



# Computer Vision 1

HC6a

## Single Shot Detection, Network Architectures

Dr. Martin Oswald, Dr. Dimitris Tzionas, Dr. Arun Mukundan,  
[m.r.oswald, d.tzionas, a.mukundan]@uva.nl

# Today's Agenda

---

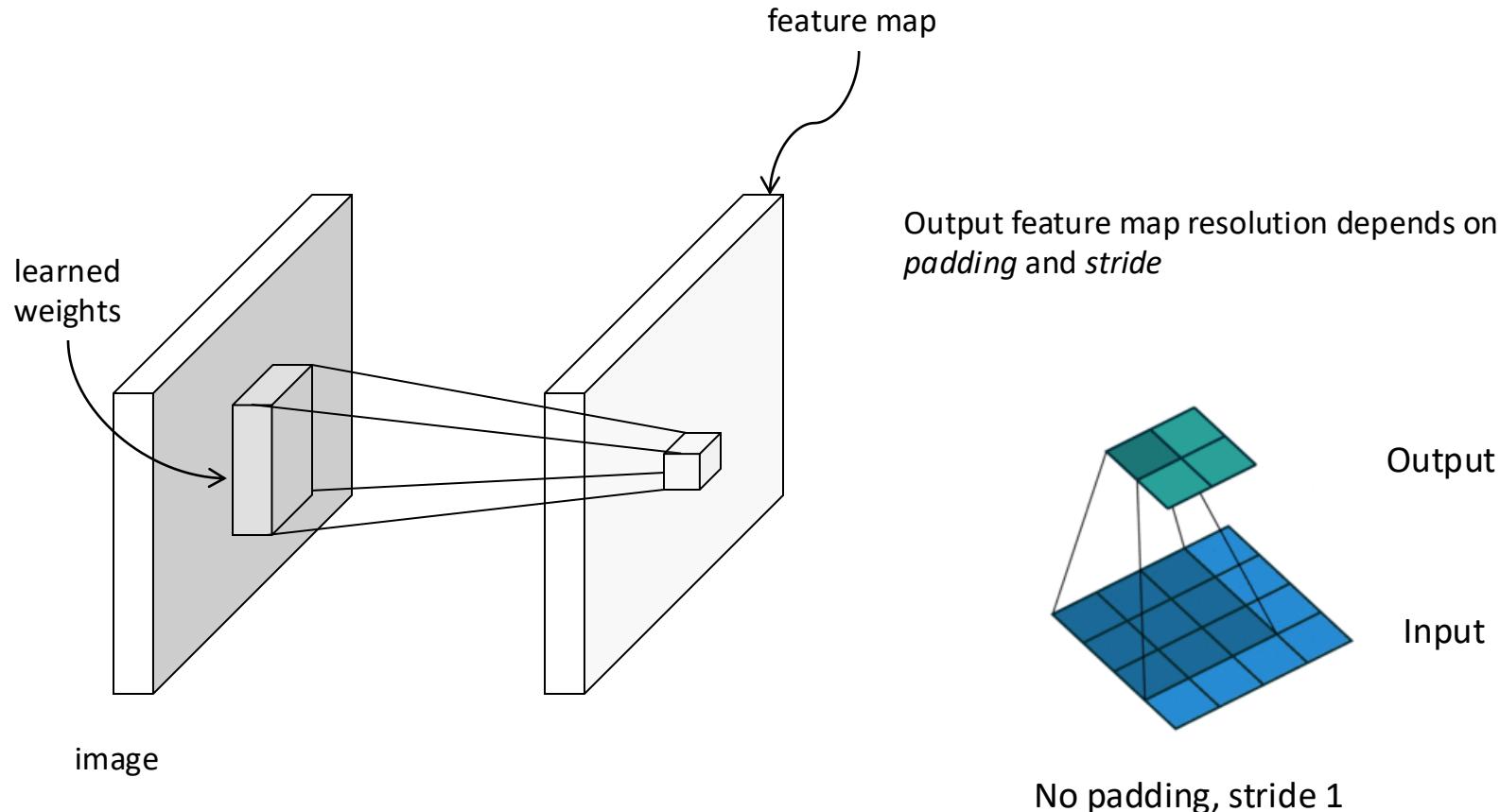
- Learning region proposals - Faster R-CNN
- Feature Pyramid Network (FPN)
- Single Shot Detectors (SSD) / Yolo
- Transfer learning
- Network architectures
  - VGG
  - ResNet
  - U-Net
  - Autoencoders

# Last Lecture

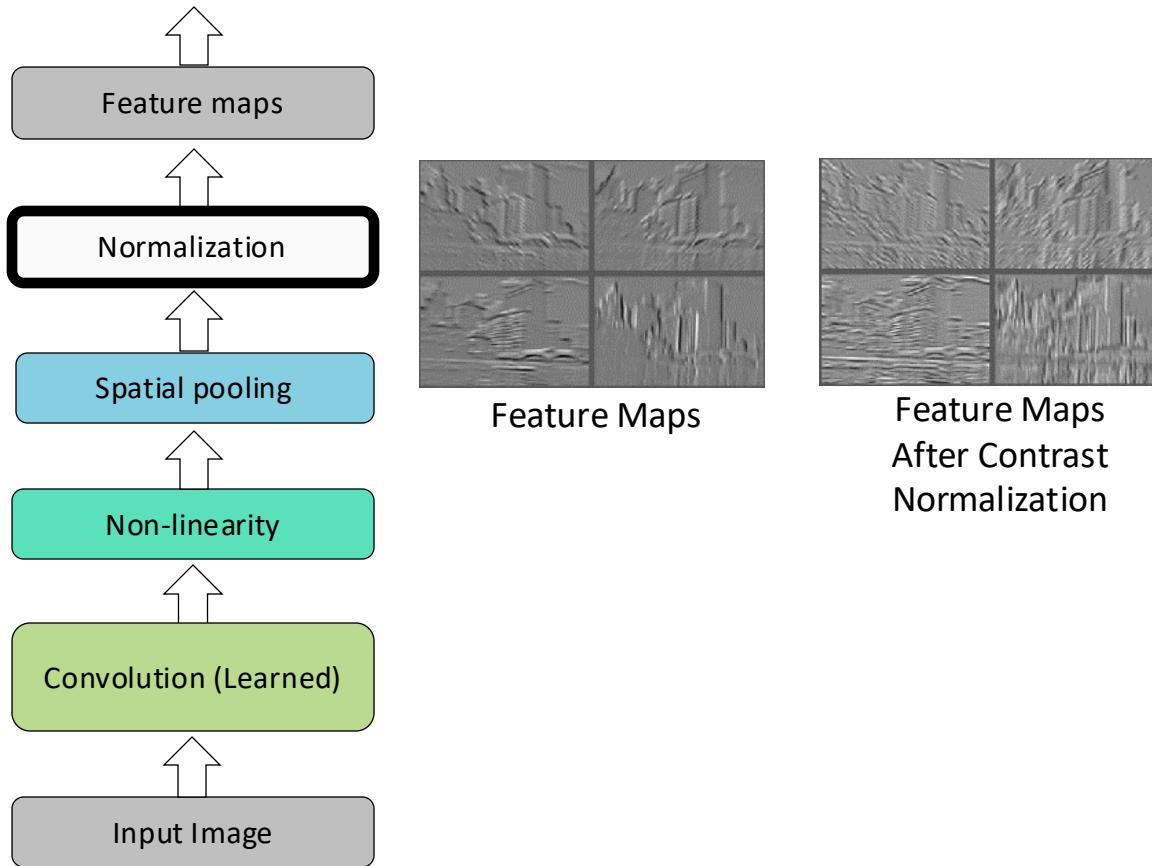
---

## Recap

# Convolutional Architecture



# Convolutional Neural Networks

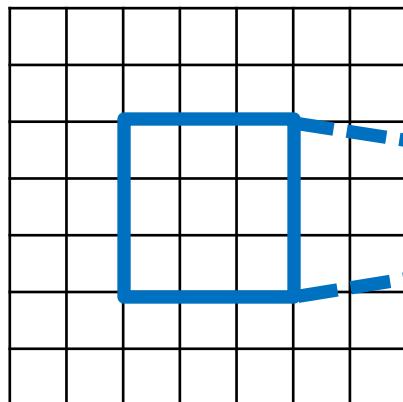


# Receptive Field

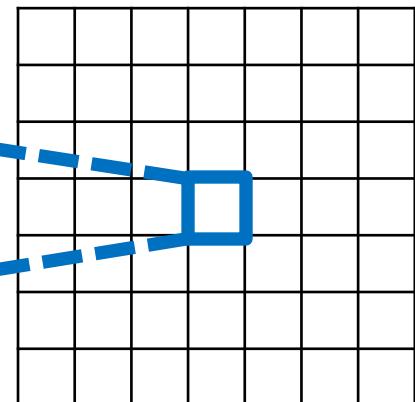
The *receptive field* of a unit is the region of the input feature map whose values contribute to the response of that unit (either in the previous layer or in the initial image)

3x3 convolutions, stride 1

Input



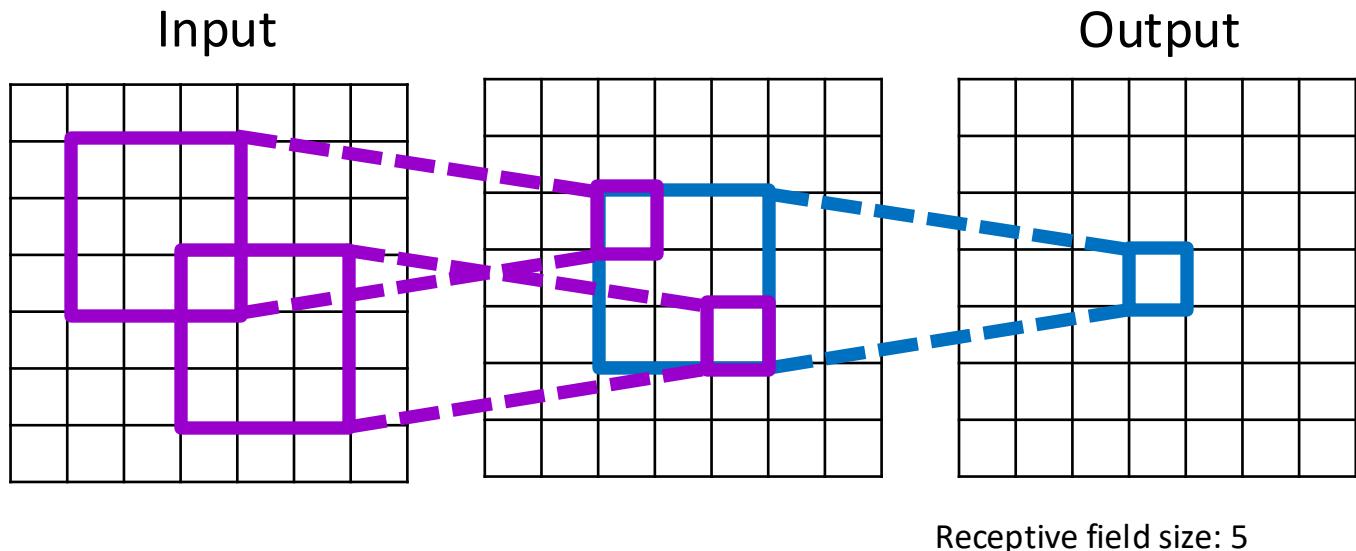
Output



Receptive field size: 3

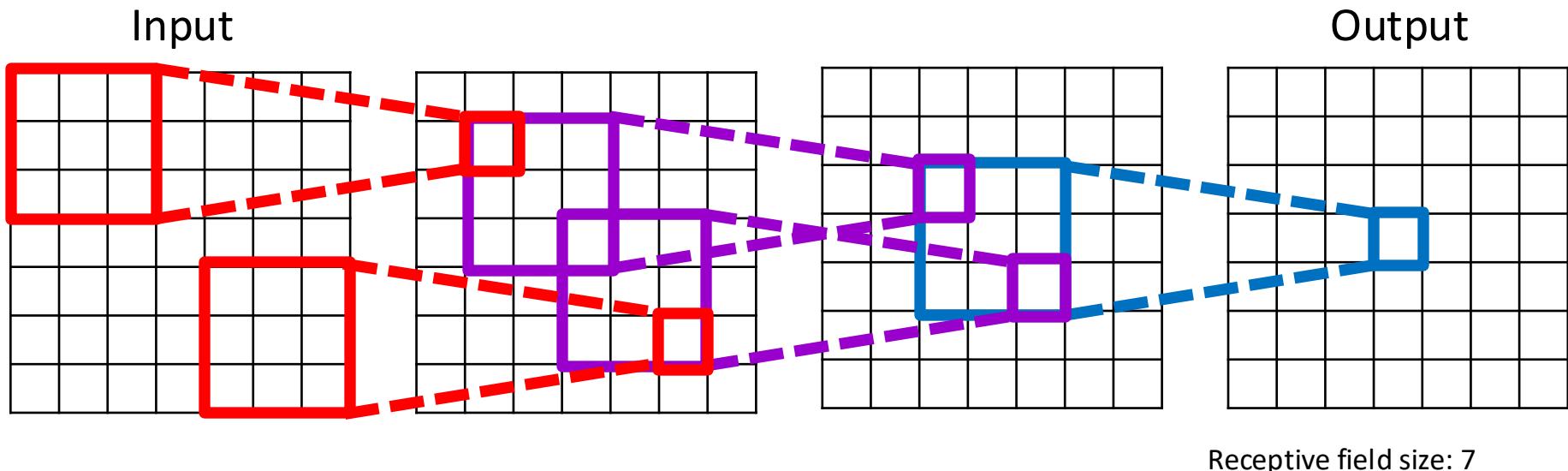
# Receptive Field

3x3 convolutions, stride 1



# Receptive Field

3x3 convolutions, stride 1

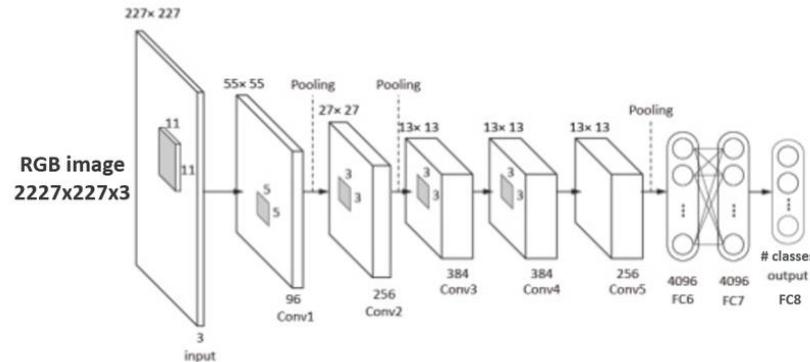


Each successive convolution adds  $F - 1$  to the receptive field size  
With  $L$  layers the receptive field size is  $1 + L * (F - 1)$

# ConvNet Architectures

## AlexNet (Alex Krizhevsky, 2012)

2012 ImageNet Large Scale Visual Recognition Challenge: 15% top 5 error



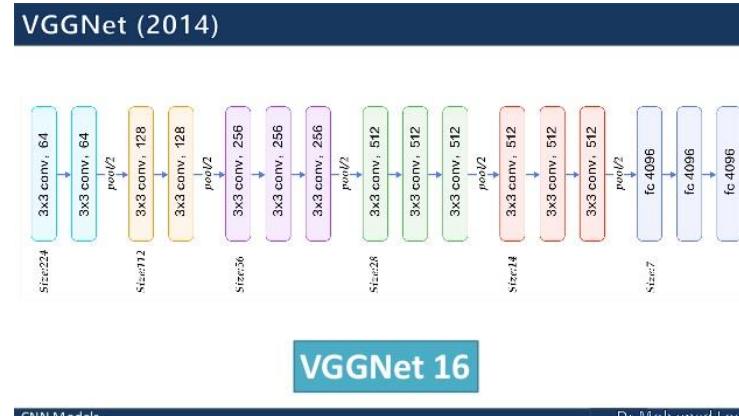
### ConvNet Architecture:

- 5 convolutional layers and 3 fully connected layers (1000 classes)
- reduces spacial dimensions: 227x227 (image) → ... → 13x13
- increases depth: 3 (image RGB depth) → ... → 256 (nr of filters)
- use of Graphical Processing Unit (GPU)

# ConvNet Architectures

## VGGNet (Zimmerman, Simonyan, 2014)

The VGGNet achieved almost 92.7% top-5 test accuracy in 2014 ImageNet

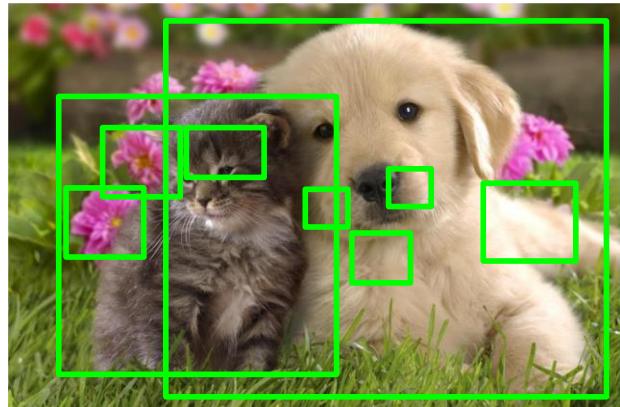


### ConvNet Architecture:

- VGG16: 13 convolutional layers and 3 fully connected layers (1000 classes)
- reduces spacial dimensions: 224x224 (image) → ... → 7x7
- increases depth: 3 (image RGB depth) → ... → 512 (nr of filters)
- Large network. Even larger version: VGG19

# Region Proposals

- Find a small set of boxes that are likely to cover all objects
- Often based on heuristics: e.g. look for “blob-like” image regions
- Relatively fast to run; e.g. Selective Search gives 2000 region proposals
- proposals in a few seconds on CPU



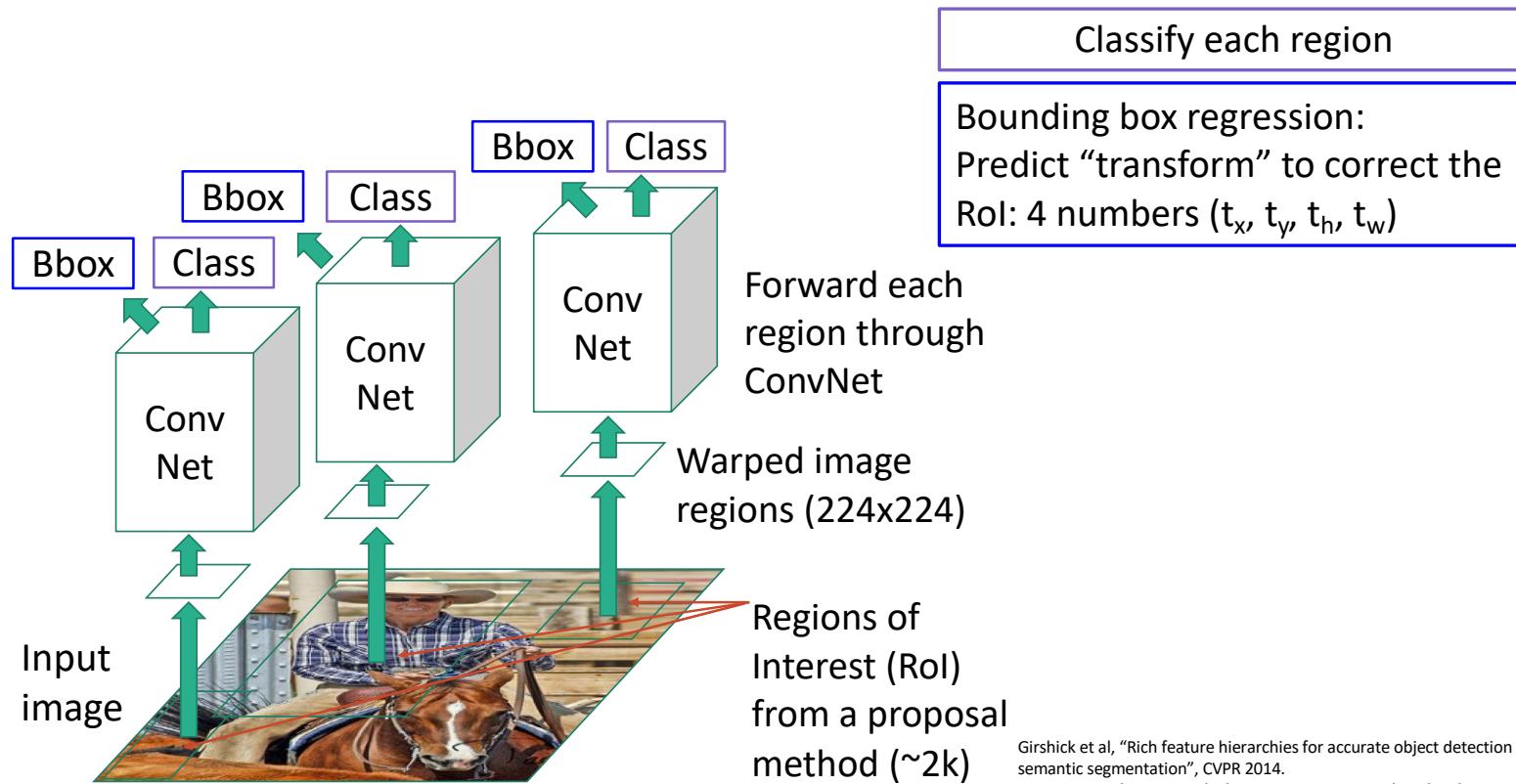
Alexe et al, “Measuring the objectness of image windows”, TPAMI 2012

Uijlings et al, “Selective Search for Object Recognition”, IJCV 2013

Cheng et al, “BING: Binarized normed gradients for objectness estimation at 300fps”, CVPR 2014

Zitnick and Dollar, “Edge boxes: Locating object proposals from edges”, ECCV 2014

# Region-based CNN (R-CNN)

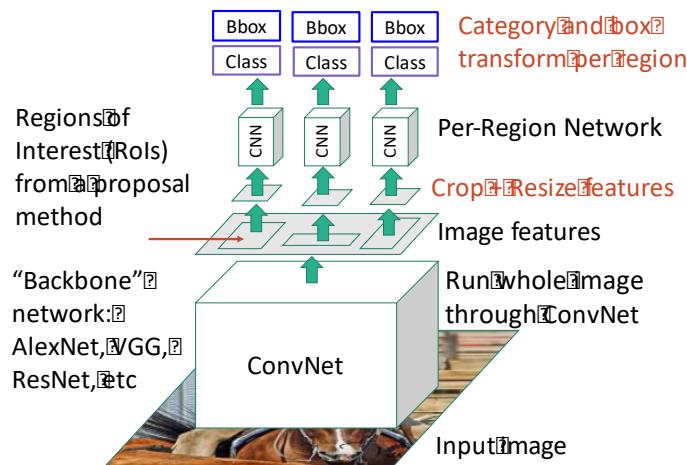


Girshick et al, "Rich feature hierarchies for accurate object detection and semantic segmentation", CVPR 2014.

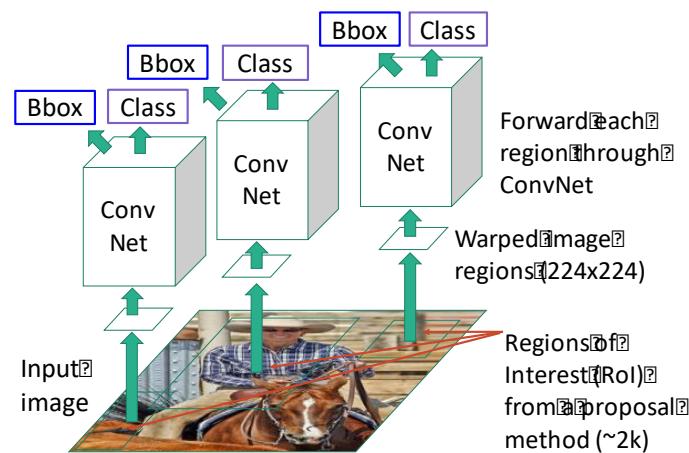
Figure copyright Ross Girshick, 2015; [source](#). Reproduced with permission.

# Fast R-CNN vs. (Slow) R-CNN

**Fast R-CNN:** Apply differentiable cropping to shared image features



**“Slow” R-CNN:** Apply differentiable cropping to shared image features

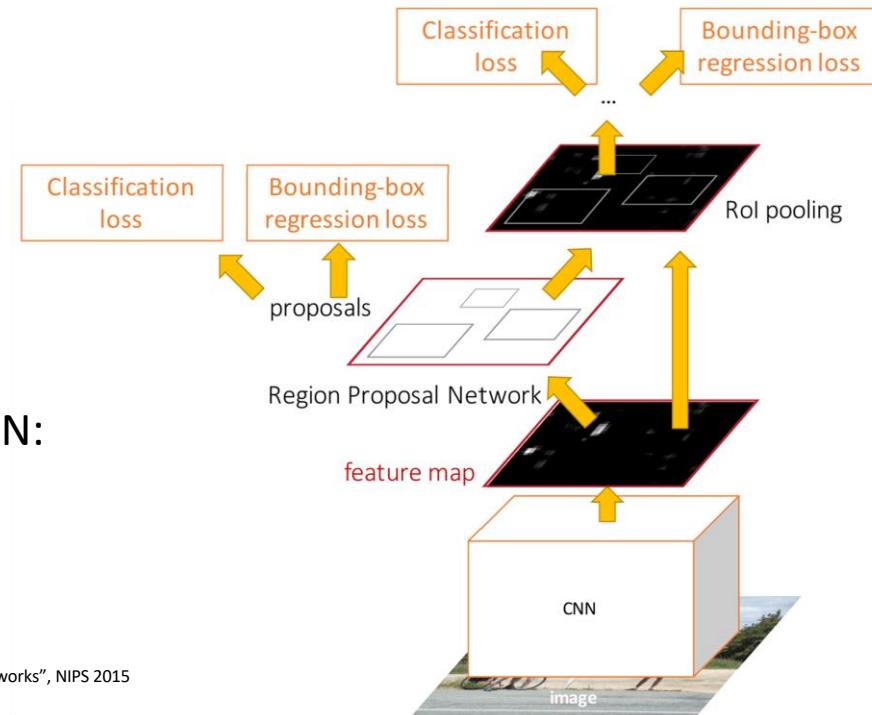


# Faster R-CNN: Learnable Region Proposals



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

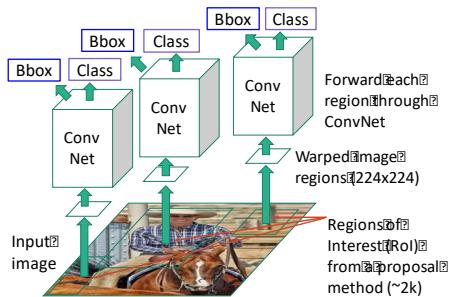
Otherwise same as Fast R-CNN:  
Crop features for each proposal, classify each one



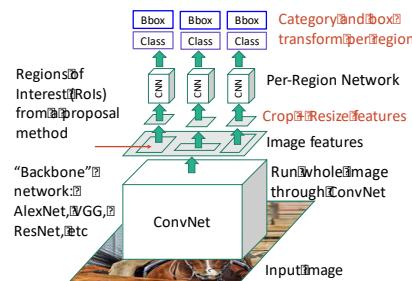
Ren et al, "Faster R-CNN: Towards Real-Time Object Detection with Region Proposal Networks", NIPS 2015  
Figure copyright 2015, Ross Girshick; reproduced with permission

# Overview: R-CNN Variants

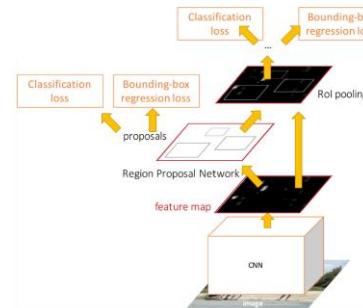
**“Slow” R-CNN:** Run CNN independently for each region



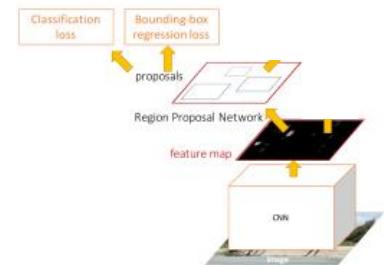
**Fast R-CNN:** Apply differentiable cropping to shared image features



**Faster R-CNN:** Compute proposals with CNN



**Single-Stage:** Fully convolutional detector



# Today's Agenda

---

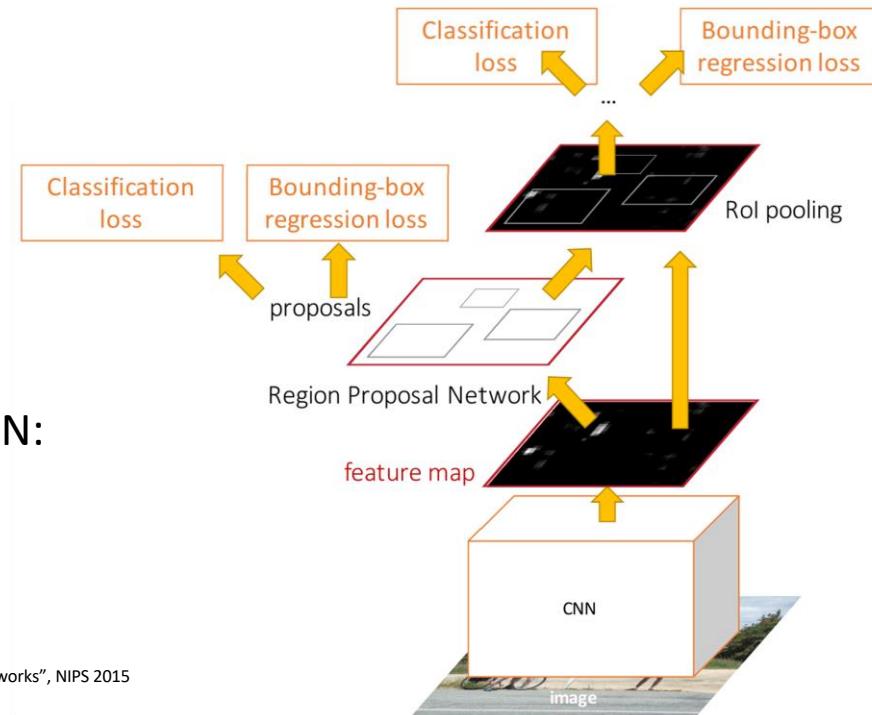
- Learning region proposals - Faster R-CNN
- Feature Pyramid Network (FPN)
- Single Shot Detectors (SSD) / Yolo
- Transfer learning
- Network architectures
  - VGG
  - ResNet
  - U-Net
  - Autoencoders

# Faster R-CNN: Learnable Region Proposals



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

Otherwise same as Fast R-CNN:  
Crop features for each proposal, classify each one



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

# Region Proposal Network (RPN)

Run backbone CNN to get  
features aligned to input image



Input Image  
(e.g.  $3 \times 640 \times 480$ )

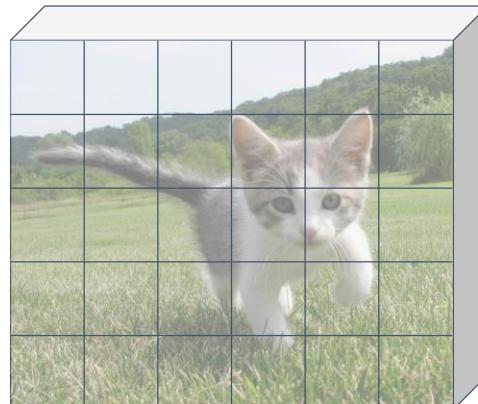


Image features  
(e.g.  $512 \times 5 \times 6$ )

Ren et al, "Faster R-CNN: Towards Real-Time Object Detection with Region Proposal Networks", NIPS 2015

# Region Proposal Network (RPN)

Run backbone CNN to get features aligned to input image



Input Image  
(e.g.  $3 \times 640 \times 480$ )



Each feature corresponds to a point in the input

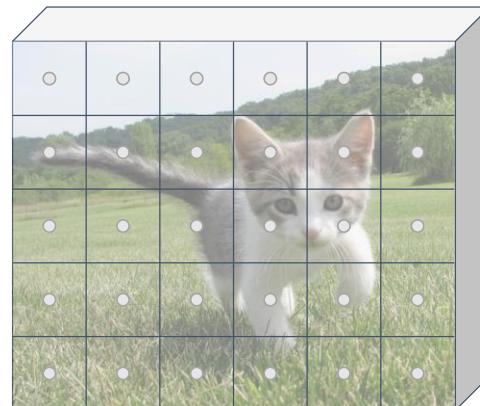
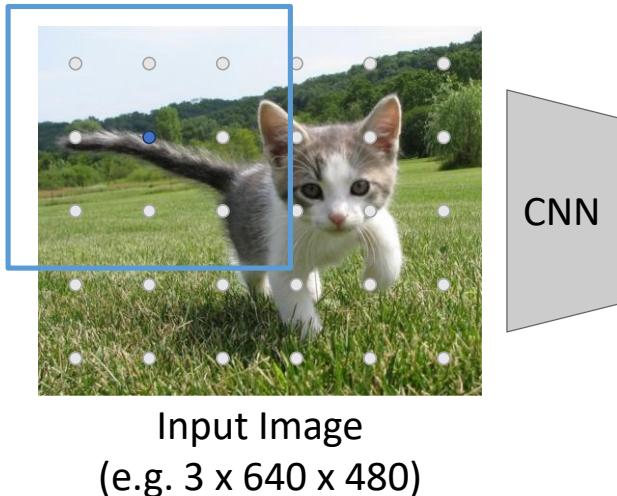


Image features  
(e.g.  $512 \times 5 \times 6$ )

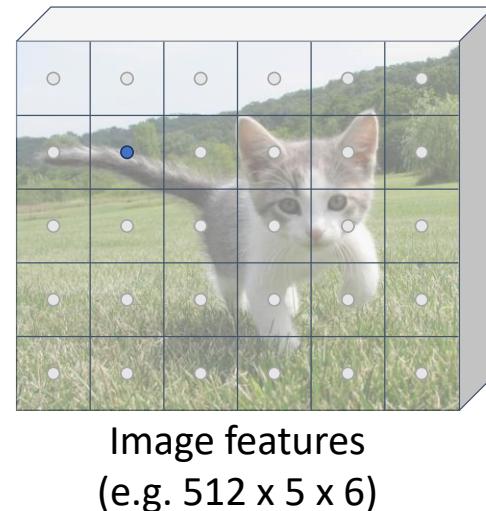
Ren et al., "Faster R-CNN: Towards Real-Time Object Detection with Region Proposal Networks", NIPS 2015

# Region Proposal Network (RPN)

Run backbone CNN to get features aligned to input image



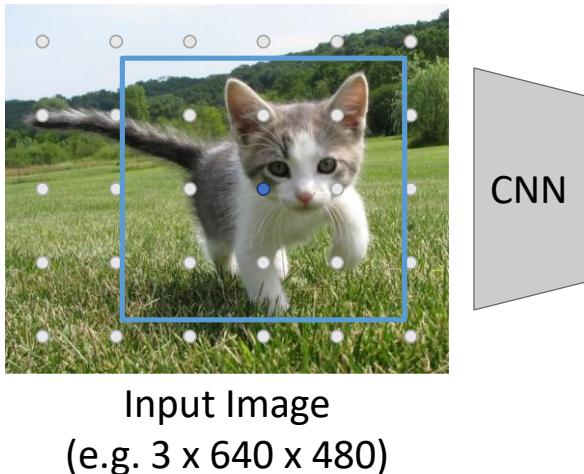
Each feature corresponds to a point in the input



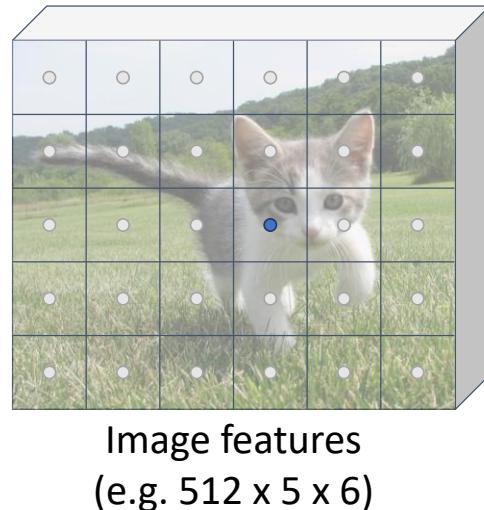
Imagine an **anchor box** of fixed size at each point in the feature map

# Region Proposal Network (RPN)

Run backbone CNN to get features aligned to input image



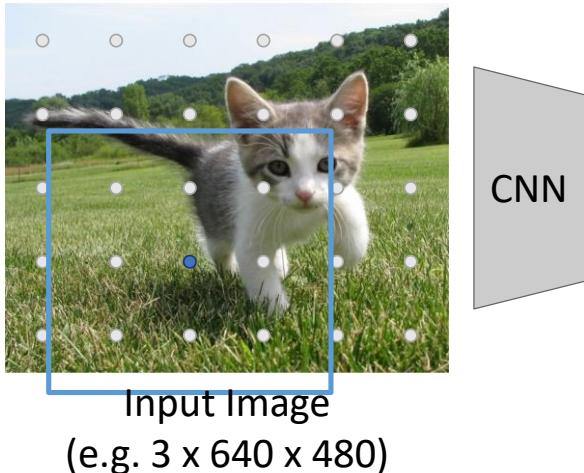
Each feature corresponds to a point in the input



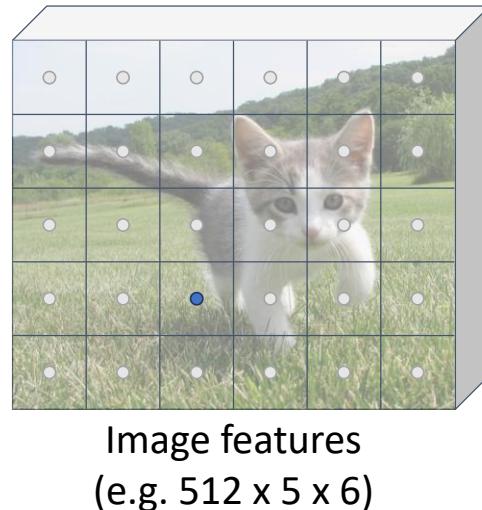
Imagine an **anchor box** of fixed size at each point in the feature map

# Region Proposal Network (RPN)

Run backbone CNN to get features aligned to input image



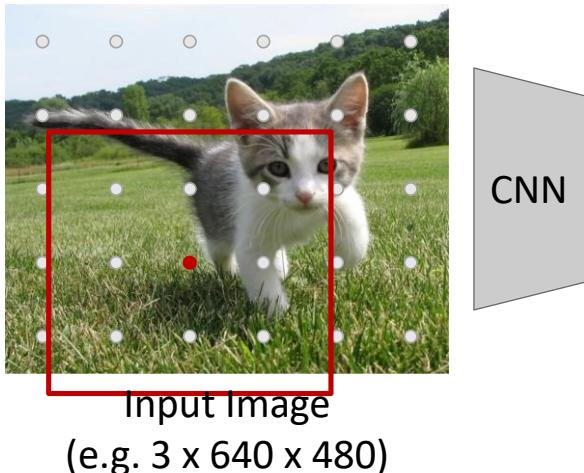
Each feature corresponds to a point in the input



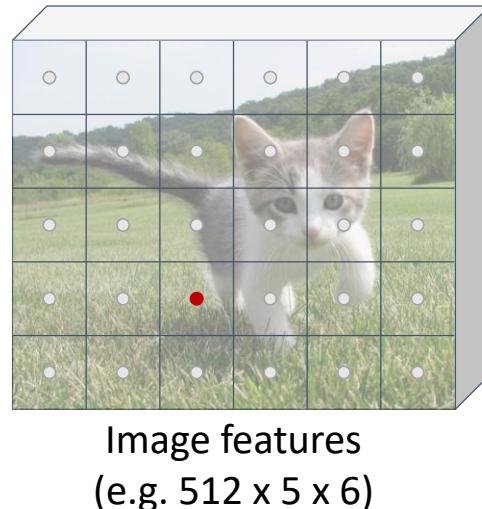
Imagine an **anchor box** of fixed size at each point in the feature map

# Region Proposal Network (RPN)

Run backbone CNN to get features aligned to input image



Each feature corresponds to a point in the input

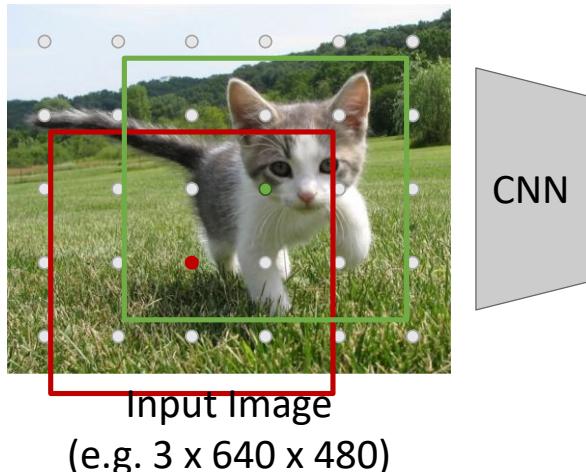


Imagine an **anchor box** of fixed size at each point in the feature map

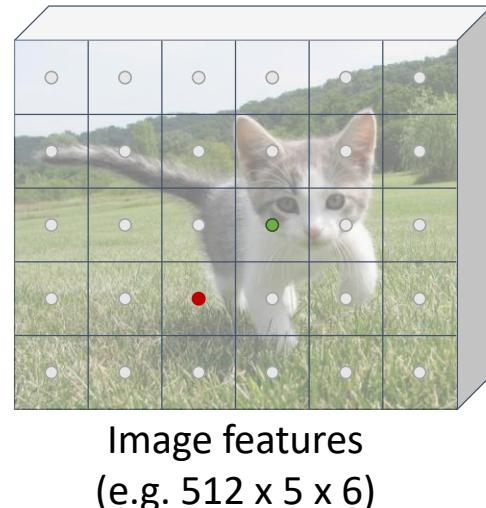
Classify each anchor as **positive (object)** or **negative (no object)**

# Region Proposal Network (RPN)

Run backbone CNN to get features aligned to input image



Each feature corresponds to a point in the input

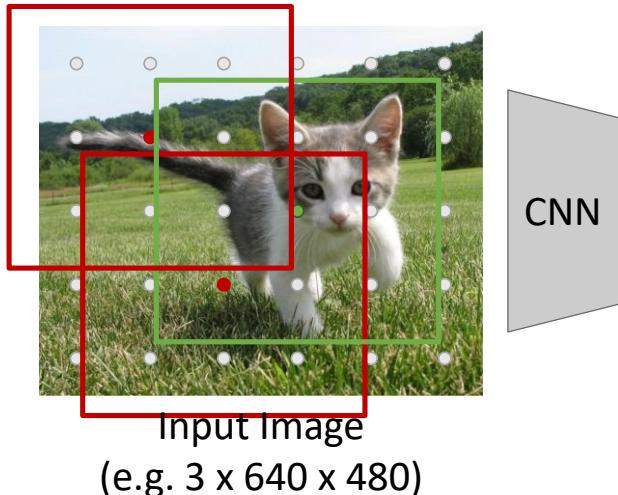


Imagine an **anchor box** of fixed size at each point in the feature map

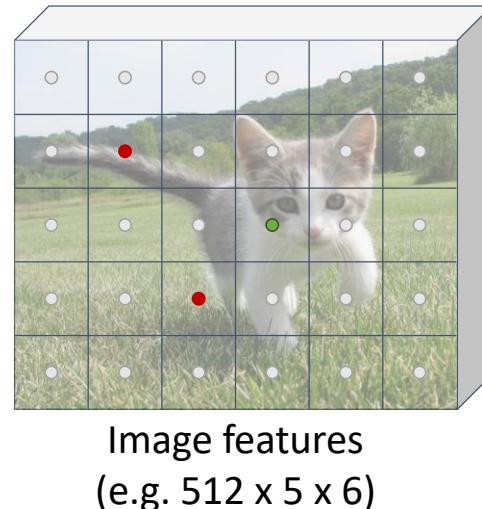
Classify each anchor as **positive (object)** or **negative (no object)**

# Region Proposal Network (RPN)

Run backbone CNN to get features aligned to input image



Each feature corresponds to a point in the input

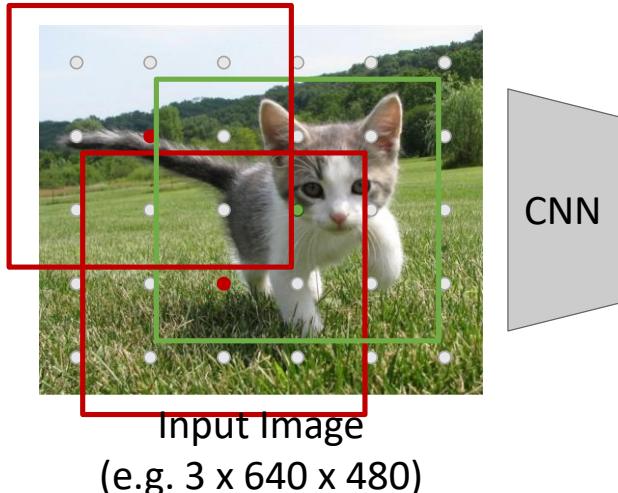


Imagine an **anchor box** of fixed size at each point in the feature map

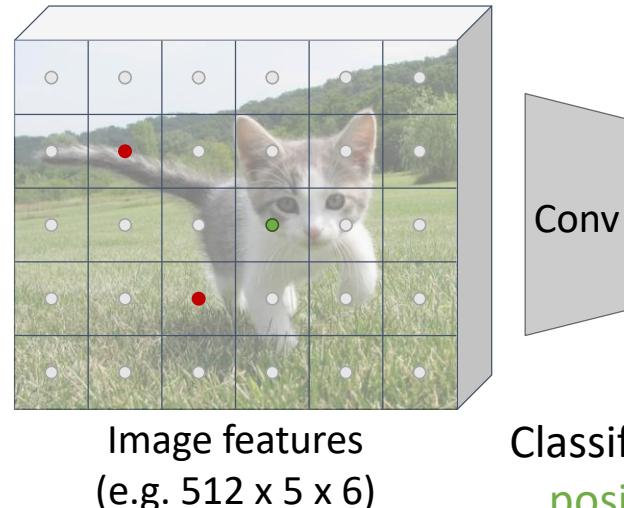
Classify each anchor as **positive (object)** or **negative (no object)**

# Region Proposal Network (RPN)

Run backbone CNN to get features aligned to input image



Each feature corresponds to a point in the input



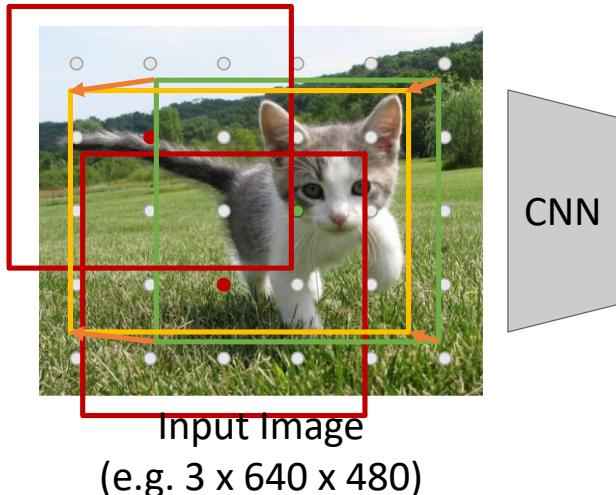
Predict object vs not object scores for all anchors with a conv layer (512 input filters, 2 output filters)

Anchor is object?  
 $2 \times 5 \times 6$

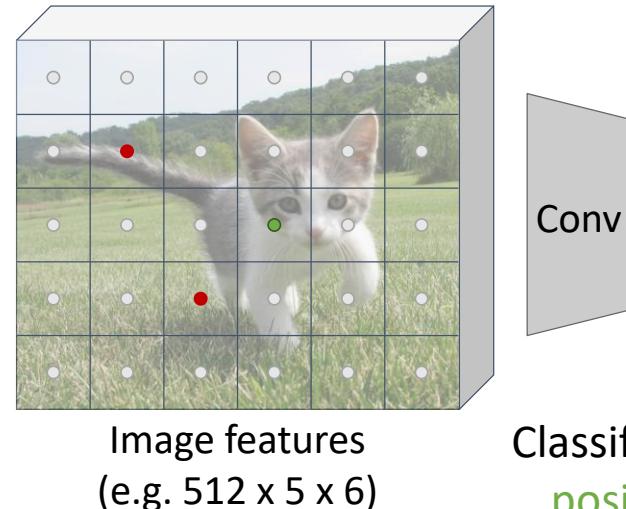
Classify each anchor as positive (object) or negative (no object)

# Region Proposal Network (RPN)

Run backbone CNN to get features aligned to input image



Each feature corresponds to a point in the input



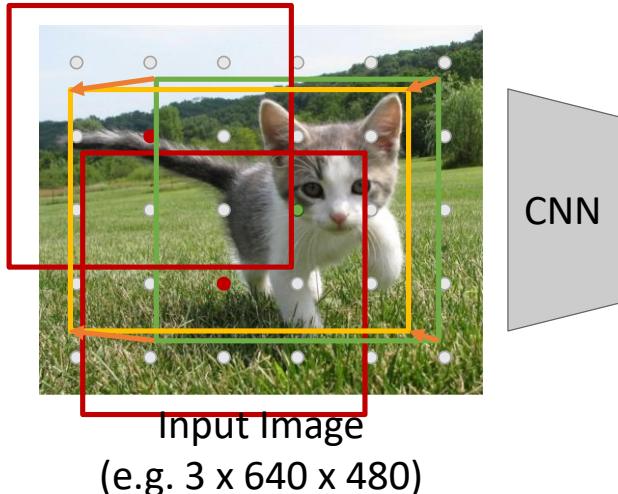
For **positive anchors**, also predict a **transform** that converting the anchor to the **GT box** (like R-CNN)

Anchor is object?  
 $2 \times 5 \times 6$

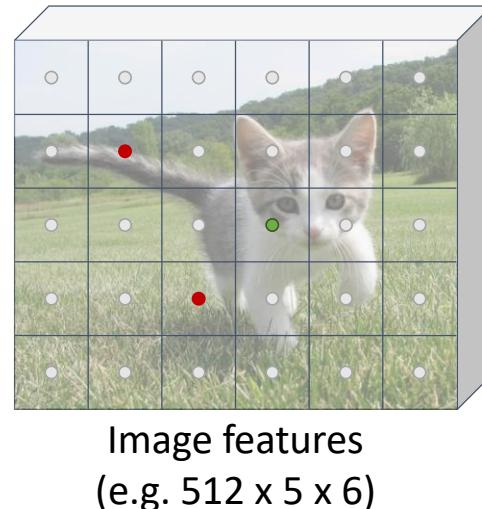
Classify each anchor as **positive (object)** or **negative (no object)**

# Region Proposal Network (RPN)

Run backbone CNN to get features aligned to input image



Each feature corresponds to a point in the input



For **positive anchors**, also predict a **transform** that converting the anchor to the **GT box** (like R-CNN)  
Predict transforms with conv

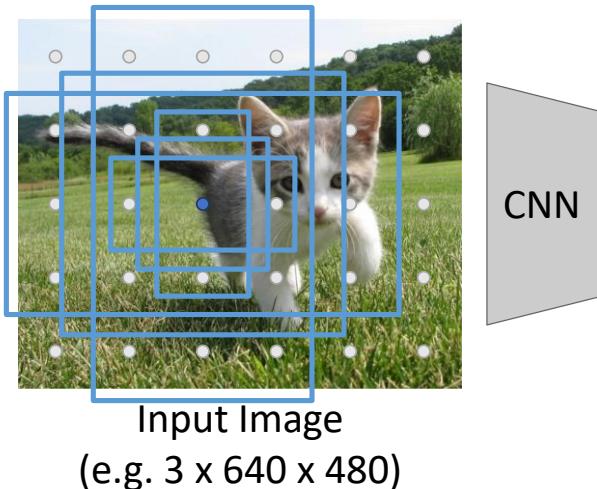
Anchor is object?  
 $2 \times 5 \times 6$

Anchor transforms  
 $4 \times 5 \times 6$

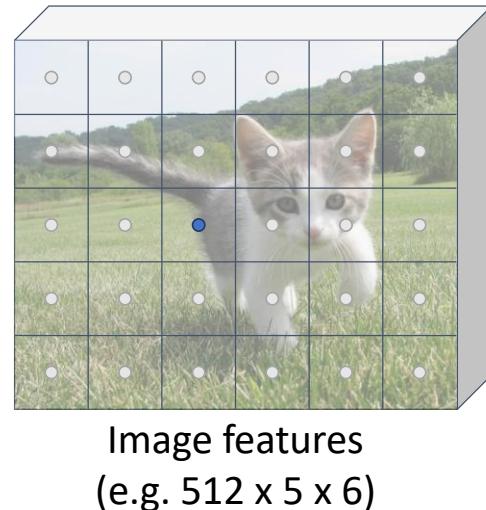
Classify each anchor as **positive (object)** or **negative (no object)**

# Region Proposal Network (RPN)

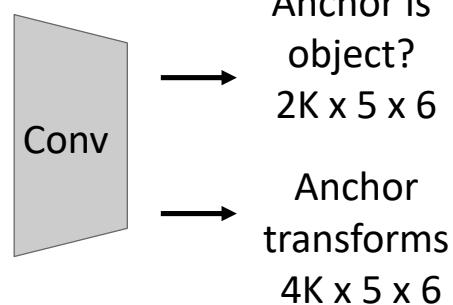
Run backbone CNN to get features aligned to input image



Each feature corresponds to a point in the input

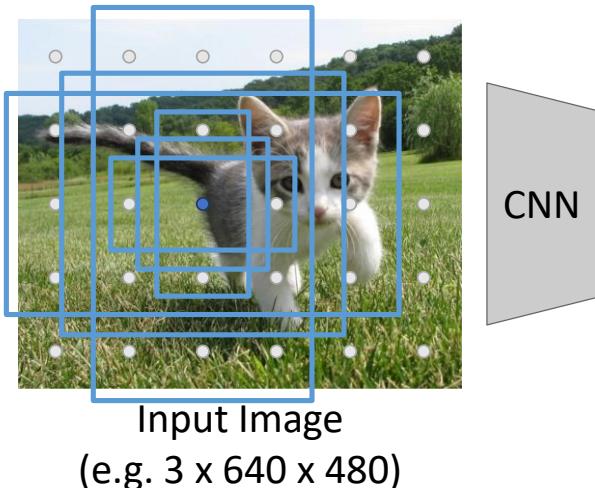


In practice: Rather than using one anchor per point, instead consider K different anchors with different size and scale (here  $K = 6$ )

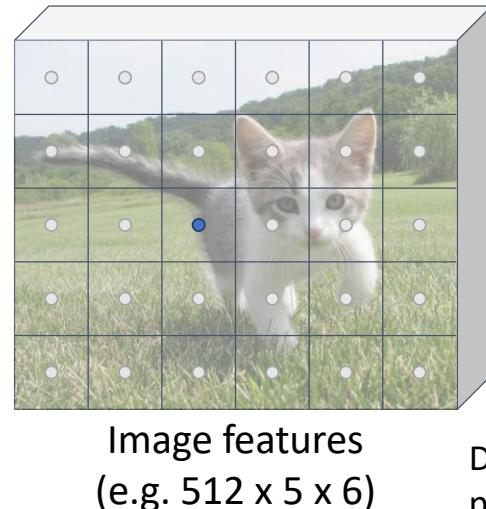


# Region Proposal Network (RPN)

Run backbone CNN to get features aligned to input image



Each feature corresponds to a point in the input



In practice: Rather than using one anchor per point, instead consider K different anchors with different size and scale (here  $K = 6$ )

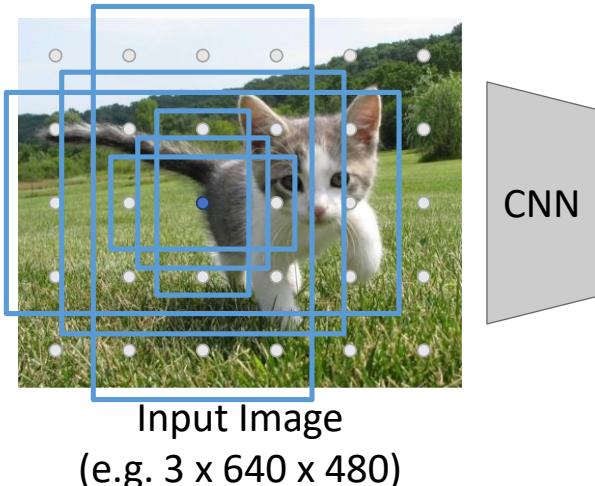
Anchor is object?  
 $2K \times 5 \times 6$

Anchor transforms  
 $4K \times 5 \times 6$

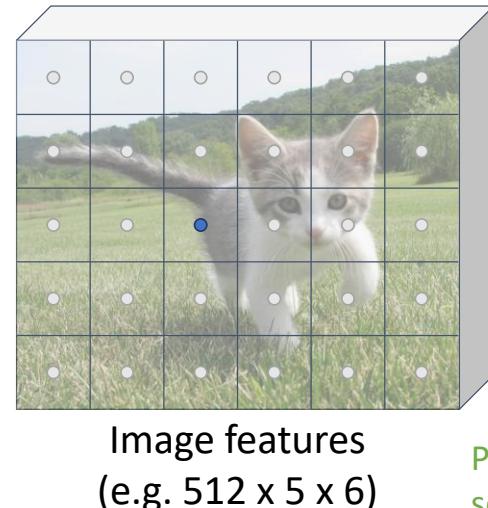
During training, supervised positive / negative anchors and box transforms like R-CNN

# Region Proposal Network (RPN)

Run backbone CNN to get features aligned to input image



Each feature corresponds to a point in the input



In practice: Rather than using one anchor per point, instead consider K different anchors with different size and scale (here  $K = 6$ )

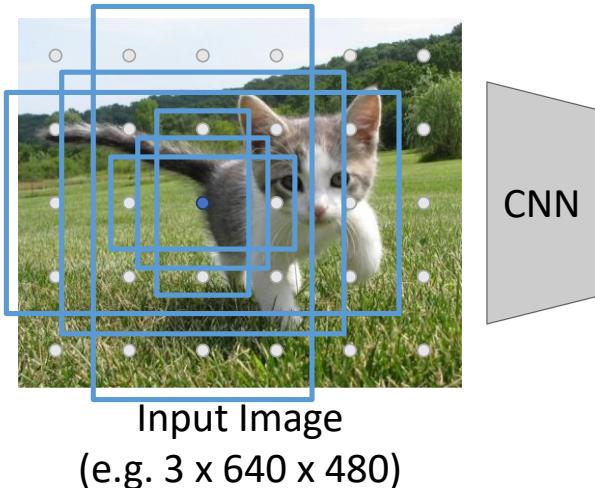
Anchor is object?  
 $2K \times 5 \times 6$

Anchor transforms  
 $4K \times 5 \times 6$

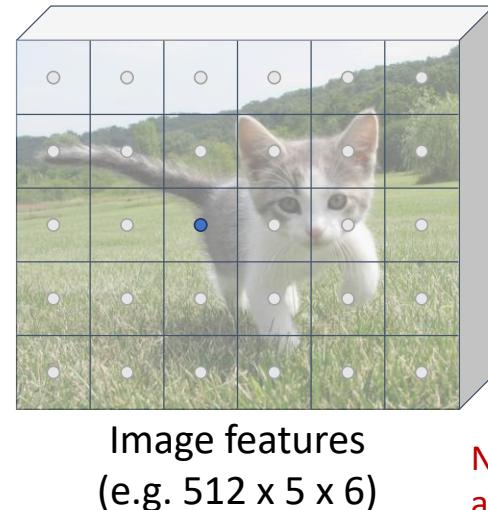
Positive anchors:  $\geq 0.7$  IoU with some GT box (plus highest IoU to each GT)

# Region Proposal Network (RPN)

Run backbone CNN to get features aligned to input image



Each feature corresponds to a point in the input



In practice: Rather than using one anchor per point, instead consider K different anchors with different size and scale (here  $K = 6$ )

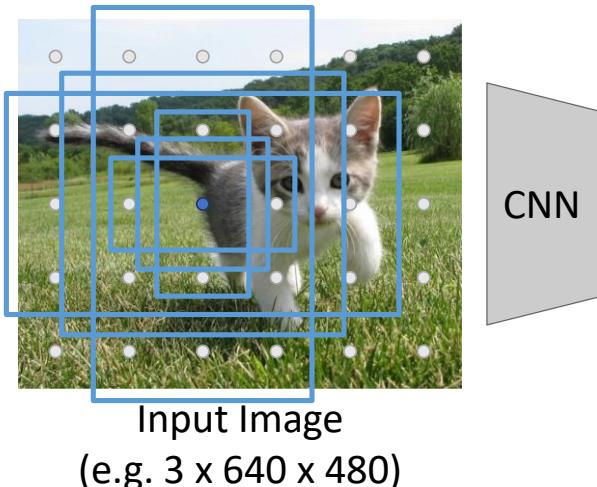
Anchor is object?  
 $2K \times 5 \times 6$

Anchor transforms  
 $4K \times 5 \times 6$

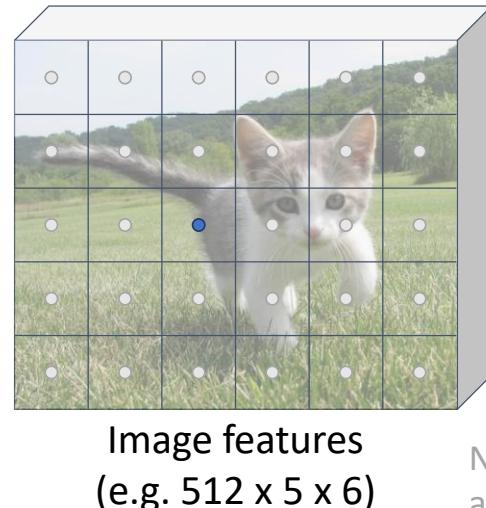
Negative anchors:  $< 0.3$  IoU with all GT boxes. Don't supervise transforms for negative boxes.

# Region Proposal Network (RPN)

Run backbone CNN to get features aligned to input image



Each feature corresponds to a point in the input



In practice: Rather than using one anchor per point, instead consider K different anchors with different size and scale (here K = 6)

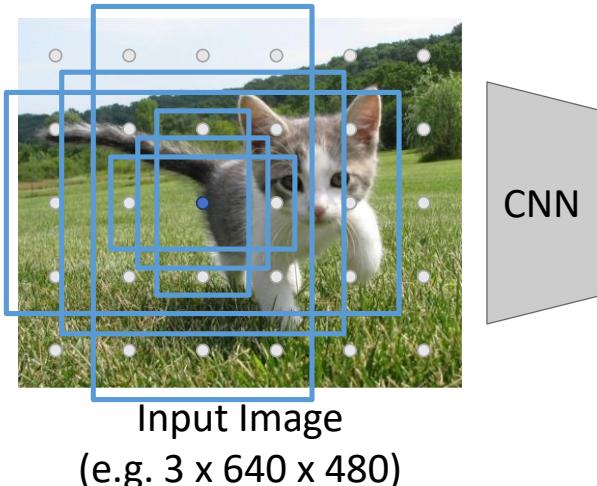
Anchor is object?  
2K x 5 x 6

Anchor transforms  
4K x 5 x 6

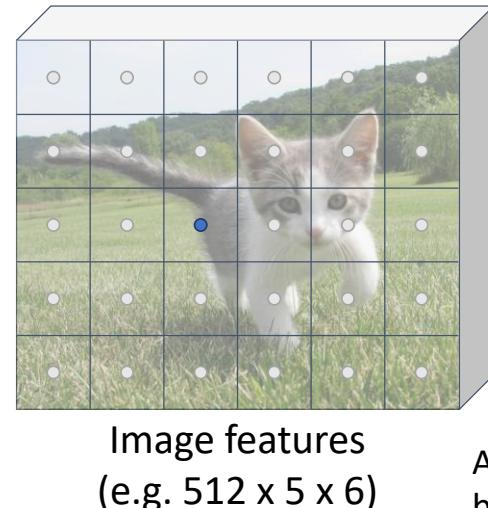
Neutral anchors: between 0.3 and 0.7 IoU with all GT boxes; ignored during training

# Region Proposal Network (RPN)

Run backbone CNN to get features aligned to input image



Each feature corresponds to a point in the input



In practice: Rather than using one anchor per point, instead consider K different anchors with different size and scale (here  $K = 6$ )

→ Anchor is object?  $2K \times 5 \times 6$   
→ Anchor transforms  $4K \times 5 \times 6$

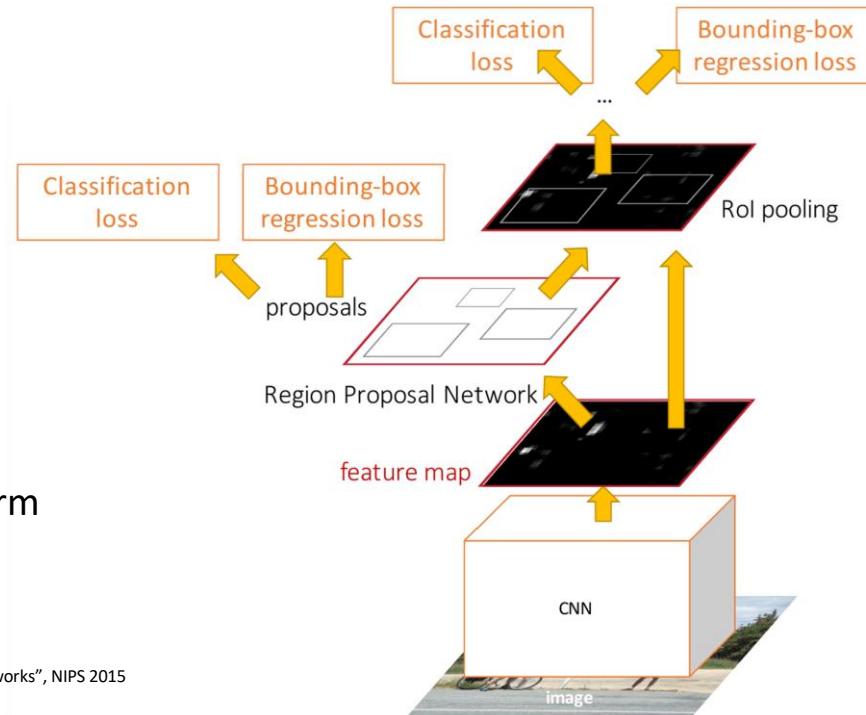
At test-time, sort all  $K \times 5 \times 6$  boxes by their positive score, take top 300 as our region proposals

# Faster R-CNN: Learnable Region Proposals



Jointly train with 4 losses:

1. **RPN classification**: anchor box is object / not an object
2. **RPN regression**: predict transform from anchor box to proposal box
3. **Object classification**: classify proposals as background / object class
4. **Object regression**: predict transform from proposal box to object box



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

# Faster R-CNN: Learnable Region Proposals

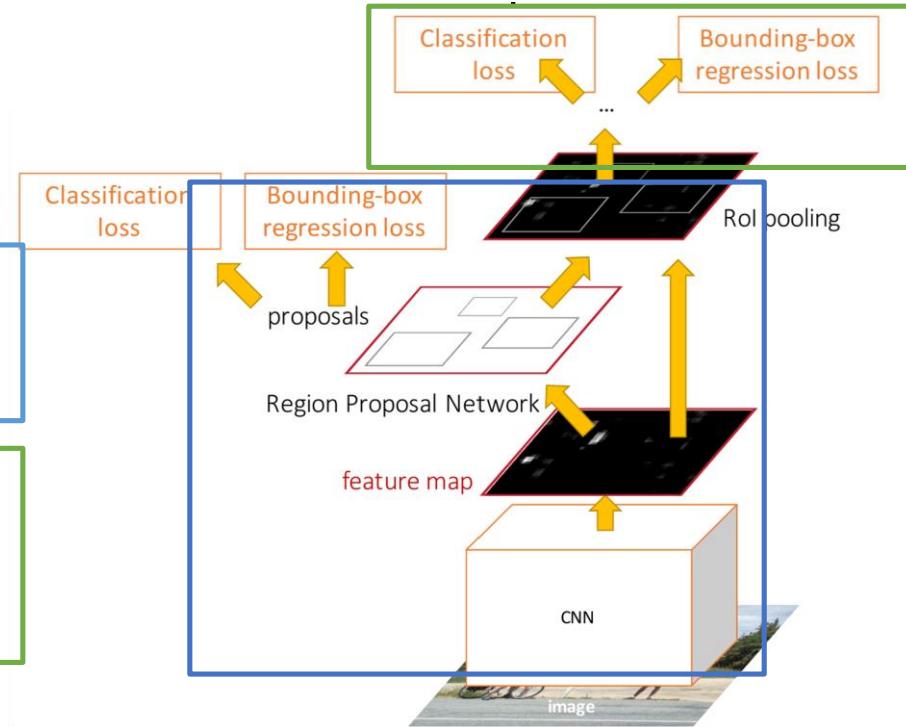
Faster R-CNN is a  
**Two-stage object detector**

First stage: Run once per image

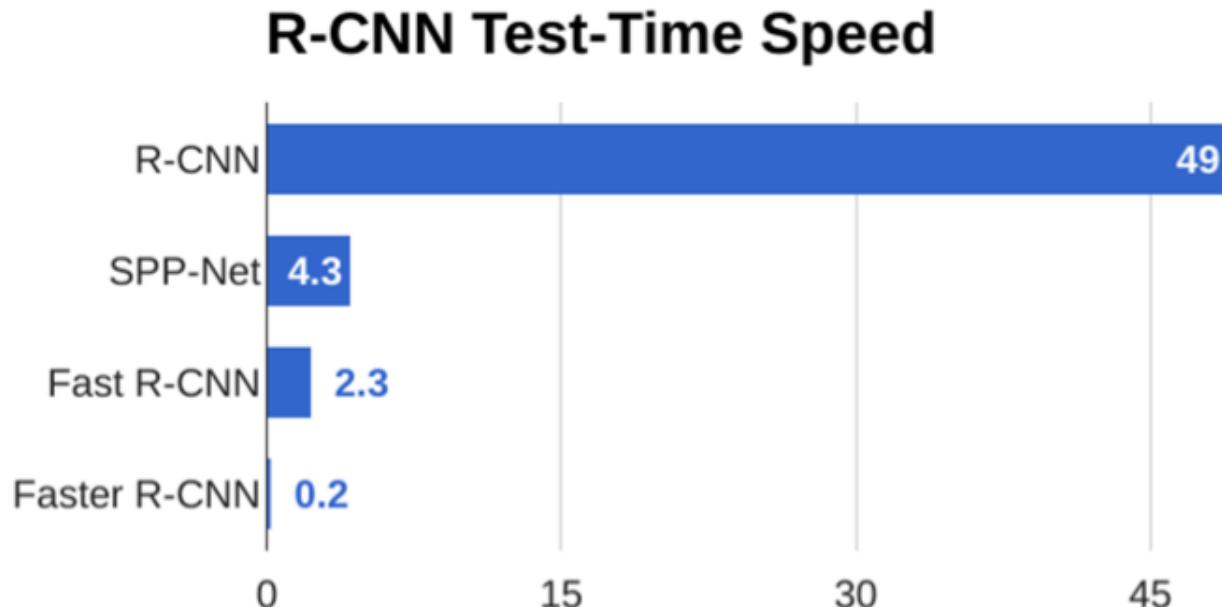
- Backbone network
- Region proposal network

Second stage: Run once per region

- Crop features: RoI pool / align
- Predict object class
- Prediction bbox offset



# Faster R-CNN: Learnable Region Proposals



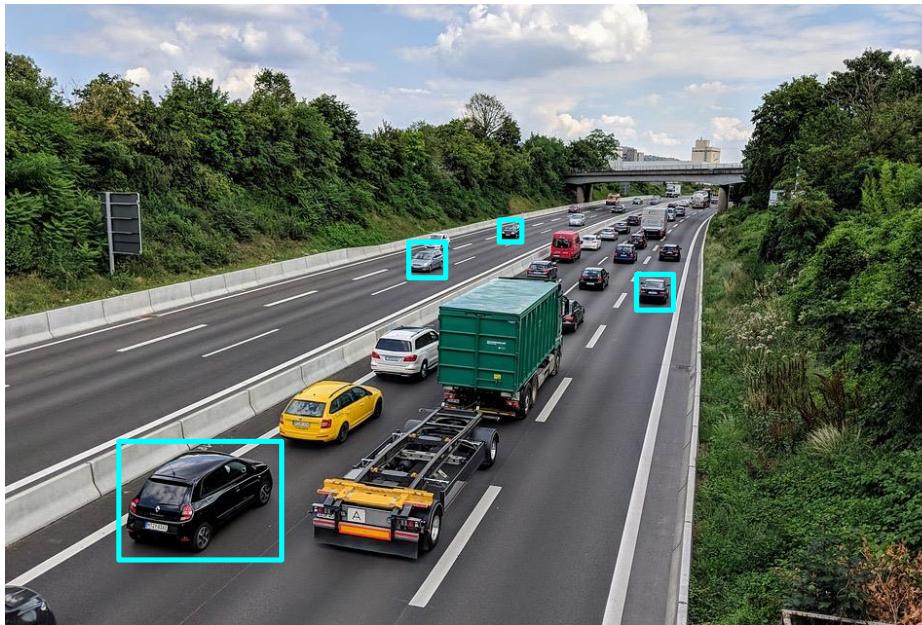
# Today's Agenda

---

- Learning region proposals - Faster R-CNN
- Feature Pyramid Network (FPN)
- Single Shot Detectors (SSD) / Yolo
- Transfer learning
- Network architectures
  - VGG
  - ResNet
  - U-Net
  - Autoencoders

# Dealing with Scale Variations

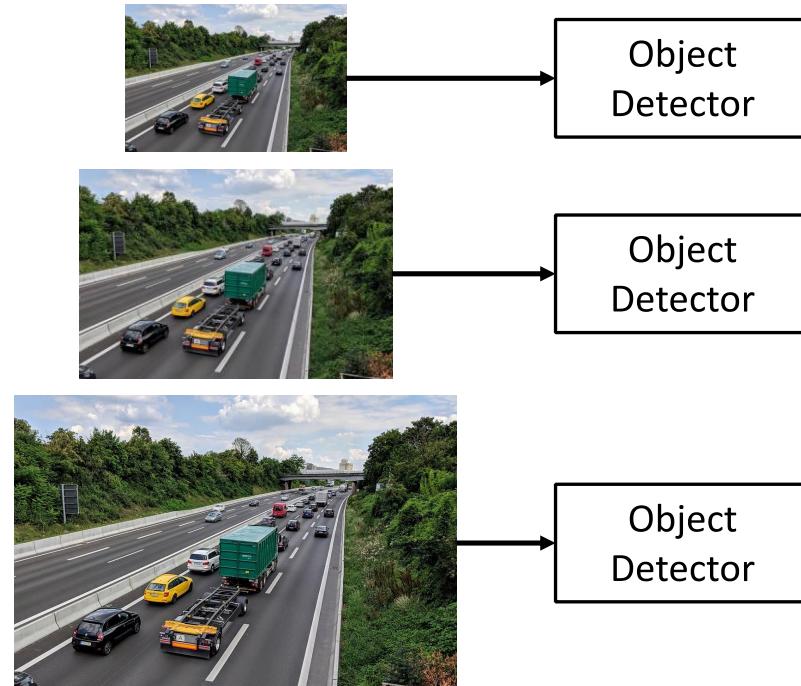
We need to detect objects of many different scales.  
How to improve *scale invariance* of the detector?



This image is free for commercial  
use under the [Pixabay license](#)

# Scale Variations: Image Pyramid

Classic idea: build an *image pyramid* by resizing the image to different scales, then process each image scale independently.



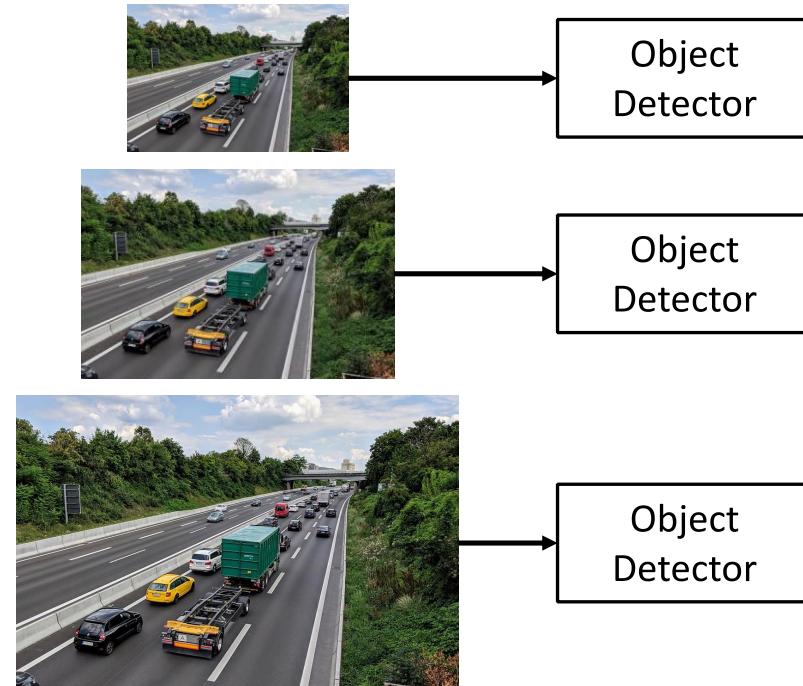
Lin et al, "Feature Pyramid Networks for Object Detection", ICCV 2017

# Scale Variations: Image Pyramid

Classic idea: build an *image pyramid* by resizing the image to different scales, then process each image scale independently.

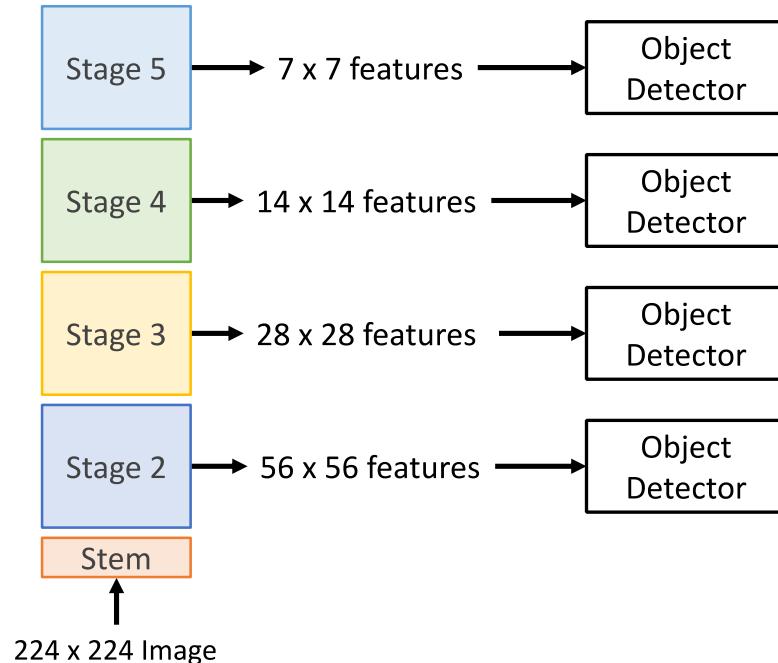
Problem: Expensive! Don't share any computation between scales

Lin et al, "Feature Pyramid Networks for Object Detection", ICCV 2017



# Scale Variations: Multiscale Features

CNNs have multiple *stages* that operate at different resolutions. Attach an independent detector to the features at each level



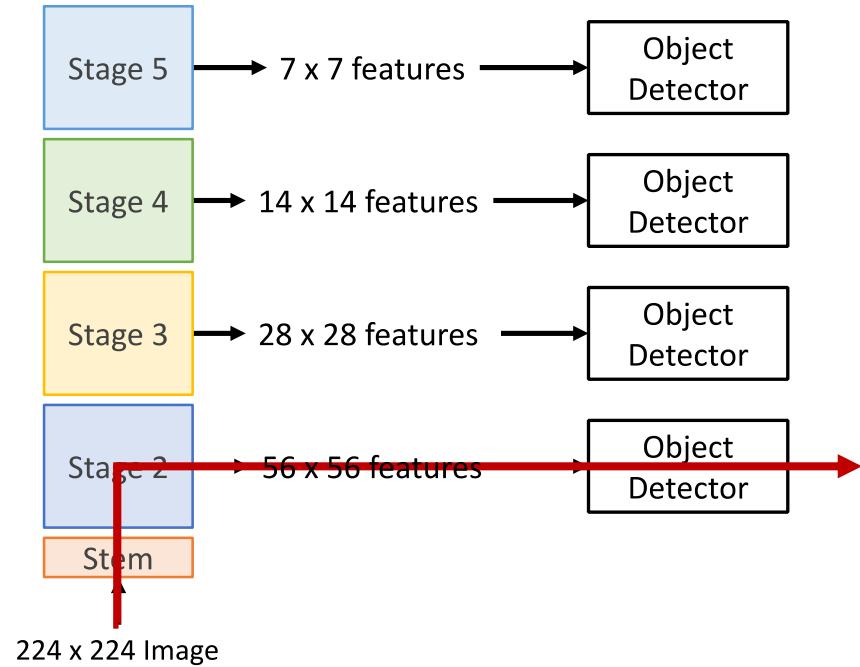
Lin et al, "Feature Pyramid Networks for Object Detection", ICCV 2017

# Scale Variations: Multiscale Features

CNNs have multiple *stages* that operate at different resolutions. Attach an independent detector to the features at each level

Problem: detector on early features doesn't make use of the entire backbone; doesn't get access to high-level features

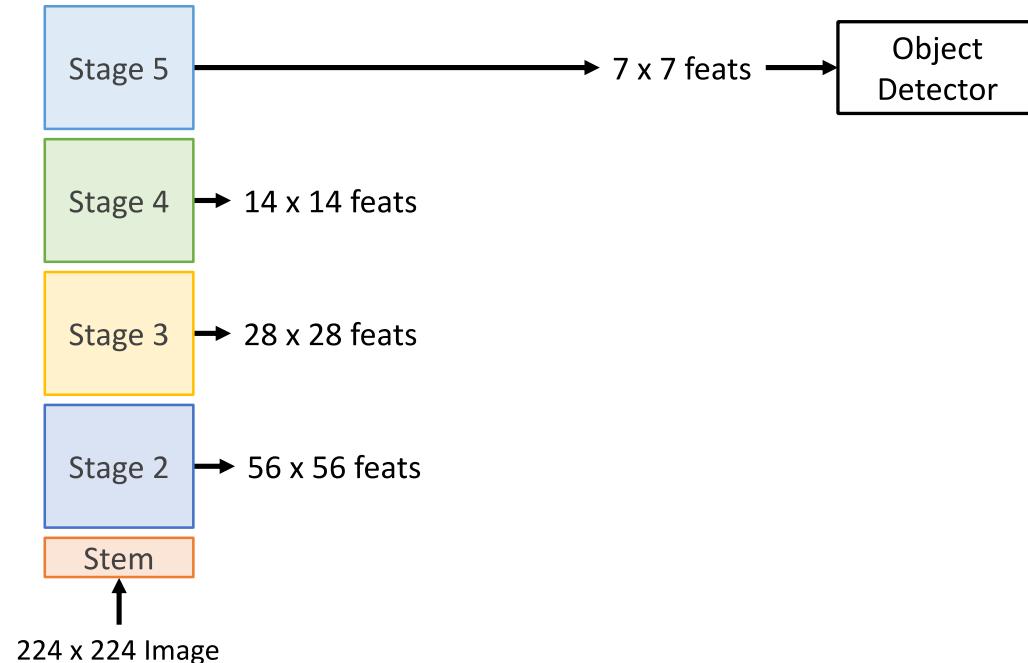
Lin et al, "Feature Pyramid Networks for Object Detection", ICCV 2017



# Scale Variations: Feature Pyramid Network



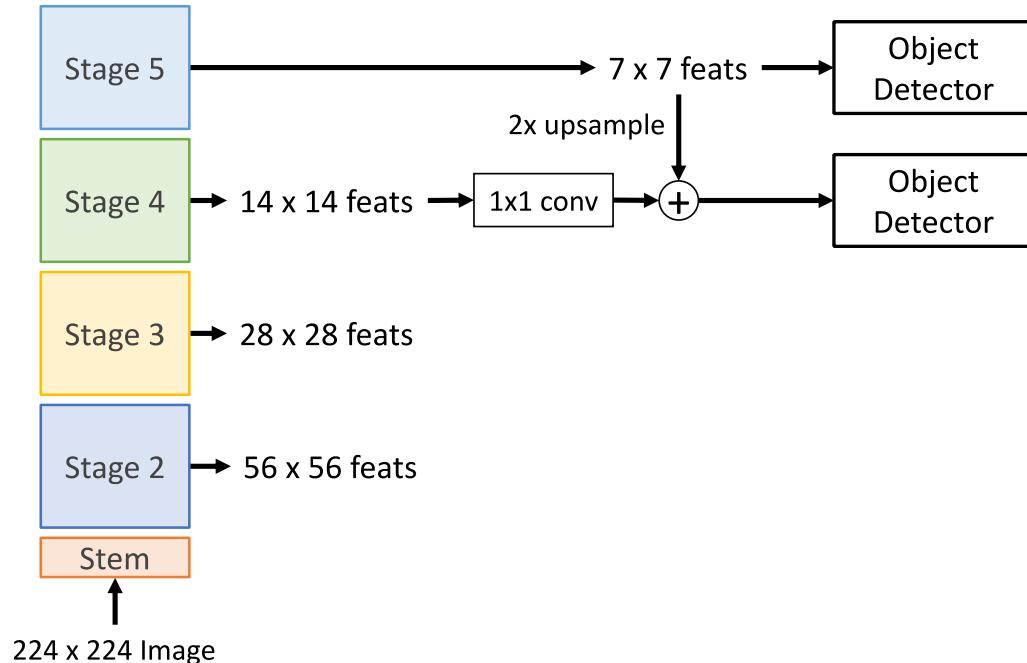
Add *top down connections* that feed information from high level features back down to lower level features



Lin et al, "Feature Pyramid Networks for Object Detection", ICCV 2017

# Scale Variations: Feature Pyramid Network

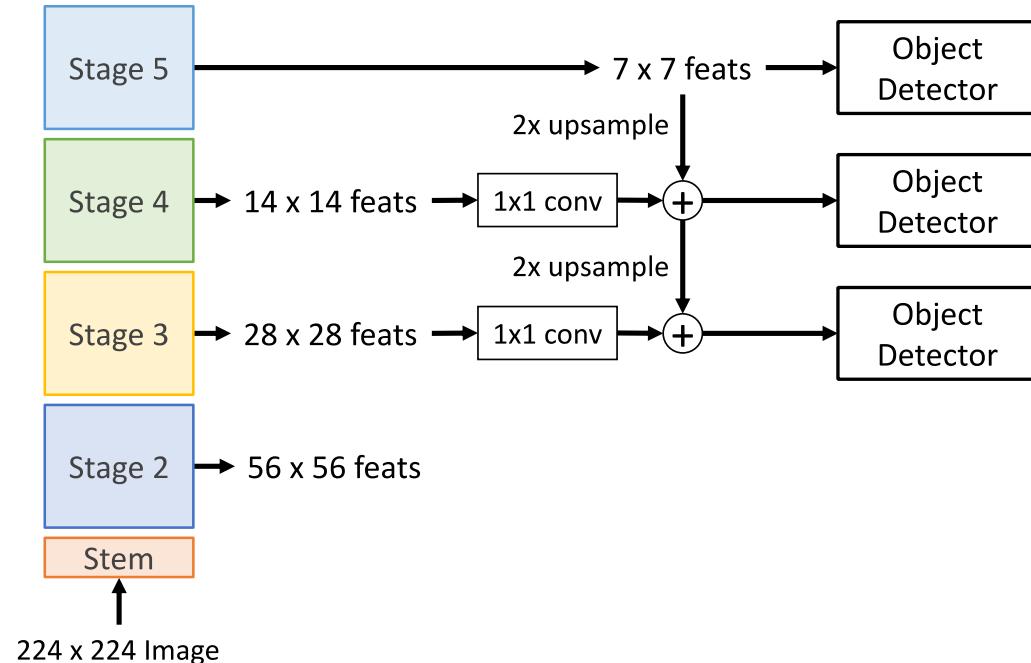
Add *top down connections* that feed information from high level features back down to lower level features



Lin et al, "Feature Pyramid Networks for Object Detection", ICCV 2017

# Scale Variations: Feature Pyramid Network

Add *top down connections* that feed information from high level features back down to lower level features

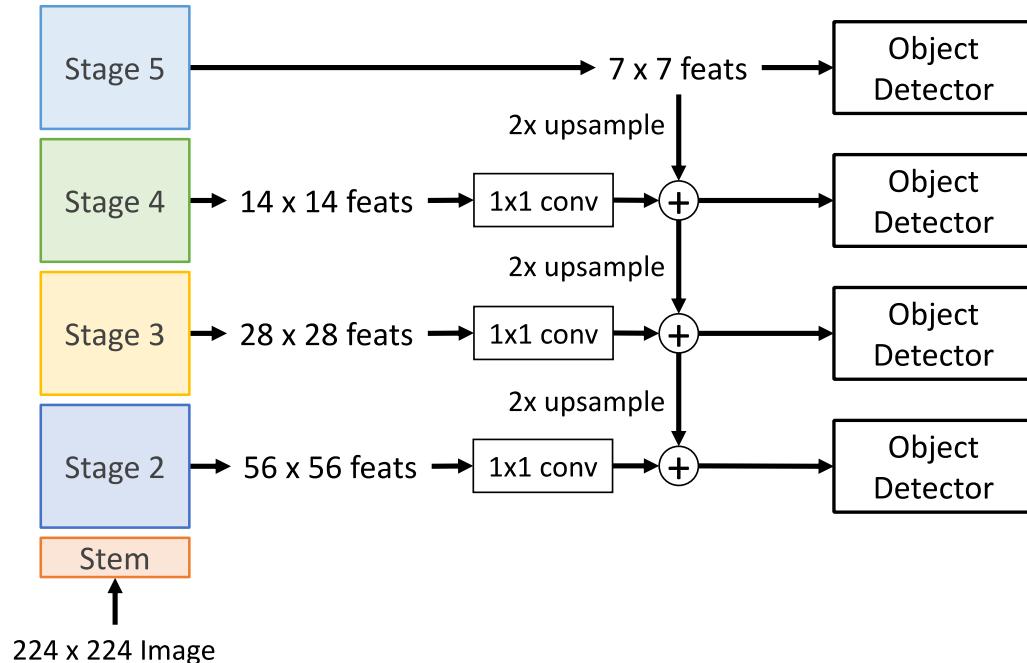


Lin et al, "Feature Pyramid Networks for Object Detection", ICCV 2017

# Scale Variations: Feature Pyramid Network



Add *top down connections* that feed information from high level features back down to lower level features



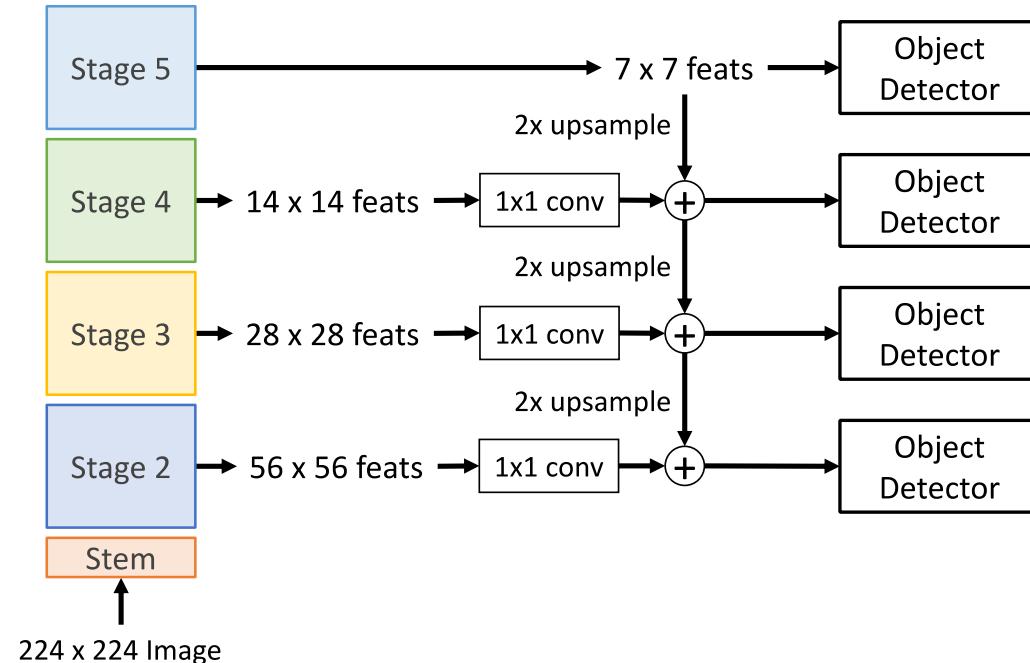
Lin et al, "Feature Pyramid Networks for Object Detection", ICCV 2017

# Scale Variations: Feature Pyramid Network



Add *top down connections* that feed information from high level features back down to lower level features

Efficient multiscale features where all levels benefit from the whole backbone! Widely used in practice



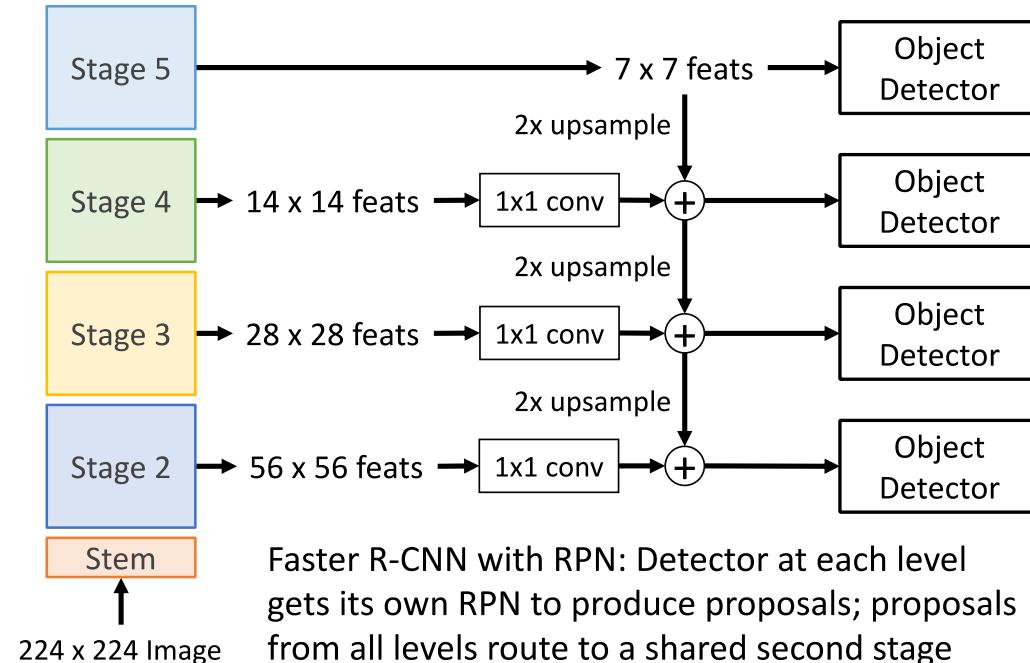
Lin et al, "Feature Pyramid Networks for Object Detection", ICCV 2017

# Scale Variations: Feature Pyramid Network



Add *top down connections* that feed information from high level features back down to lower level features

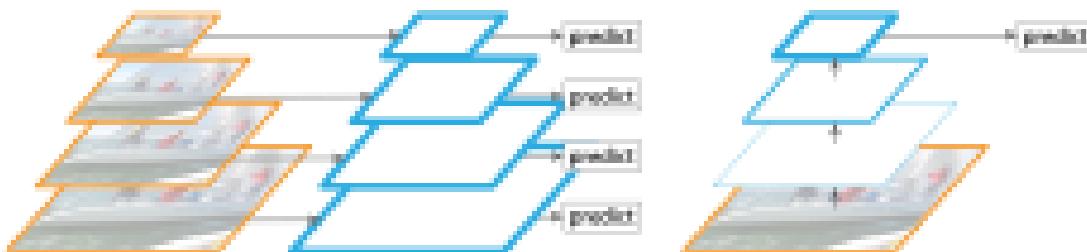
Efficient multiscale features where all levels benefit from the whole backbone! Widely used in practice



Lin et al, "Feature Pyramid Networks for Object Detection", ICCV 2017

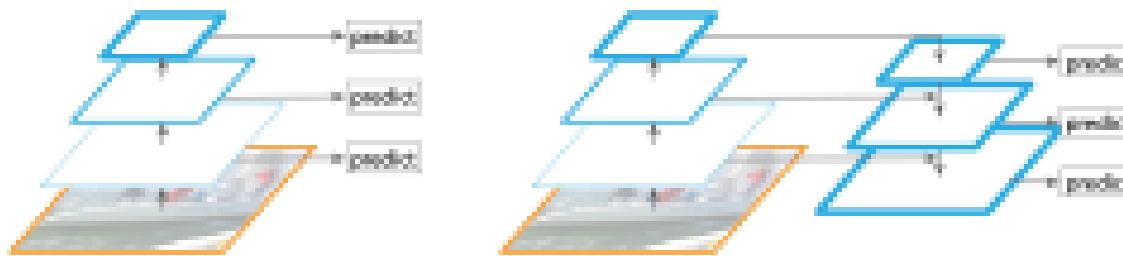
# Scale Variations: Feature Pyramid Network

CV



(a) Featureized image pyramid

(b) Single feature map



(c) Pyramidal feature hierarchy

(d) Feature Pyramid Network

# Faster R-CNN: Learnable Region Proposals

Faster R-CNN is a  
**Two-stage object detector**

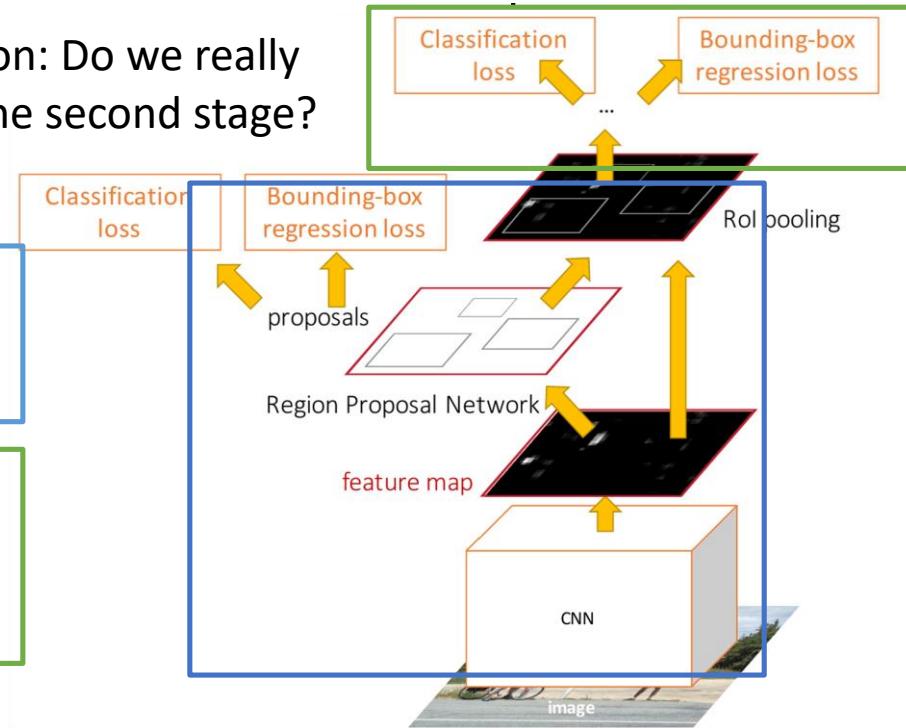
First stage: Run once per image

- Backbone network
- Region proposal network

Second stage: Run once per region

- Crop features: RoI pool / align
- Predict object class
- Prediction bbox offset

Question: Do we really  
need the second stage?



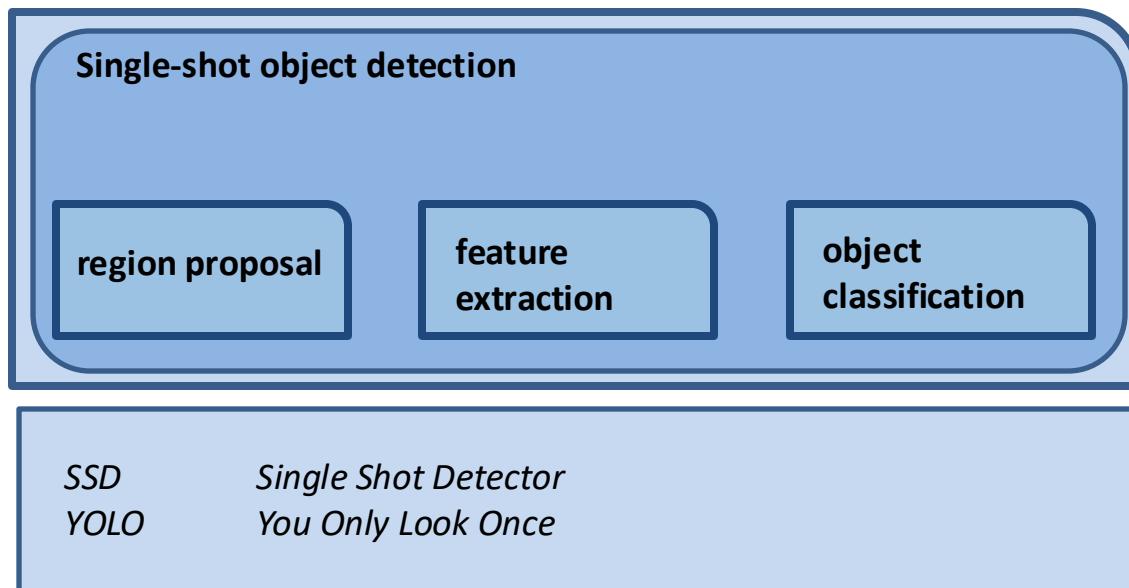
# Today's Agenda

---

- Learning region proposals - Faster R-CNN
- Feature Pyramid Network (FPN)
- Single Shot Detectors (SSD) / Yolo
- Transfer learning
- Network architectures
  - VGG
  - ResNet
  - U-Net
  - Autoencoders

# Single-Shot Detection

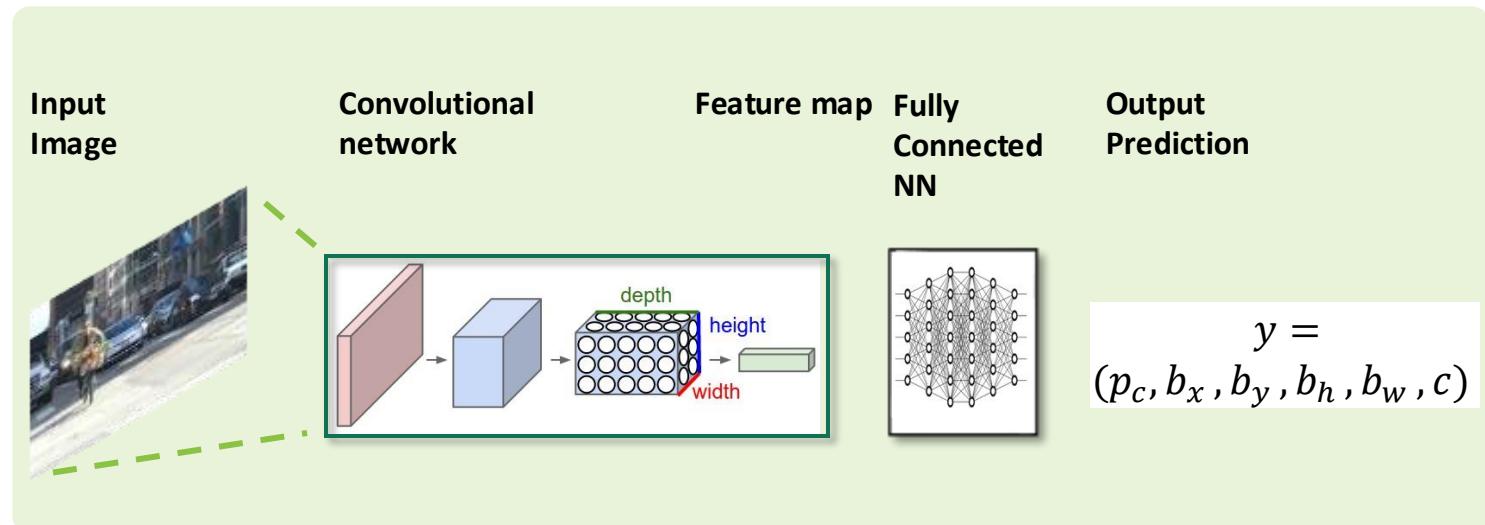
Single-Shot Detection methods do the region proposal, feature extraction and classification steps in one go



# Single-Shot Detection

## Single-Shot Detectors:

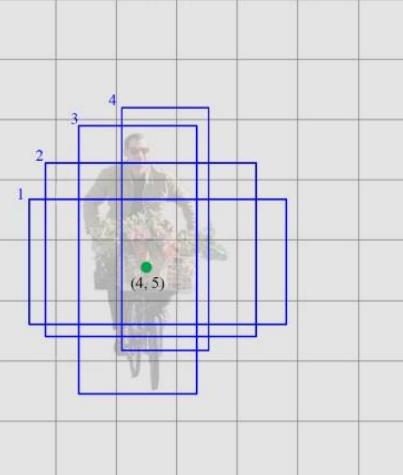
- Feature map is created for the whole image
- Prediction is a combination of bounding box coordinates and class confidence.



# Single-Shot Detection

Anchor-based single-shot detectors work with following steps:

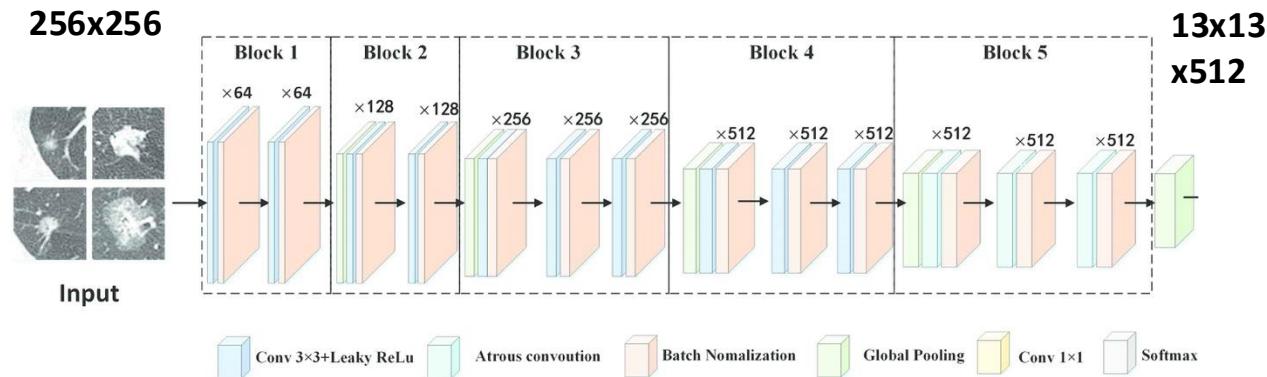
- A. Divide image in grid cells
- B. Pre-define anchor-boxes
- C. Refine box coordinates
- D. Train class confidence
- E. Filter best predictions



# Single-Shot Detection

## A. Divide image in grid cells

Use a convolutional neural network to divide image in grid cells and extract convolutional features for each grid cell

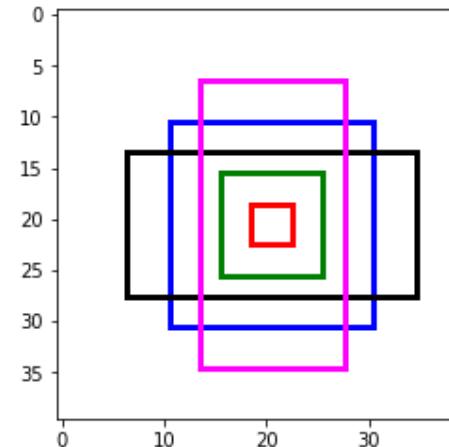
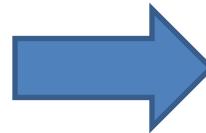
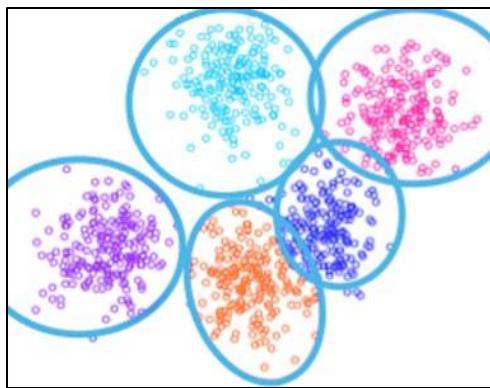


CNN typically reduces dimension in each convolutional block (by pooling)  
This example CNN leads to an Feature Map with 512 features (filter outcomes) for each of  $13 \times 13$  grid cells.

# Single-Shot Detection

## B. Pre-define anchor boxes

Define default anchor boxes with typical object sizes and aspect ratio's



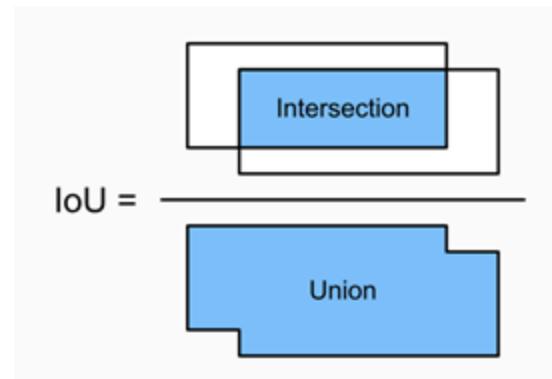
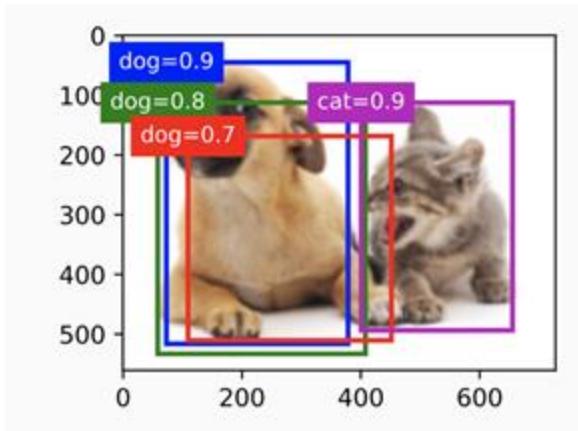
5 default anchor boxes are defined manually or based on k-means clustering on an existing image dataset. This step is usually done off-line (not included in CNN training)

# Single-Shot Detection

## C. Refine box coordinates

For each grid cell, a number of box sizes can be learned by adjusting the pre-defined anchor boxes.

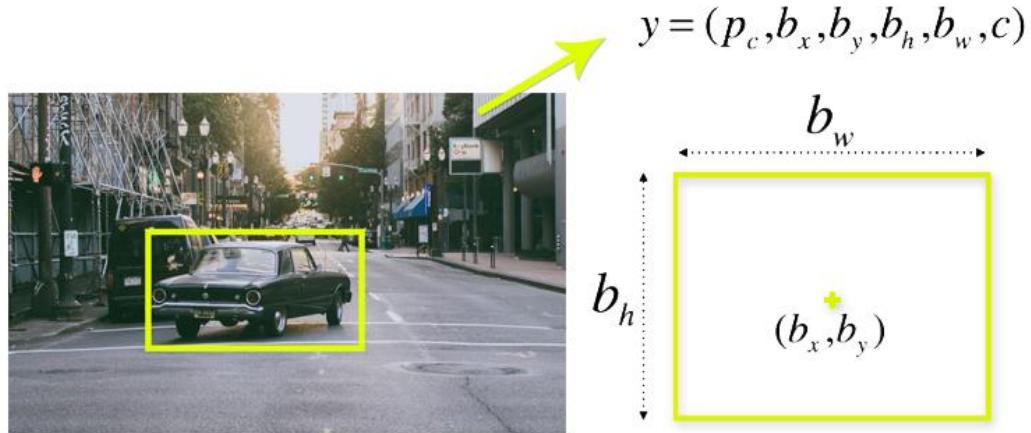
Learning is by minimizing Loss (Intersection over Union) between prediction and ground truth.



# Single-Shot Detection

## D. Train class confidence (jointly with step C)

Based on convolutional features for each box with an object, a class confidence can be learned.



Yolo example: class confidence is combination of  $c$  and  $p_c$

$b_x, b_y, b_h, b_w$  box coordinates

$c$  box probability that there is an object in box

$p_c$  conditional class probability for each class  $c$

# Single-Shot Detection

## D. Train class confidence (jointly with step C)

Location and class probability are trained simultaneously.  
 Loss function (regression) combined both.

$$\begin{aligned}
 & \lambda_{\text{coord}} \sum_{i=0}^{S^2} \sum_{j=0}^B \mathbb{1}_{ij}^{\text{obj}} \left[ (x_i - \hat{x}_i)^2 + (y_i - \hat{y}_i)^2 \right] \\
 & + \lambda_{\text{coord}} \sum_{i=0}^{S^2} \sum_{j=0}^B \mathbb{1}_{ij}^{\text{obj}} \left[ (\sqrt{w_i} - \sqrt{\hat{w}_i})^2 + (\sqrt{h_i} - \sqrt{\hat{h}_i})^2 \right] \\
 & + \sum_{i=0}^{S^2} \sum_{j=0}^B \mathbb{1}_{ij}^{\text{obj}} (C_i - \hat{C}_i)^2 \\
 & + \lambda_{\text{noobj}} \sum_{i=0}^{S^2} \sum_{j=0}^B \mathbb{1}_{ij}^{\text{noobj}} (C_i - \hat{C}_i)^2 \\
 & + \sum_{i=0}^{S^2} \mathbb{1}_i^{\text{obj}} \sum_{c \in \text{classes}} (p_i(c) - \hat{p}_i(c))^2
 \end{aligned}$$

Prediction      Ground truth

**regulators**

Box location

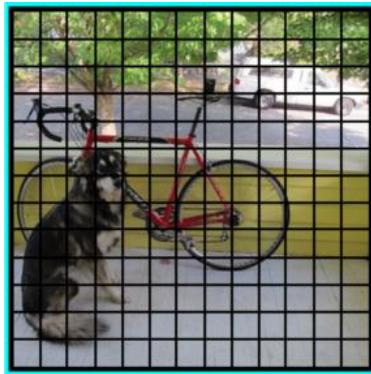
Box Confidence score (IoU)

conditional class probability  
for class c

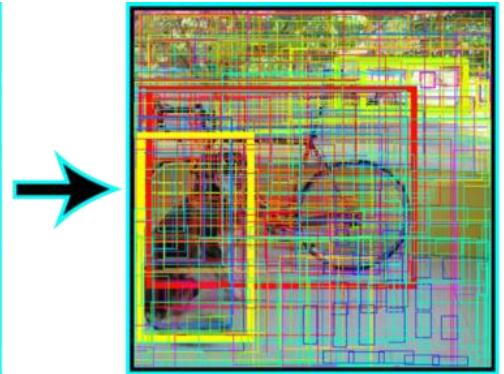
# Single-Shot Detection

## E. Filter best predictions

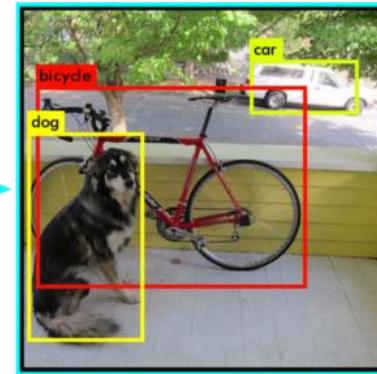
Non-maximum suppression filters out predictions with low class confidence. Filter out boxes with confidence below threshold.



**Image grid**  
 $S \times S$   
 $S = 13$



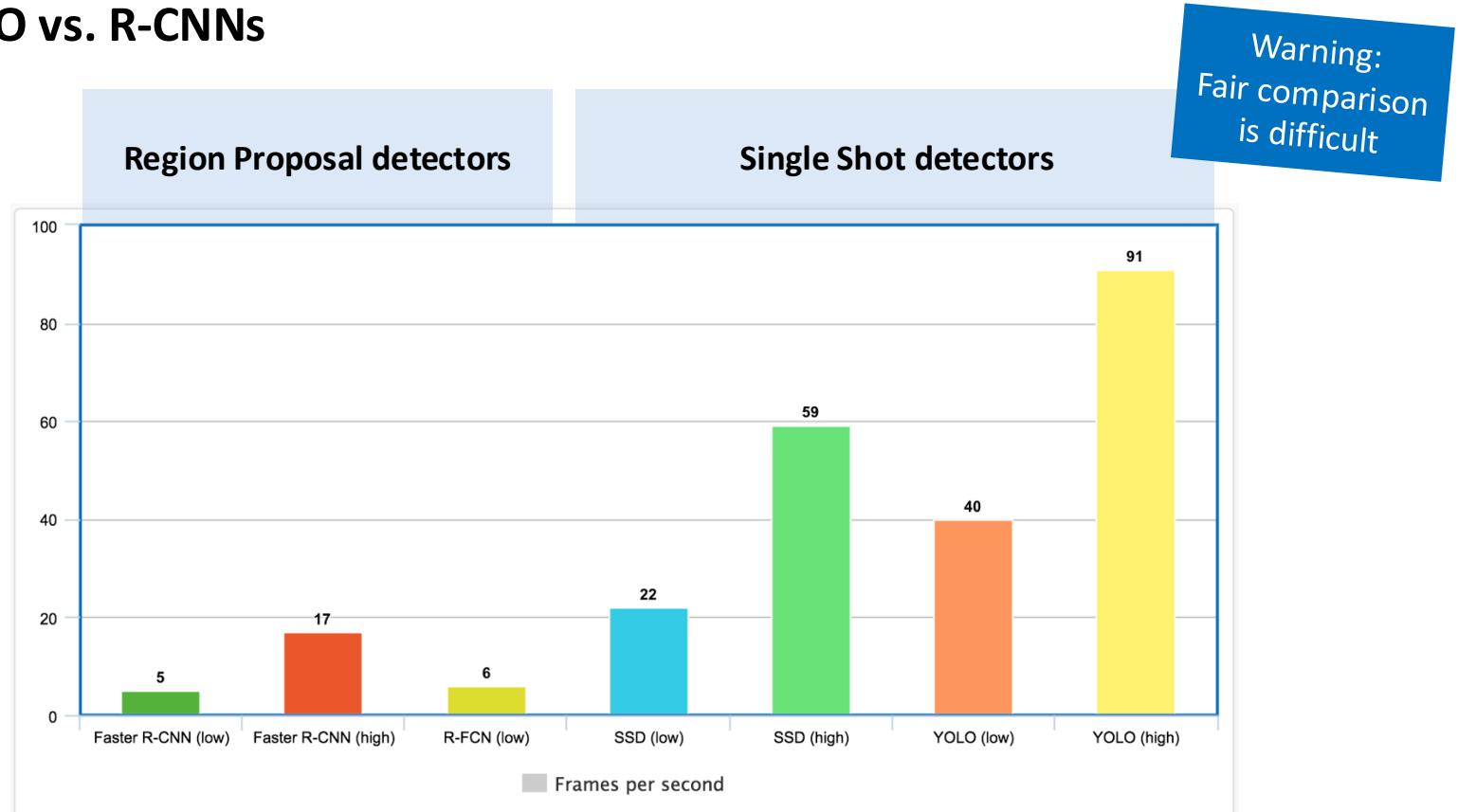
**Prediction**  
Every grid cell B boxes  
 $B = 5$   
  
So:  
 $13 \times 13 \times 5 = 845$  boxes



**Thresholding:**  
Box confidence score  $> 0.3$   
  
**Non-Maximum suppression**  
Class predictions with highest class confidence score

# Single-Shot Detection

## SSD / YOLO vs. R-CNNs

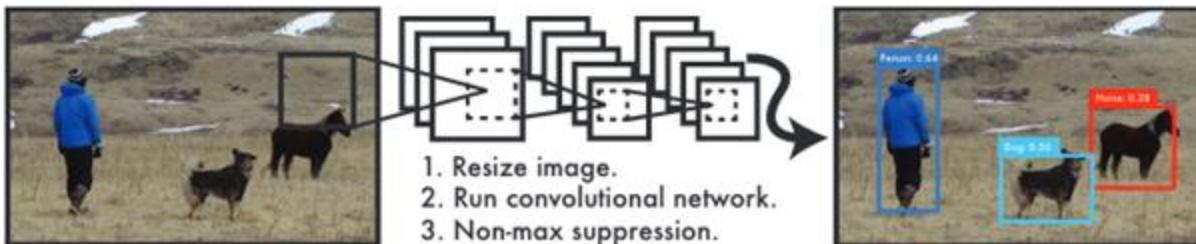


# You Only Look Once (YOLO)

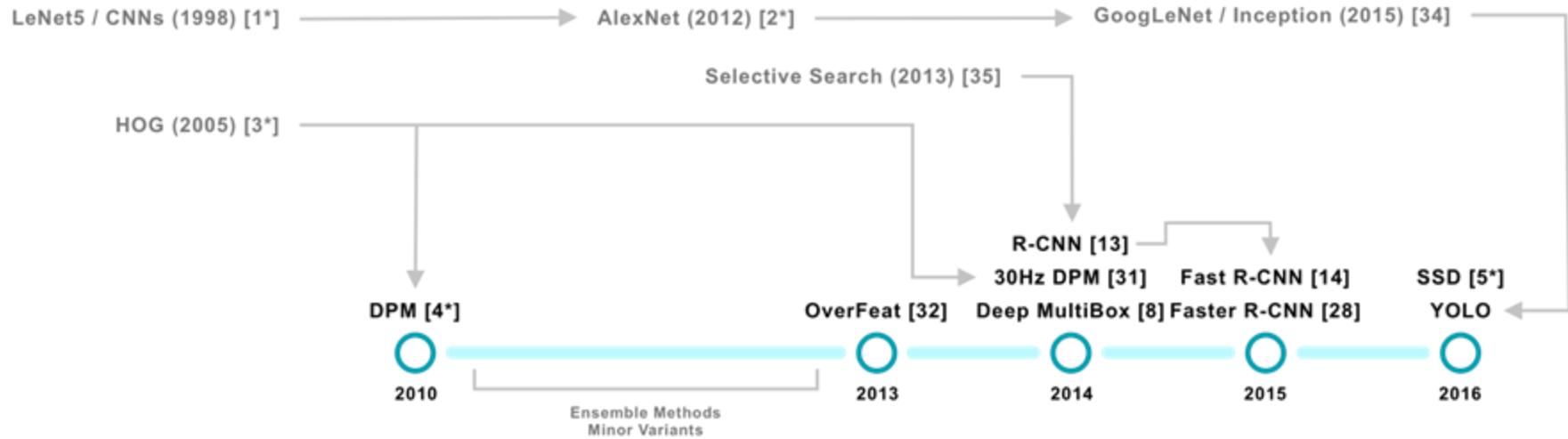
CV



- Simple architecture
- Single regression problem, from pixels to bounding box coordinates and class probabilities:  
**simpler to optimize**
- **FAST:** 45 FPS for YOLO and 150 FPS for YOLO-Fast
- Sees the global context, not just patches



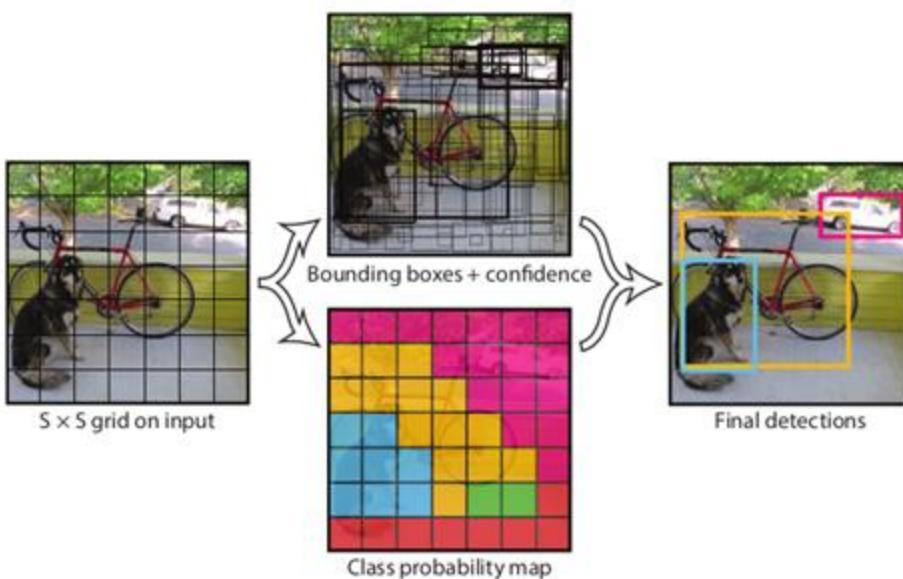
# You Only Look Once (YOLO)



Developed and released a novel object detection approach that:

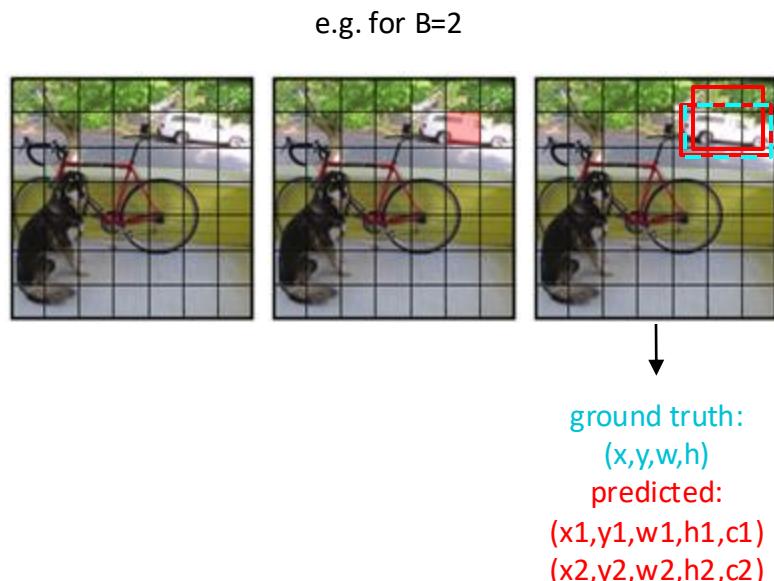
- Is 2.5 times faster than competing methods and 6.5 times faster than SoTA while sacrificing ~10% mAP
- Only consists of a single neural network which is easy to train
- Considers the global context

# YOLO – Unified Detection



- Divide image into  $S \times S$  grid.
- Each cell predicts:
  - **B** bounding boxes and confidence scores (*Localization*)
  - **C** conditional class probabilities (*Classification*).for the object whose center it roughly contains.

# YOLO - Localization



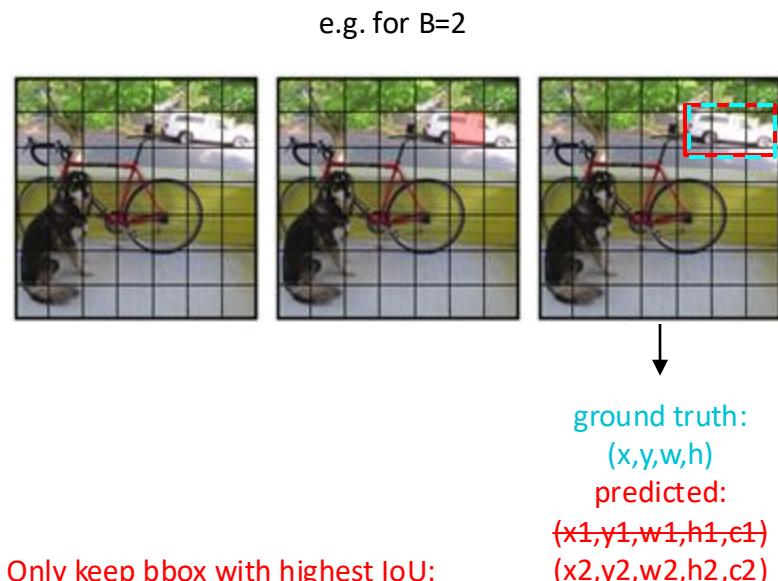
**B** bounding boxes and confidence scores.

$(x, y, w, h, c) \rightarrow B \times 5$  outputs

- $x, y$ : center of the b. box
- $w, h$ : dimensions of b. box
- $c$ : confidence

$$\text{Pr}(\text{Object}) * \text{IOU}_{\text{pred}}^{\text{truth}}$$

# YOLO - Localization



B bounding boxes and confidence scores.

$(x, y, w, h, c) \rightarrow B \times 5$  outputs

- $x, y$ : center of the b. box
- $w, h$ : dimensions of b. box
- $c$ : confidence

$$\text{Pr(Object)} * \text{IOU}_{\text{pred}}^{\text{truth}}$$

# YOLO - Localization

e.g. for B=2



ground truth:

-  
predicted:  
 $(x_1, y_1, w_1, h_1, 0)$

Ideally:

$$\text{confidence} = \begin{cases} \text{conf}=0 & \text{if object present} \\ \text{conf}=\text{IoU} & \text{if no object present} \end{cases}$$

**B** bounding boxes and confidence scores.

$(x, y, w, h, c) \rightarrow B \times 5$  outputs

- **x,y**: center of the b. box
- **w,h**: dimensions of b. box
- **c**: confidence

$\text{Pr}(\text{Object}) * \text{IOU}_{\text{pred}}^{\text{truth}}$

# YOLO - Classification



$$(p_{\text{dog}}, p_{\text{bike}}, p_{\text{car}}) = (0.1, 0.2, 0.8)$$

C conditional class probabilities.

$$(p_1, \dots, p_C) \rightarrow C \text{ outputs}$$

- $p_i$ :

$$\Pr(\text{Class}_i | \text{Object})$$

# YOLO - Confidence Score

$$\underbrace{\Pr(\text{Class}_i | \text{Object}) * \Pr(\text{Object})}_{\text{Classification output}} * \underbrace{\text{IOU}_{\text{pred}}^{\text{truth}}}_{\text{Localization output}} = \underbrace{\Pr(\text{Class}_i) * \text{IOU}_{\text{pred}}^{\text{truth}}}_{\text{Class-specific confidence score}}$$

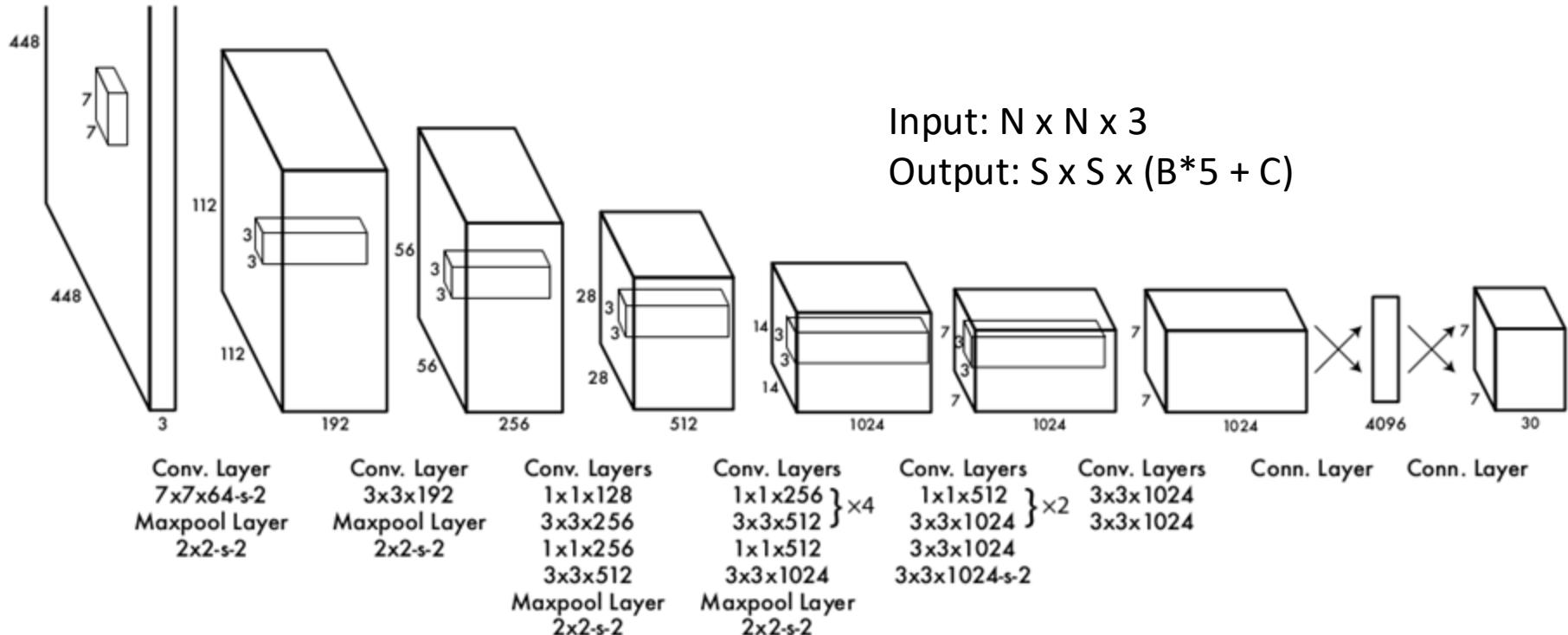
Classification  
output

Localization  
output

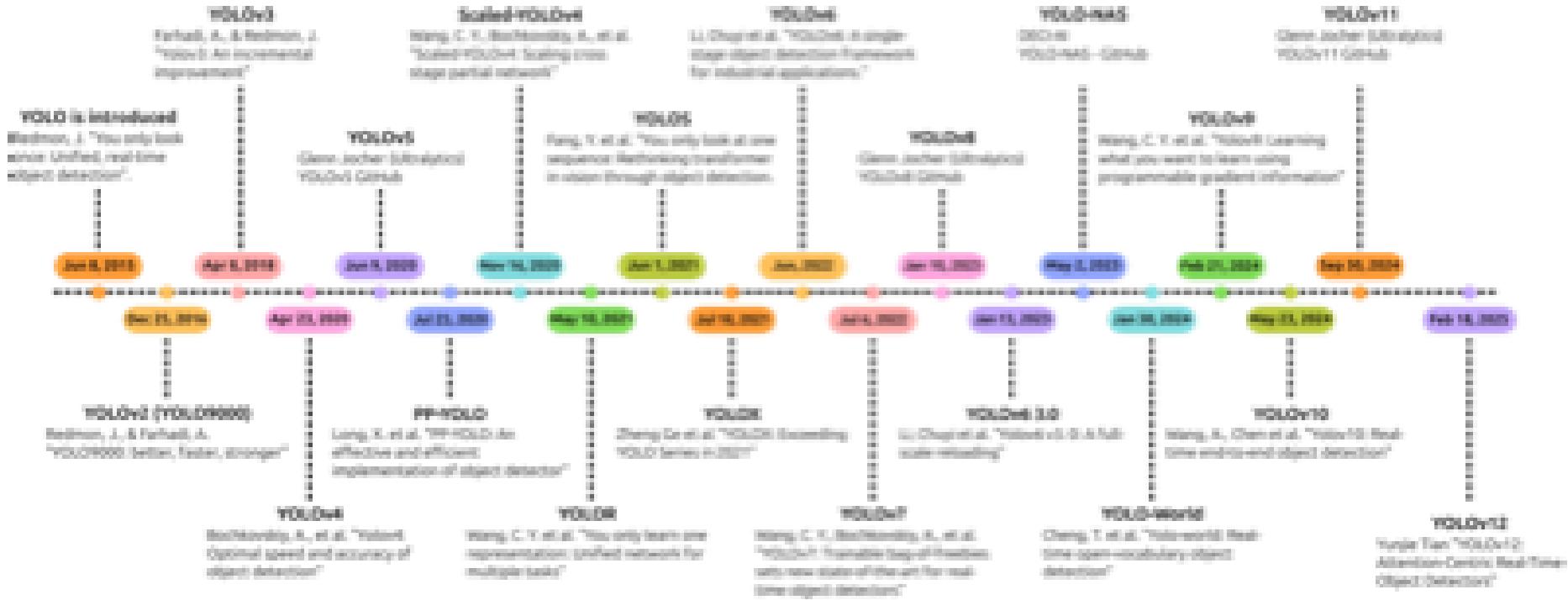
Class-specific  
confidence score



# YOLO - Architecture



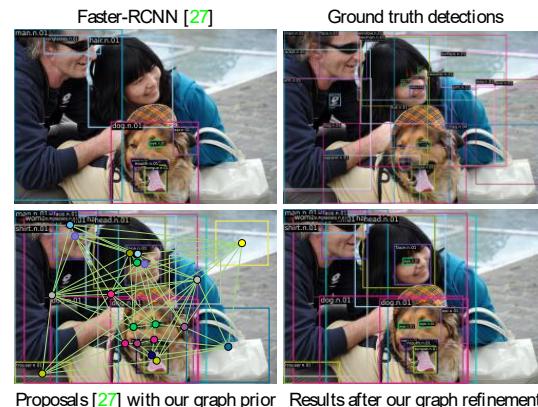
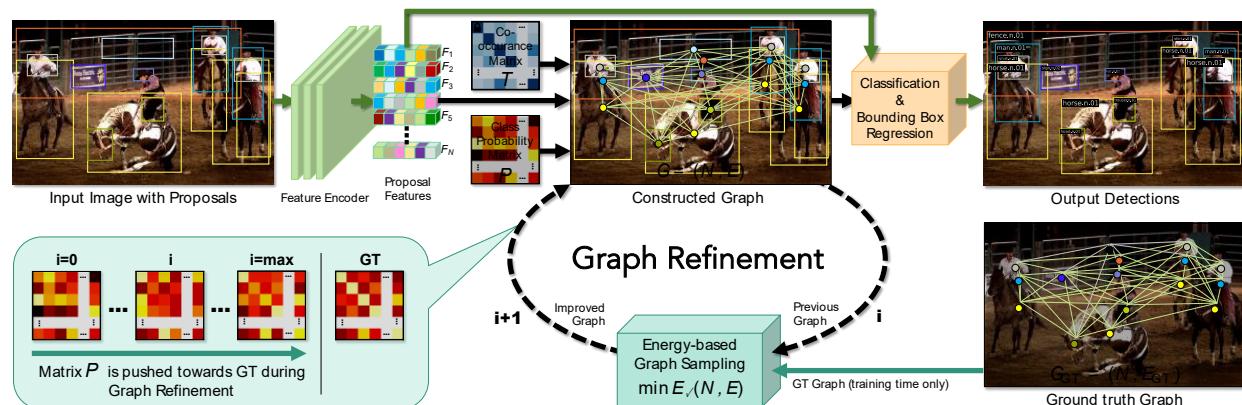
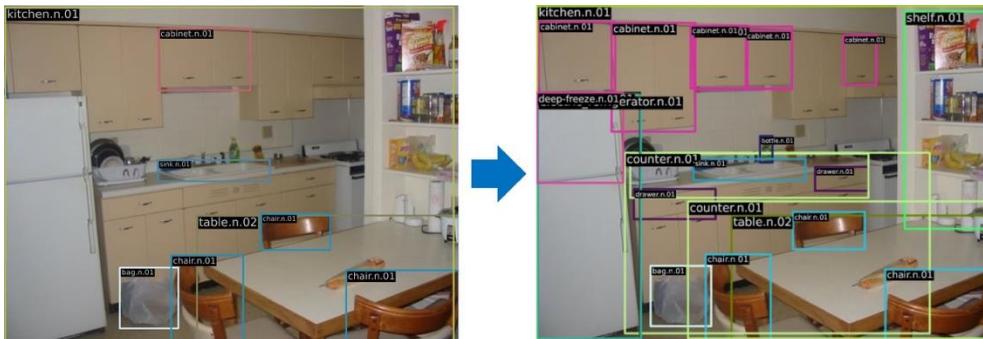
# YOLO - Architecture Timeline



# Recent Work: Object Detection @UvA

## Key ideas:

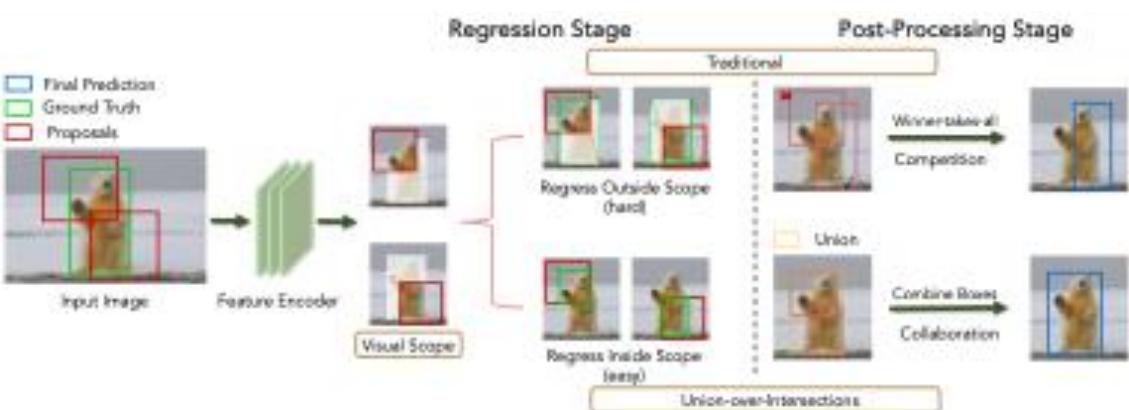
- Leverage context information about surrounding objects
- Context encoded via graph network of object proposals
- Learn & use co-occurrence statistics



# Recent Work: Object Detection @UvA

## Key ideas:

- Regress bbox size only within scope
- Combine resized bboxes for final result rather than selecting only one with NMS



Consistent improvements for various two stage detectors:

Method	Backbone	mAP $\uparrow$	AP $_{50}\uparrow$	AP $_{75}\uparrow$	AP $_{S}\uparrow$	AP $_{M}\uparrow$	AP $_{L}\uparrow$
Faster R-CNN (Ren et al., 2015)	R-50-fpn	37.4	58.1	40.4	21.2	41.0	48.1
Faster R-CNN w/ <i>IoU</i>	R-50-fpn	<b>38.1</b>	<b>58.7</b>	<b>40.9</b>	<b>21.7</b>	<b>41.8</b>	<b>49.5</b>
Faster R-CNN (Ren et al., 2015)	R-101-fpn	39.4	60.1	43.1	22.4	43.7	51.1
Faster R-CNN w/ <i>IoU</i>	R-101-fpn	<b>40.3</b>	<b>60.8</b>	<b>43.5</b>	<b>23.1</b>	<b>44.5</b>	<b>52.8</b>
Mask R-CNN (He et al., 2017)	R-50-fpn	38.2	58.8	41.4	21.9	40.9	49.5
Mask R-CNN w/ <i>IoU</i>	R-50-fpn	<b>38.8</b>	<b>59.6</b>	<b>41.8</b>	<b>22.2</b>	<b>41.6</b>	<b>50.9</b>
Mask R-CNN (He et al., 2017)	R-101-fpn	39.8	60.3	43.4	23.1	43.8	52.5
Mask R-CNN w/ <i>IoU</i>	R-101-fpn	<b>40.9</b>	<b>61.1</b>	<b>44.0</b>	<b>23.5</b>	<b>44.7</b>	<b>54.6</b>
Cascade R-CNN (Cai & Vasconcelos, 2018)	R-101-fpn	42.5	60.7	46.4	23.5	46.5	56.4
Cascade R-CNN w/ <i>IoU</i>	R-101-fpn	<b>43.1</b>	<b>61.9</b>	<b>46.8</b>	<b>24.0</b>	<b>47.2</b>	<b>57.3</b>
YOLOv3 (Redmon & Farhadji, 2018)	DarkNet-53	33.7	56.6	35.3	19.4	36.8	44.3
YOLOv3 w/ <i>IoU</i>	DarkNet-53	<b>34.5</b>	<b>57.5</b>	<b>35.9</b>	<b>19.9</b>	<b>37.3</b>	<b>45.2</b>
Def-DETR (Zhu et al., 2021)	R-50-fpn	44.3	63.2	48.6	26.8	47.7	58.8
Def-DETR w/ <i>IoU</i>	R-50-fpn	<b>44.8</b>	<b>63.9</b>	<b>49.1</b>	<b>27.2</b>	<b>48.3</b>	<b>59.8</b>



# Today's Agenda

---

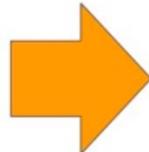
- Learning region proposals - Faster R-CNN
- Feature Pyramid Network (FPN)
- Single Shot Detectors (SSD) / Yolo
- **Transfer learning**
- Network architectures
  - VGG
  - ResNet
  - U-Net
  - Autoencoders

# Classification Model Size

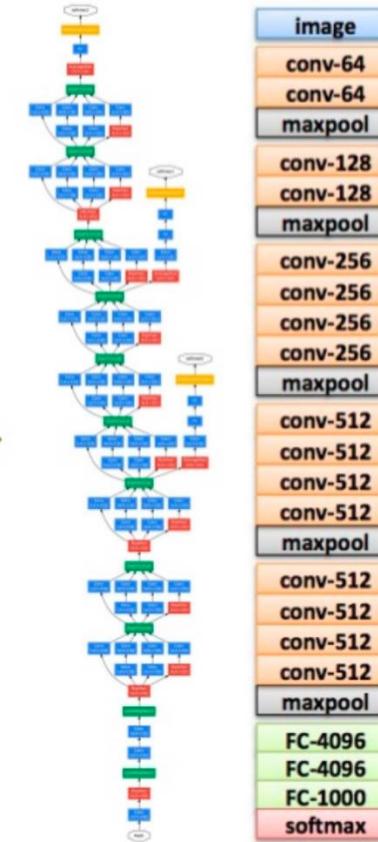
## ImageNet Challenge: 2014

AlexNet

image
conv-64
conv-192
conv-384
conv-256
conv-256
FC-4096
FC-4096
FC-1000



GoogLeNet

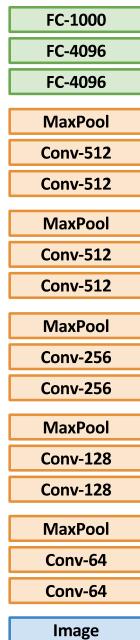


- Model size steadily increases
- Expensive training on image net
- Idea: Re-use training results

# Transfer Learning

## Transfer Learning

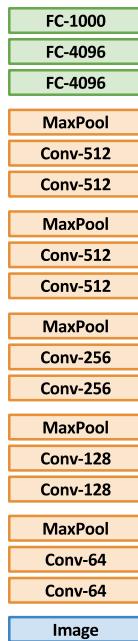
### 1. Train on ImageNet



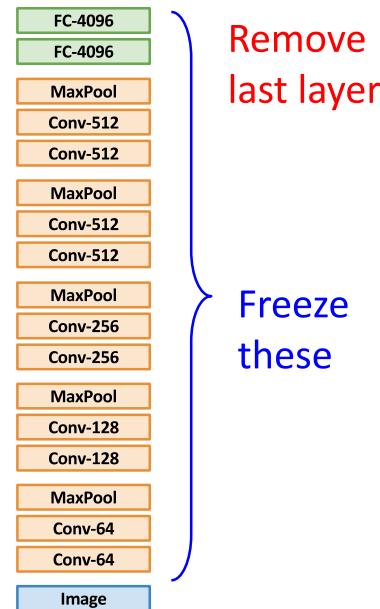
# Transfer Learning

## Transfer Learning: Feature Extraction

1. Train on ImageNet



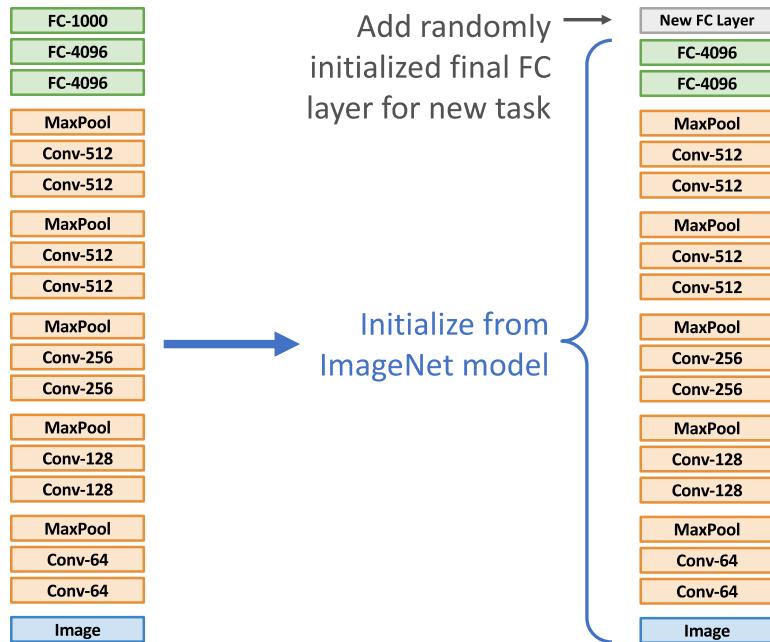
2. Extract features with CNN, train linear model



# Transfer Learning

## Transfer Learning: Fine-Tuning

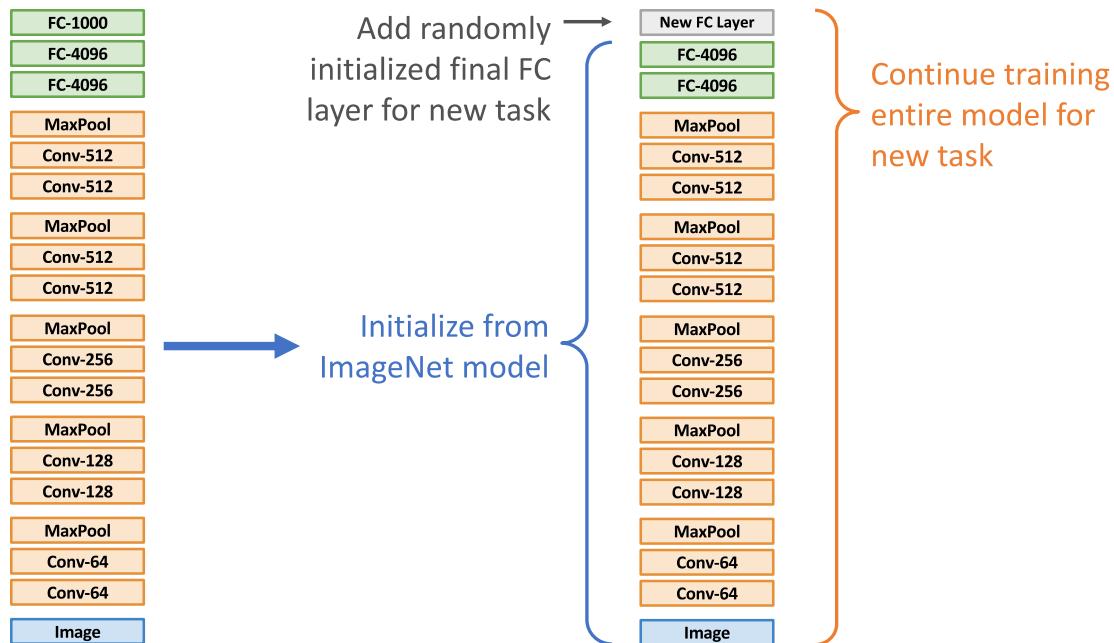
### 1. Train on ImageNet



# Transfer Learning

## Transfer Learning: Fine-Tuning

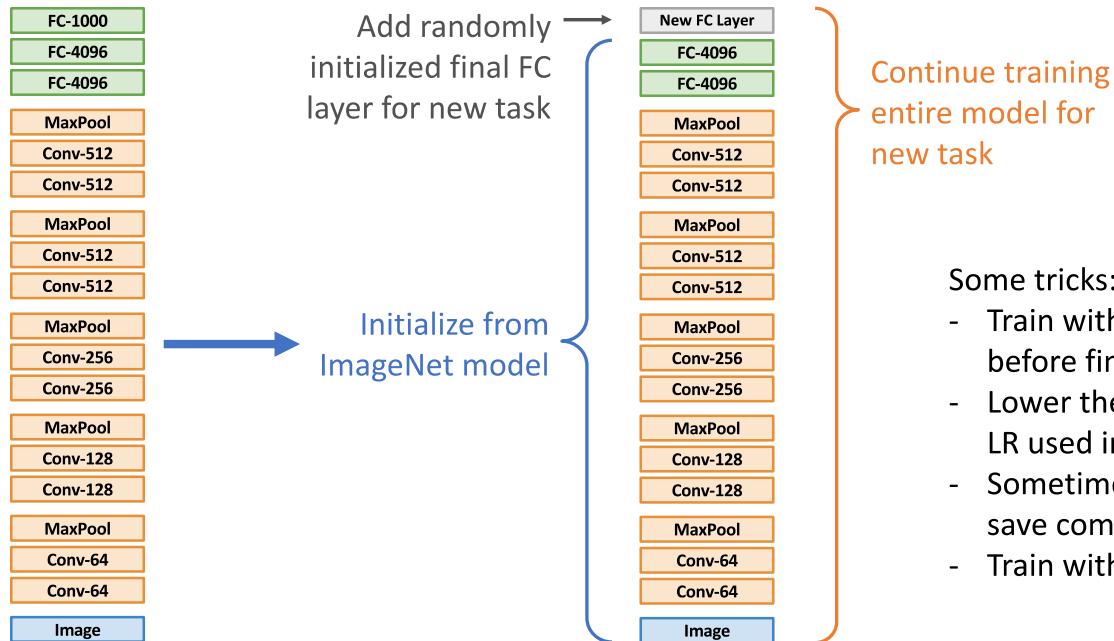
### 1. Train on ImageNet



# Transfer Learning

## Transfer Learning: Fine-Tuning

### 1. Train on ImageNet



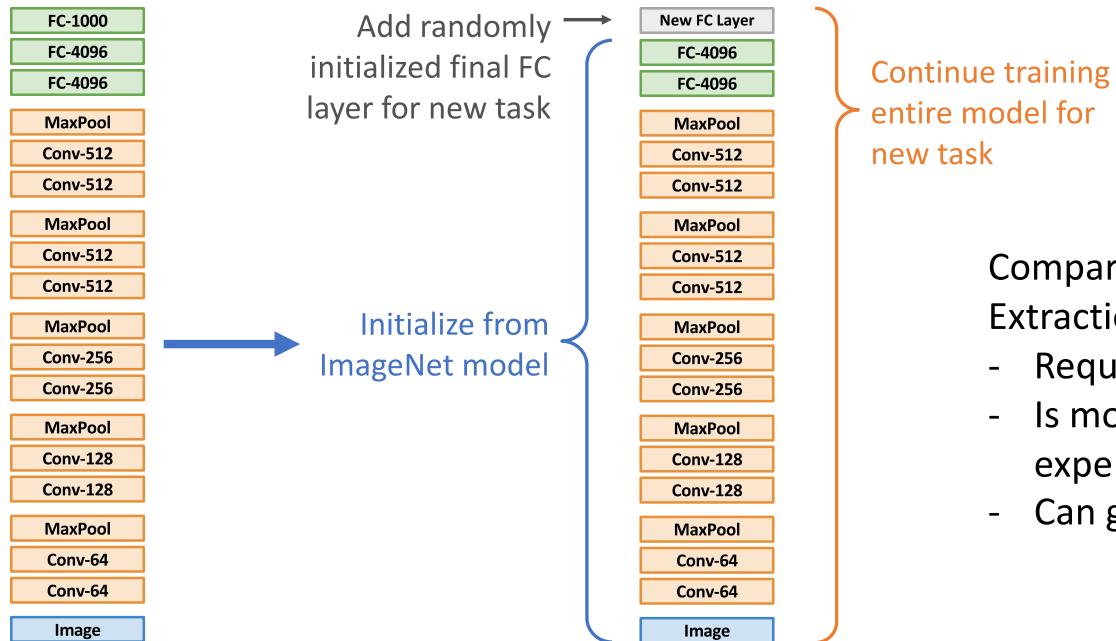
Some tricks:

- Train with feature extraction first before fine-tuning
- Lower the learning rate: use  $\sim 1/10$  of LR used in original training
- Sometimes freeze lower layers to save computation
- Train with BatchNorm in “test” mode

# Transfer Learning

## Transfer Learning: Fine-Tuning

### 1. Train on ImageNet

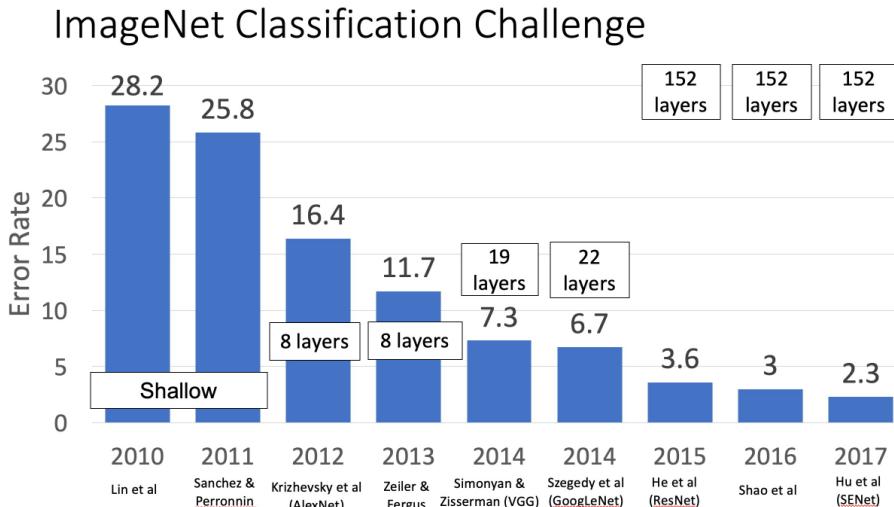


Compared with Feature Extraction, Fine-Tuning:

- Requires more data
- Is more computationally expensive
- Can give higher accuracies

# Transfer Learning

## Transfer Learning: Architecture Matters!

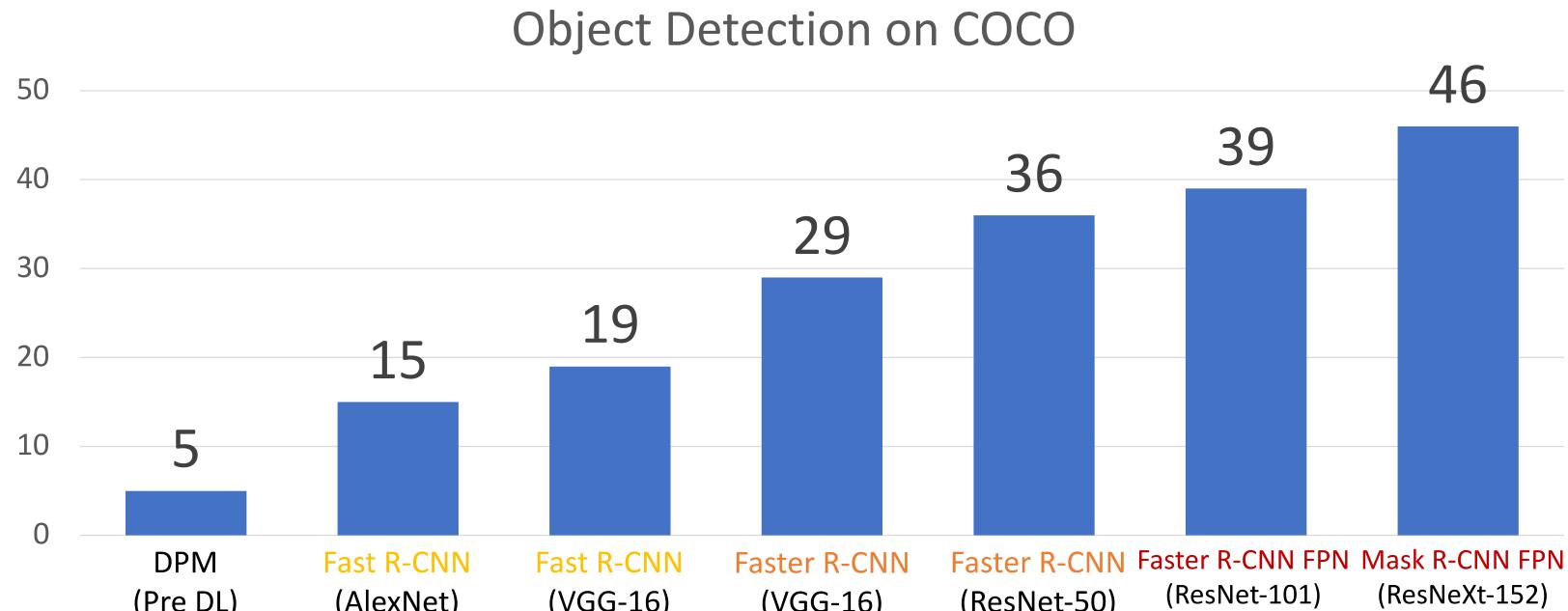


Improvements in CNN architectures lead to improvements in many downstream tasks thanks to transfer learning!



# Transfer Learning

Transfer Learning: Architecture Matters!

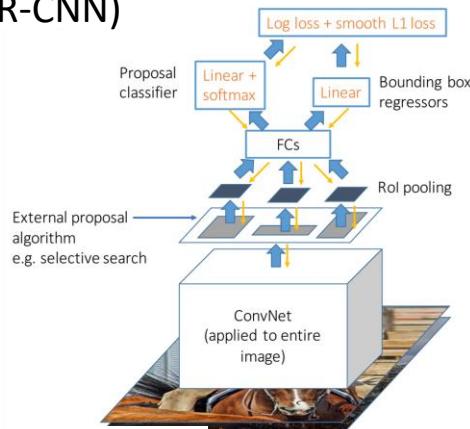


Ross Girshick, "The Generalized R-CNN Framework for Object Detection", ICCV 2017 Tutorial on Instance-Level Visual Recognition

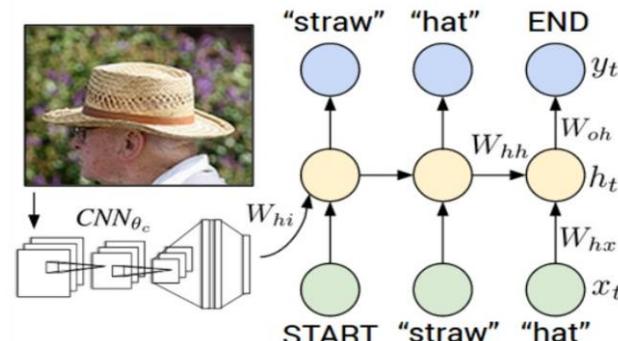
# Transfer Learning

Transfer learning is pervasive!  
It's the norm, not the exception

## Object Detection (Fast R-CNN)



Girshick, "Fast R-CNN", ICCV 2015  
Figure copyright Ross Girshick, 2015. Reproduced with permission.

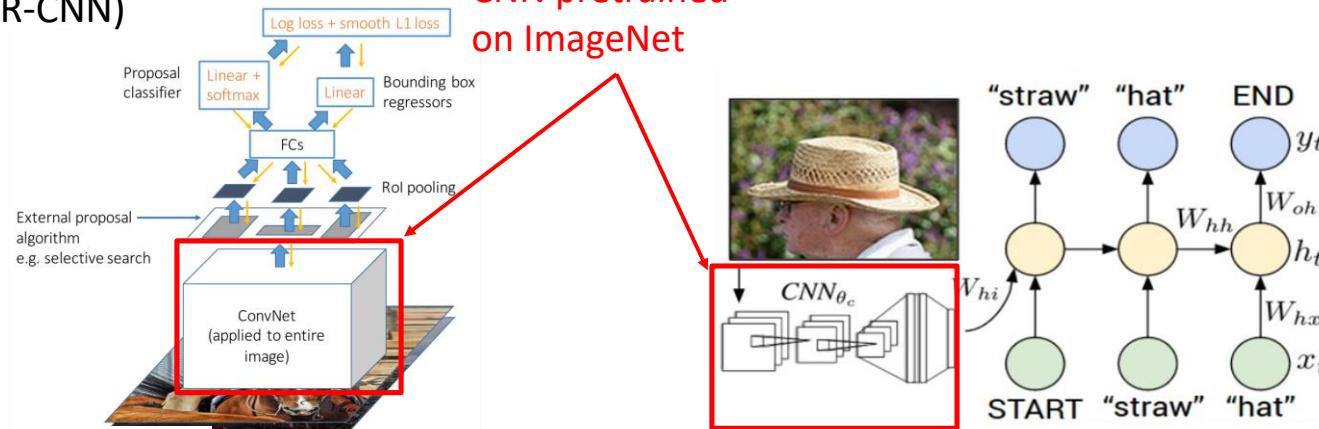


Karpathy and Fei-Fei, "Deep Visual-Semantic Alignments for Generating Image Descriptions", CVPR 2015

# Transfer Learning

Transfer learning is pervasive!  
It's the norm, not the exception

Object Detection  
(Fast R-CNN)



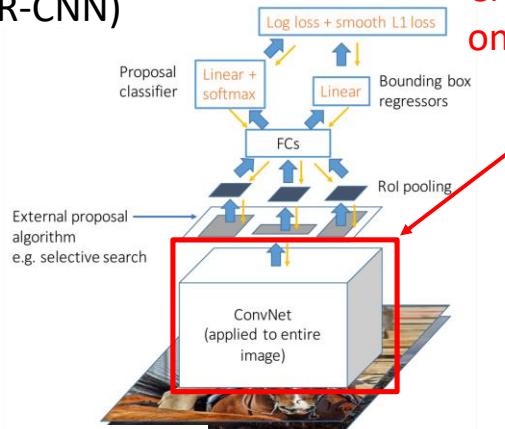
Girshick, "Fast R-CNN", ICCV 2015  
Figure copyright Ross Girshick, 2015. Reproduced with permission.

Karpathy and Fei-Fei, "Deep Visual-Semantic Alignments for Generating Image Descriptions", CVPR 2015

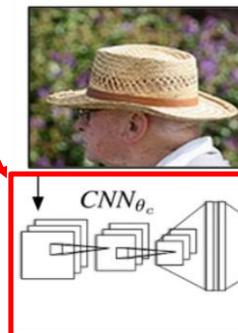
# Transfer Learning

Transfer learning is pervasive!  
It's the norm, not the exception

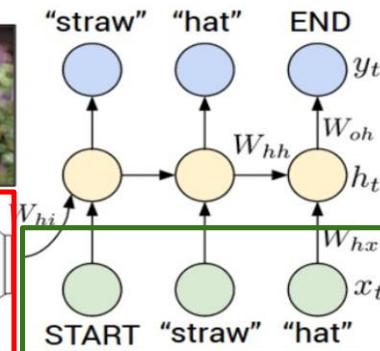
Object Detection  
(Fast R-CNN)



CNN pretrained  
on ImageNet



Word vectors pretrained  
with word2vec



Karpathy and Fei-Fei, "Deep Visual-Semantic Alignments  
for Generating Image Descriptions", CVPR 2015

Girshick, "Fast R-CNN", ICCV 2015  
Figure copyright Ross Girshick, 2015. Reproduced with permission.

# Today's Agenda

---

- Learning region proposals - Faster R-CNN
- Feature Pyramid Network (FPN)
- Single Shot Detectors (SSD) / Yolo
- Transfer learning
- Network architectures
  - VGG
  - ResNet
  - U-Net
  - Autoencoders

## VGG: Deeper Networks, Regular Design

### VGG Design rules:

All conv are 3x3 stride 1 pad 1

All max pool are 2x2 stride 2

After pool, double #channels

Network has 5 convolutional stages:

Stage 1: conv-conv-pool

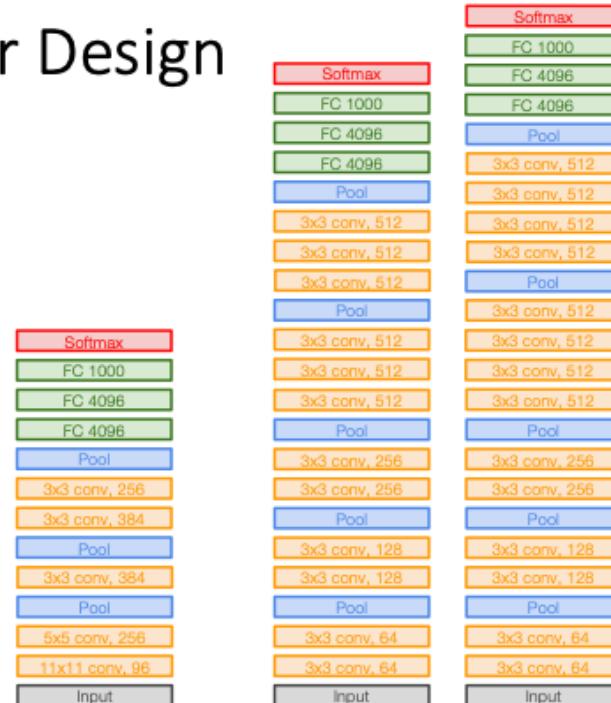
Stage 2: conv-conv-pool

Stage 3: conv-conv-pool

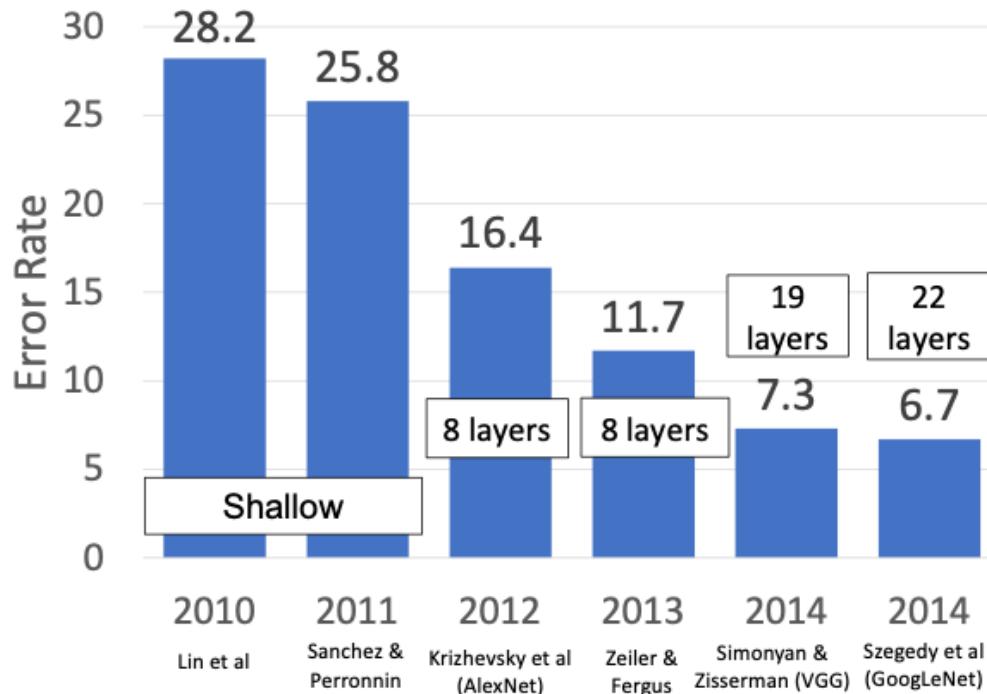
Stage 4: conv-conv-conv-[conv]-pool

Stage 5: conv-conv-conv-[conv]-pool

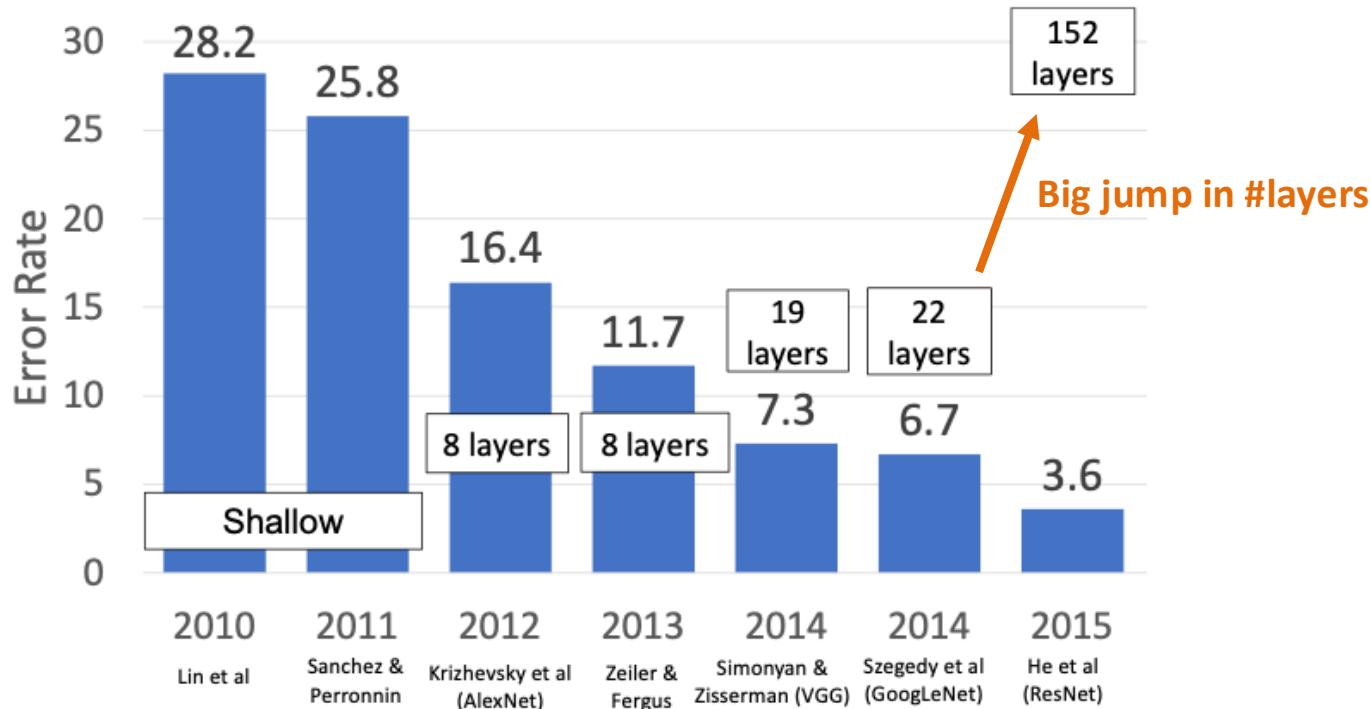
(VGG-19 has 4 conv in stages 4 and 5)



# ImageNet Classification Challenge



# ImageNet Classification Challenge



# ResNet (ILSVRC 2015 winner)

AlexNet, 8 layers  
(ILSVRC 2012)



VGG, 19 layers  
(ILSVRC 2014)



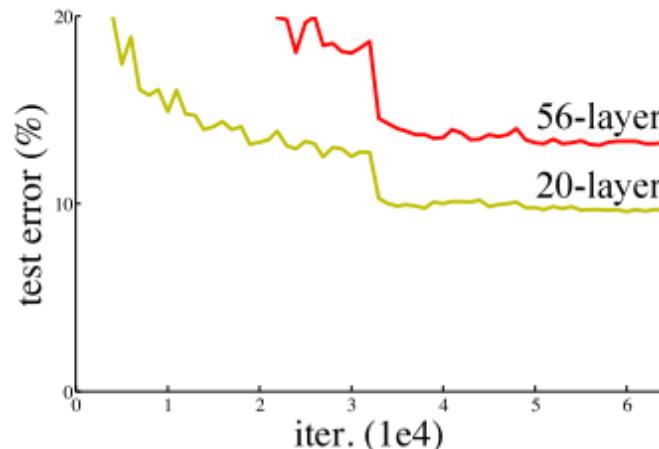
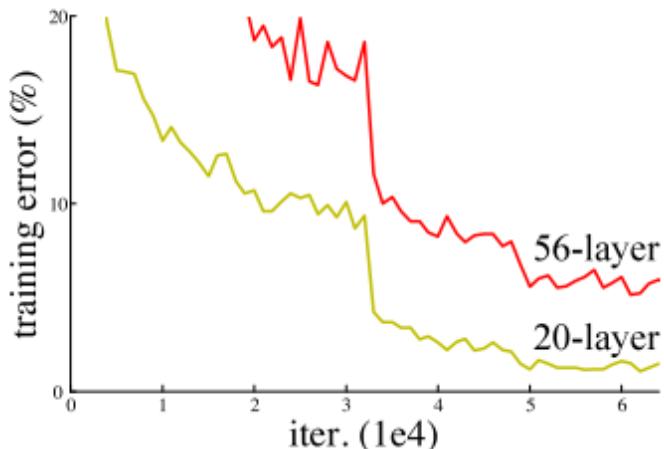
ResNet, **152 layers**  
(ILSVRC 2015)



# ResNet – Motivation: Vanishing Gradients

Deeper networks tend to saturate. Problem comes from vanishing gradients.

Deeper networks were underfitting rather than overfitting due to limitations of the optimization.



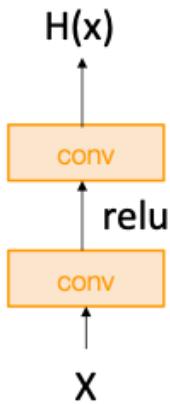
- A deeper model can emulate a shallower model: copy layers from shallower model, set extra layers to identity.

Thus deeper models should do at least as good as shallow models.

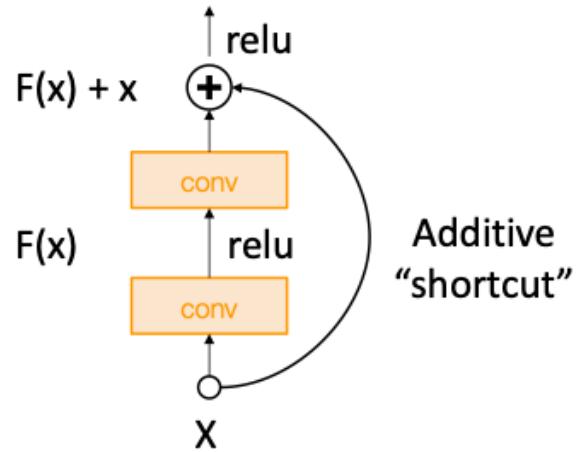
- **Hypothesis:** This is an optimization problem. Deeper models are harder to optimize, and in particular don't learn identity functions to emulate shallow models.
- **Solution:** Change the network so learning identity functions with extra layers is easy!

# ResNet

**Solution:** Change the network so learning identity functions with extra layers is easy!



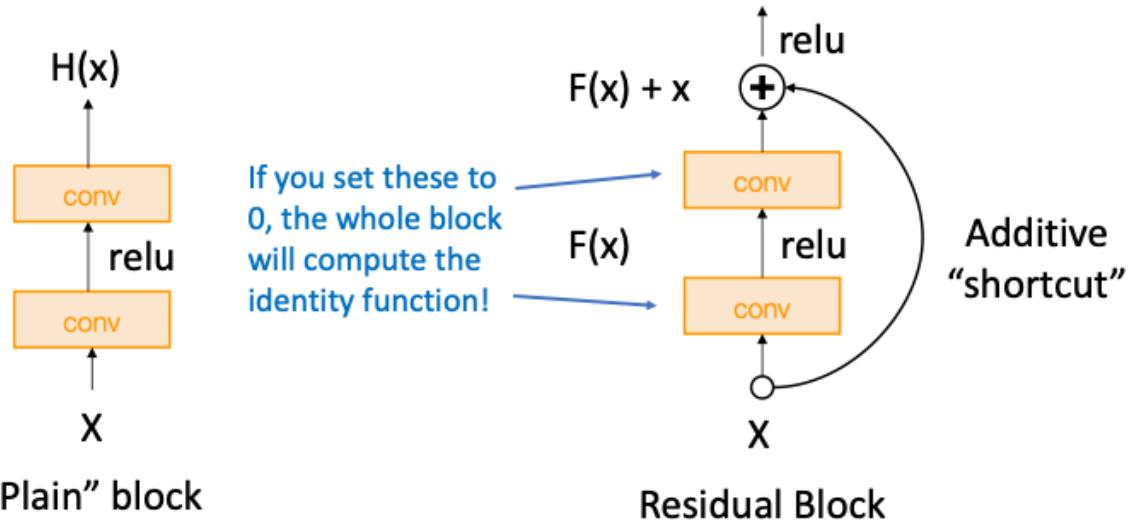
“Plain” block



Residual Block

# ResNet

**Solution:** Change the network so learning identity functions with extra layers is easy!

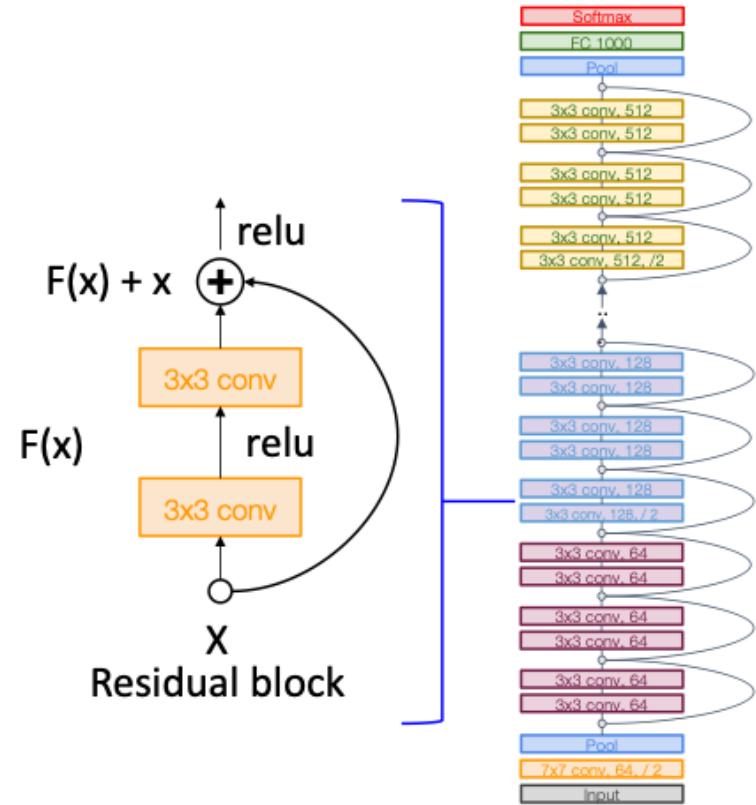


# ResNet

A residual network is a stack of many residual blocks

Regular design, like VGG: each residual block has two 3x3 conv

Network is divided into **stages**: the first block of each stage halves the resolution (with stride-2 conv) and doubles the number of channels



# ResNet

## ResNet-18:

Stem: 1 conv layer

Stage 1 (C=64): 2 res. block = 4 conv

Stage 2 (C=128): 2 res. block = 4 conv

Stage 3 (C=256): 2 res. block = 4 conv

Stage 4 (C=512): 2 res. block = 4 conv

Linear

ImageNet top-5 error: 10.92

GFLOP: 1.8

## ResNet-34:

Stem: 1 conv layer

Stage 1: 3 res. block = 6 conv

Stage 2: 4 res. block = 8 conv

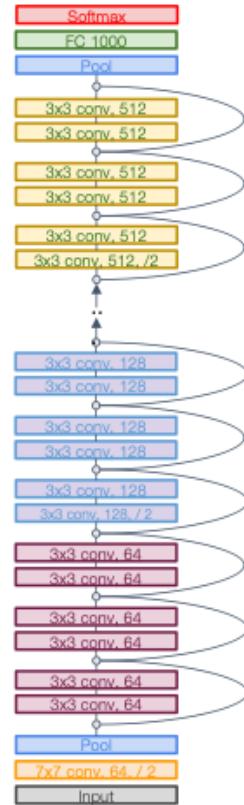
Stage 3: 6 res. block = 12 conv

Stage 4: 3 res. block = 6 conv

Linear

ImageNet top-5 error: 8.58

GFLOP: 3.6



- ResNets are collections of many paths of different length, and shorter paths predominantly contribute to training

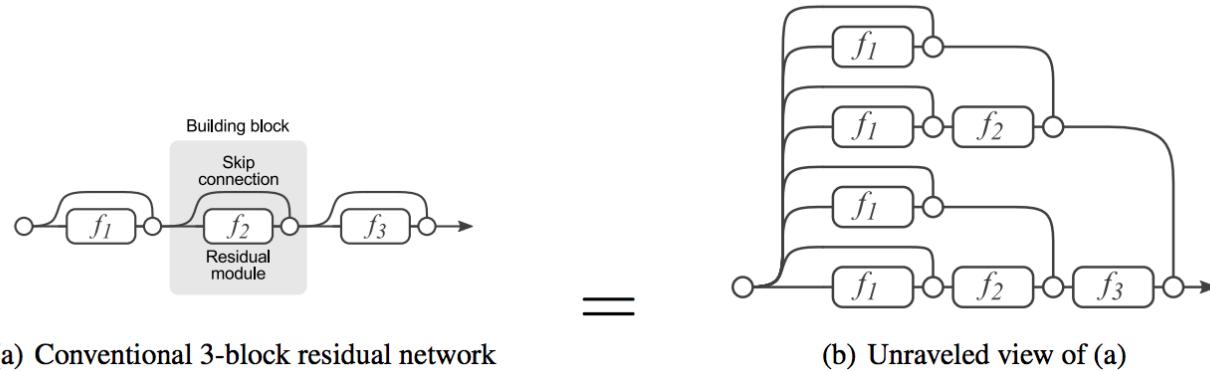
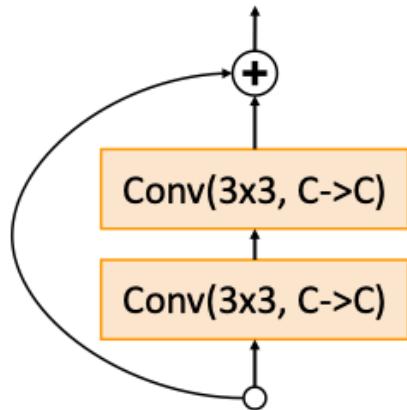


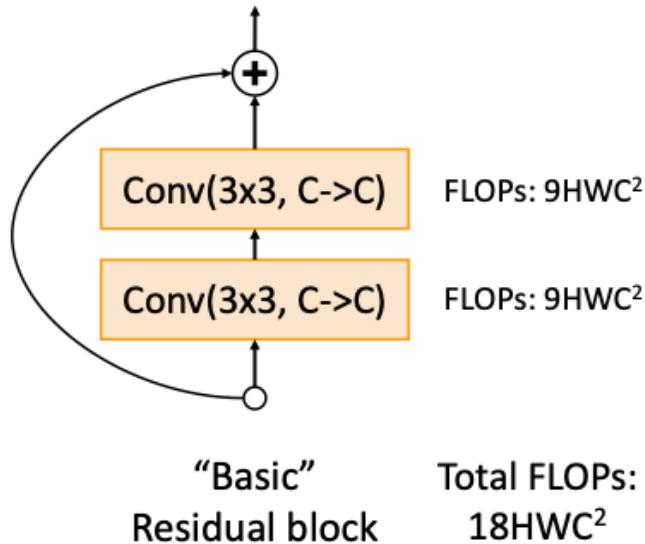
Figure 1: Residual Networks are conventionally shown as (a), which is a natural representation of Equation (1). When we expand this formulation to Equation (6), we obtain an *unraveled view* of a 3-block residual network (b). Circular nodes represent additions. From this view, it is apparent that residual networks have  $O(2^n)$  implicit paths connecting input and output and that adding a block doubles the number of paths.

# ResNet

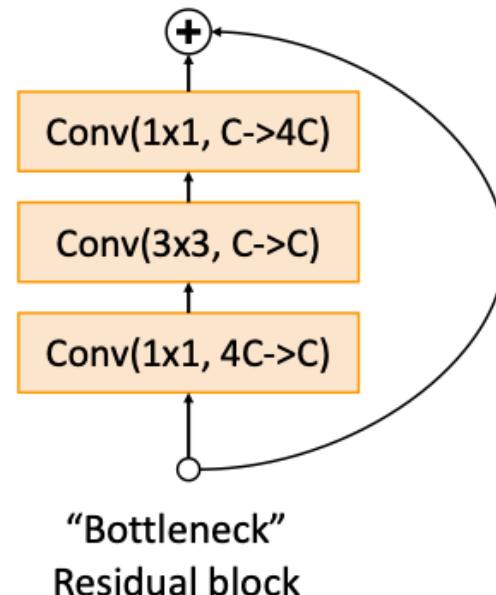
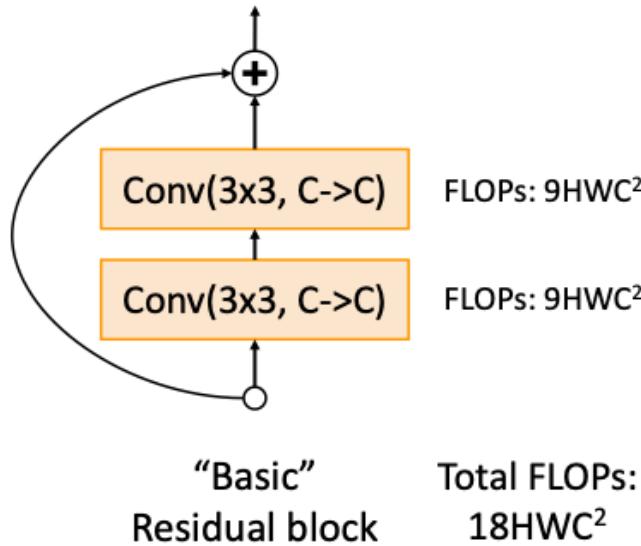


“Basic”  
Residual block

# ResNet

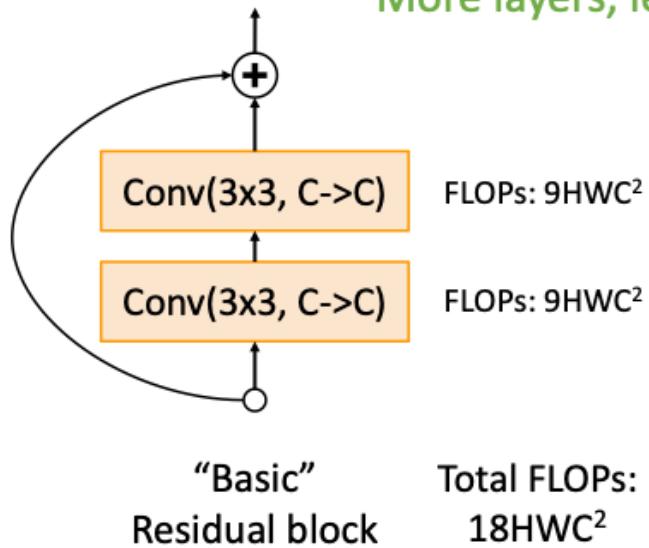


# ResNet



# ResNet

More layers, less computational cost!

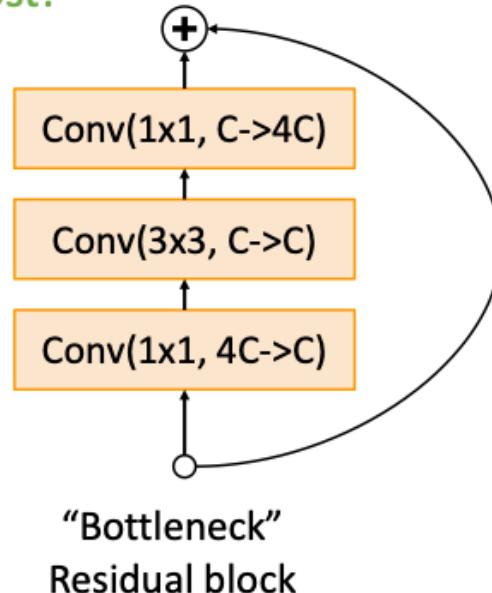


FLOPs:  $4HWC^2$

FLOPs:  $9HWC^2$

FLOPs:  $4HWC^2$

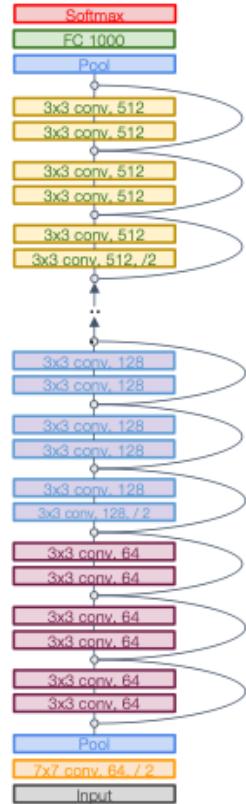
Total FLOPs:  
 $17HWC^2$



## ResNet

Deeper ResNet-101 and ResNet-152 models are more accurate, but also more computationally heavy

			Stage 1		Stage 2		Stage 3		Stage 4				
	Block type	Stem layers	Blocks	Layers	Blocks	Layers	Blocks	Layers	Blocks	Layers	FC layers	GFLOP	ImageNet top-5 error
ResNet-18	Basic	1	2	4	2	4	2	4	2	4	1	1.8	10.92
ResNet-34	Basic	1	3	6	4	8	6	12	3	6	1	3.6	8.58
ResNet-50	Bottle	1	3	9	4	12	6	18	3	9	1	3.8	7.13
ResNet-101	Bottle	1	3	9	4	12	23	69	3	9	1	7.6	6.44
ResNet-152	Bottle	1	3	9	8	24	36	108	3	9	1	11.3	5.94



# ResNet

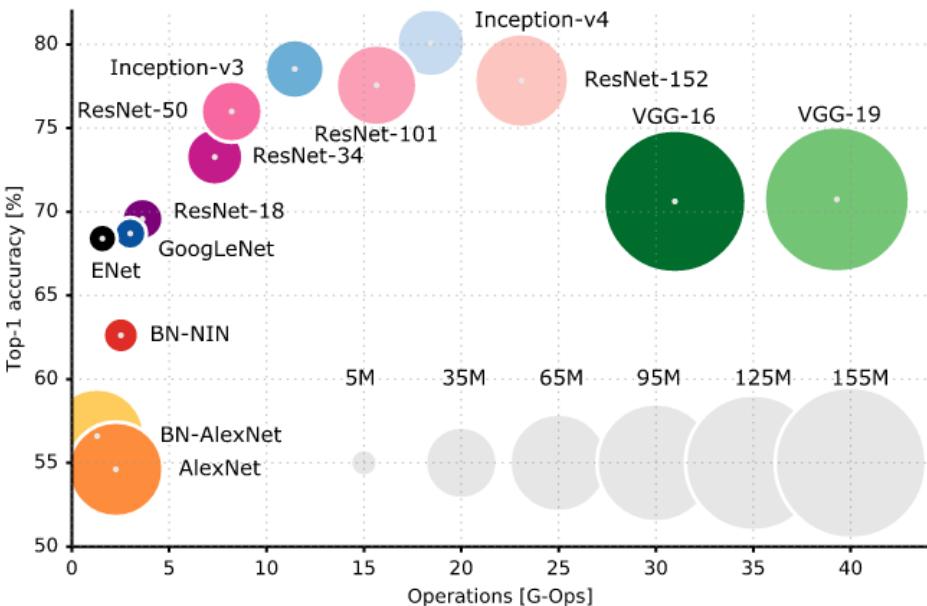
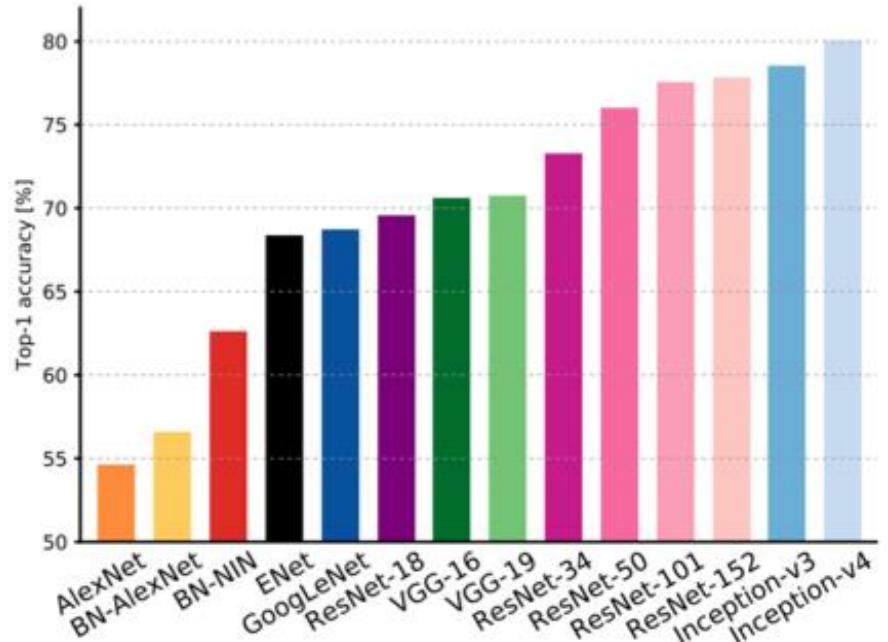
- Able to train very deep networks
- Deeper networks do better than shallow networks (as expected)
- Swept 1st place in all ILSVRC and COCO 2015 competitions
- Still widely used today!

## 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

# Comparing Architectures



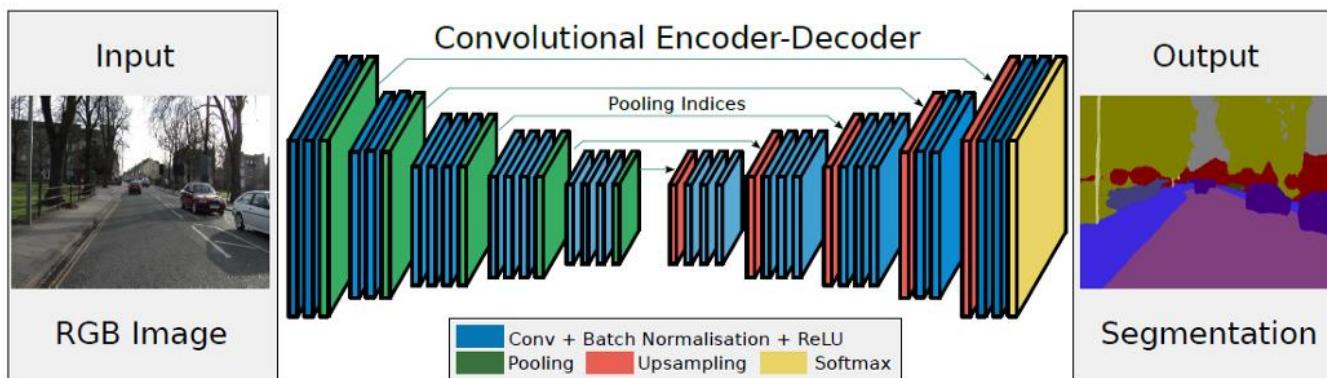
# Encoder-Decoder Architecture

Two common use-cases:

- Pixel wise predictions - output is of same dimension of input
  - semantic segmentation (e.g. Segnet, U-Net)
- Auto-encoder for
  - unsupervised feature learning
  - unsupervised pre-training

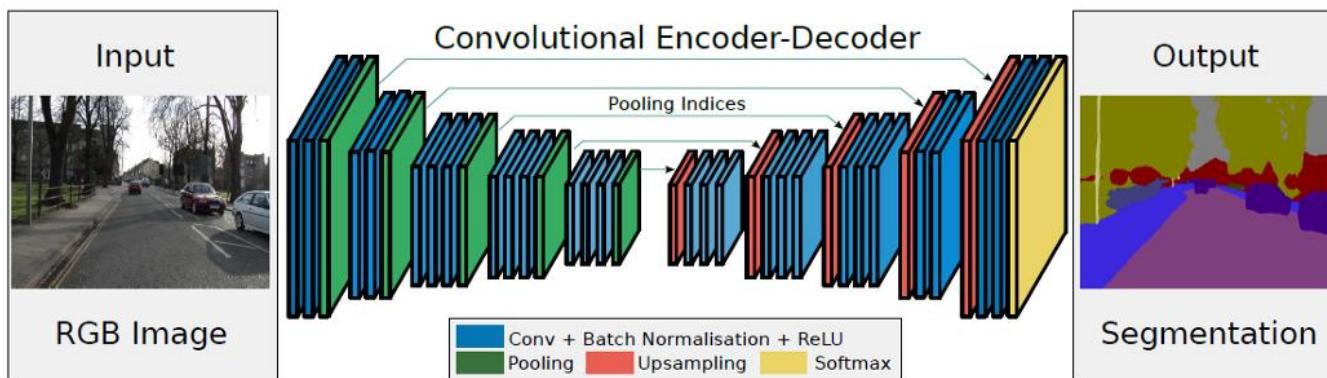
# Encoder-Decoder Architecture - SegNet

Segnet uses an encoder-decoder architecture for segmantic segmentation.



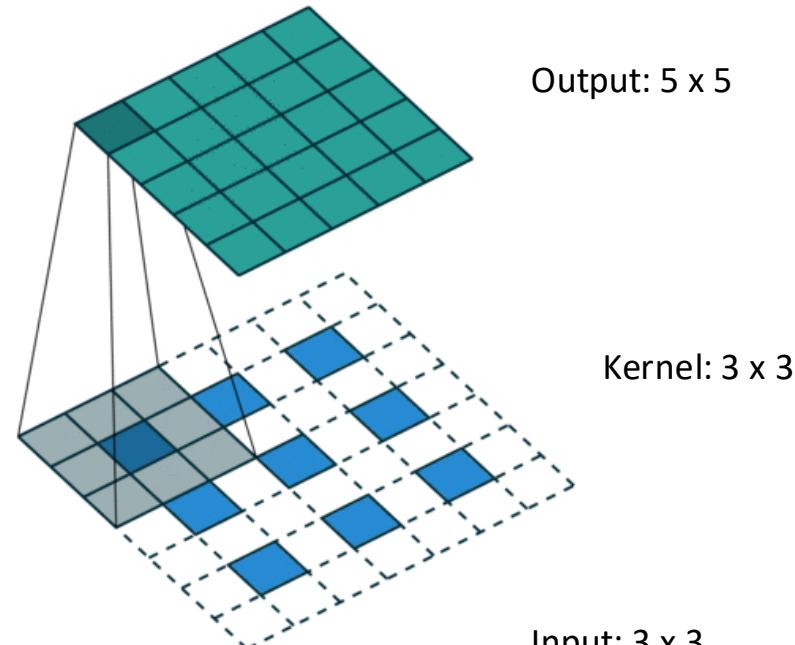
# Encoder-Decoder Architecture - SegNet

- Encoder is ‘normal’ convolution + pooling.
- Decoder uses upsampling + convolution
- (Learned) convolution is to refine the upsampling



# SegNet: Transposed Convolution

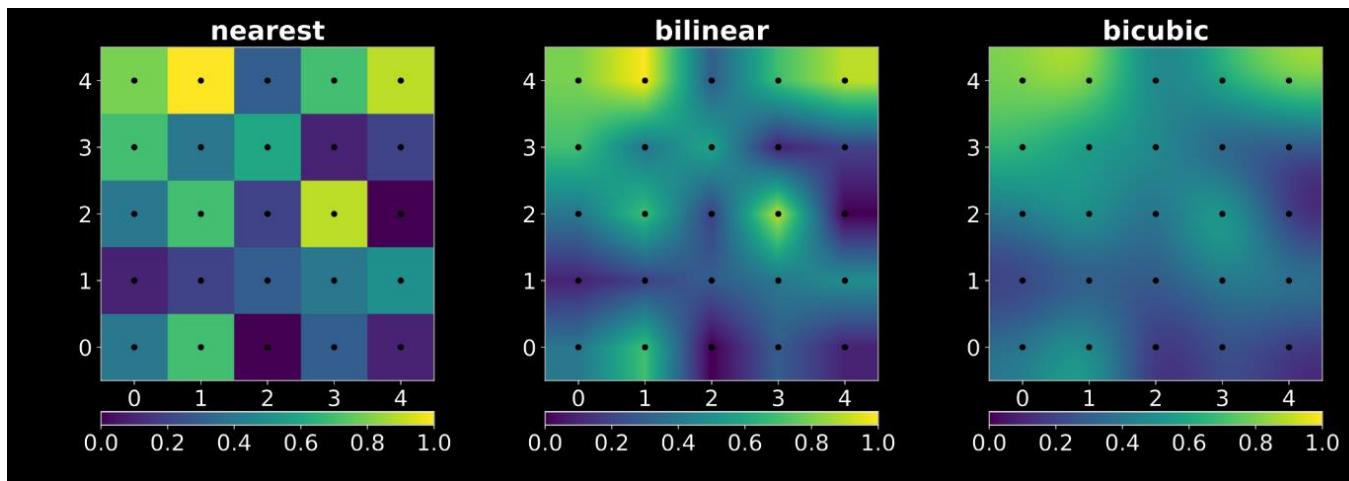
- upsampling: unpooling
- convolution



also: up-convolution,  
not ~~deconvolution~~

# Upsampling Methods

Recall: as alternative for convolution, different interpolation methods exist.



# Upsampling Methods

Use learned network – kernels and feature map to ‘deconvolve’ to original pixel space.

**Output**


**Kernel**

0	1
2	3

**Input Feature map**

0	1
2	3



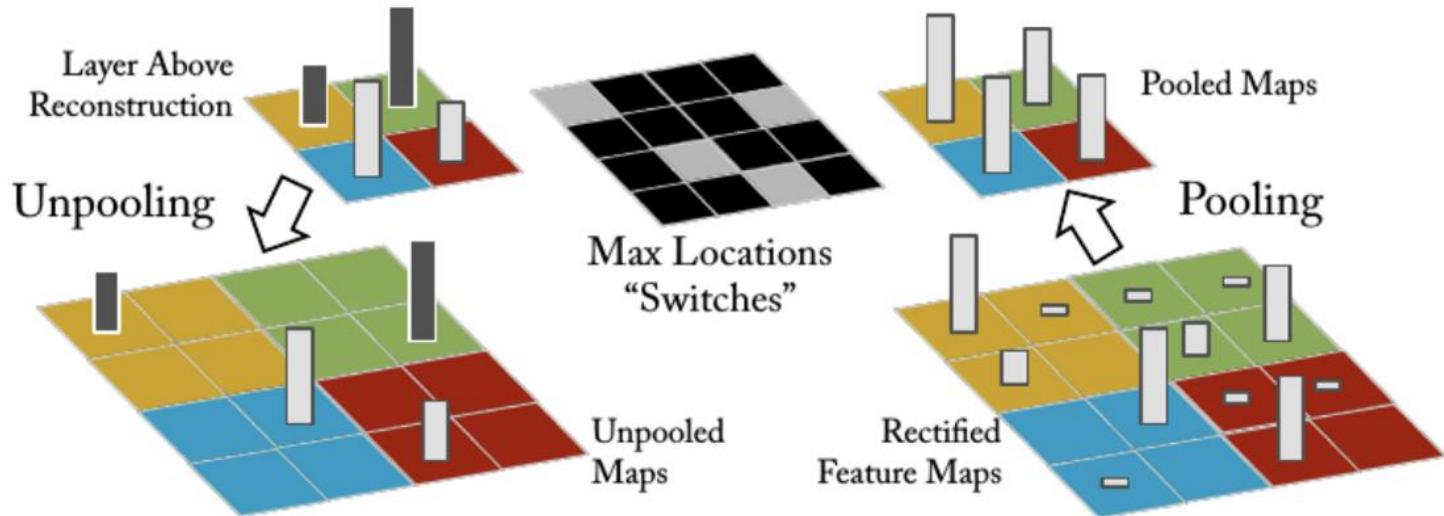
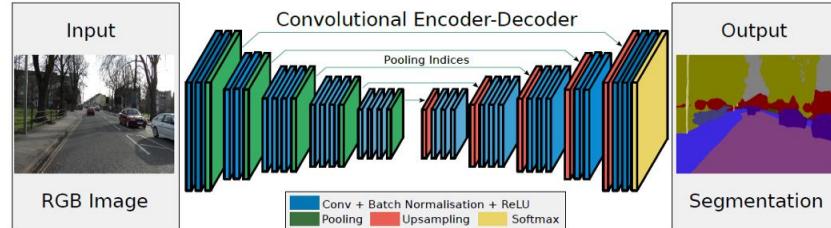
0	0	1
0	4	6
4	12	9



$$\begin{array}{c} \begin{array}{|c|c|c|} \hline 0 & 0 & \\ \hline 0 & 0 & \\ \hline & & \\ \hline \end{array} + \begin{array}{|c|c|c|} \hline & 0 & 1 \\ \hline & 2 & 3 \\ \hline & & \\ \hline \end{array} + \begin{array}{|c|c|c|} \hline & & \\ \hline 0 & 2 & \\ \hline 4 & 6 & \\ \hline \end{array} + \begin{array}{|c|c|c|} \hline & & \\ \hline 0 & 3 & \\ \hline 6 & 9 & \\ \hline \end{array} \end{array}$$

# Segnet: Alternative Unpooling operator

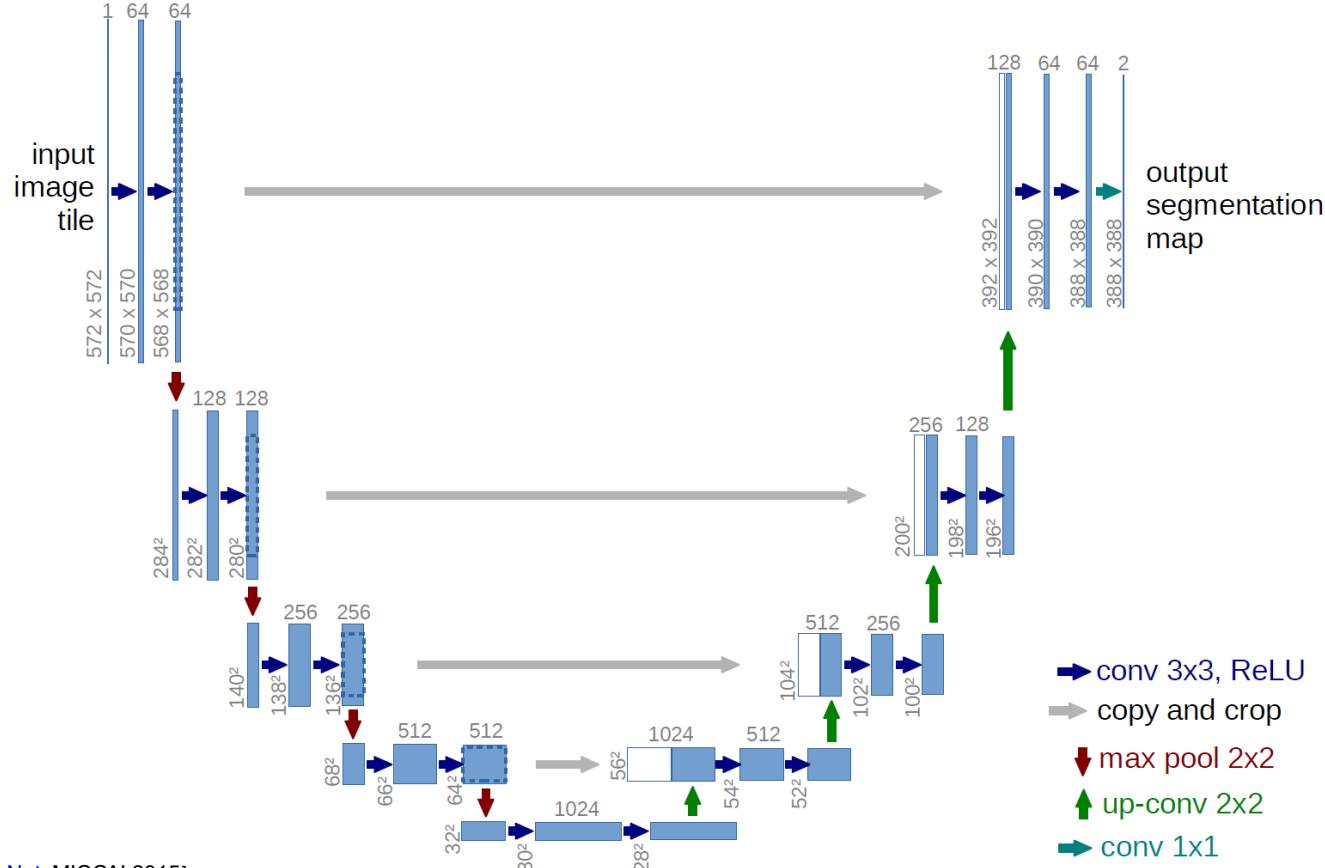
The unpooling uses the max-locations ('switches') that have been stored during pooling.



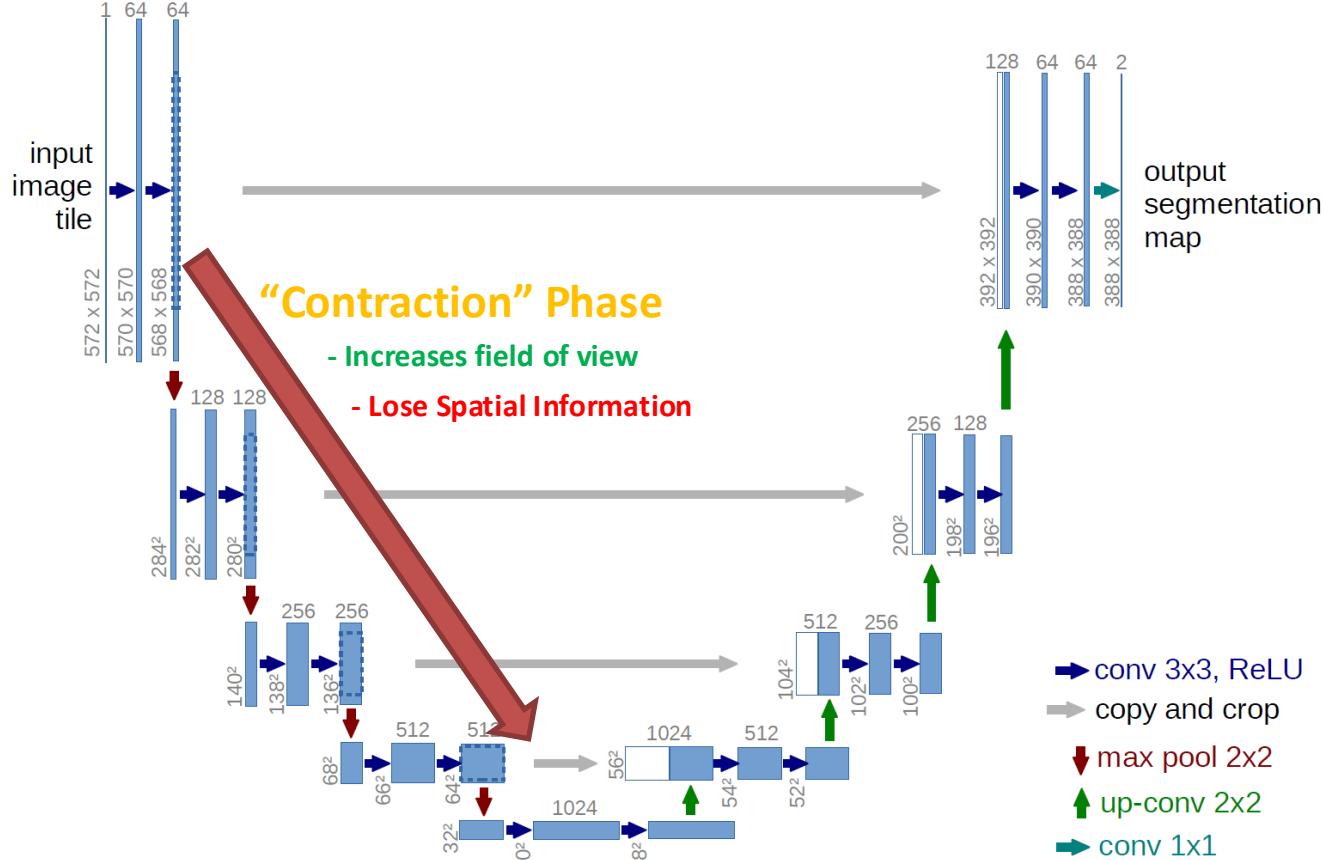
# Segnet - Results



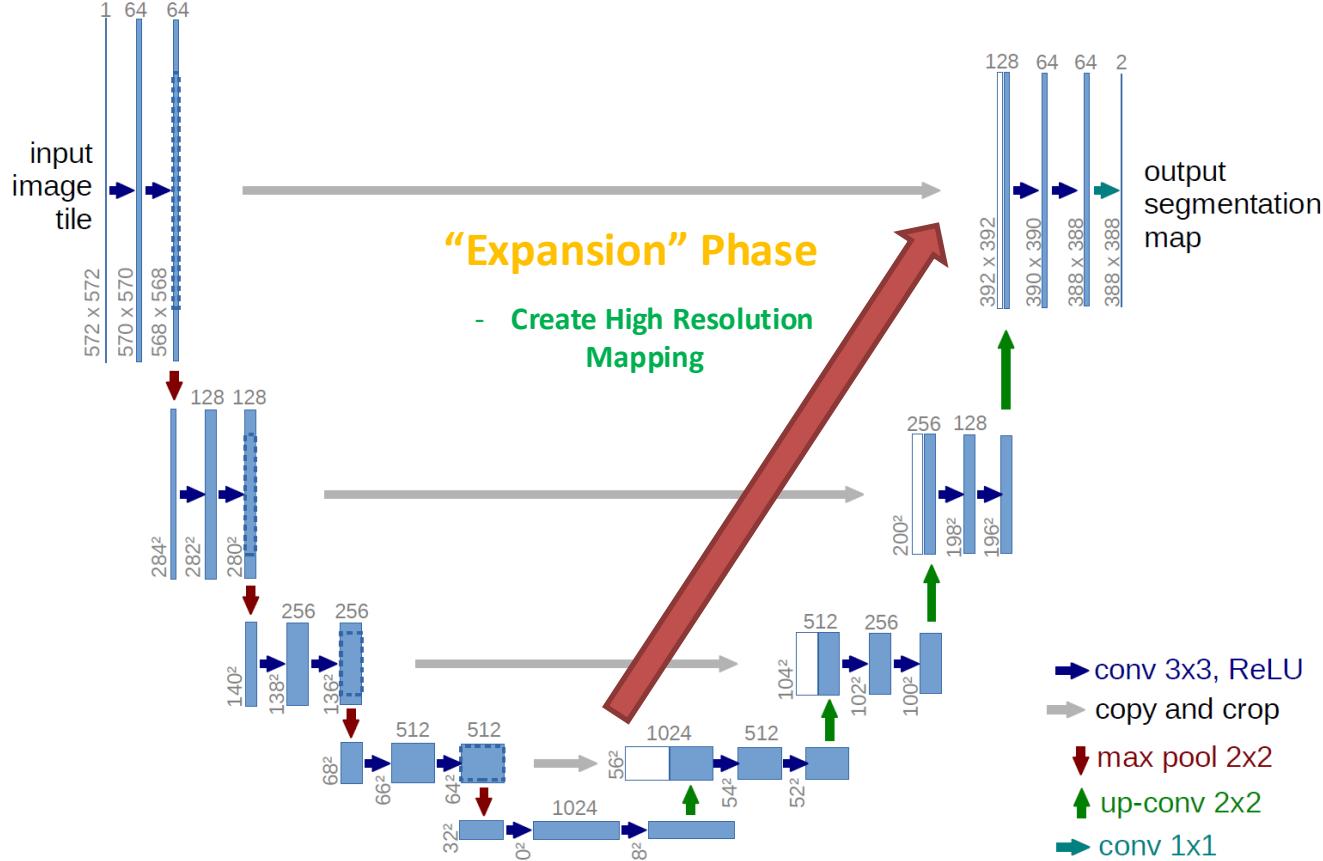
# U-Net Architecture



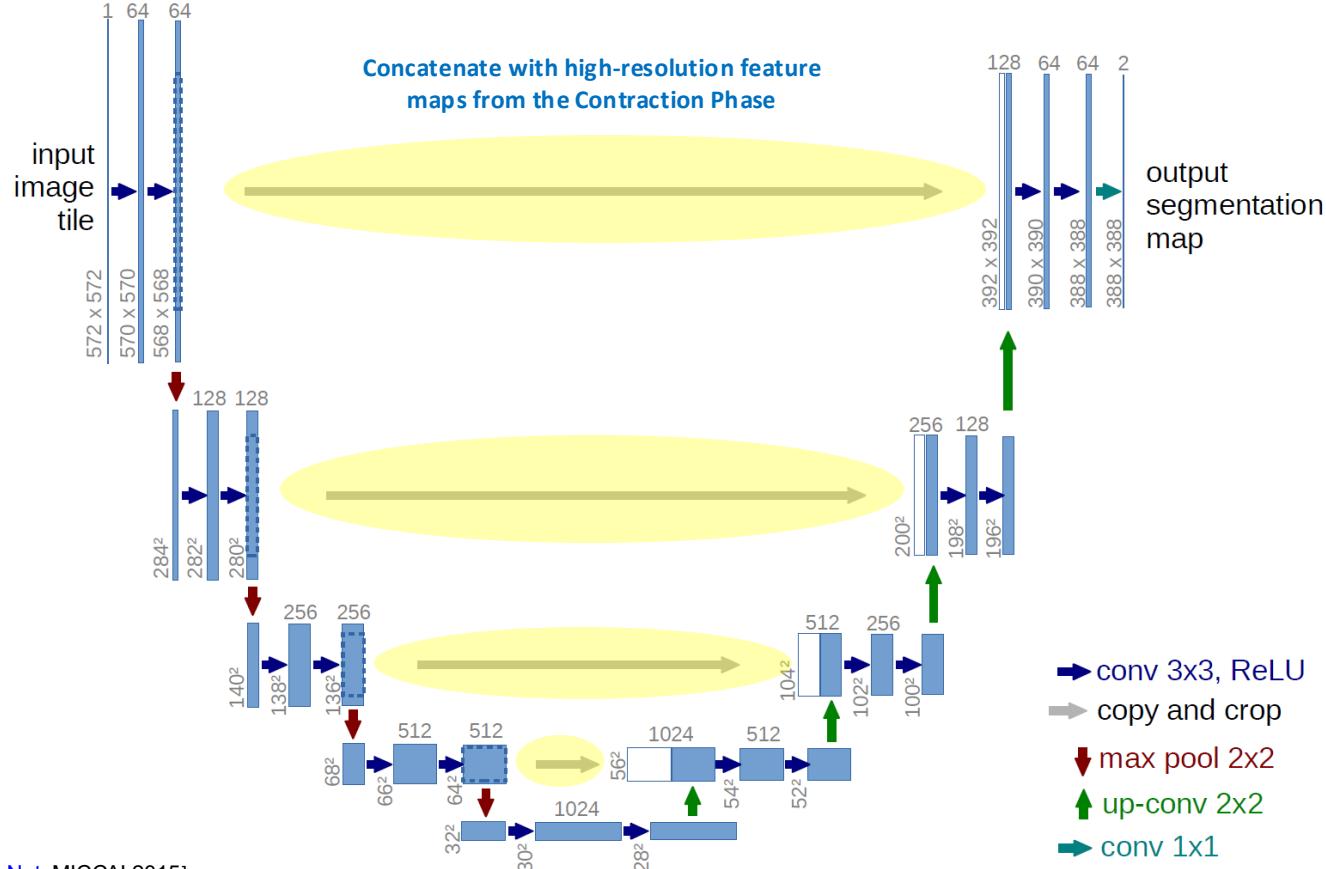
# U-Net Architecture



# U-Net Architecture

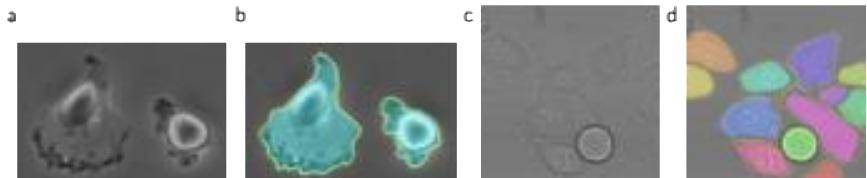
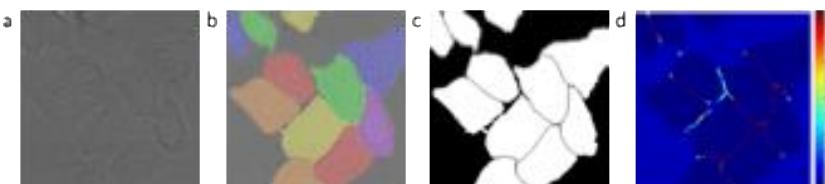


# U-Net Architecture



# U-Net Architecture

- ConvNet architecture: inductive bias for translation equivariance
- Image-to-Image predictions: input resolution = output resolution
  - Segmentation
  - Superresolution
  - Denoising
  - Inpainting
  - Style transfer
  - etc.
- Encoder-decoder architecture: contracting & expanding path
- Combines the idea of ResNet with coarse-to-fine / pyramid networks: skip connections across hierarchy levels
- widely used today:
- originally developed for medical image segmentation:



# Autoencoders

---

**So far:**

Supervised learning:

- training with large number of labeled images

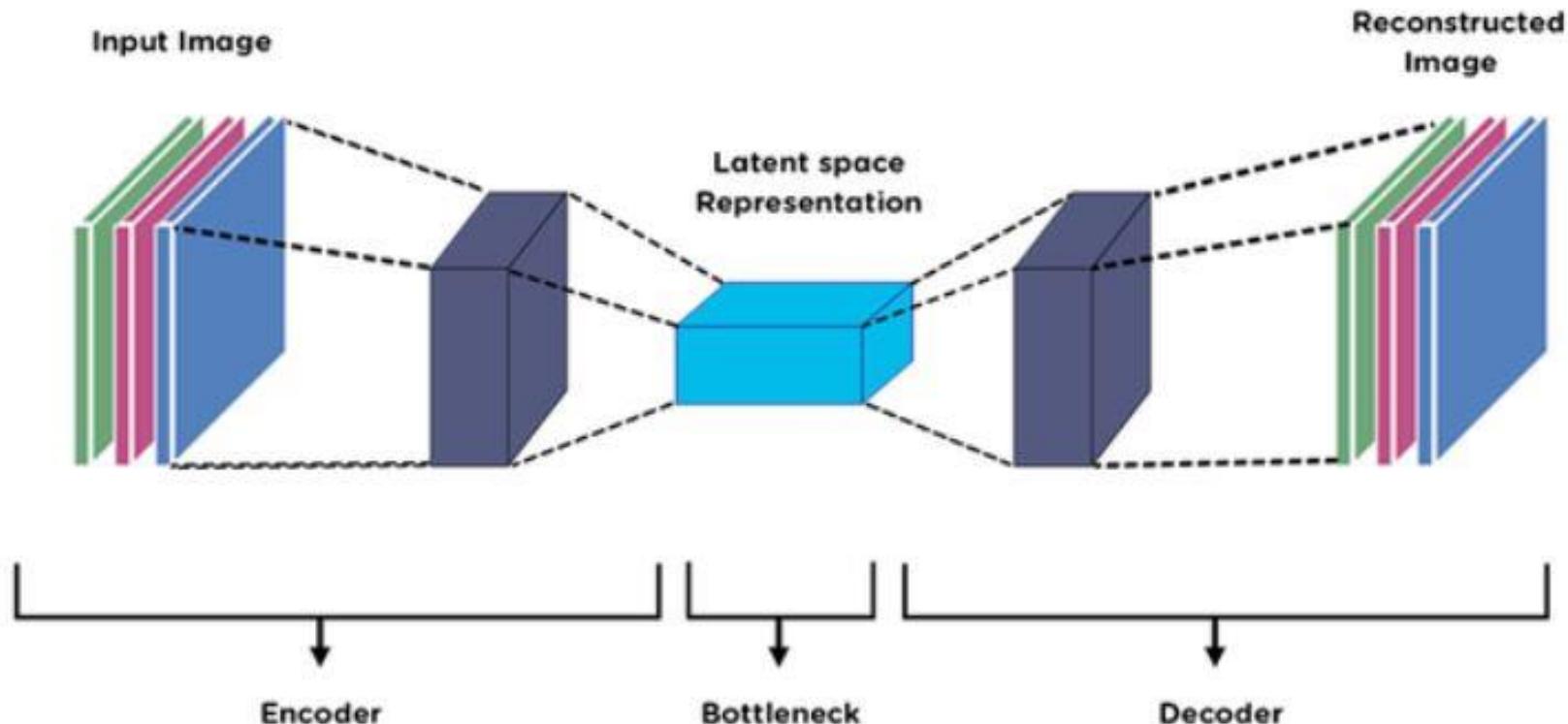
**Now:**

Unsupervised approach:

- to learn lower dimensional feature representations from unlabeled training data.

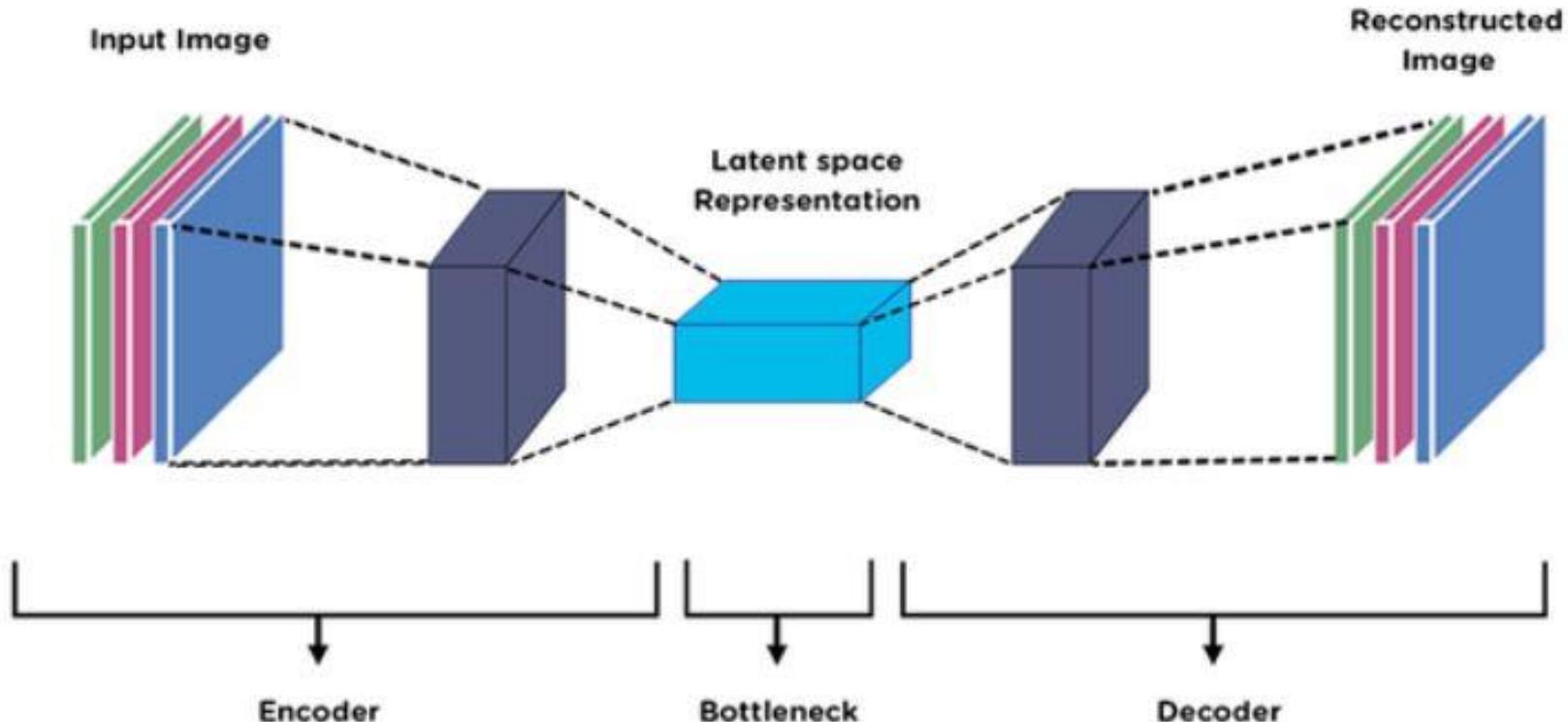
# Autoencoders

Encoder – decoder. Network is trained to minimize difference between Reconstructed Image and Input Image (e.g. LSE loss).



# Autoencoders

Latent space is of lower dimension than input.  
So reconstruction typical not perfect.

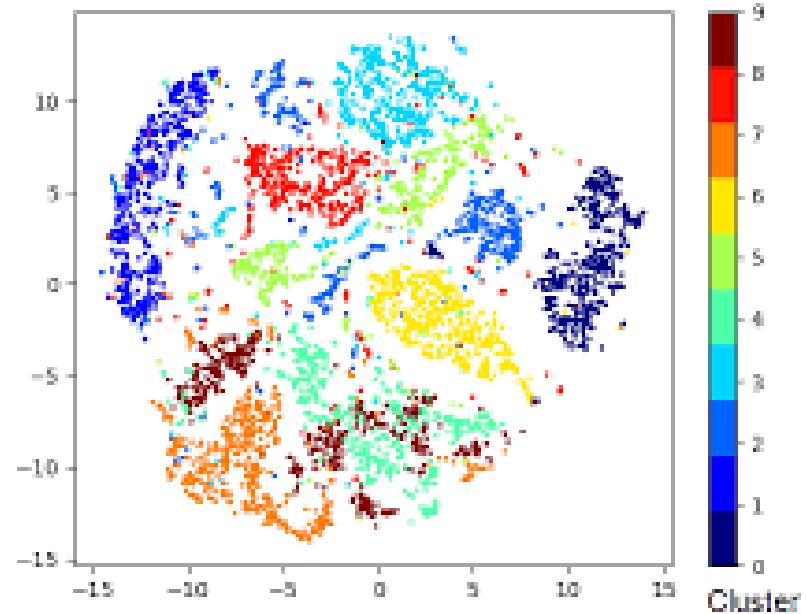


# Autoencoders – Use case: Clustering

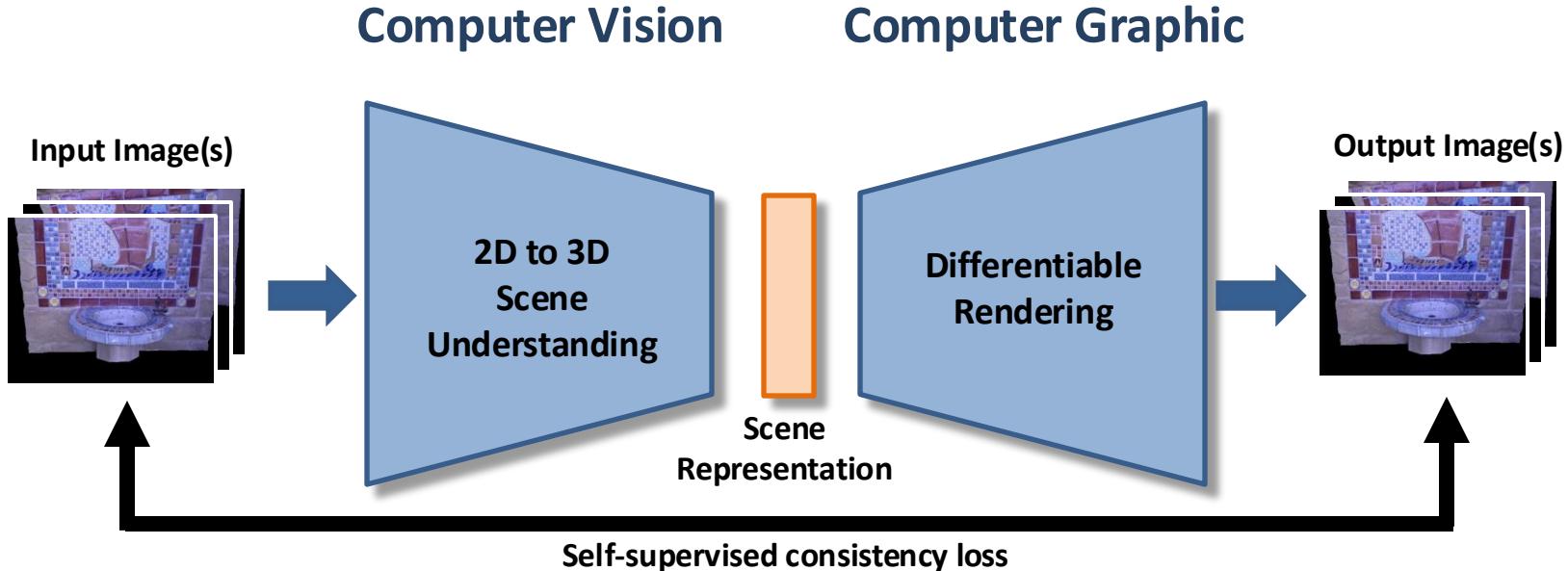
Embedding of MNIST dataset



MNIST

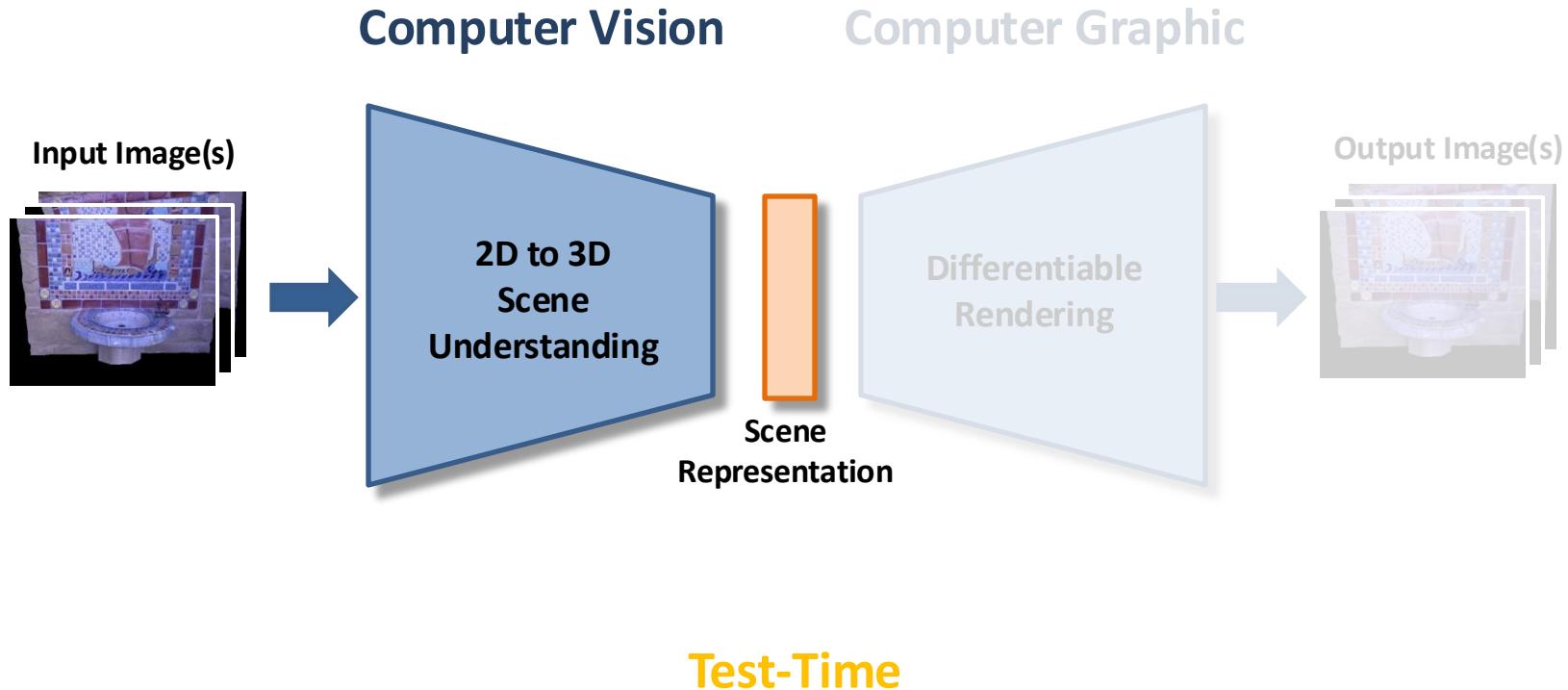


# Autoencoders – Self-Supervised Learning

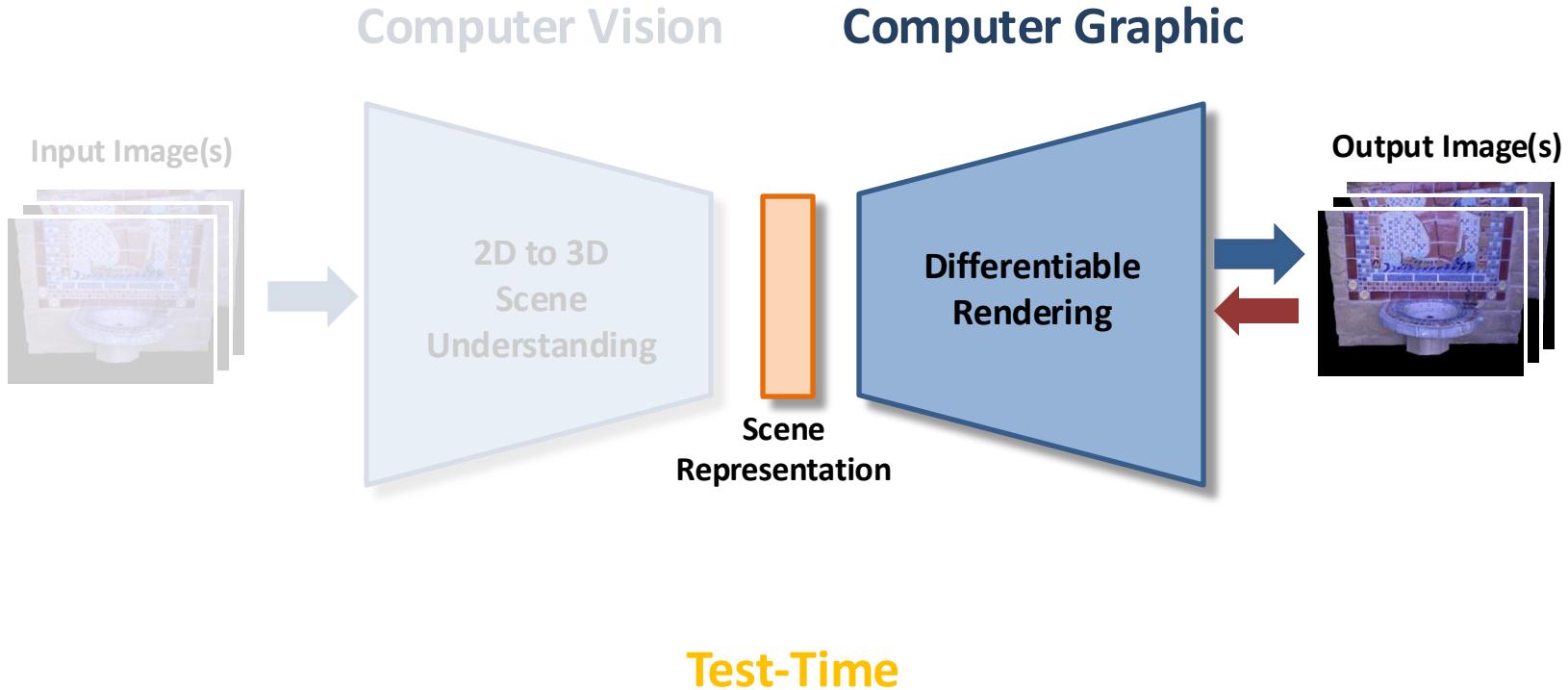


Training-Time

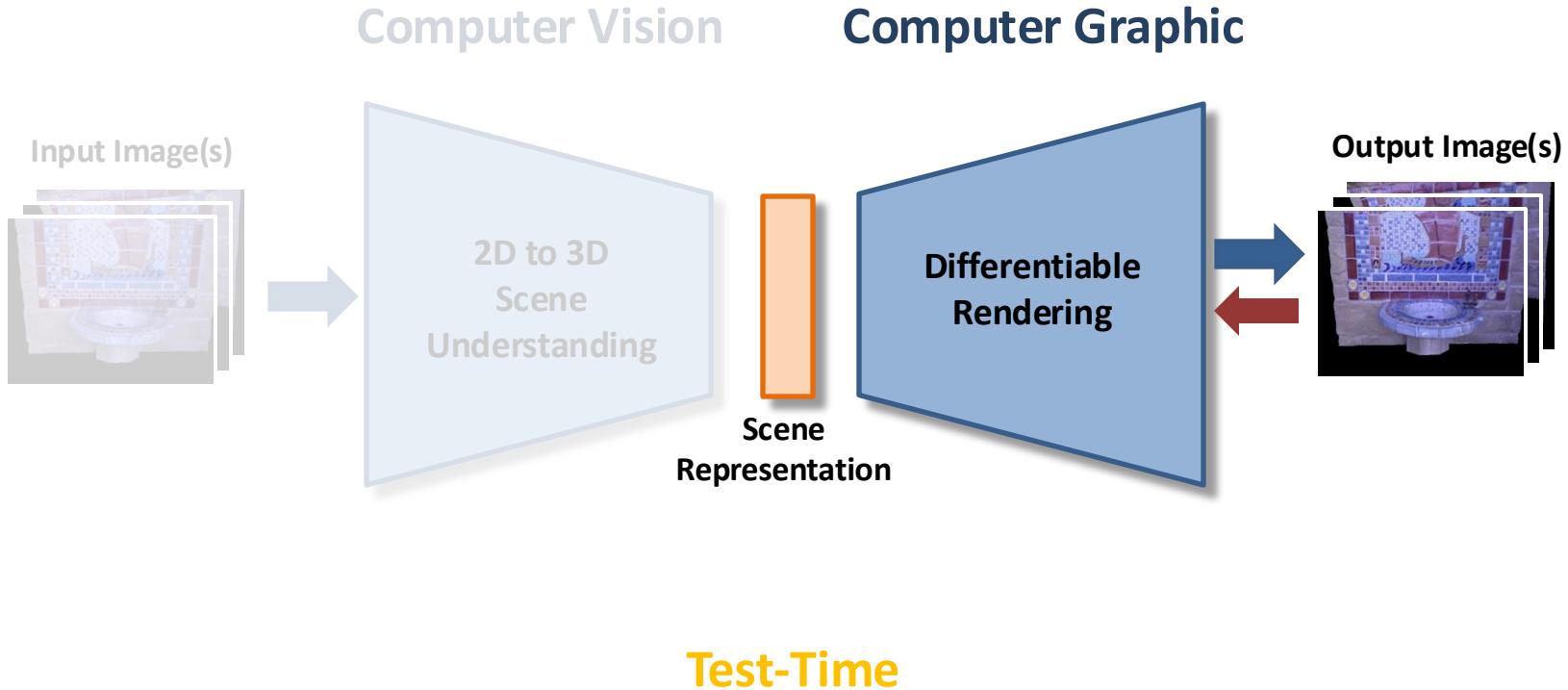
# Autoencoders – Self-Supervised Learning



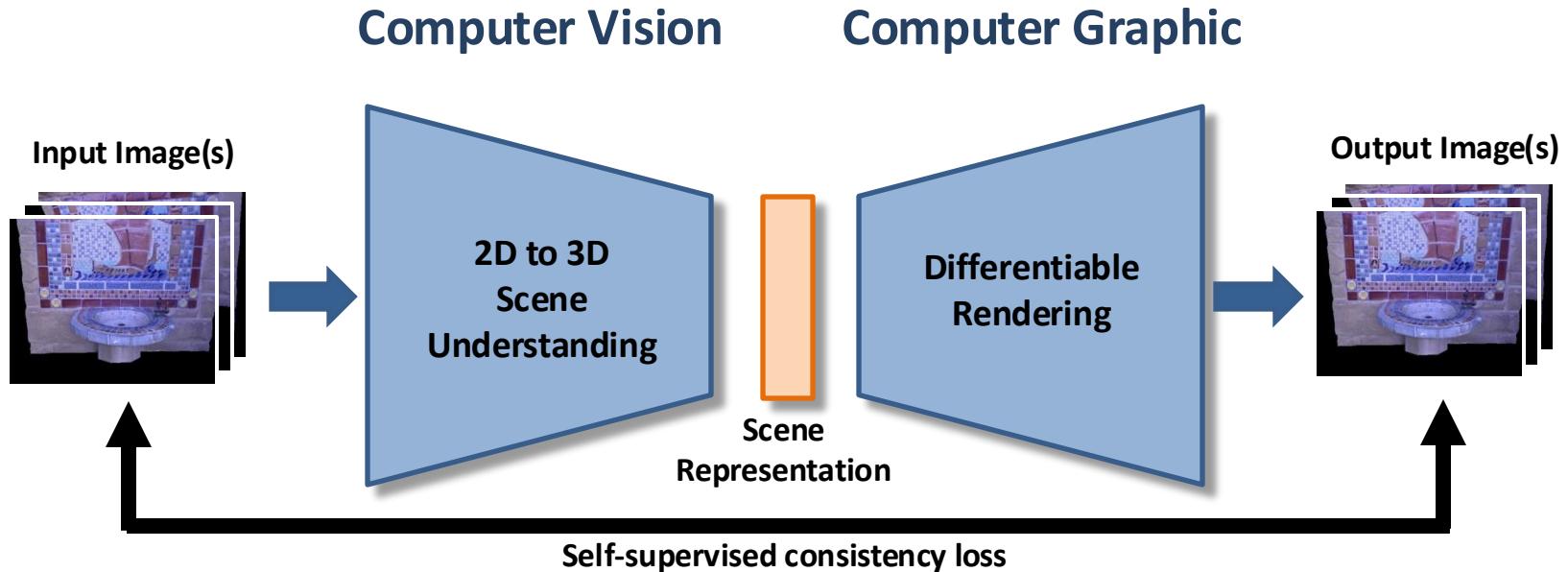
# Autoencoders – Self-Supervised Learning



# Autoencoders – Self-Supervised Learning

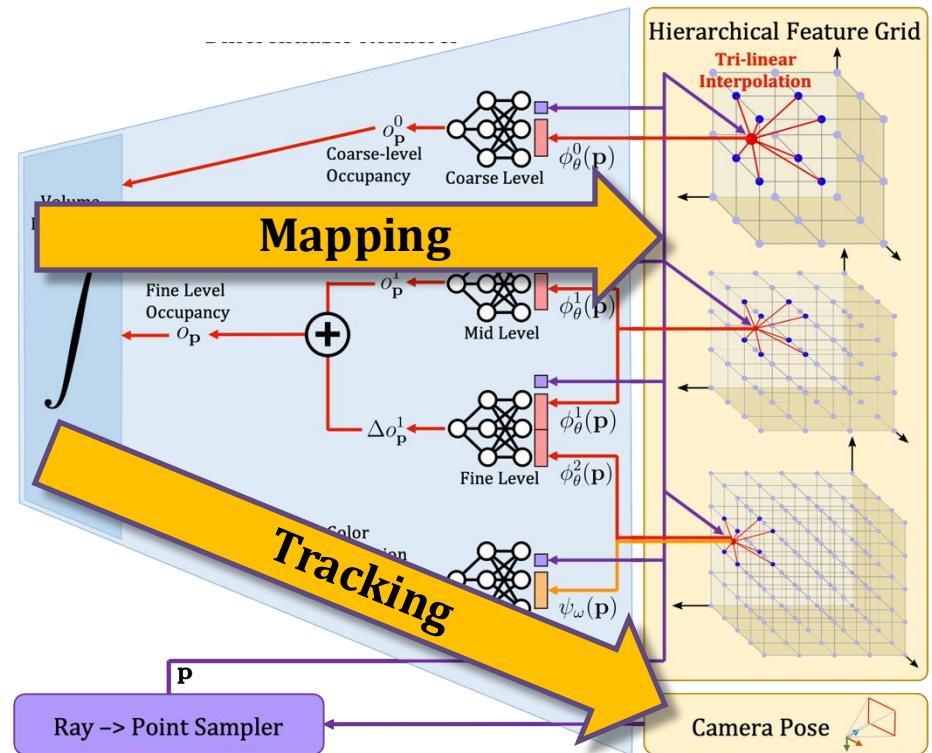
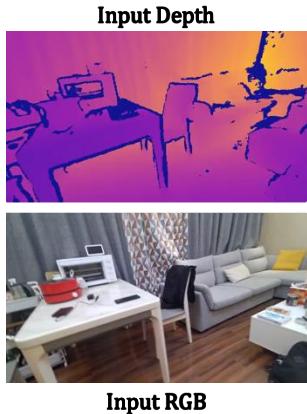


# Autoencoders – Self-Supervised Learning



- This principle is the foundation of many recent works, e.g. Neural radiance fields (NeRF, Gaussian Splatting, etc.), 3D reconstruction, 3D localization (Dust3R, Mast3R, etc.), Human pose estimation, ...
- Enables scalable learning without human labeling: only images/videos required

# Example: NICE-SLAM



---

# Questions?

# Disclaimer

---

Many of the slides used here are obtained from online resources (including many open lecture materials) without appropriate acknowledgement. They are used here for the sole purpose of classroom teaching. All the credit and all the copyrights belong to the original authors. You should not copy it, redistribute it, put it online, or use it for any other purposes than for this course.