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

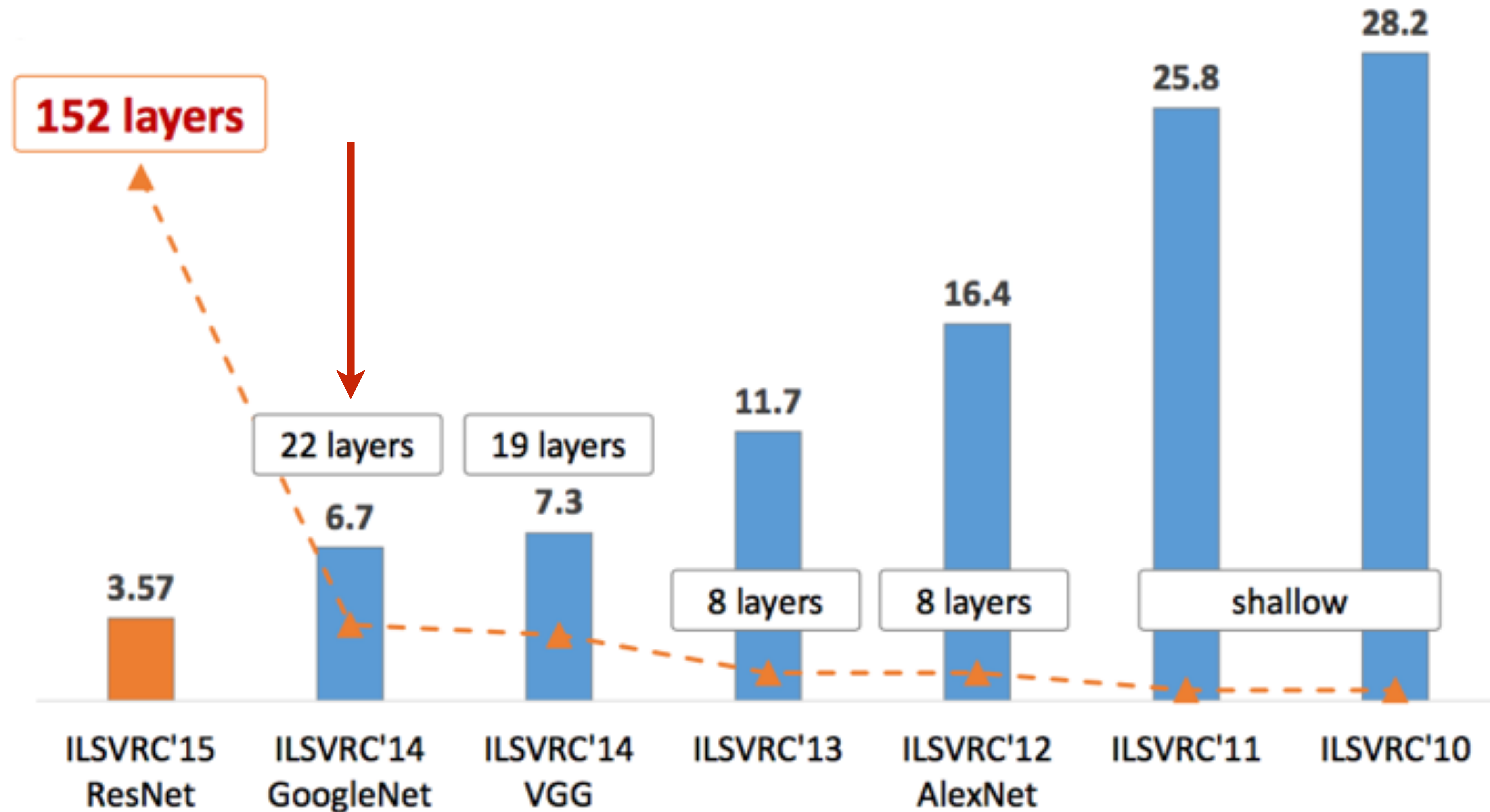**Lecture 6: Convolutional Neural Networks (Part 3)**

# **Logistics**:

**Assignment 2** is due on **Wednsday, 11:59pm**

Groups + idea for project by **Thursday, January 31st**
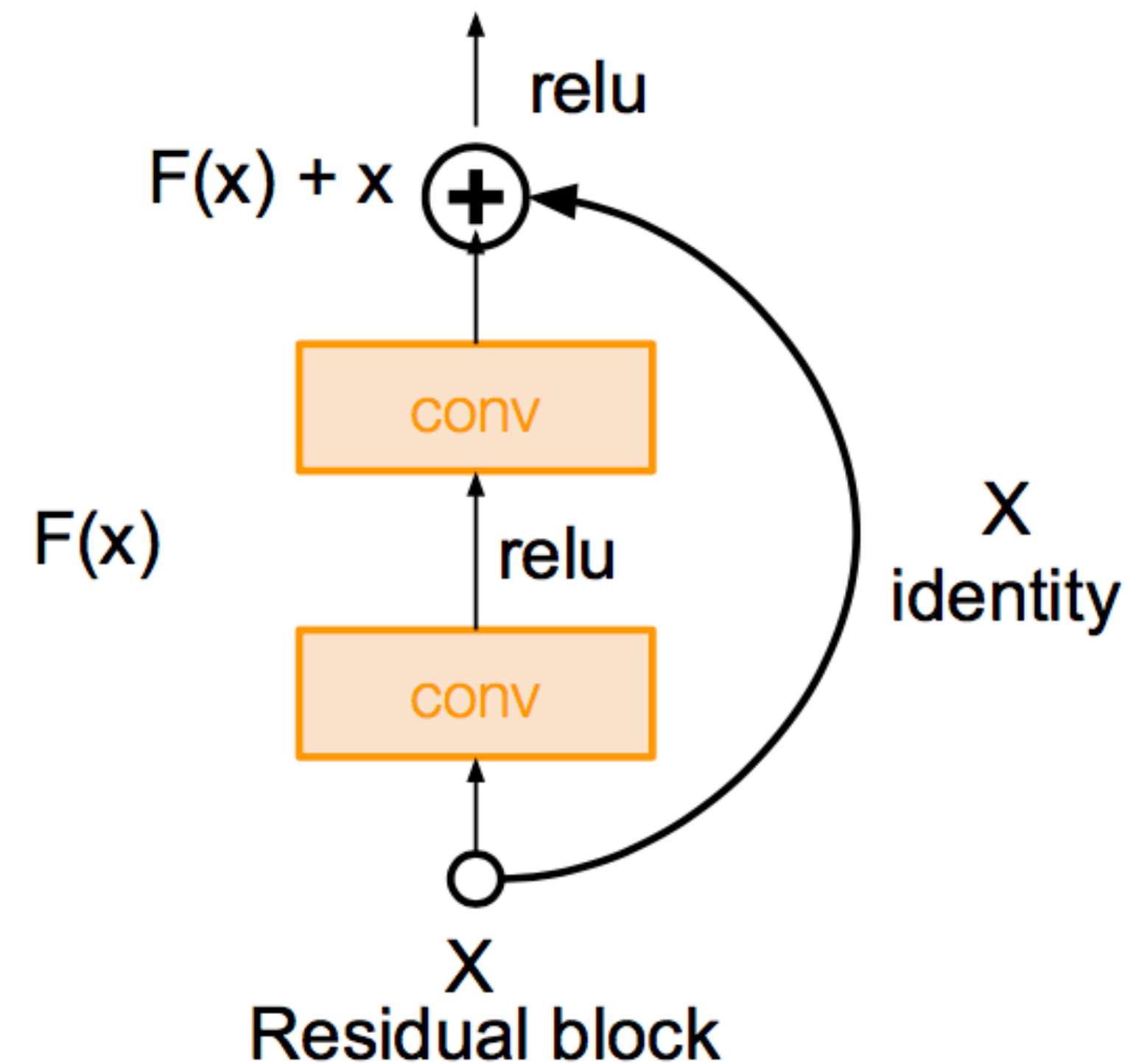
**Paper list** to be posted by **Monday, January 28th**
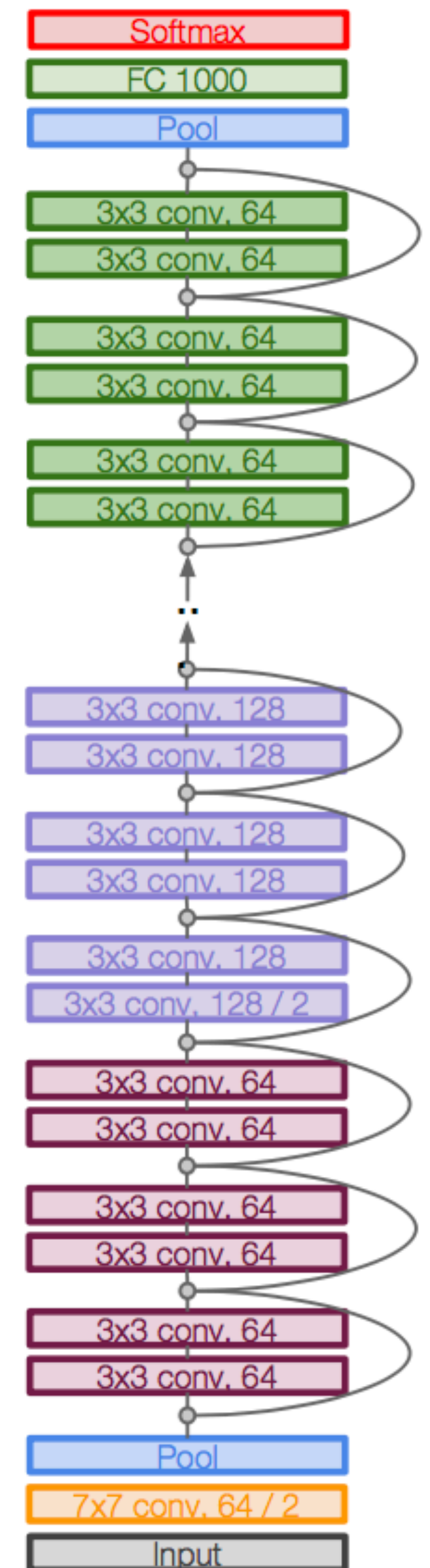
# **ILSVRC** winner 2012



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

# 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
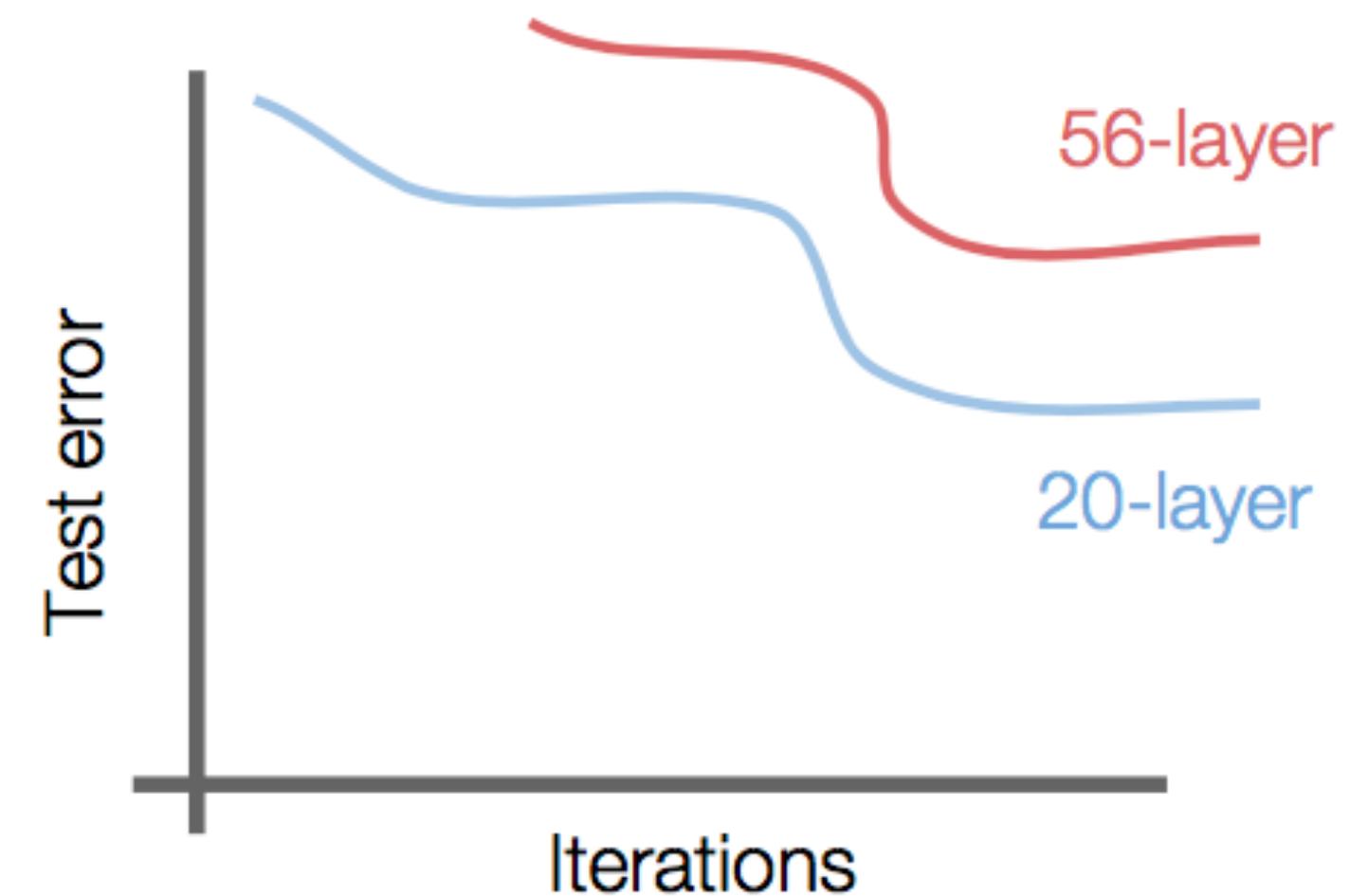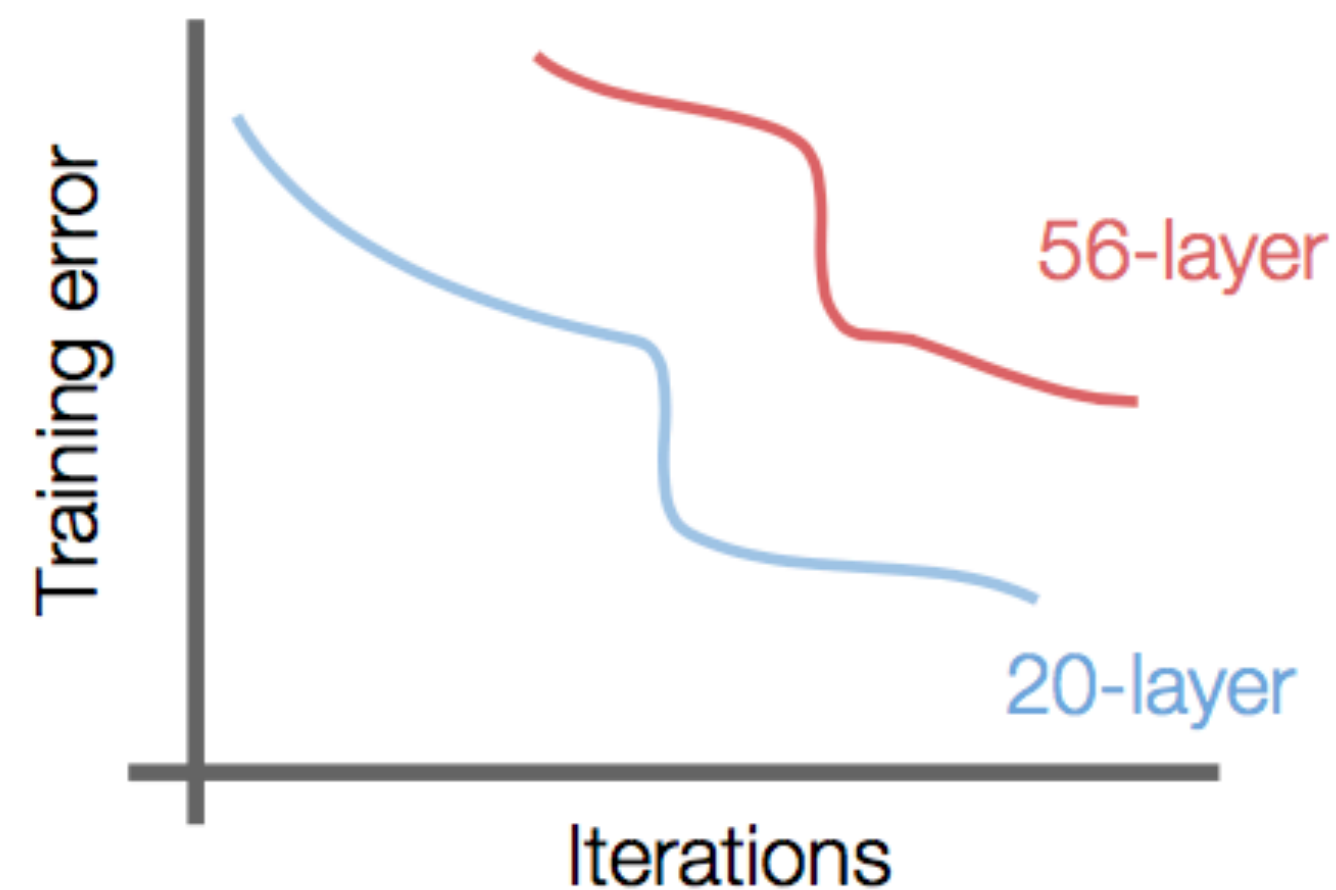
# ResNet: Motivation

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



Whats the **problem**?

# ResNet: Motivation

[ He et al., 2015 ]

**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)
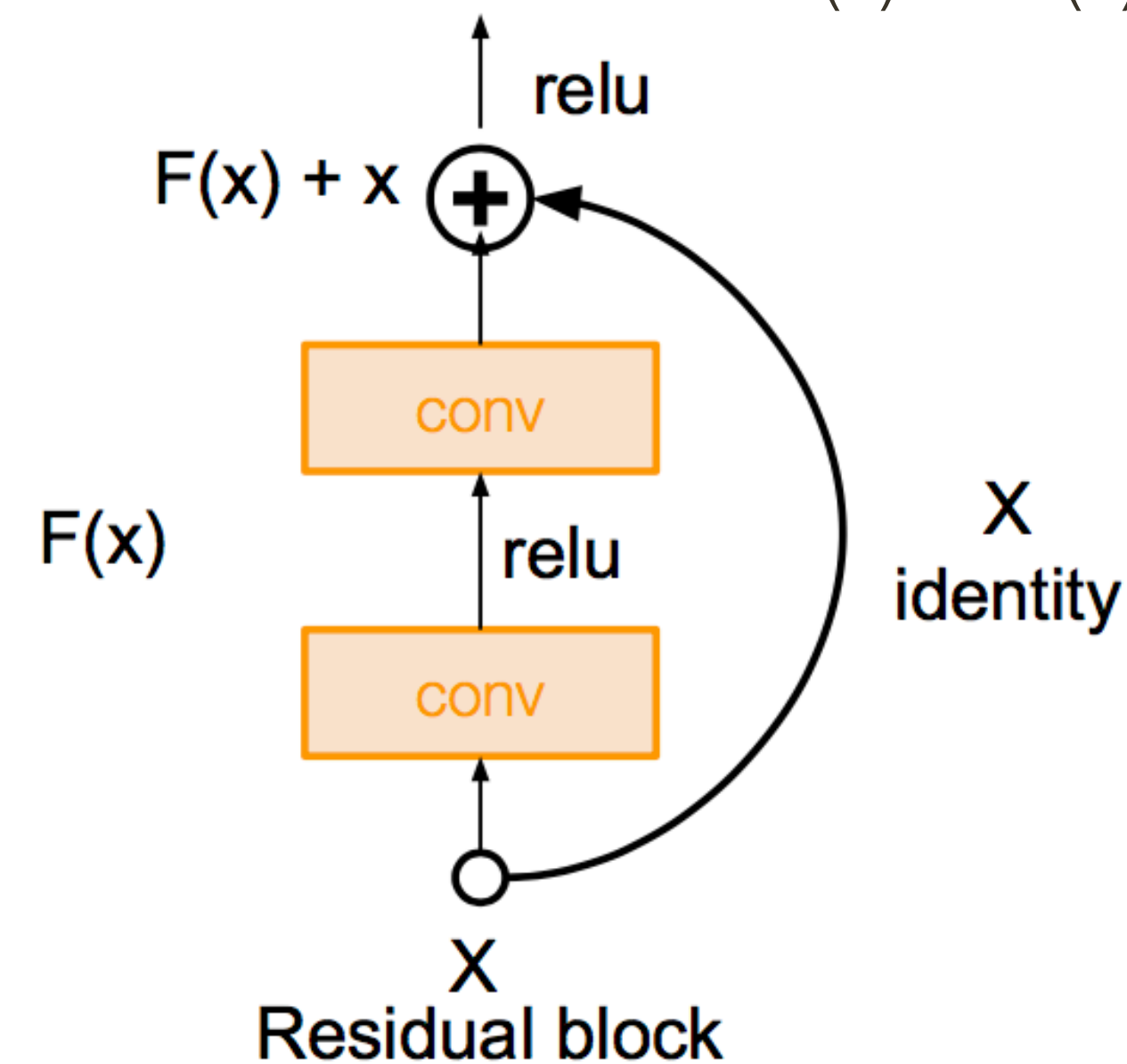
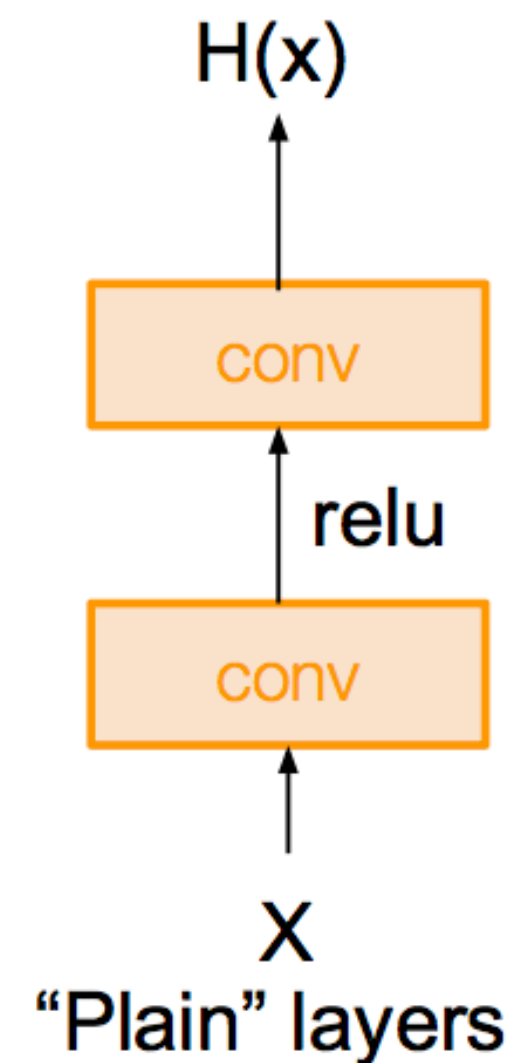How do we implement this idea in practice

# 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**
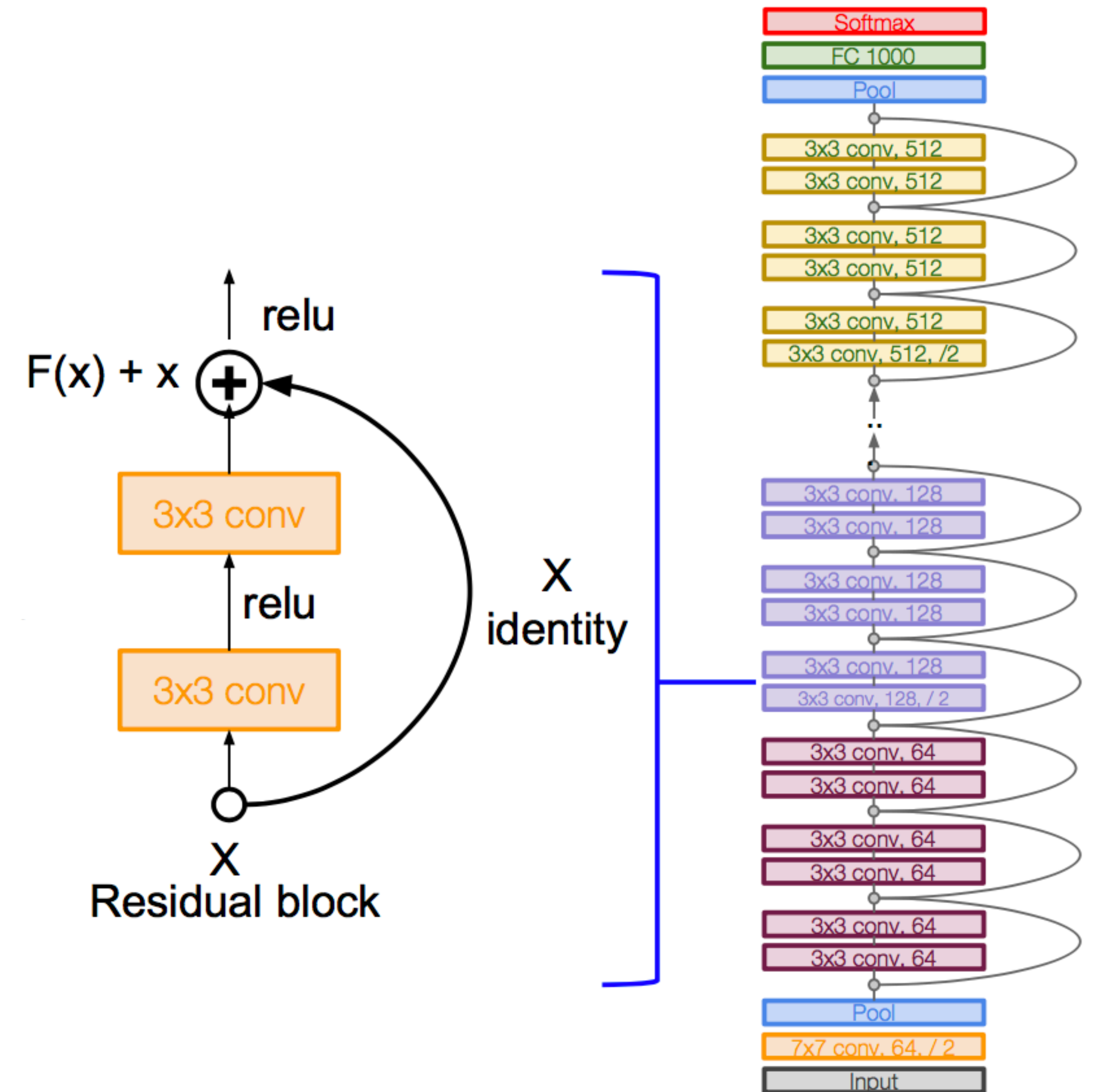$F(x) = H(x) - X$ instead of $H(x)$ directly
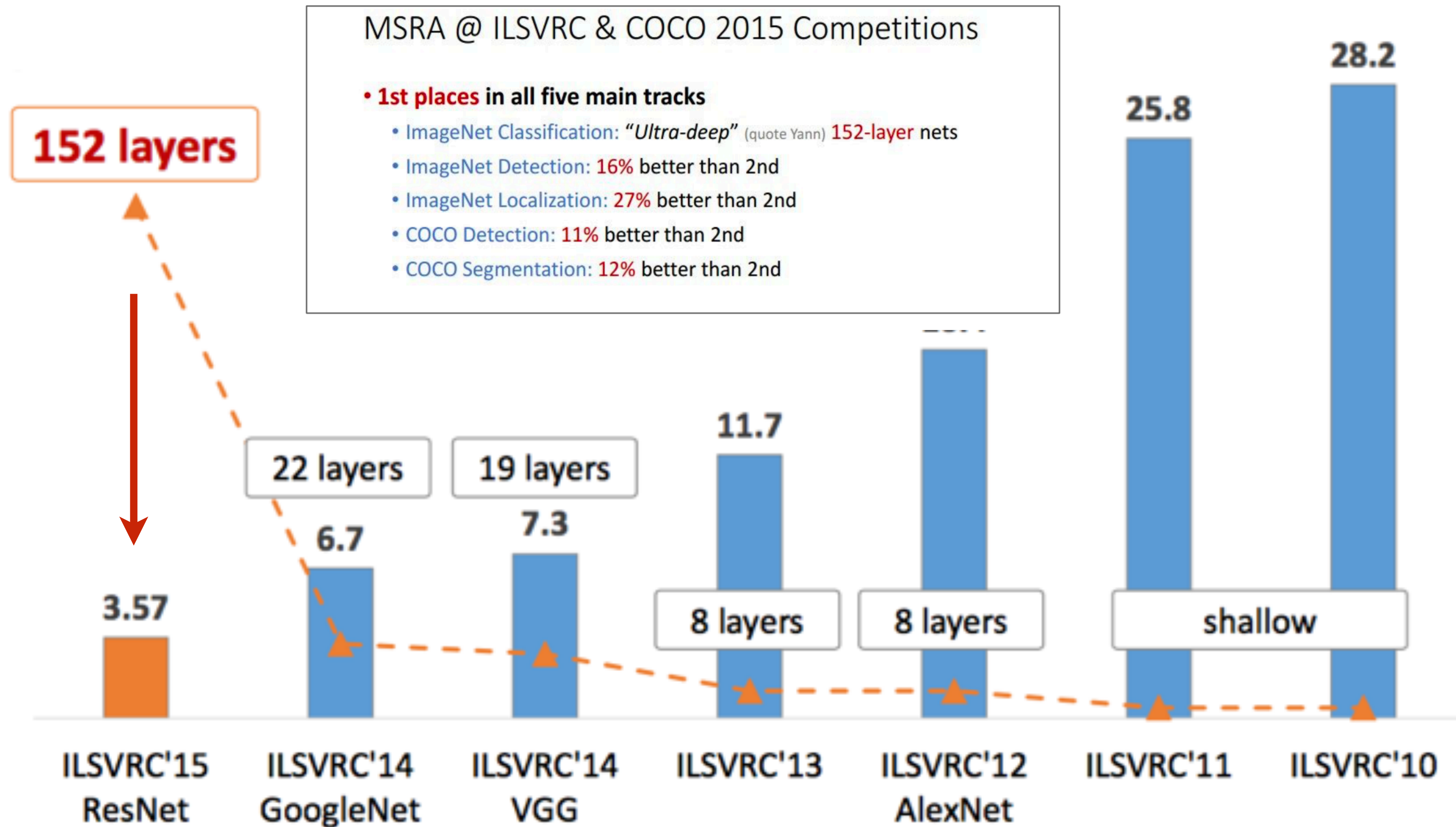


"Plain" layers

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)



Residual block

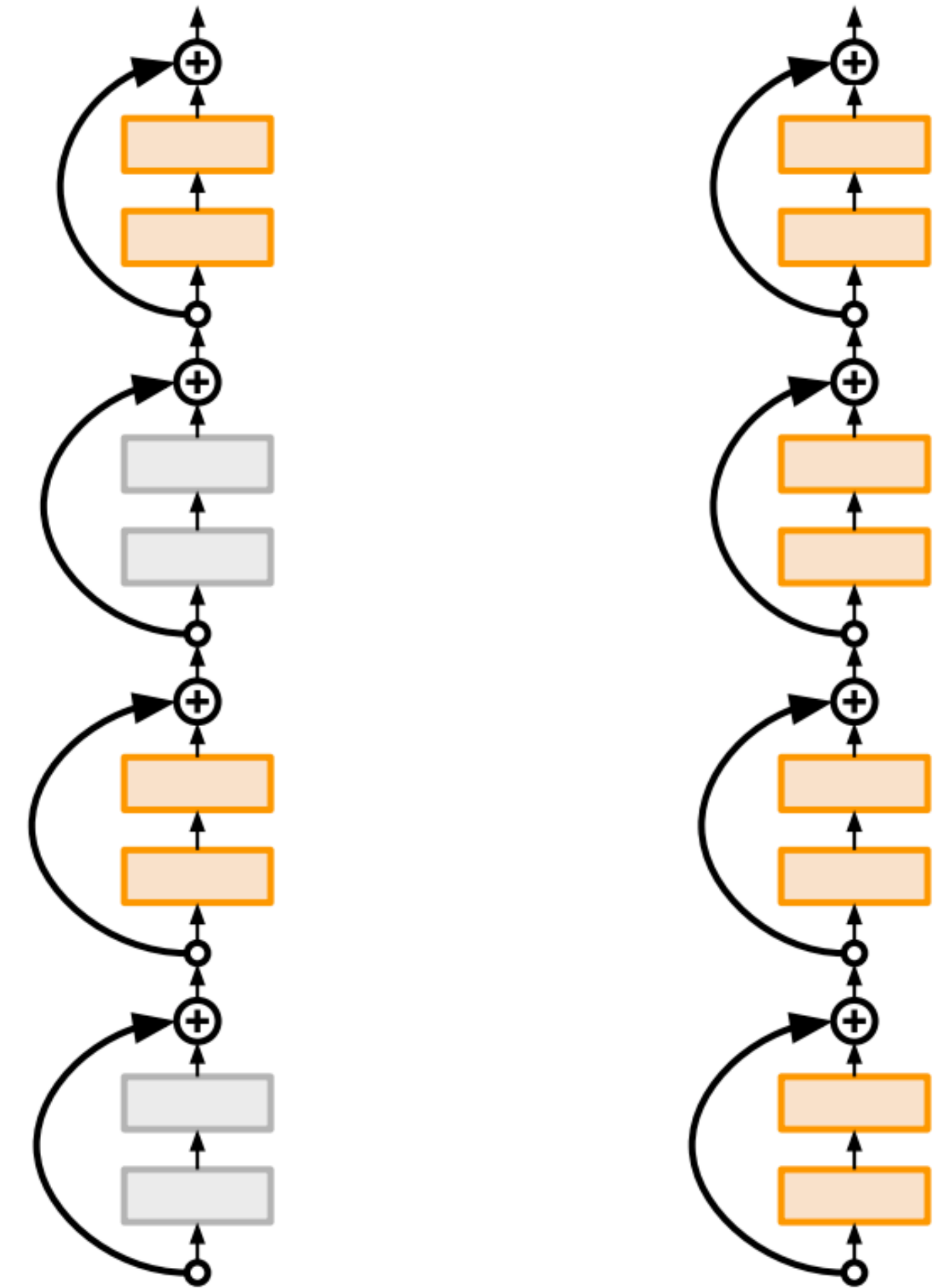X
identity

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



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

11.7

8 layers    8 layers    shallow

25.8    28.2

6.7    7.3

3.57

ILSVRC'15 ResNet · ILSVRC'14 GoogleNet · ILSVRC'14 VGG · ILSVRC'13 · ILSVRC'12 AlexNet · ILSVRC'11 · ILSVRC'10
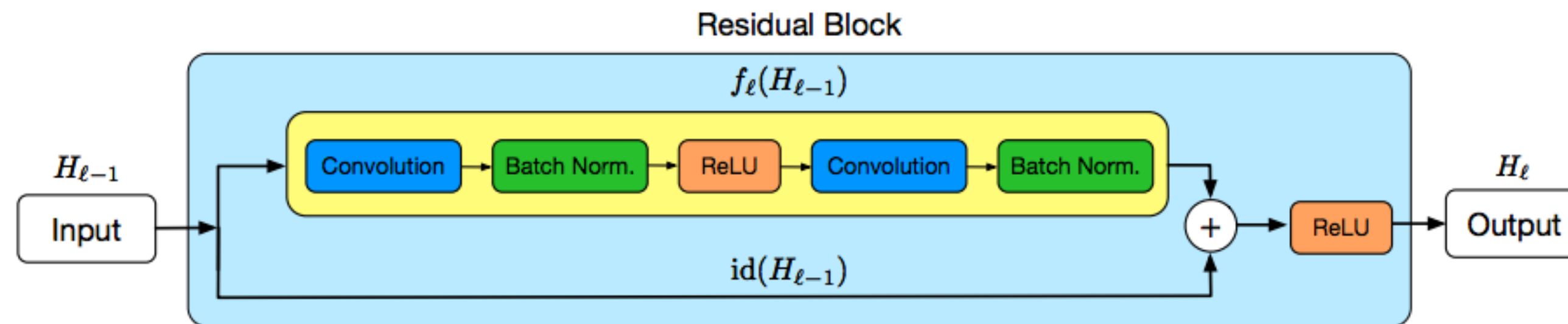
# Regularization: Stochastic Depth

Effectively "dropout" but for layers

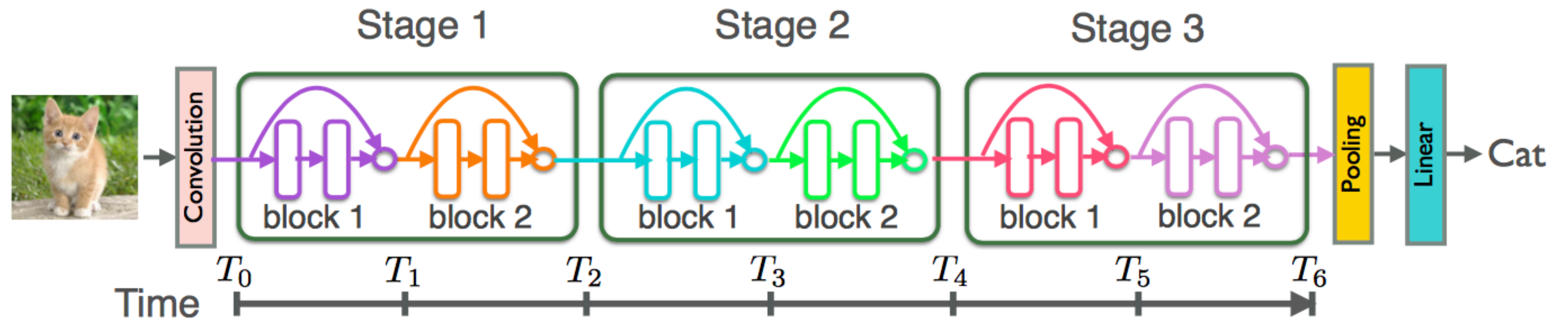Stochastically with some probability **turn off some layer** (for each batch)

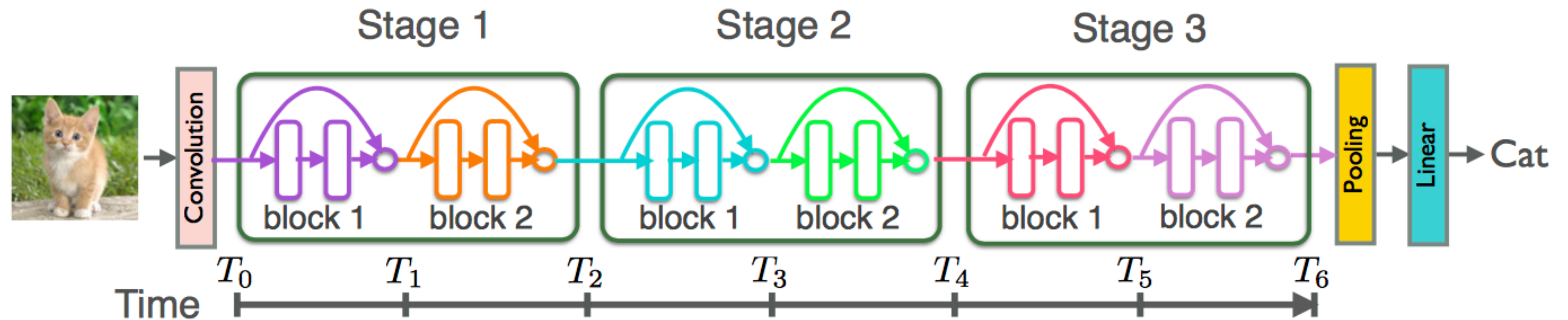Effectively trains a collection of neural networks

# ResNet: A little theory

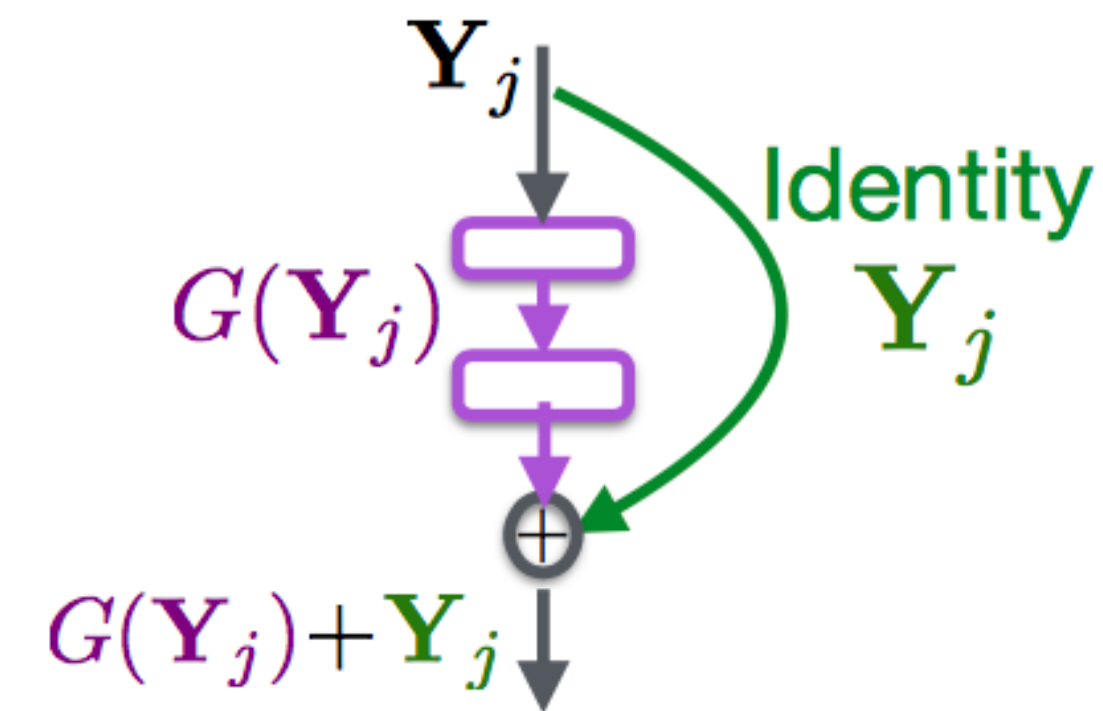One can view a sequence of outputs from residual layers as a **Dynamical System**



[ Cheng et al., ICLR 2018 ]

# **ResNet**: A little theory

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



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

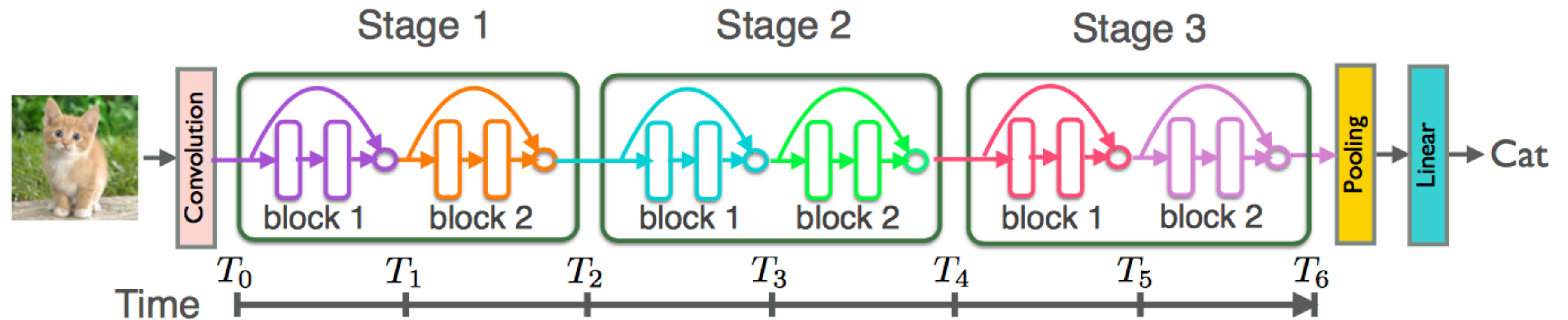[ Cheng et al., ICLR 2018 ]

# **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

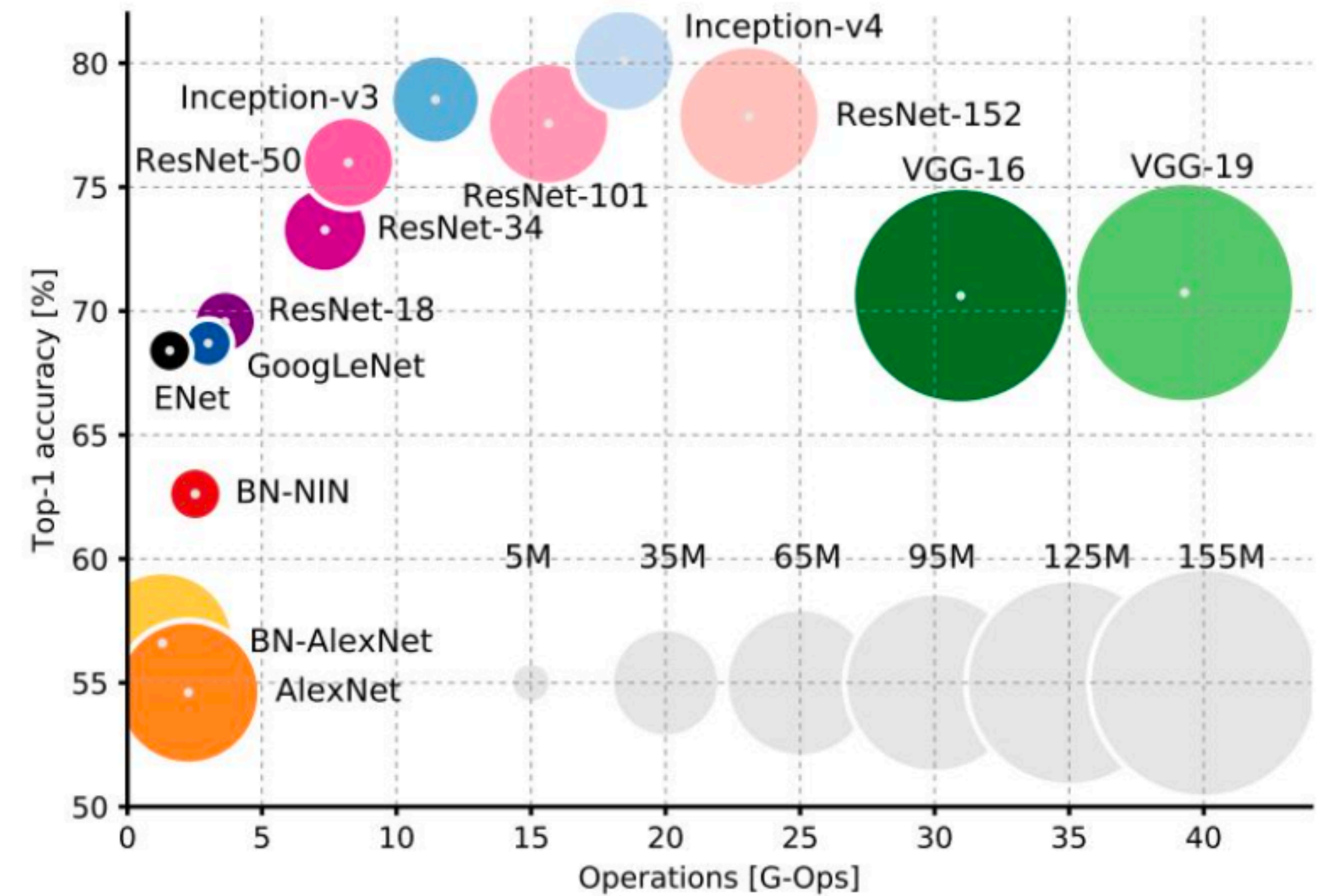One can view a sequence of outputs from residual layers as a **Dynamical System**
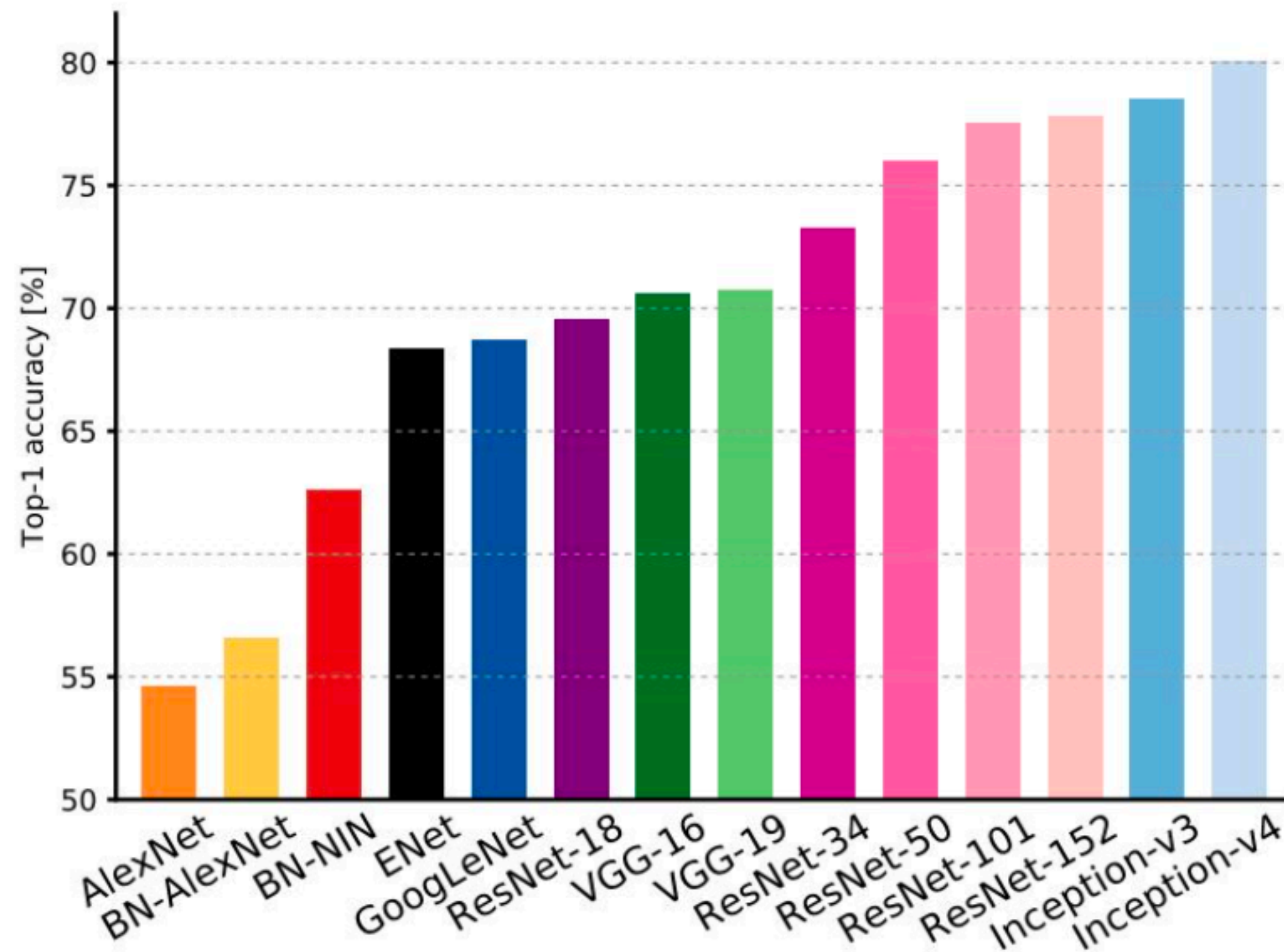


What happens if you take more layers and take smaller steps?

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

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

# Comparing **Complexity**



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

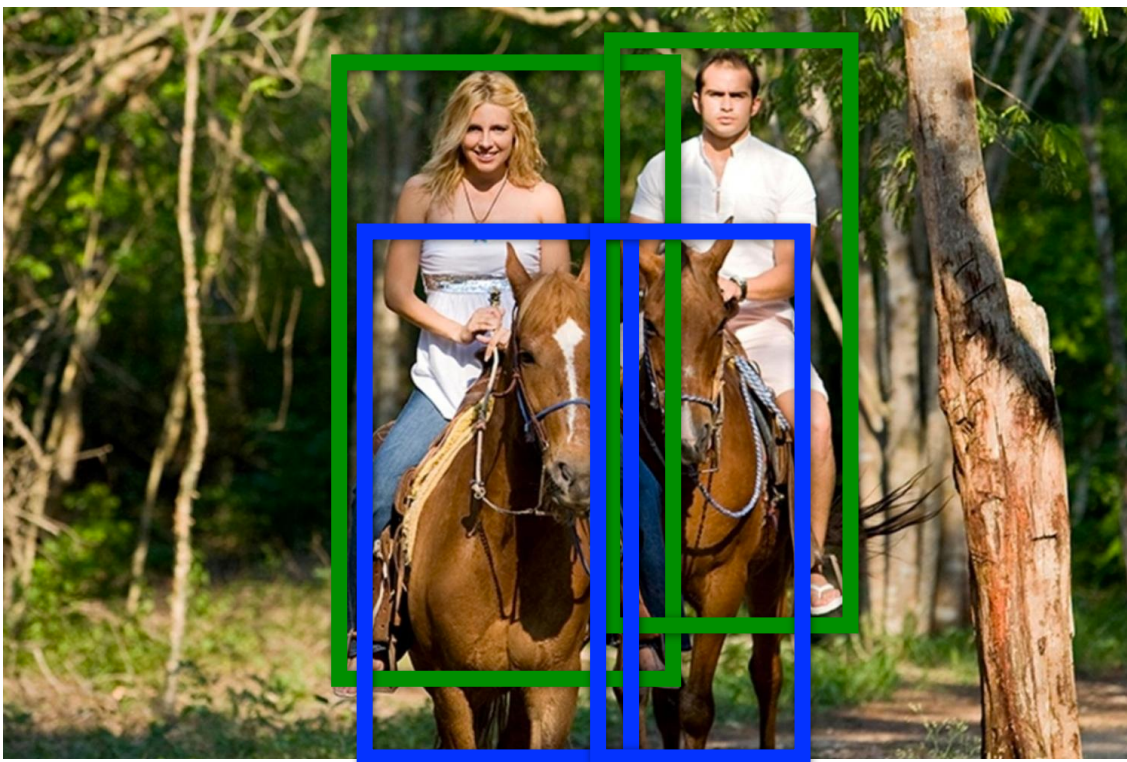# Computer **Vision Problems** (no language for now)

## Categorization



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

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
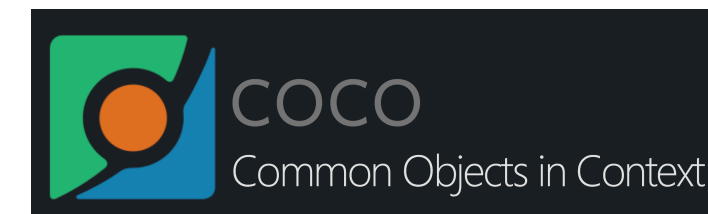Common Objects in Context

## Instance Segmentation



Horse1
Horse2
Person1
Person2

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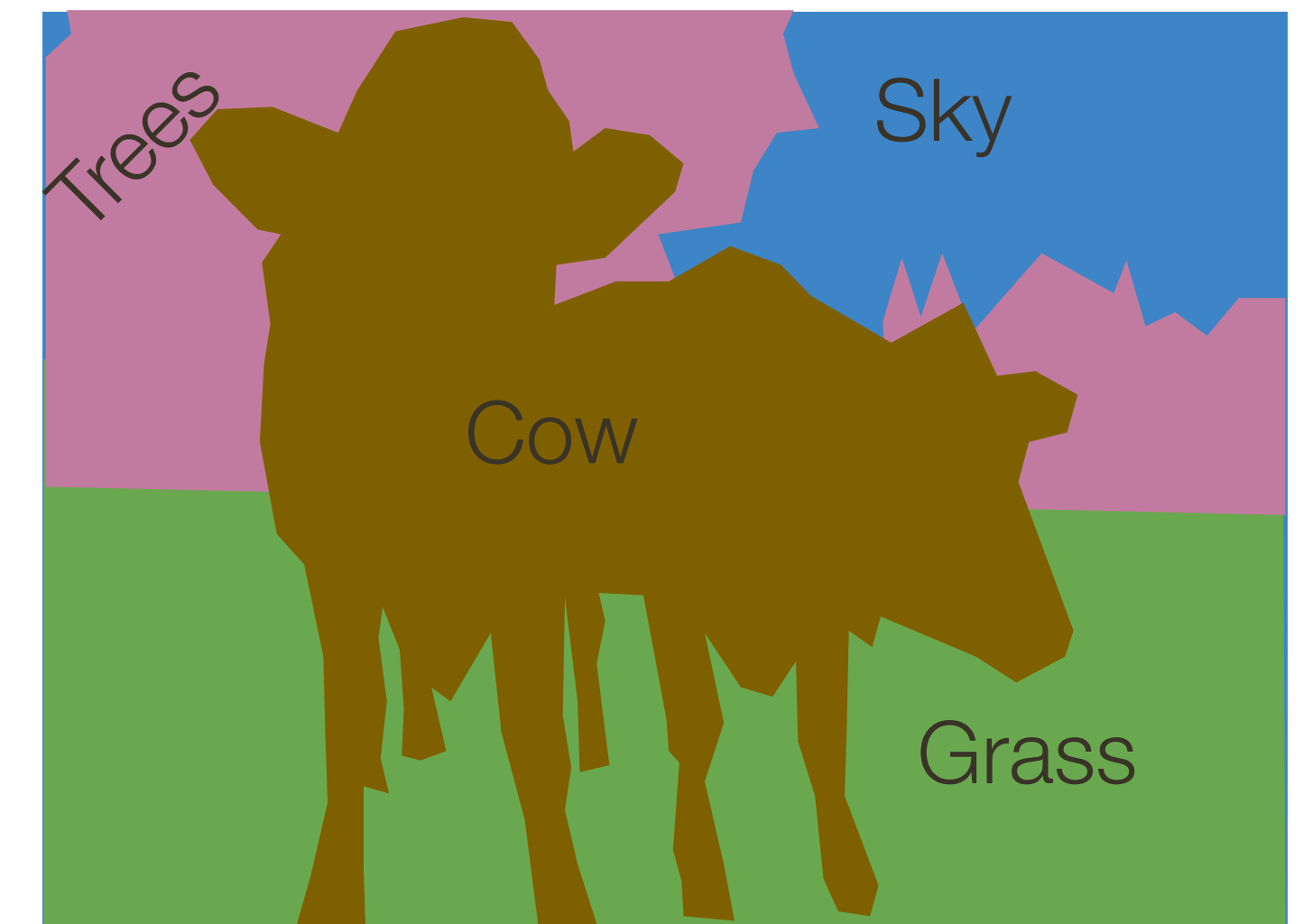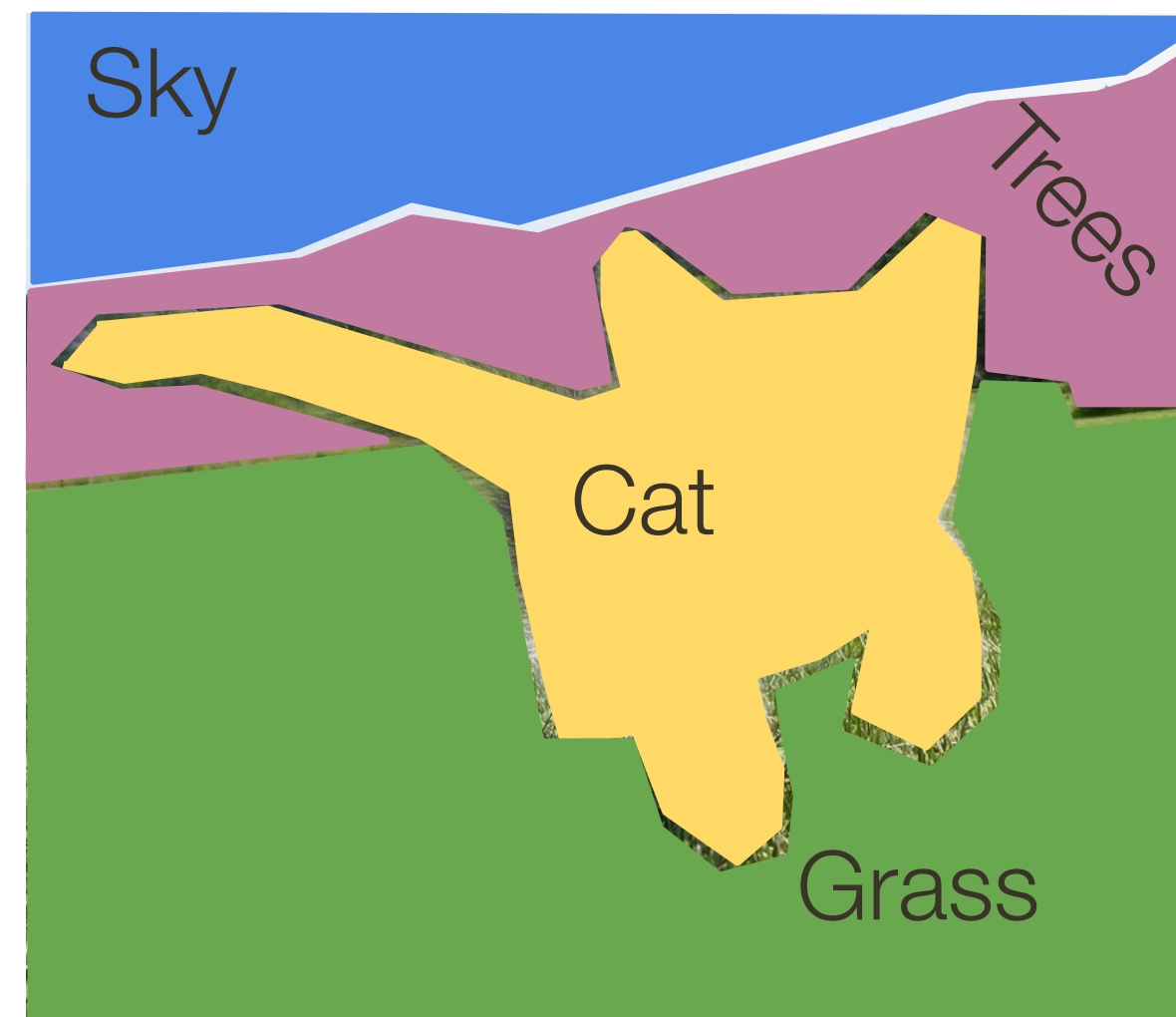
## Segmentation



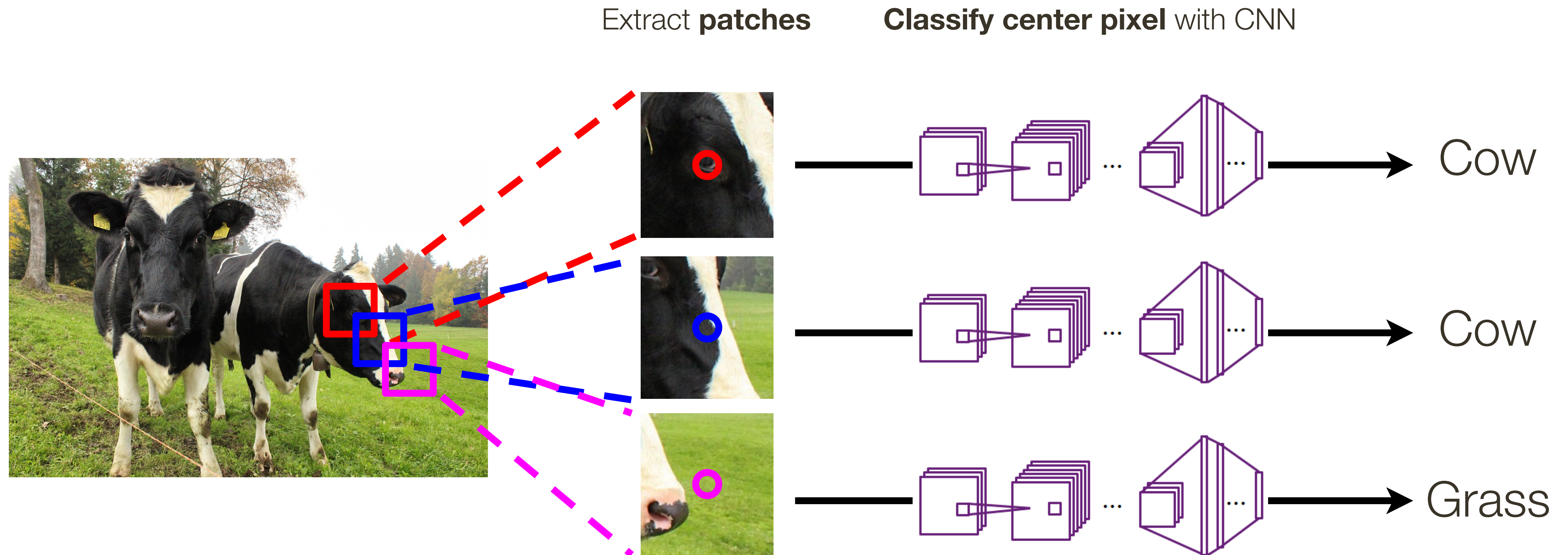Horse
Person

COCO
Common Objects in Context

# Semantic **Segmentation**

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

# Semantic **Segmentation:** Sliding Window

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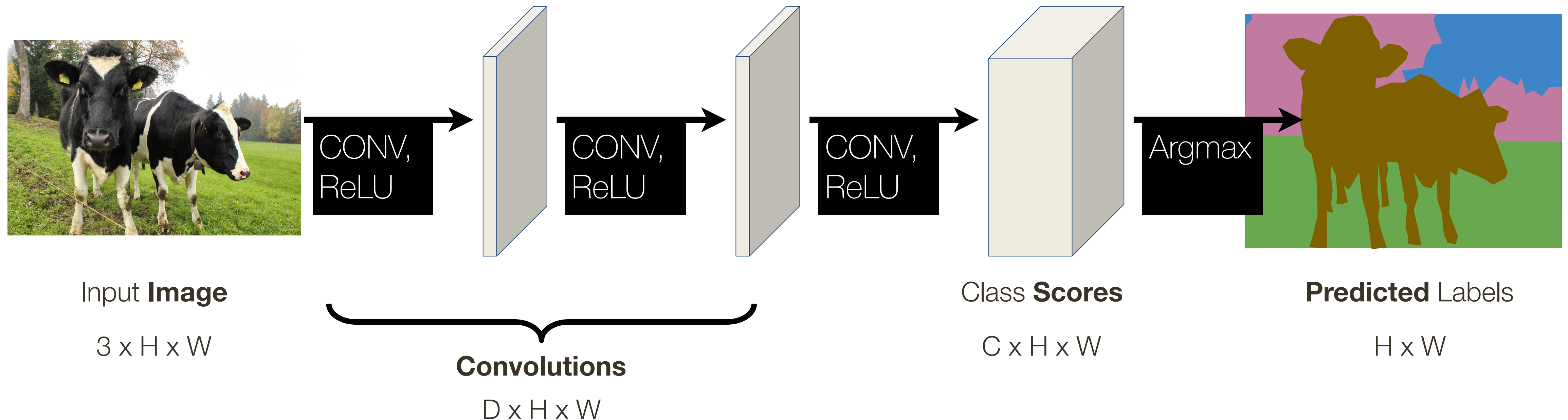


Cow

Cow

Grass

**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

**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
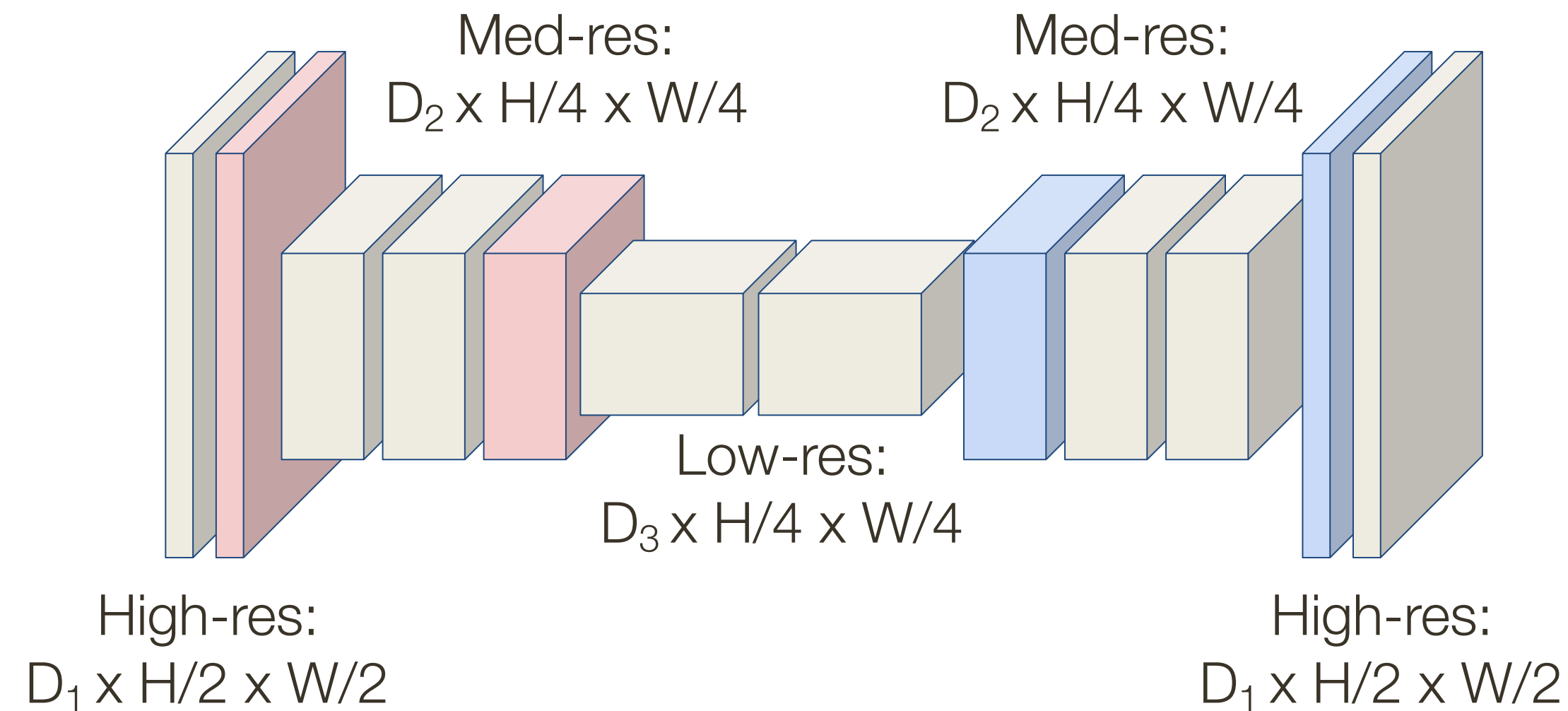
# Semantic **Segmentation:** Fully Convolutional CNNs

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



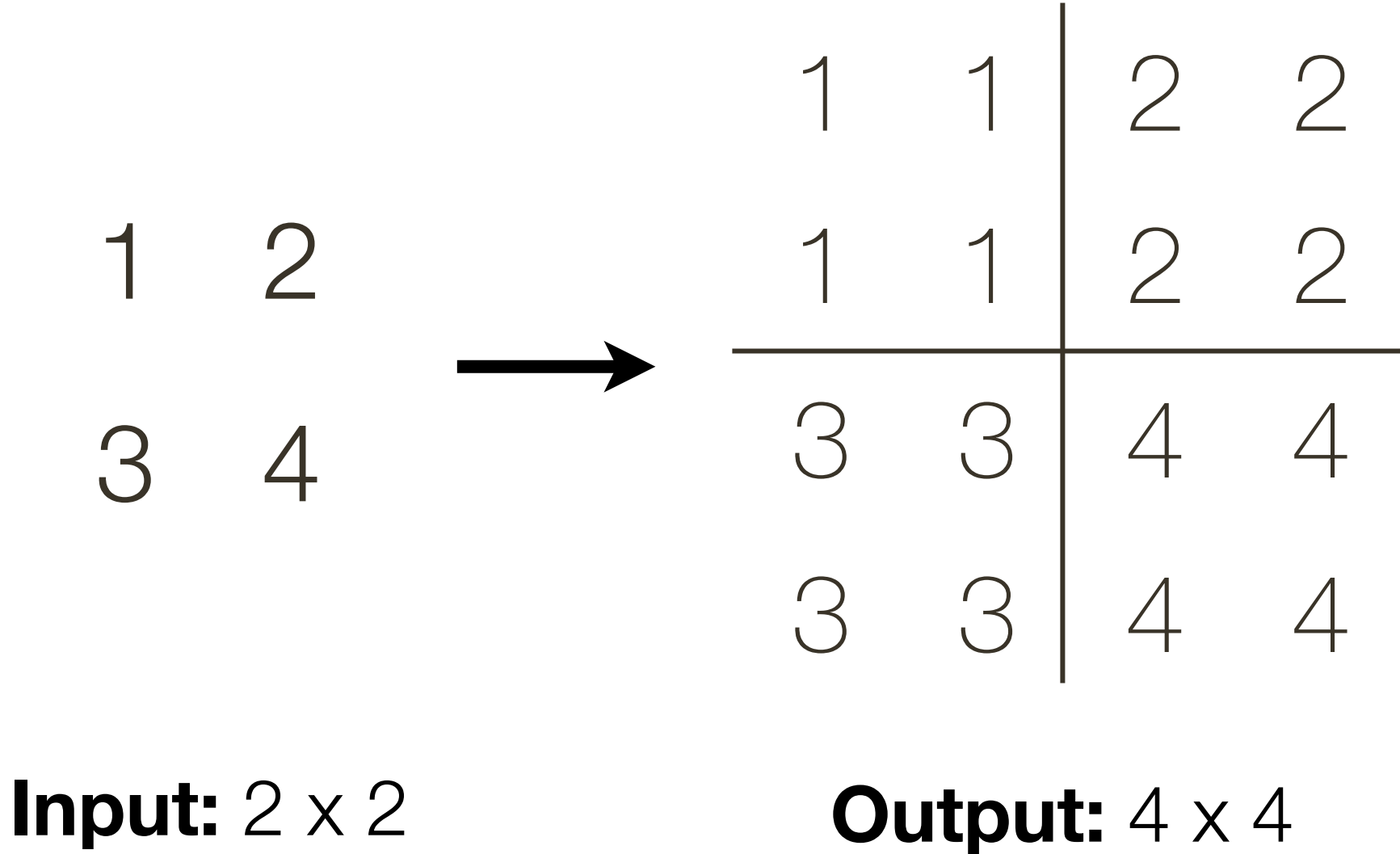Input **Image**

$3 \times H \times W$

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**

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

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

**"Bed of Nails"**

$$
\begin{bmatrix} 1 & 2 \\ 3 & 4 \end{bmatrix} \rightarrow
\begin{array}{cc|cc}
1 & 0 & 2 & 0 \\
0 & 0 & 0 & 0 \\
\hline
3 & 0 & 4 & 0 \\
0 & 0 & 0 & 0
\end{array}
$$

**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

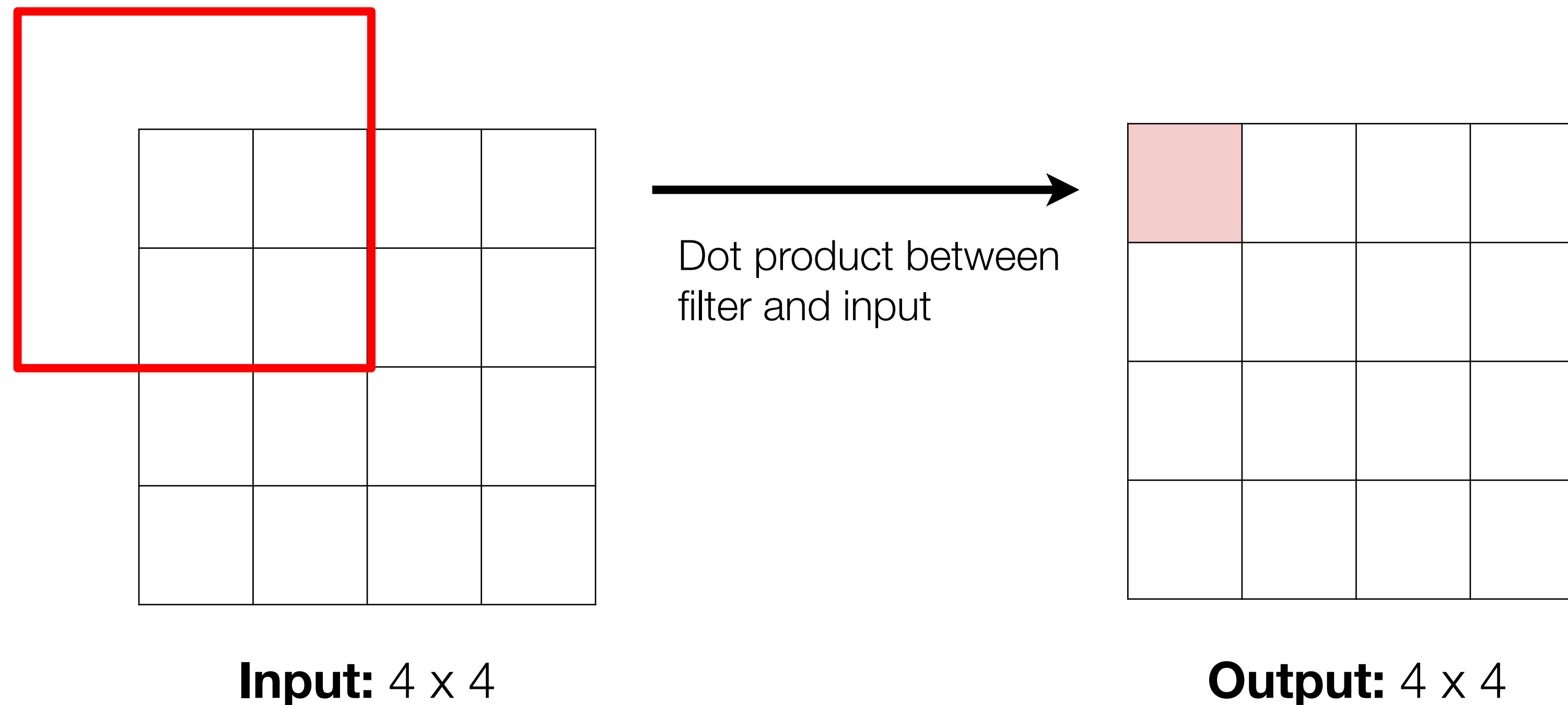**Output:** 2 x 2

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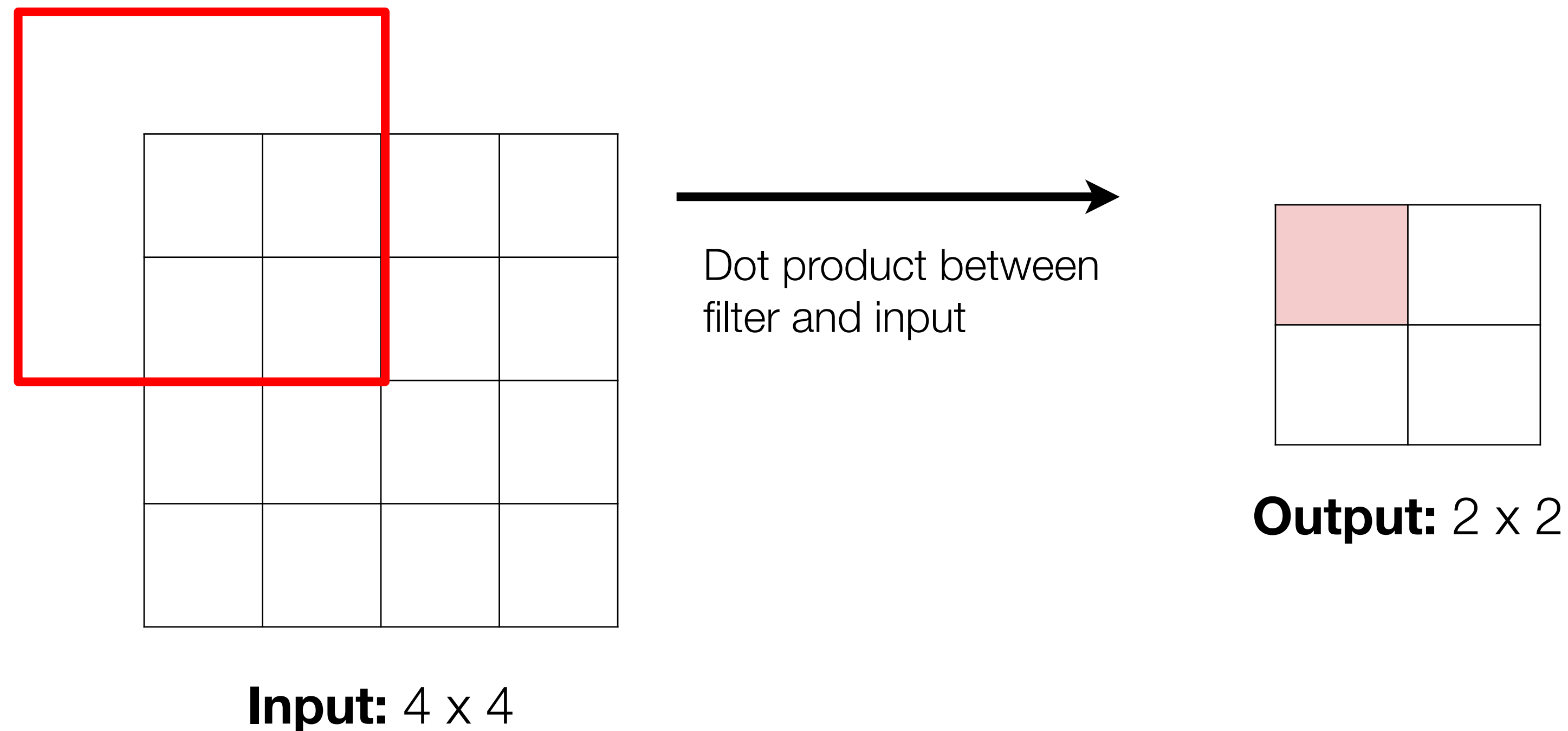
**Recall:** Normal 3 x 3 convolution, stride 1 pad 1



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

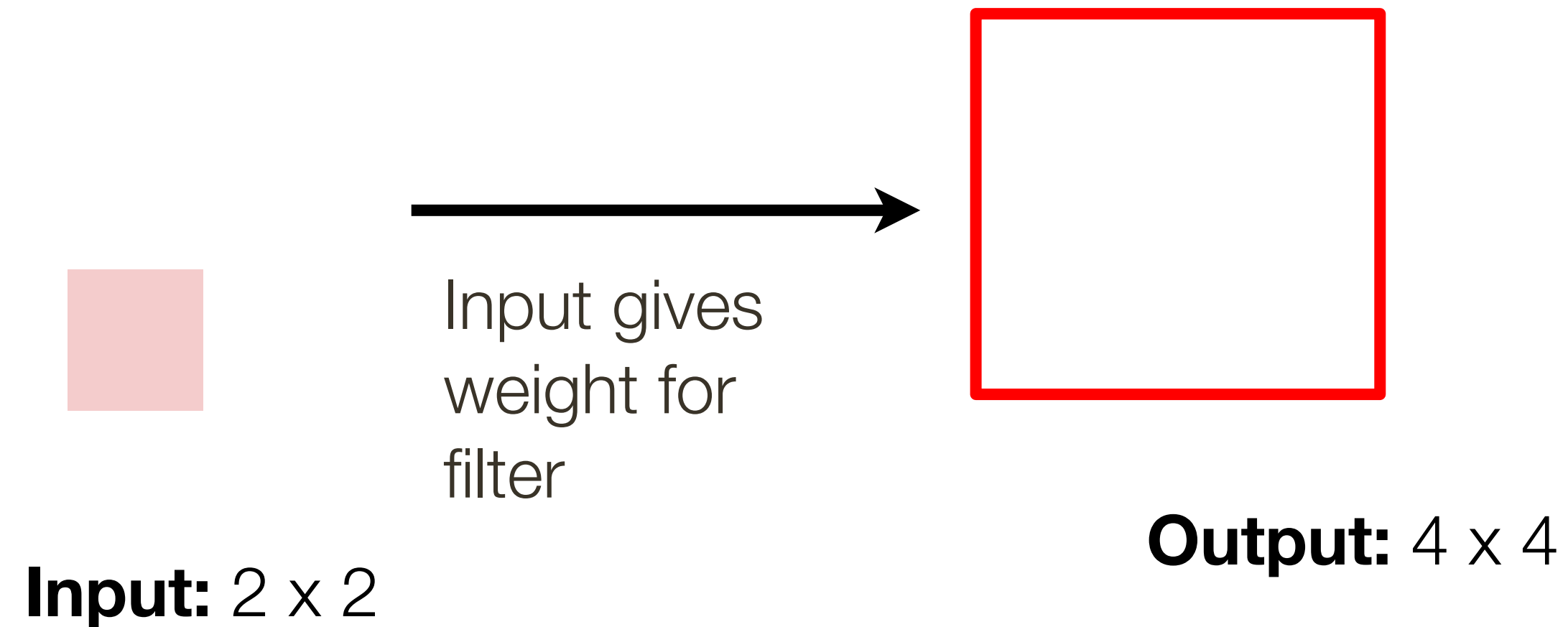3 x 3 **transpose** convolution, stride 2 pad 1
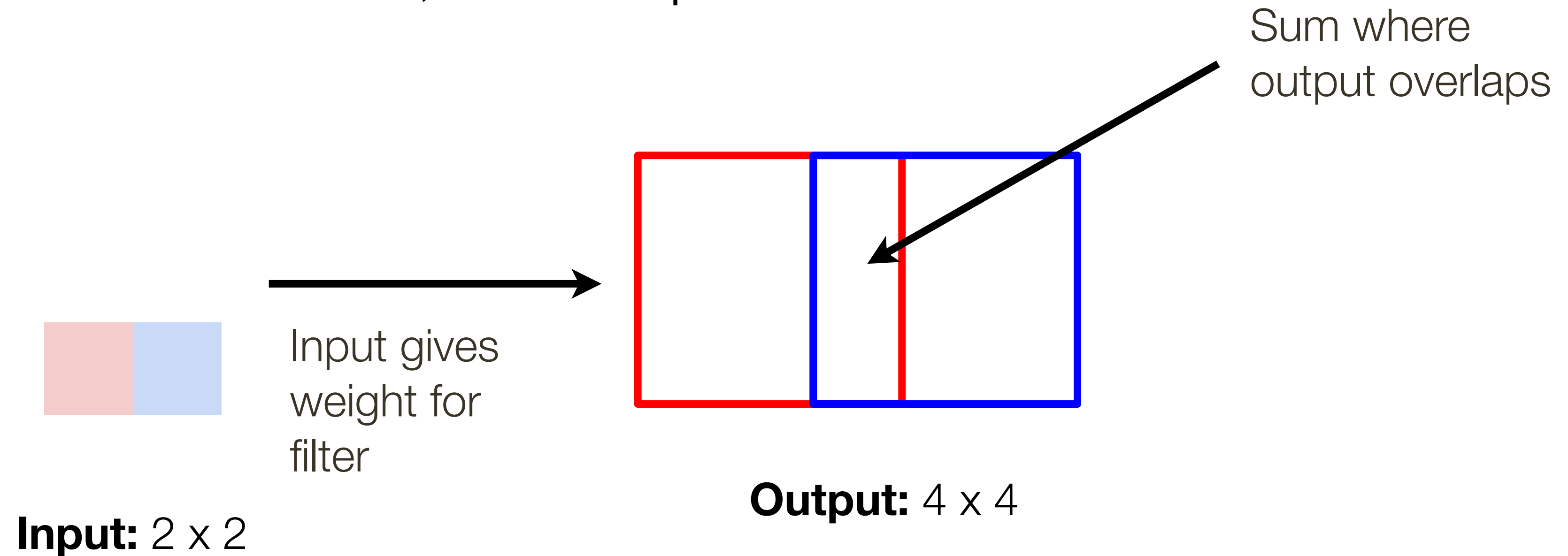
**Input:** 2 x 2

**Output:** 4 x 4

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

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



Input gives weight for filter

**Input:** 2 x 2

**Output:** 4 x 4

# 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**
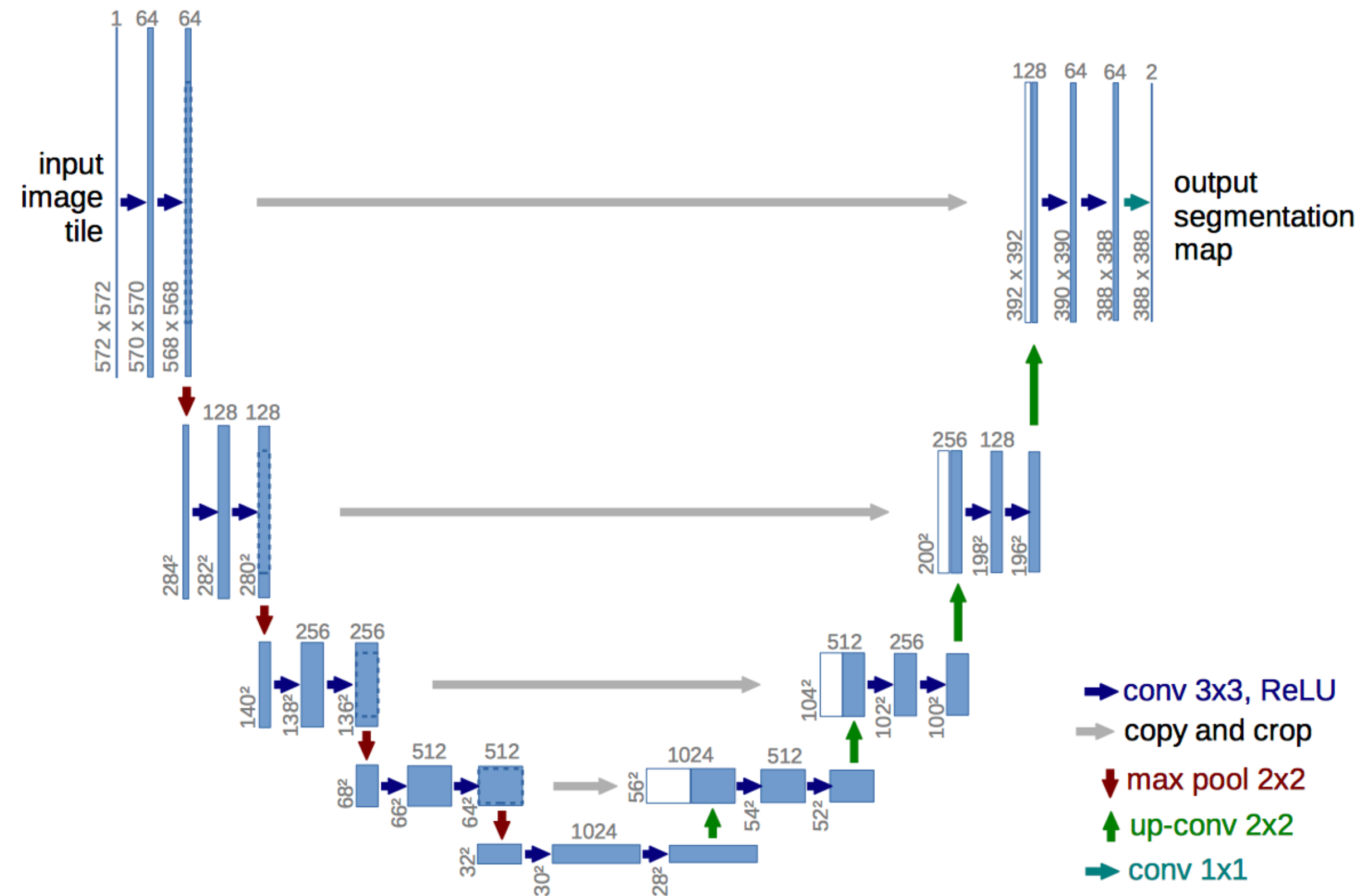
Stride gives ratio in movement in output vs input

# **Transpose Convolution**: 1-D Example

**Input**

**Filter**

**Output**



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



[ Ronneberger et al, CVPR 2015 ]

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

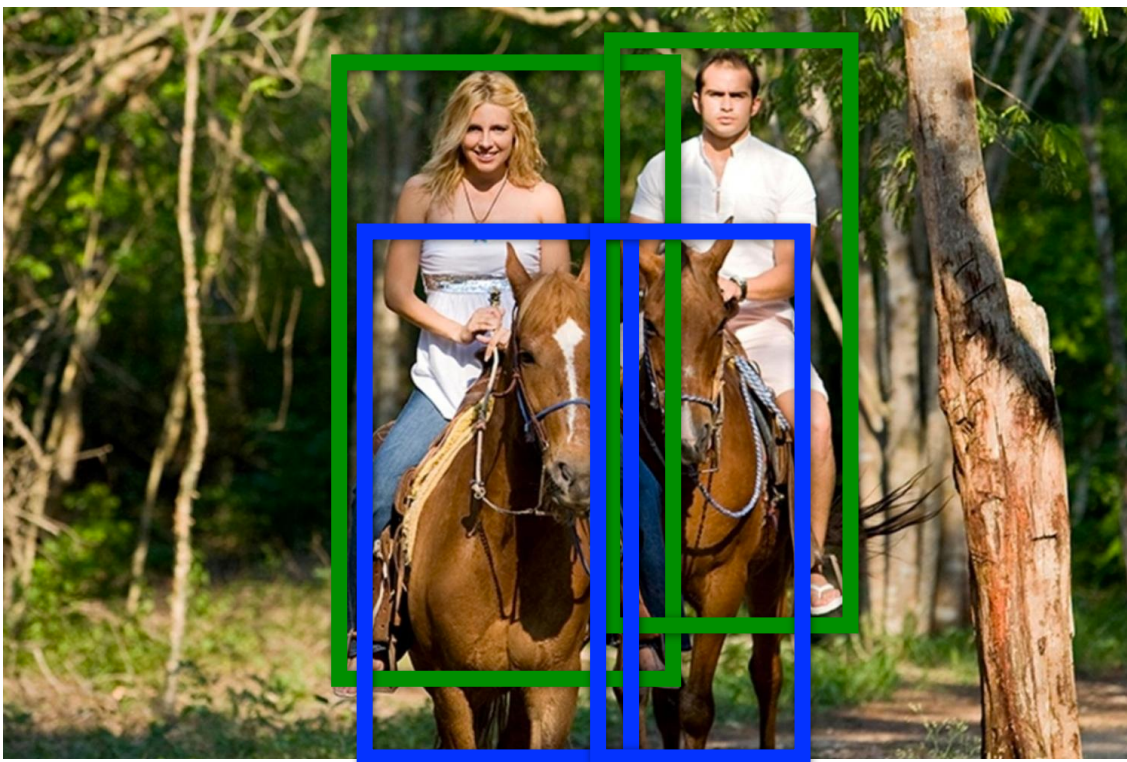| Categorization | Detection | Segmentation | Instance Segmentation |
|---|---|---|---|



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

IMAGENET

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

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
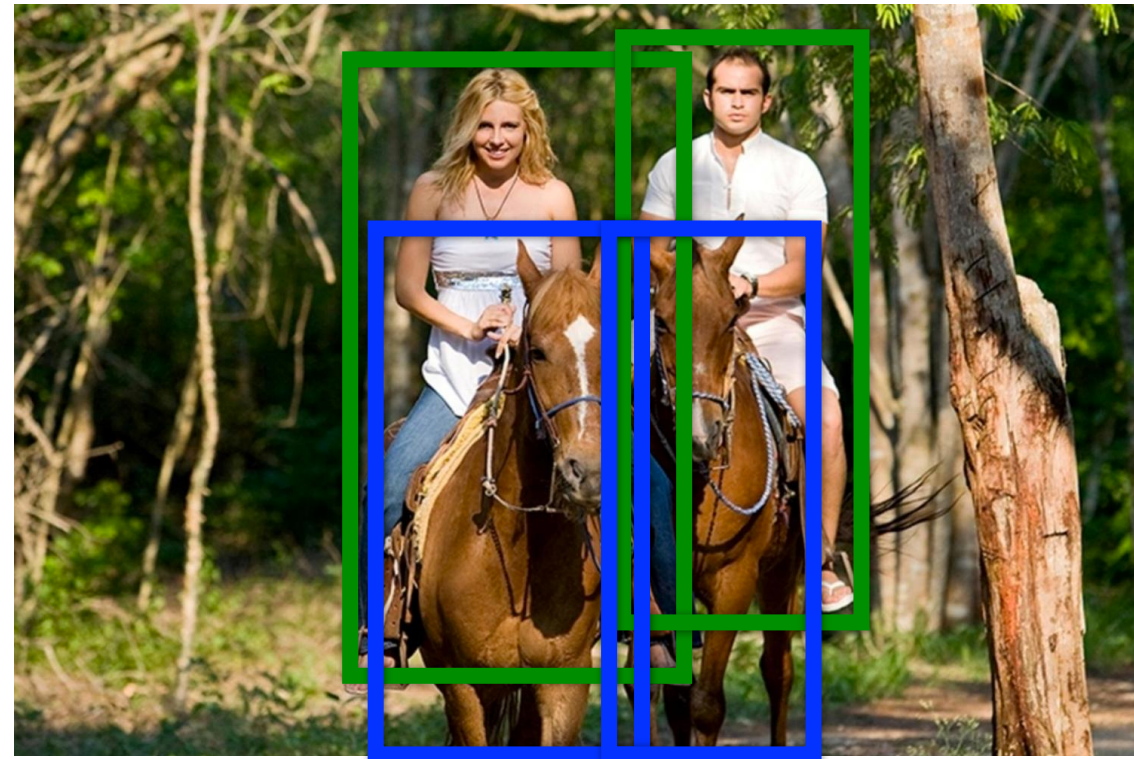
Horse
Person

COCO
Common Objects in Context

Horse1
Horse2
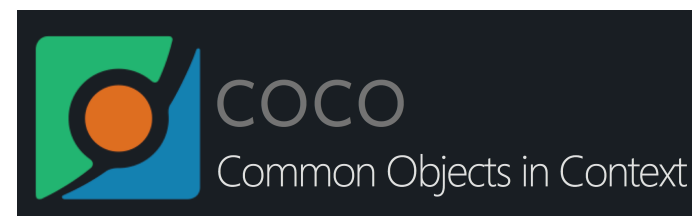Person1
Person2

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

Detection

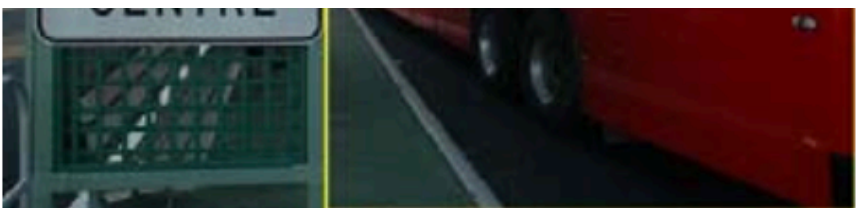

Horse (x, y, w, h)
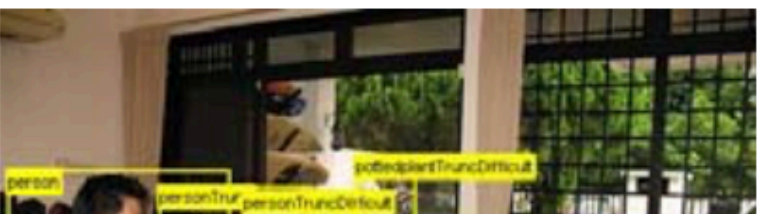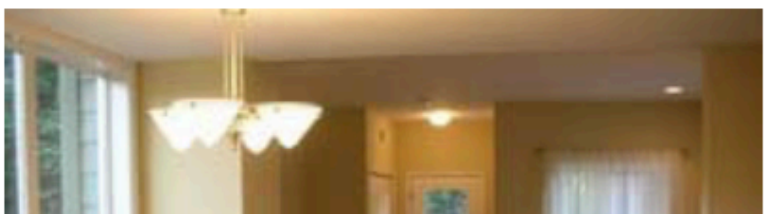Horse (x, y, w, h)
Person (x, y, w, h)
Person (x, y, w, h)

# Datasets: Pascal VOC

20 classes: aeroplane, bicycle, boat, bottle, bus, car, cat, chair, cow, dining table, dog, horse, motorbike, person, potted plant, sheep, train, TV



|  | Training | Testing |
|---|---|---|
| **Images** | 10,103 | 9,637 |
| **Objects** | 23,374 | 22,992 |

Real images downloaded from flickr, not filtered for "quality"
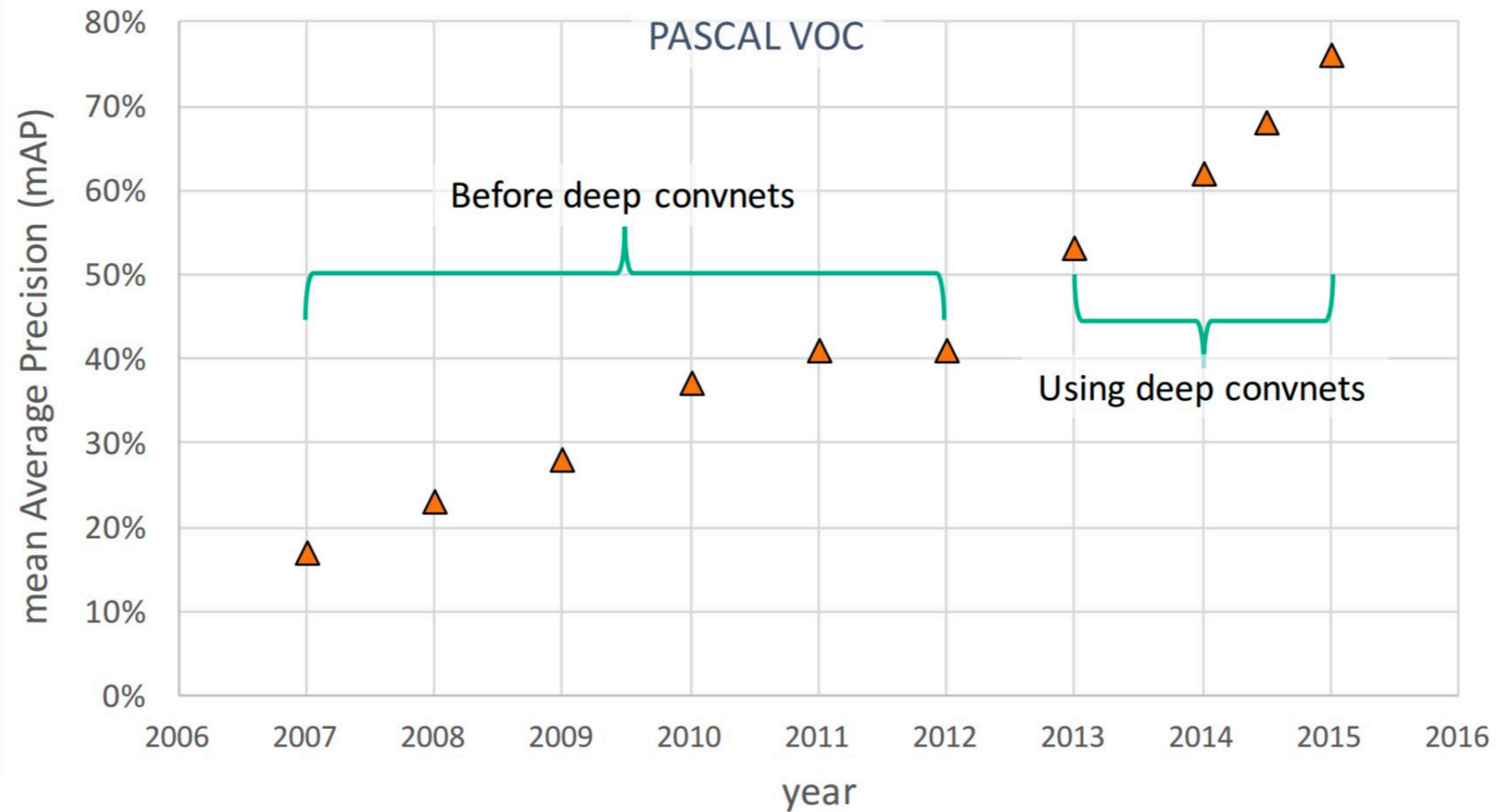
# **Datasets:** COCO



- ✔ Object segmentation
- ✔ Recognition in context
- ✔ Superpixel stuff segmentation
- ✔ 330K images (>200K labeled)
- ✔ 1.5 million object instances
- ✔ 80 object categories
- ✔ 91 stuff categories
- ✔ 5 captions per image
- ✔ 250,000 people with keypoints

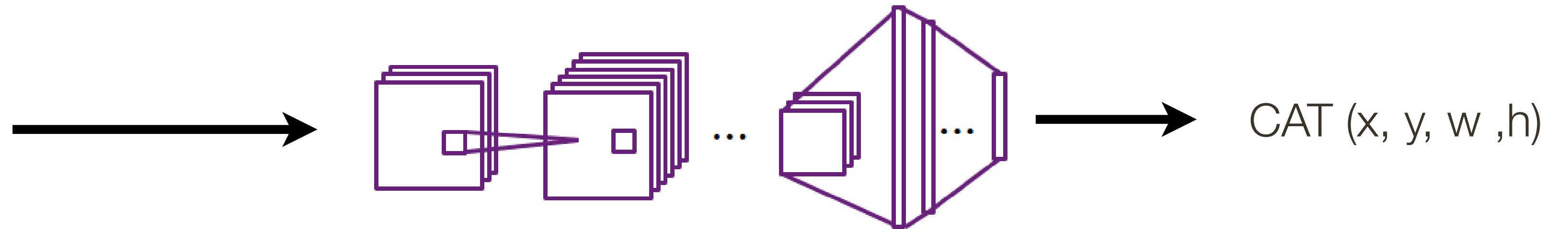# Object **Detection**



* plot from Ross Girshick, 2015

# Object **Detection** as Regression Problem

CAT (x, y, w ,h)

**Problem:** each image needs a different number of outputs

DUCK (x, y, w ,h)
DUCK (x, y, w ,h)
DUCK (x, y, w ,h)
DUCK (x, y, w ,h)
DUCK (x, y, w ,h)
DUCK (x, y, w ,h)
DUCK (x, y, w ,h)
DUCK (x, y, w ,h)
…

# Object **Detection** as Classification Problem

| Category | Prediction |
|---|---|
| Dog | No |
| Cat | No |
| Couch | No |
| Flowers | No |
| Background | **Yes** |
| … | … |

Apply CNN to many different crops in the image and (classification) CNN classifies each patch as object or background

# Object **Detection** as Classification Problem

| Category | Prediction |
|---|---|
| Dog | No |
| Cat | No |
| Couch | No |
| Flowers | No |
| Background | **Yes** |
| … | … |

Apply CNN to many different crops in the image and (classification) CNN classifies each patch as object or background

# Object **Detection** as Classification Problem



| Category | Prediction |
|---|---|
| Dog | **Yes** |
| Cat | No |
| Couch | No |
| Flowers | No |
| Background | No |
| … | … |

Apply CNN to many different crops in the image and (classification) CNN classifies each patch as object or background

# Object **Detection** as Classification Problem



| Category | Prediction |
|---|---|
| Dog | **Yes** |
| Cat | No |
| Couch | No |
| Flowers | No |
| Background | No |
| … | … |

Apply CNN to many different crops in the image and (classification) CNN classifies each patch as object or background

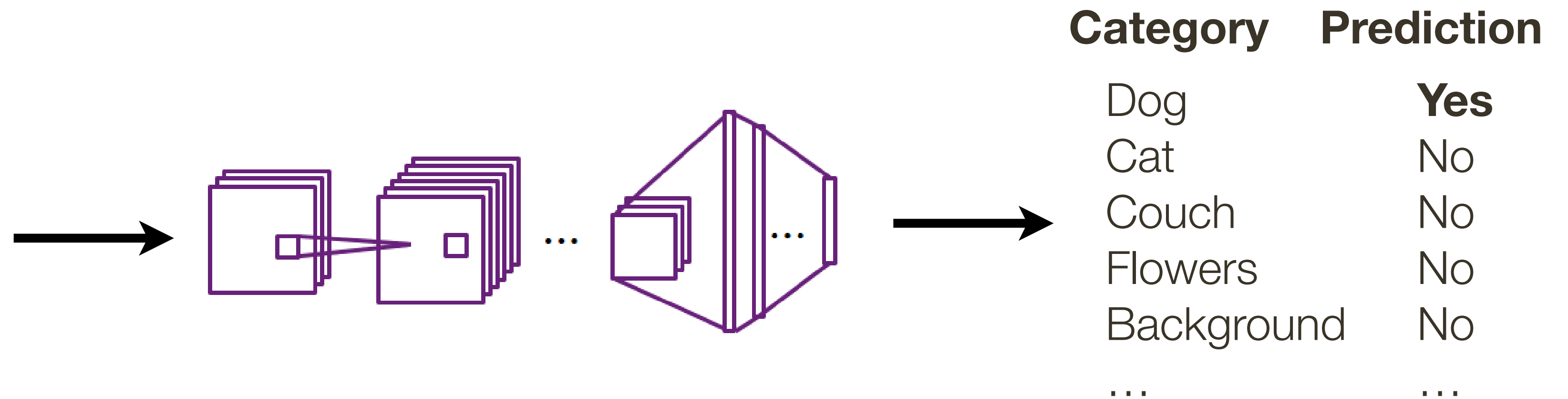# Object **Detection** as Classification Problem

**Problem:** Need to apply CNN to **many** patches in each image



| Category | Prediction |
|---|---|
| Dog | No |
| Cat | **Yes** |
| Couch | No |
| Flowers | No |
| Background | No |
| … | … |

Apply CNN to many different crops in the image and (classification) CNN classifies each patch as object or background

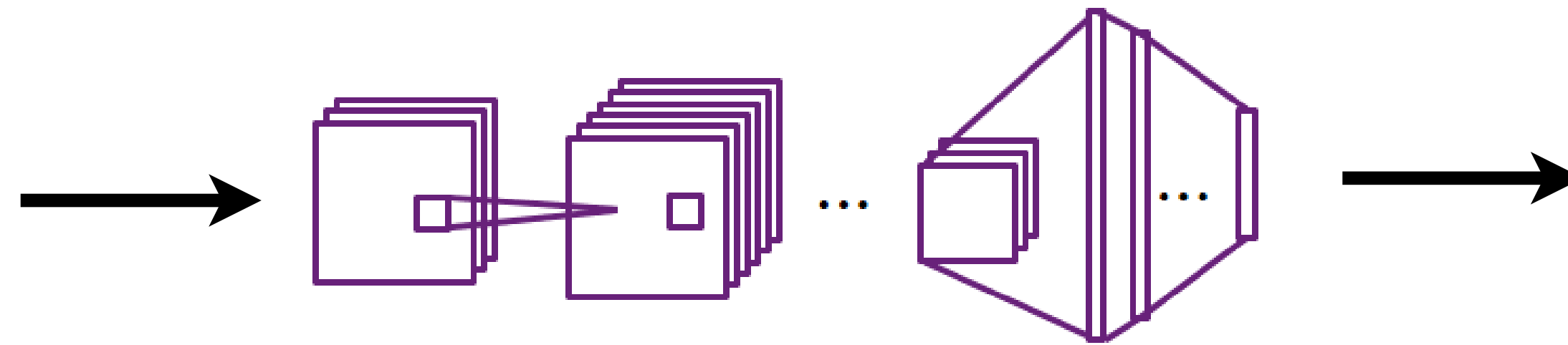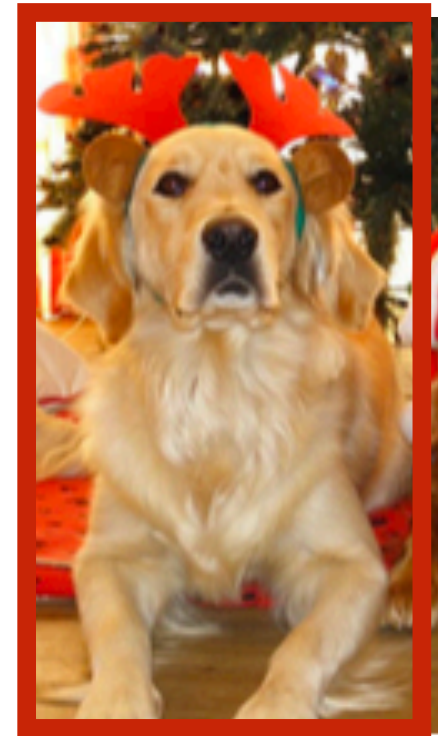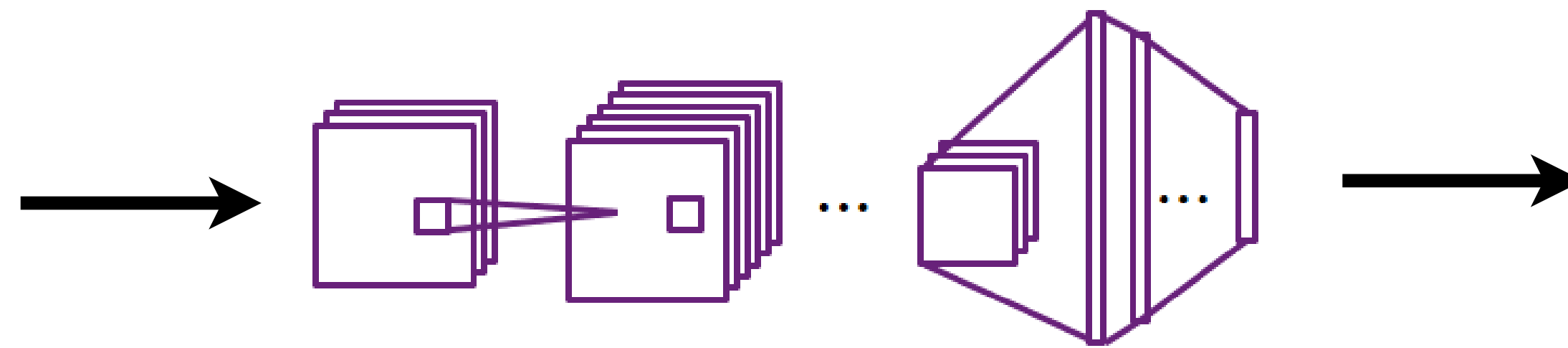# **Region** Proposals (older idea in vision)

[ Alexe et al, TPAMI 2012 ]
[ Uijkings et al, IJCV 2013 ]
[ Cheng et al, CVPR 2014 ]
[ Zitnick and Dollar, ECCV 2014 ]

Find image **regions that are likely contain objects** (any object at all)

- typically works by looking at histogram distributions, region aspect ratio, closed contours, coherent color

Relatively **fast to run** (Selective Search gives 1000 region proposals in a few seconds on a CPU)



**Goal:** Get "true" object regions to be in as few top K proposals as possible

# R-CNN

Input **Image**

* image from Ross Girshick

# R-CNN

**Regions of Interest** from
a proposal method (~2k)

Input **Image**

* image from Ross Girshick

# R-CNN

**Warped** image regions

**Regions of Interest** from a proposal method (~2k)

Input **Image**

* image from Ross Girshick

# R-CNN

Forward each region through a **CNN**

**Warped** image regions

**Regions of Interest** from a proposal method (~2k)

Input **Image**

* image from Ross Girshick

# R-CNN

**Classify** regions with SVM

Forward each region through a **CNN**

**Warped** image regions

**Regions of Interest** from a proposal method (~2k)

Input **Image**

* image from Ross Girshick

# R-CNN

**Linear Regression** for bounding box offsets

**Classify** regions with SVM

Forward each region through a **CNN**

**Warped** image regions

**Regions of Interest** from a proposal method (~2k)

Input **Image**

* image from Ross Girshick

# R-CNN: Training

**Fine-tuning ImageNet** CNN on object proposal patches

— > 50% Intersection-over-Union overlap with GT considered "object" others "background"

— batches of 128 (**32 positives, 96 negatives**)



* image from Ross Girshick

# **R-CNN:** Issues

## **Ad-hoc** training objectives

— Fine-tune network with softmax objective (**log** loss)

— Train post-hoc linear SVM (**hinge** loss)

— Train post-hoc bounding-box regression (**least squares**)

## **Training** is slow

— 84 hours and takes a lot of disk space

## Inference / **Detection is slow**

— 47 sec / image with  VGG16 [ Simonyan et al, ICLR 2015 ]

* image from Ross Girshick

# R-CNN vs. SPP

[ He et al, ECCV 2014 ]



**R-CNN**
2000 nets on image regions

**SPP-net**
**1 net on full image**

# Fast **R-CNN**

Input **Image**

* image from Ross Girshick

# Fast **R-CNN**

"conv5" feat

Forward wl

ConvNet

Input **Image**

# Fast **R-CNN**

"conv5" feature map

Forward prop the **whole image** through CNN

ConvNet

Input **Image**

* image from Ross Girshick

# Fast **R-CNN**

**Regions of Interest** from the proposal method

"conv5" feature map

Forward prop the **whole image** through CNN

ConvNet

Input **Image**

# Fast **R-CNN**

**Regions of Interest** from the proposal method

"RoI Pooling" layer

"conv5" feature map

Forward prop the **whole image** through CNN

ConvNet

Input **Image**

Girshick, "Fast R-C
Figure copyright R

# Fast **R-CNN**

Log loss + Smooth L1 loss — Multi-task loss

Object classification

Linear + softmax

Linear — Bounding box regression

FCs

**Regions of Interest** from the proposal method

"RoI Pooling" layer

"conv5" feature map

ConvNet — Forward prop the **whole image** through CNN

Input **Image**

* image from Ross Girshick

# Fast **R-CNN:** Training



Log loss + Smooth L1 loss — Multi-task loss

[ Girshick et al, ICCV 2015 ]

Object classification

Linear + softmax

Linear — Bounding box regression

FCs

**Regions of Interest** from the proposal method

"RoI Pooling" layer

"conv5" feature map

ConvNet

Forward prop the **whole image** through CNN

Input **Image**

* image from Ross Girshick

# R-CNN vs. SPP vs. Fast R-CNN

[ Girshick et al, CVPR 2014 ]
[ Girshick et al, ICCV 2015 ]
[ He et al, ECCV 2014 ]



**Training time (Hours)**

- R-CNN: 84
- SPP-Net: 25.5
- Fast R-CNN: 8.75

**Test time (seconds)**

Including Region propos... / Excluding Region Propo...

- R-CNN: 49 / 47
- SPP-Net: 4.3 / 2.3
- Fast R-CNN: 2.3 / 0.32

**Observation:** Performance dominated by the region proposals at this point!
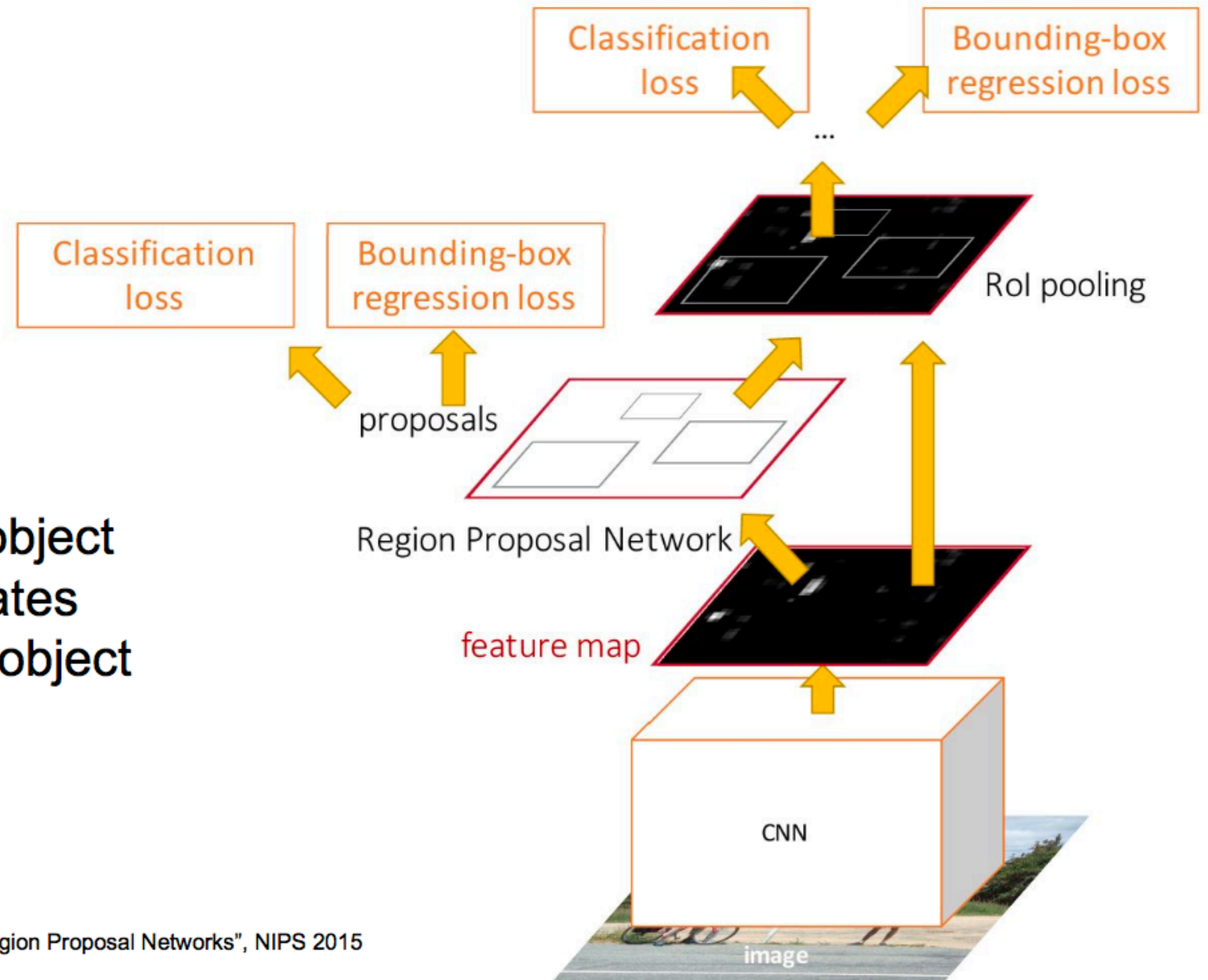
# Faster **R-CNN**

Make CNN do proposals!

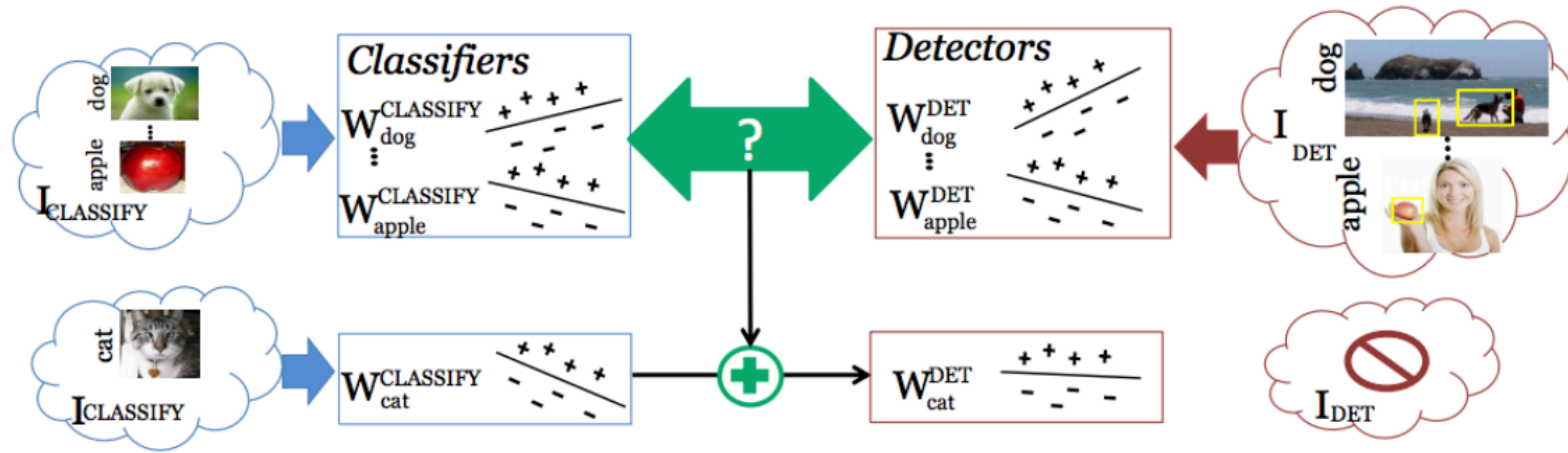Insert **Region Proposal Network (RPN)** to predict proposals from features

Jointly train with 4 losses:
1. RPN classify object / not object
2. RPN regress box coordinates
3. Final classification score (object classes)
4. Final box coordinates



Classification loss

Bounding-box regression loss

RoI pooling

Classification loss

Bounding-box regression loss

proposals

Region Proposal Network

feature map

CNN

image

Ren et al, "Faster R-CNN: Towards Real-Time Object Detection with Region Proposal Networks", NIPS 2015
Figure copyright 2015, Ross Girshick; reproduced with permission

# LSDA: Large Scale Detection through Adaptation



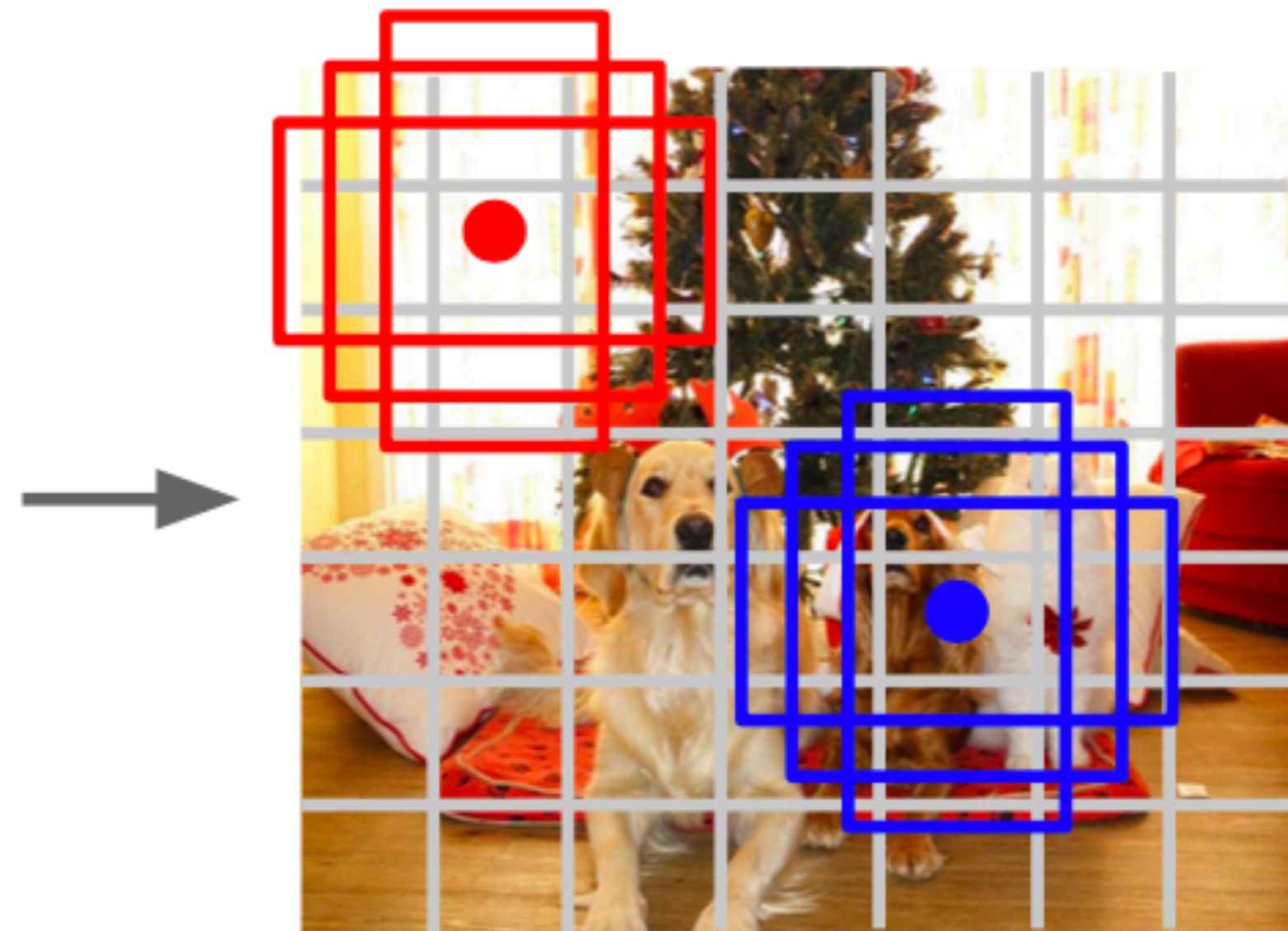$$\mathbf{W}_{cat}^{DETECT} = \mathbf{W}_{cat}^{CLASIFY} + \delta\mathbf{W}_{cat}$$

[ Hoffman et al, NIPS 2014 ]

# YOLO: You Only Look Once

Input image
3 x H x W

Divide image into grid
7 x 7

Image a set of **base boxes**
centered at each grid cell
Here B = 3

Within each grid cell:
- Regress from each of the B base boxes to a final box with 5 numbers:
  (dx, dy, dh, dw, confidence)
- Predict scores for each of C classes (including background as a class)
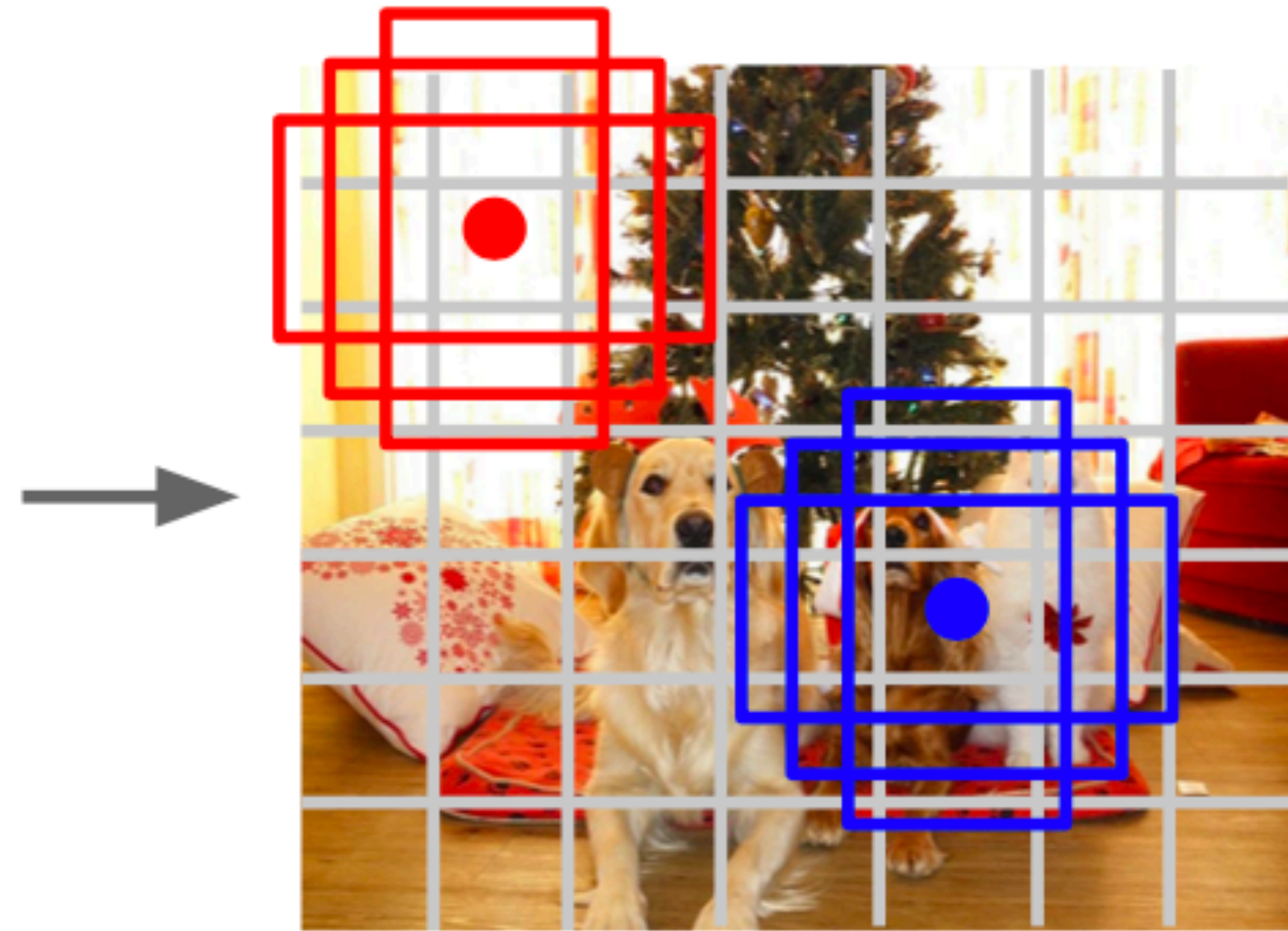
Output:
7 x 7 x (5 * B + C)

# YOLO: You Only Look Once

[ Redmon et al, CVPR 2016 ]



Input image
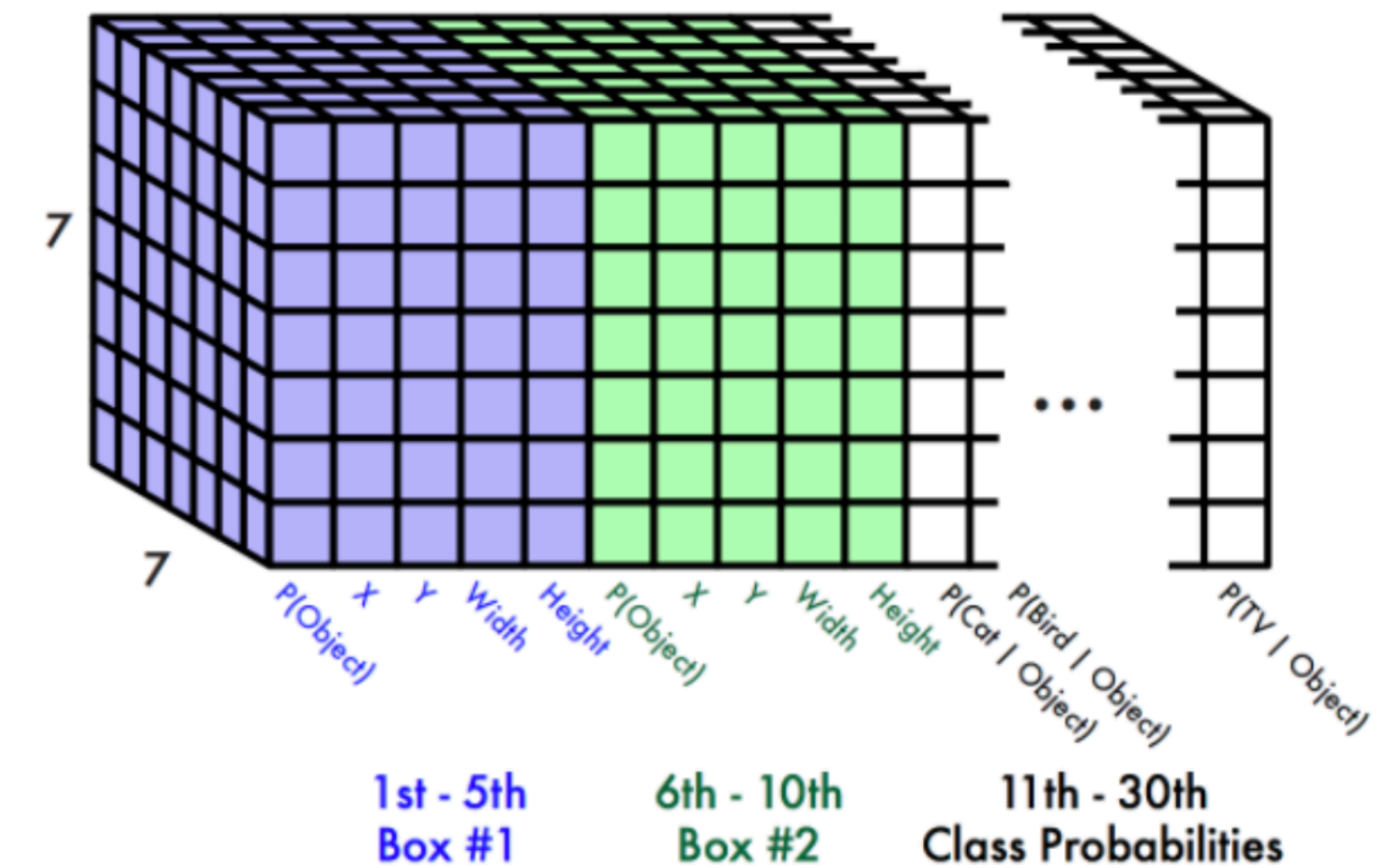3 x H x W

Divide image into grid
7 x 7

Image a set of **base boxes**
centered at each grid cell
Here B = 3

1st - 5th
Box #1
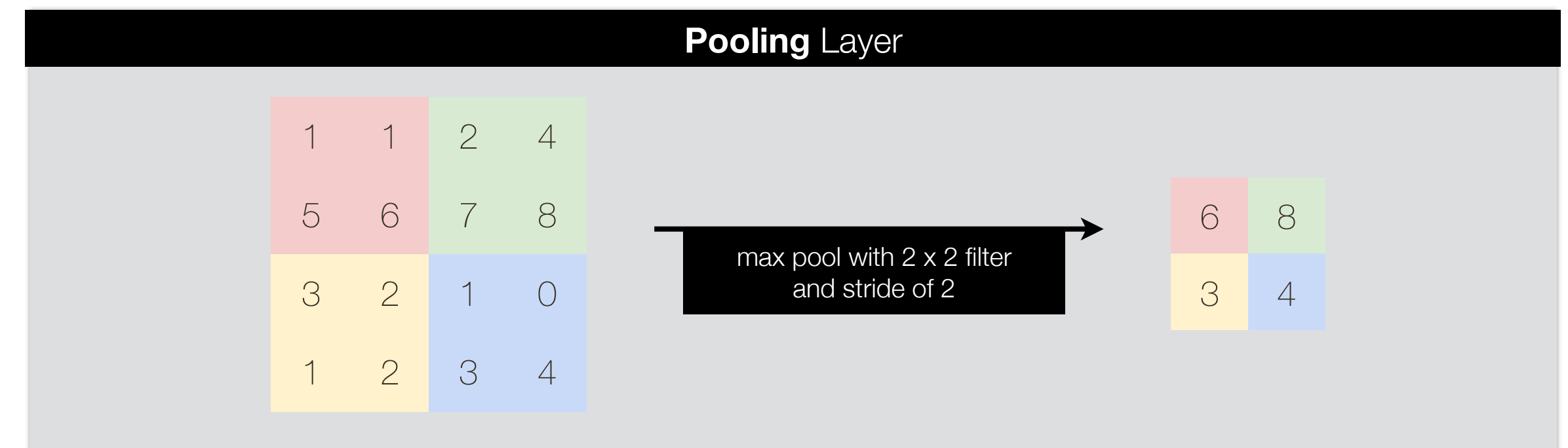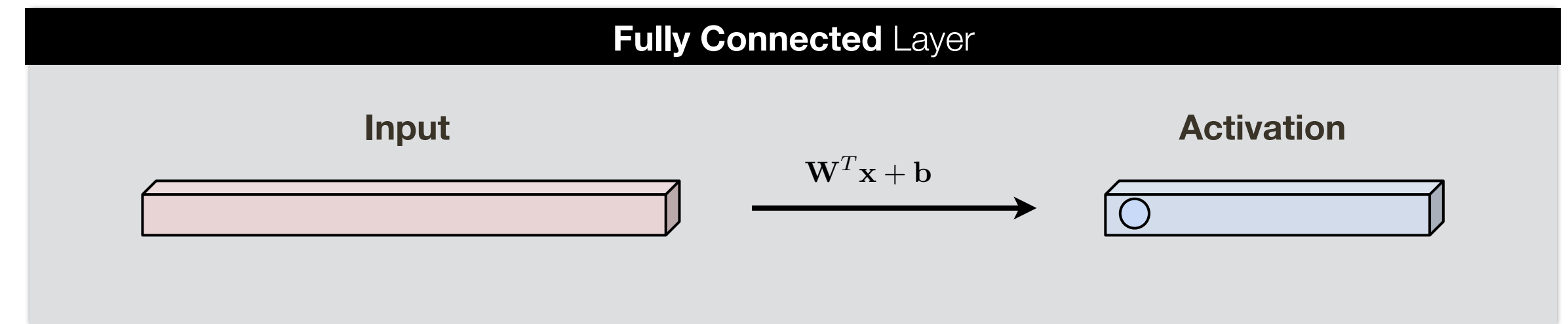
6th - 10th
Box #2

11th - 30th
Class Probabilities

# YOLO v2

http://pjreddie.com/yolo

# Review of **CNNs**



**Convolutional** Layer

activation map

**Fully Connected** Layer

Input    $\mathbf{W}^T\mathbf{x} + \mathbf{b}$    Activation

**Pooling** Layer

max pool with 2 x 2 filter and stride of 2
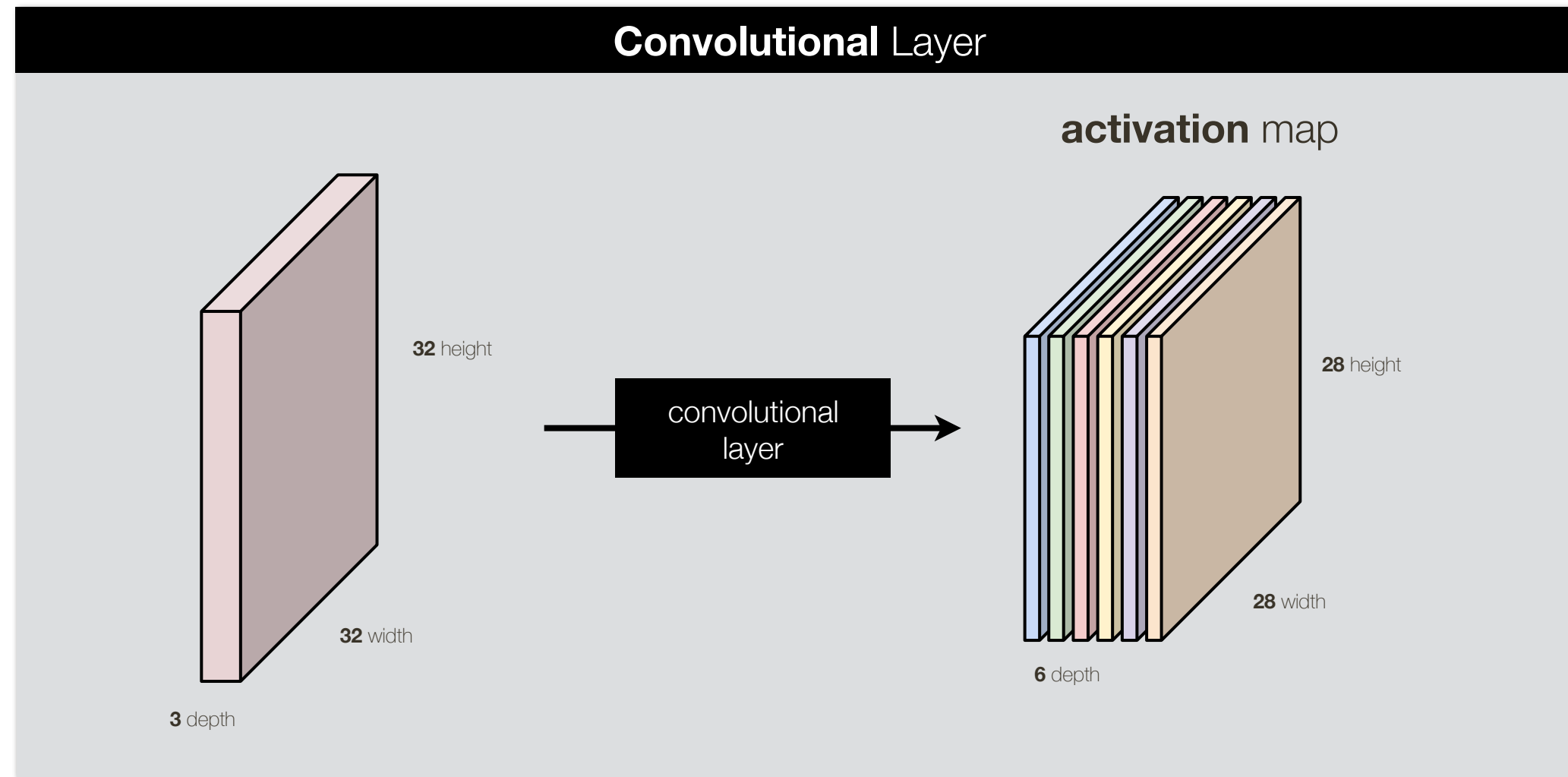
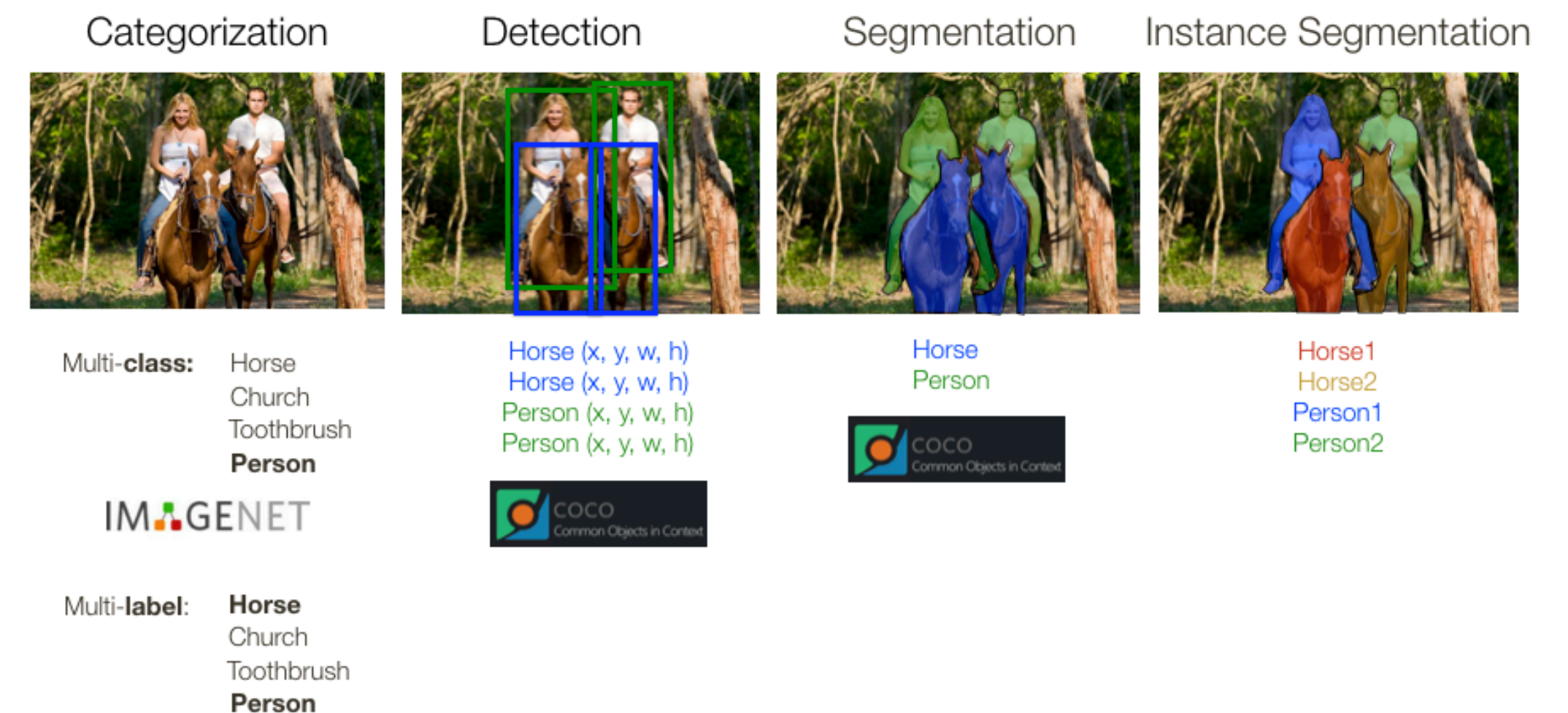## Effective Techniques for **Training**

— **Regularization:** L1, L2, data augmentation

— **Transfer Learning:** fine-tuning networks

## Vision **Applications** of CNNs

— **Classification:** AlexNet, VGG, GoogleLeNet, ResNet

— **Segmentation:** Fully convolutional CNNs

— **Detection:** R-CNN, Fast R-CNN, Faster R-CNN, YOLO

# Any **CNN** Could be **Fully Convolutional**



Image

224 x 224

VGG

1 x 1000

# Any **CNN** Could be **Fully Convolutional**



Image
**225 x 225**

VGG

**2 x 2 x 1000**

# Review of **CNNs**

**Convolutional** Layer

**activation** map

32 height
32 width
3 depth

convolutional layer

28 height
28 width
6 depth

**Fully Connected** Layer

Input

$\mathbf{W}^T\mathbf{x} + \mathbf{b}$

Activation

**Pooling** Layer

| 1 | 1 | 2 | 4 |
|---|---|---|---|
| 5 | 6 | 7 | 8 |
| 3 | 2 | 1 | 0 |
| 1 | 2 | 3 | 4 |

max pool with 2 x 2 filter and stride of 2

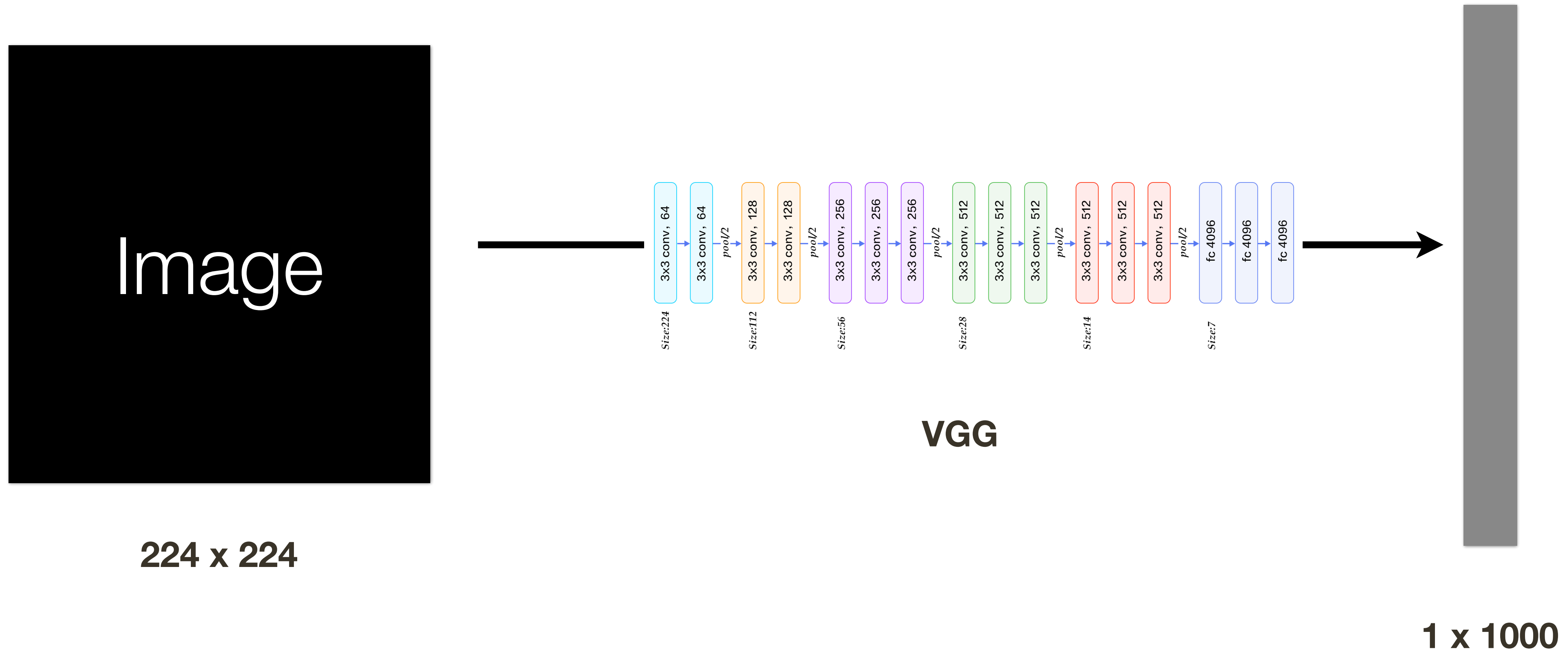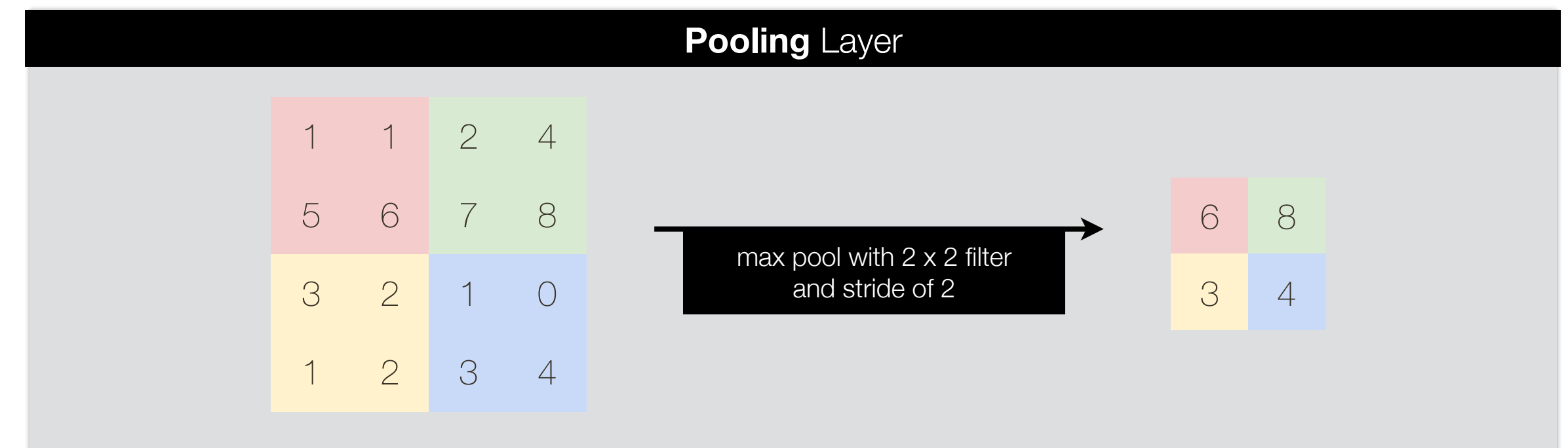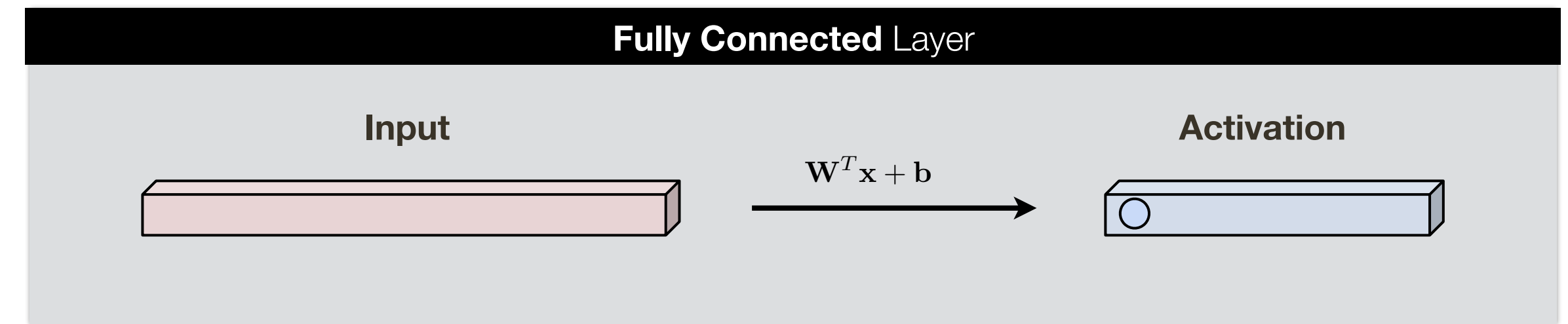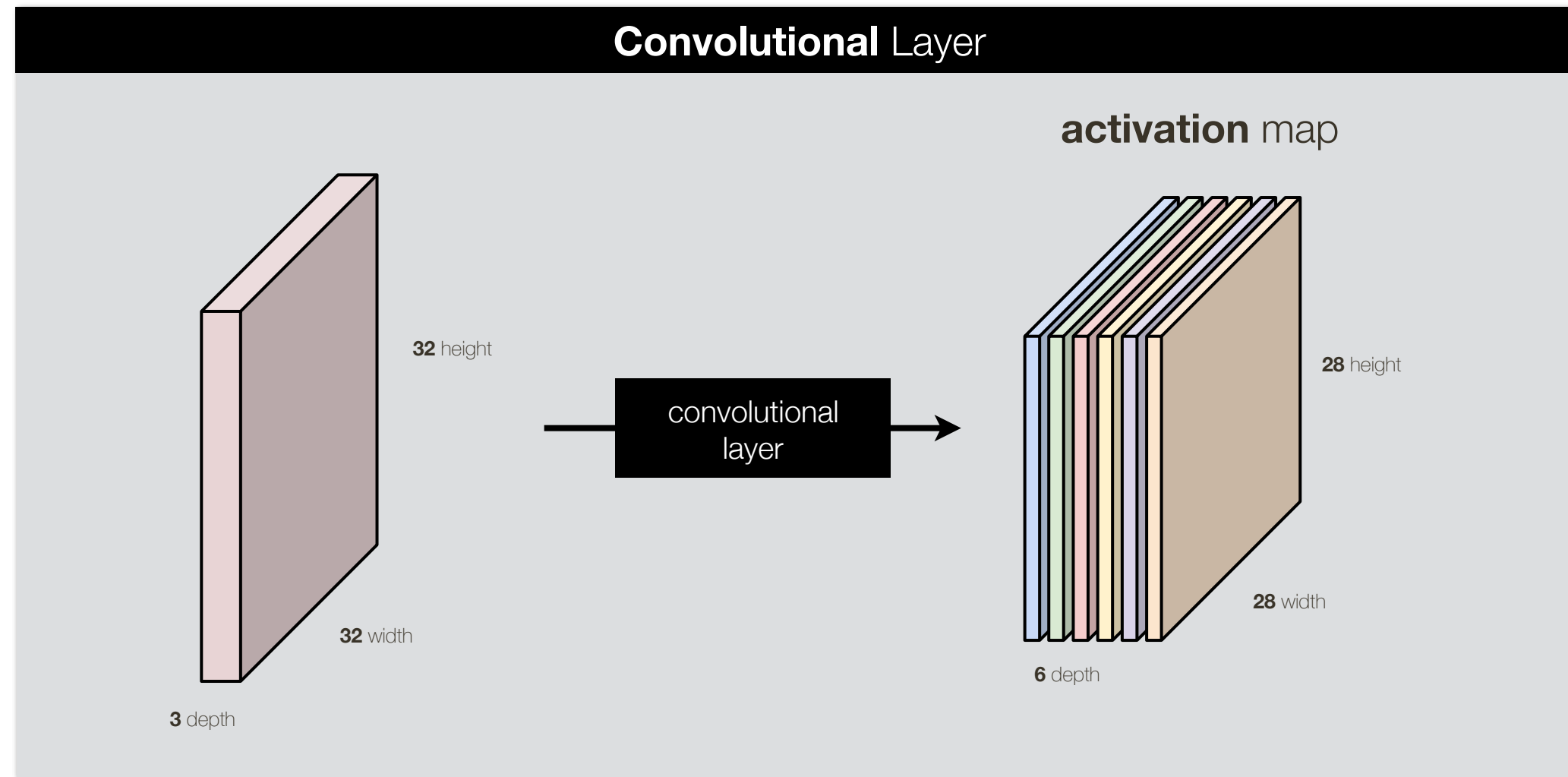| 6 | 8 |
|---|---|
| 3 | 4 |

## Effective Techniques for **Training**

— **Regularization:** L1, L2, data augmentation

— **Transfer Learning:** fine-tuning networks

## Vision **Applications** of CNNs

— **Classification:** AlexNet, VGG, GoogleLeNet, ResNet

— **Segmentation:** Fully convolutional CNNs

— **Detection:** R-CNN, Fast R-CNN, Faster R-CNN, YOLO

Categorization

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

IM·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