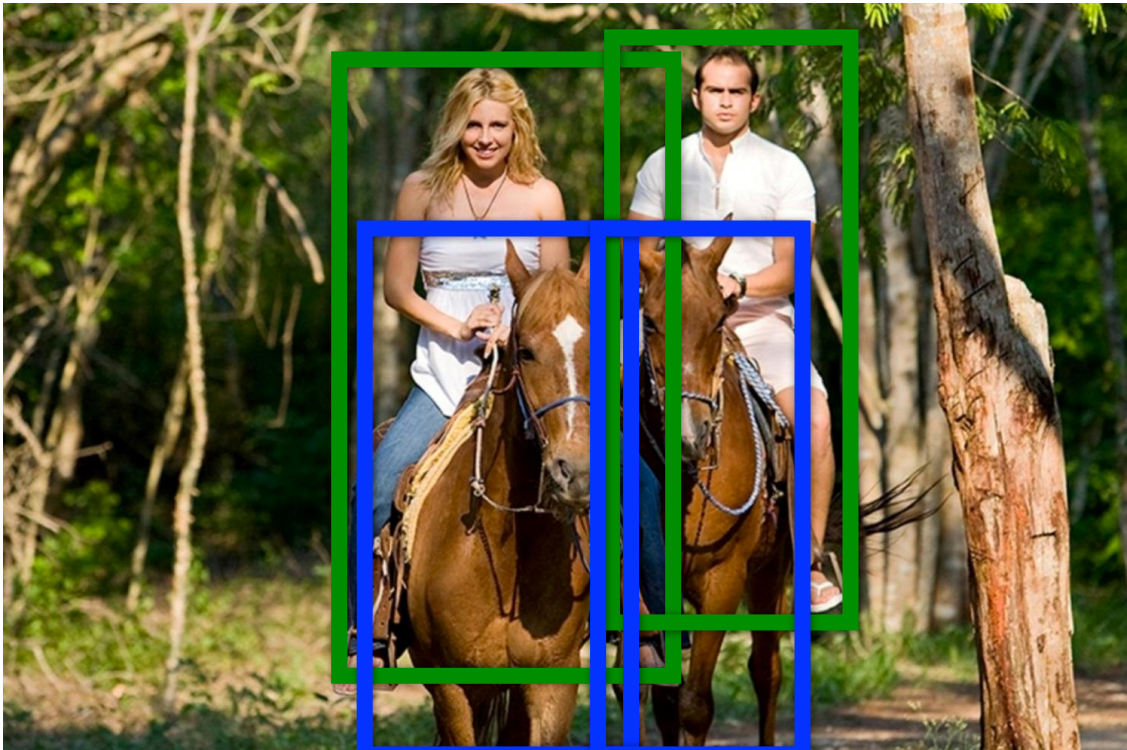# Computer **Vision Problems** (no language for now)

| Categorization | Detection | Segmentation | Instance Segmentation |
|---|---|---|---|



**Categorization**

Multi-**class:**　Horse
　　　　　　　Church
　　　　　　　Toothbrush
　　　　　　　**Person**

IMAGENET

Multi-**label**:　**Horse**
　　　　　　　Church
　　　　　　　Toothbrush
　　　　　　　**Person**

**Detection**

Horse (x, y, w, h)
Horse (x, y, w, h)
Person (x, y, w, h)
Person (x, y, w, h)

COCO
Common Objects in Context

**Segmentation**

Horse
Person

COCO
Common Objects in Context

**Instance Segmentation**

Horse1
Horse2
Person1
Person2