

Object Detection

Longbin Jin



Artificial Intelligence
& Computer Vision
Laboratory

Slide reference: Stanford University cs231n, Fei-Fei Li

The Generalized R-CNN Framework for Object Detection, Ross Girshick, ICCV 2019 tutorial

Computer Vision Tasks

Classification



CAT

No spatial extent

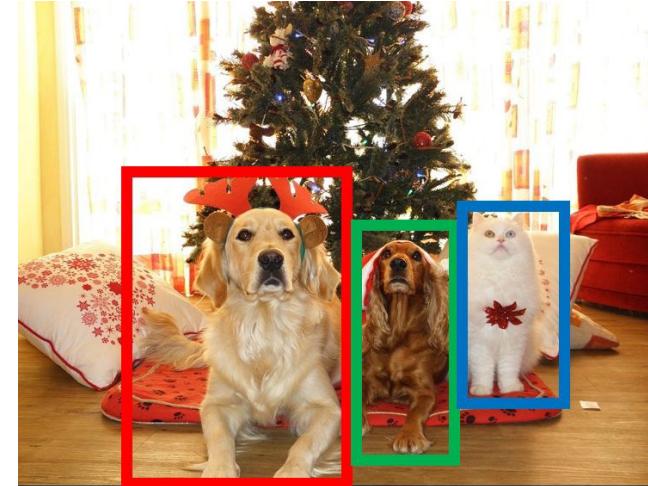
Semantic Segmentation



GRASS, CAT
TREE, SKY

No objects, just pixels

Object Detection



DOG, DOG, CAT

Multiple Object

Instance Segmentation

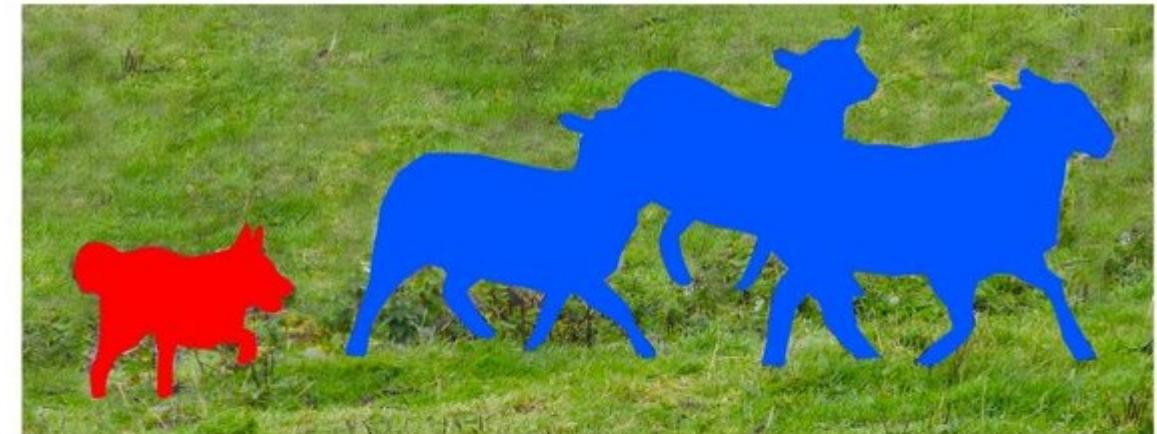


DOG, DOG, CAT

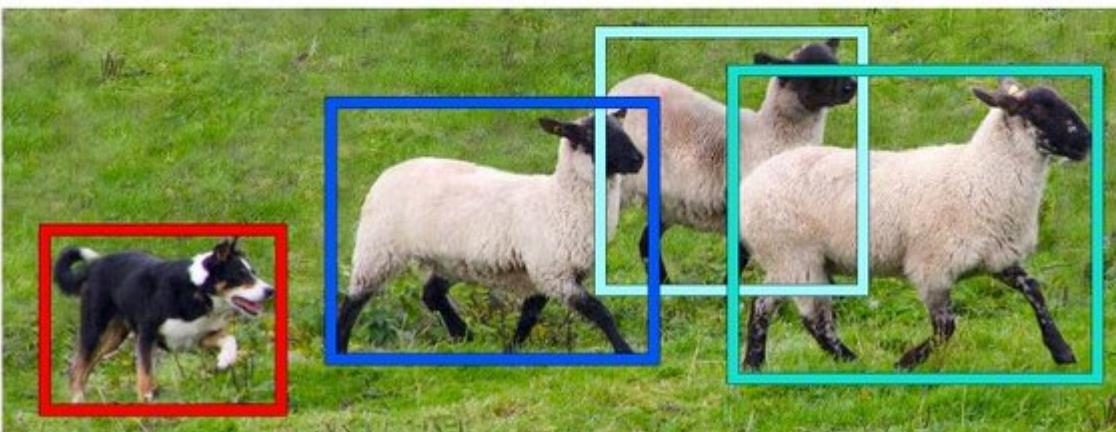
Computer Vision Tasks



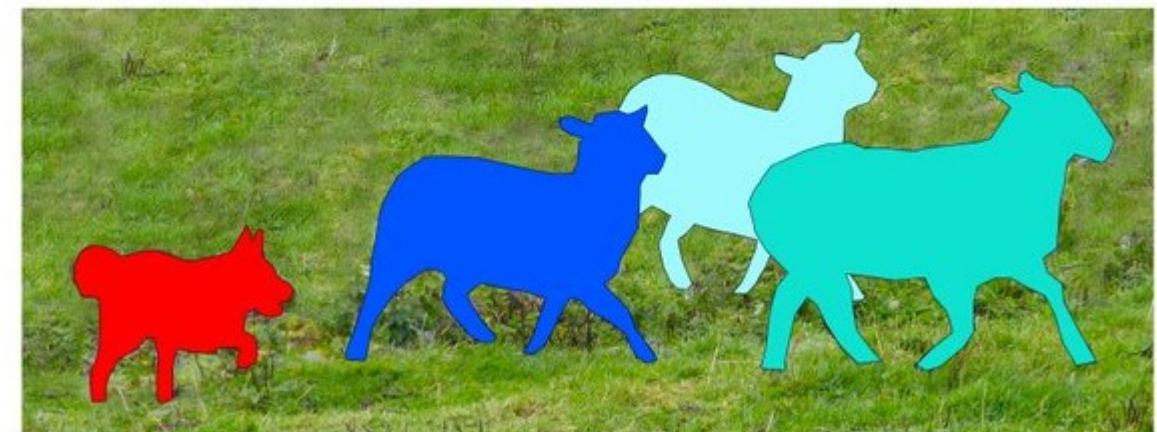
Image Recognition



Semantic Segmentation



Object Detection



Instance Segmentation

Segment Anything Model (SAM) 2023

Segment Anything

Research By Meta AI

Home Demo Dataset Blog



SA-1B Dataset Explorer

Filter by masks per image, mask area, or image id e.g. Masks Per Image > 300

Search

50000 Images Hide Annotations

50k preview of the full 11M image dataset. [Access the full dataset](#)



◀ 1 / 1000 ▶

Semantic Segmentation

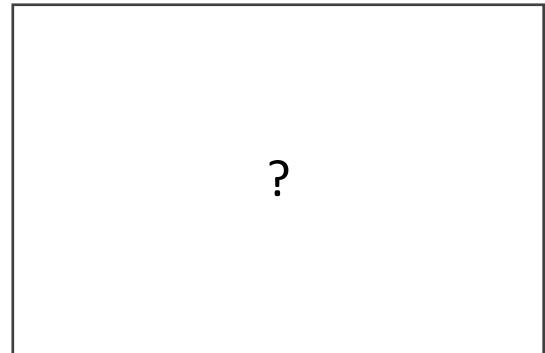
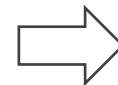
Train



GRASS, CAT
TREE, SKY

Paired training data: for each training image,
each pixel is labeled with a semantic category

Test



At test time, classify each pixel of a new image.

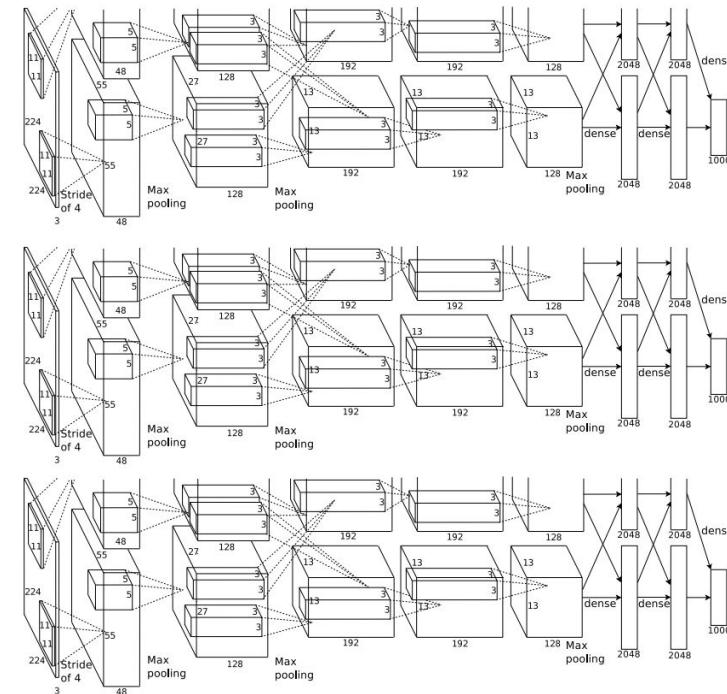
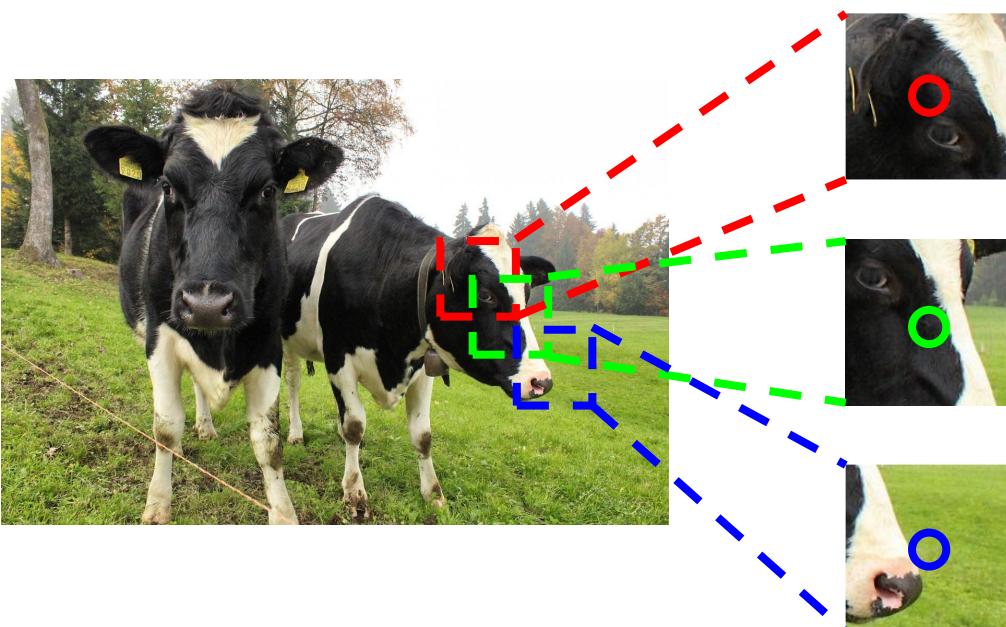
Semantic Segmentation Idea: Sliding Window



Impossible to classify without context

Q: how do we include context?

Semantic Segmentation Idea: Sliding Window



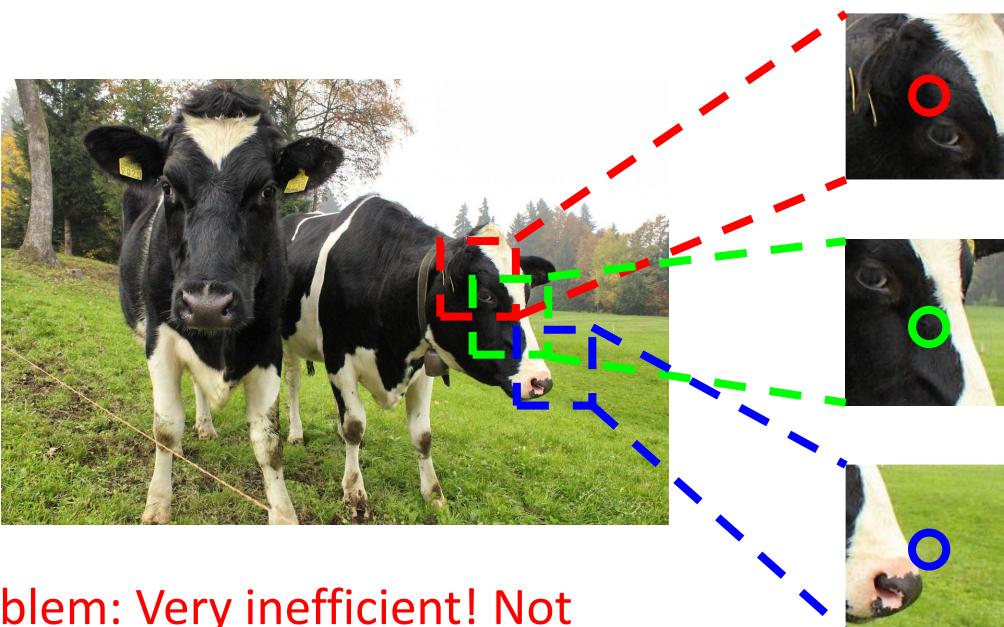
Cow

Cow

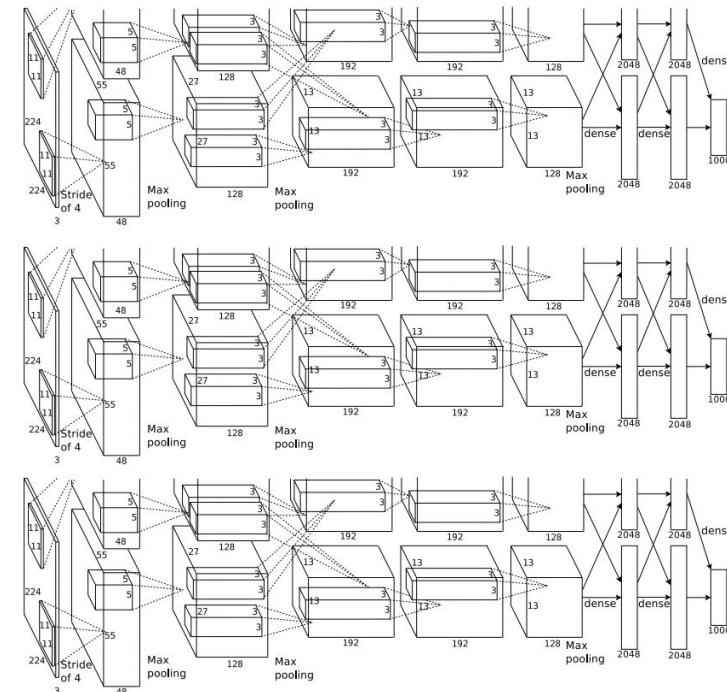
Grass

Q: how do we model this?

Semantic Segmentation Idea: Sliding Window



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

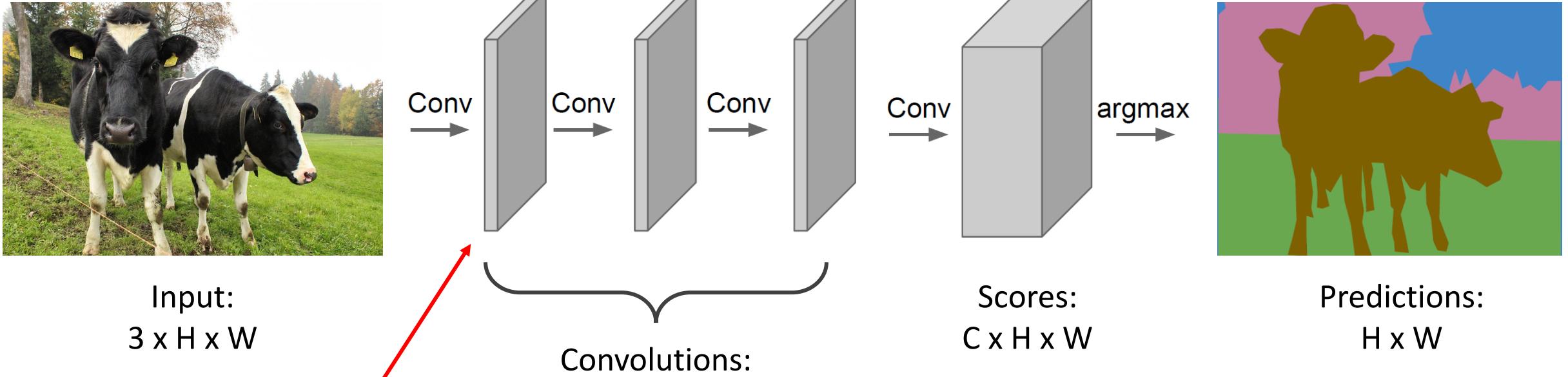


Cow

Cow

Grass

Semantic Segmentation Idea: Fully Convolutional



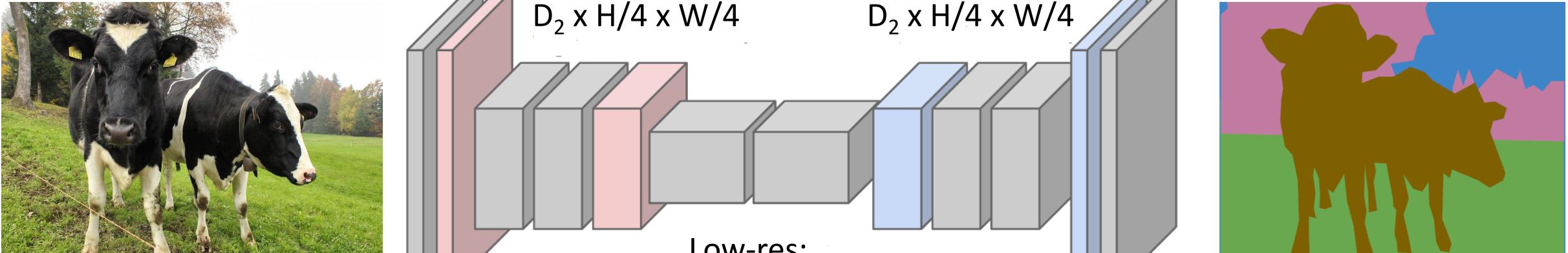
Problem: convolutions at
original image resolution will
be very expensive ...

Semantic Segmentation Idea: Fully Convolutional



Input:
 $3 \times H \times W$

Downsampling:
Pooling, strided
convolution



Predictions:
 $H \times W$

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

Upsampling:
??

In-Network upsampling: “Unpooling”

Nearest Neighbor

1	2
3	4

1	1	2	2
1	1	2	2
3	3	4	4
3	3	4	4

Input: 2 x 2

Output: 4 x 4

“Bed of Nails”

1	2
3	4

1	0	2	0
0	0	0	0
3	0	4	0
0	0	0	0

Input: 2 x 2

Output: 4 x 4



Semantic Segmentation Idea: Fully Convolutional

Max Pooling

Remember which element was max!

1	2	6	3
3	5	2	1
1	2	2	1
7	3	4	8

5	6
7	8

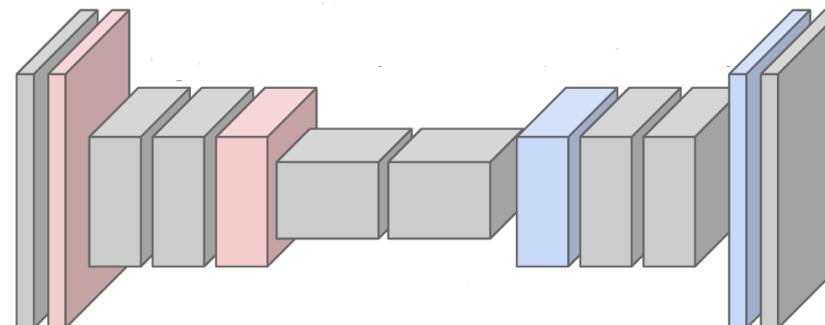
Rest of the
network

...

1	2
3	4

0	0	2	0
0	1	0	0
0	0	0	0
3	0	0	4

Corresponding pairs of
downsampling and
upsampling layers



Max Unpooling

Use positions from pooling layer

Learnable Upsampling

3 x 3 **transposed** convolution

Input: 2 x 2 Learnable filter:
$$\begin{matrix} 1 & 2 \\ 3 & 4 \end{matrix} * T \quad \begin{matrix} 1 & 1 & 1 \\ 1 & 1 & 1 \\ 1 & 1 & 1 \end{matrix}$$

Learnable filter

= SUM (

1	1	1	0
1	1	1	0
1	1	1	0
0	0	0	0

0	0	0	0
3	3	3	0
3	3	3	0
3	3	3	0

0	2	2	2
0	2	2	2
0	2	2	2
0	0	0	0

0	0	0	0
0	4	4	4
0	4	4	4
0	4	4	4

)

Learnable Upsampling

3 x 3 **transposed** convolution

Input: 2 x 2 Learnable filter:
 3 x 3

Learnable filter

1	1	1
1	1	1
1	1	1

$*^T$

=

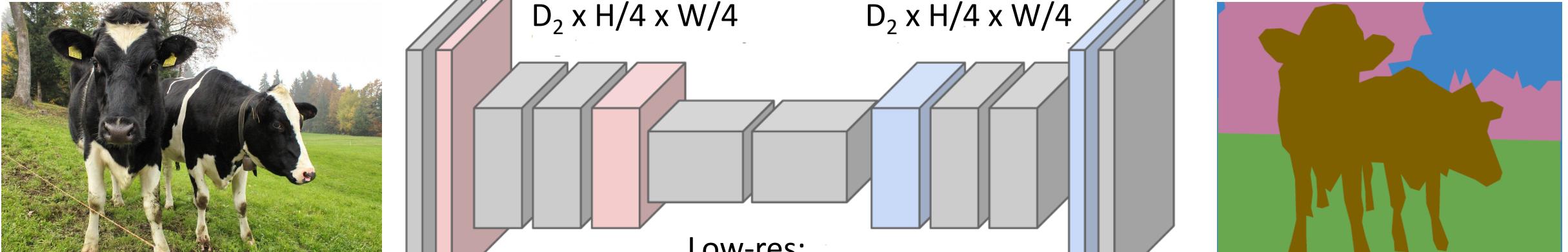
1	3	3	2
4	10	10	6
4	10	10	6
3	7	7	4

Semantic Segmentation Idea: Fully Convolutional



Input:
 $3 \times H \times W$

Downsampling:
Pooling, strided
convolution

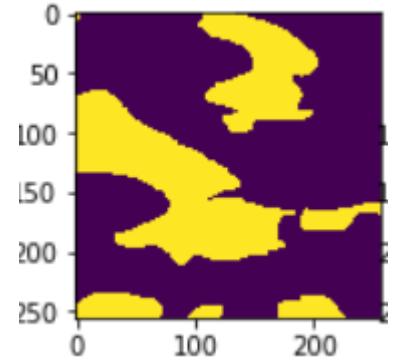
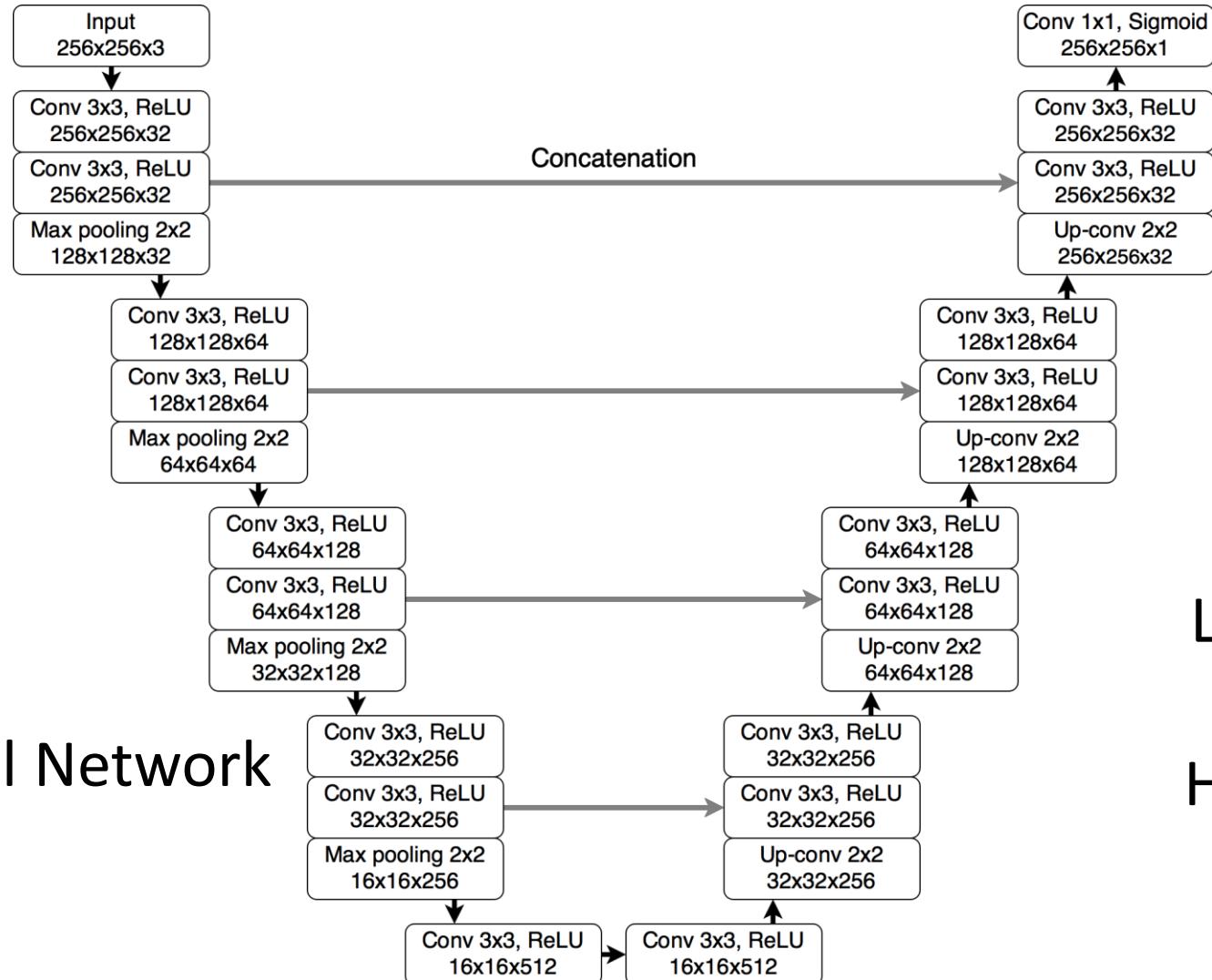
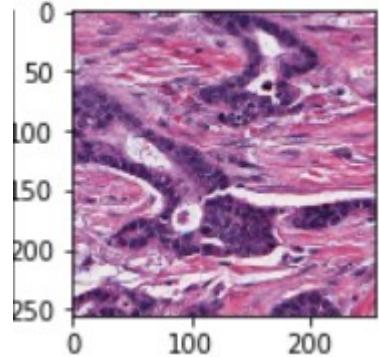


Design network as a bunch of convolutional layers, with
downsampling and **upsampling** inside the network!

Predictions:
 $H \times W$

Upsampling:
Unpooling or strided
transposed convolution

Unet



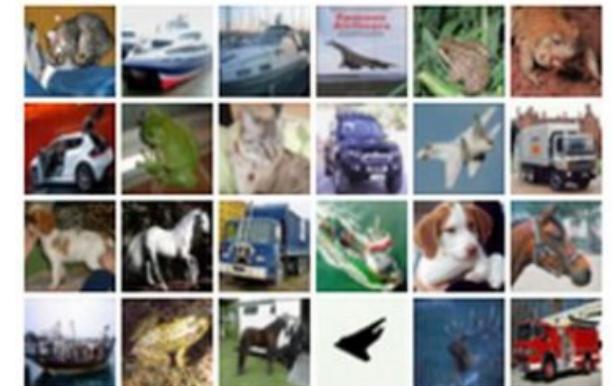
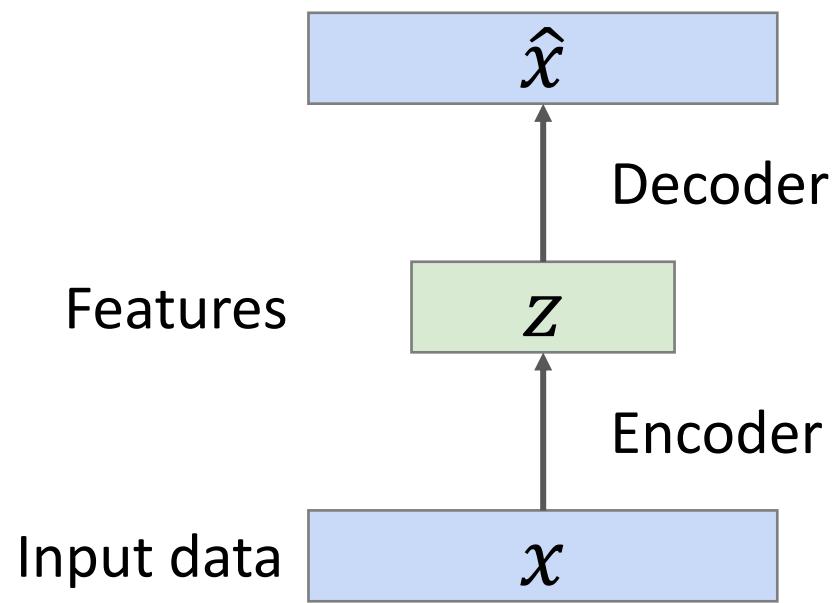
Fully Convolutional Network
FCN

Low-level feature
+
High-level feature

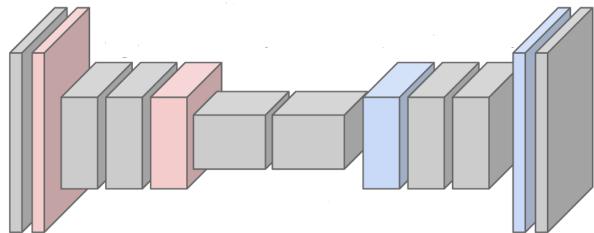
AutoEncoder (AE)

- Unsupervised approach for learning a lower-dimensional feature representation from unlabeled training data

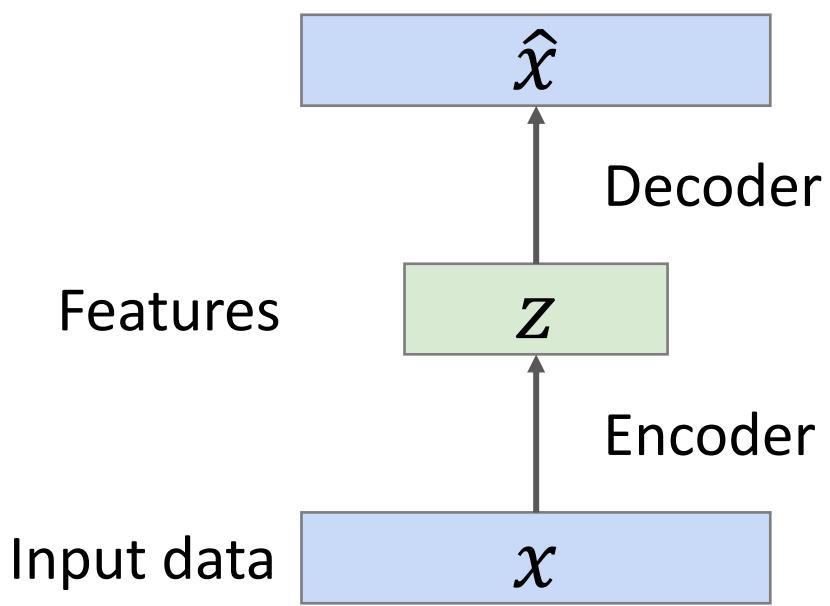
Want features to capture meaningful factors of variation in data automatically



AutoEncoder (AE)

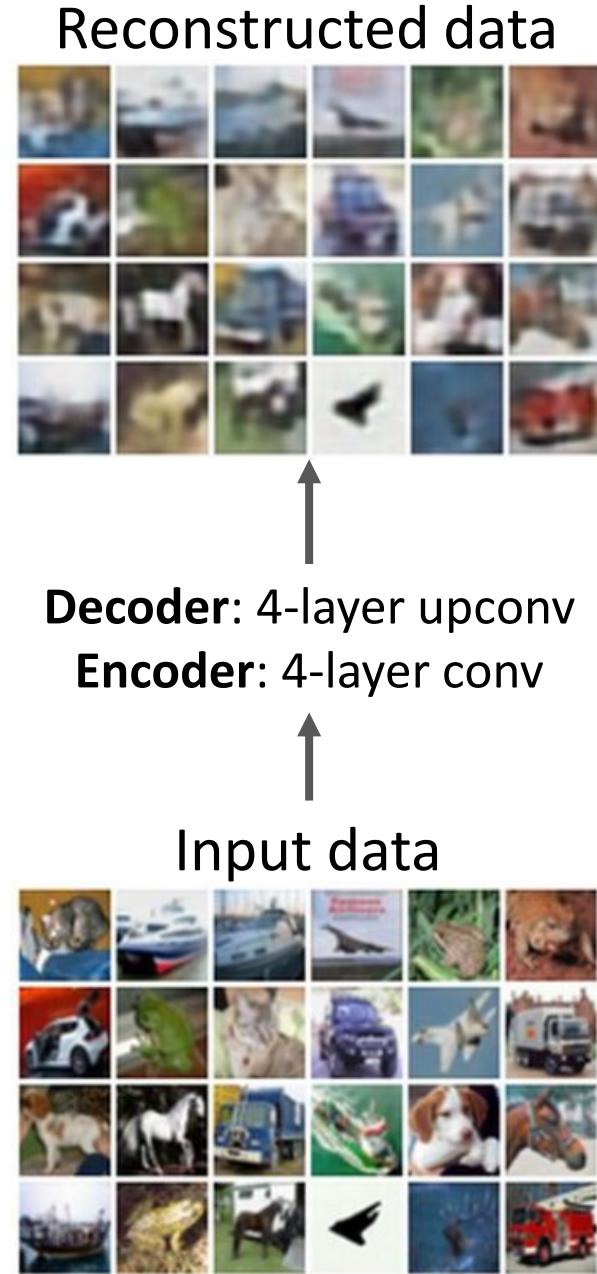


Reconstructed
input data



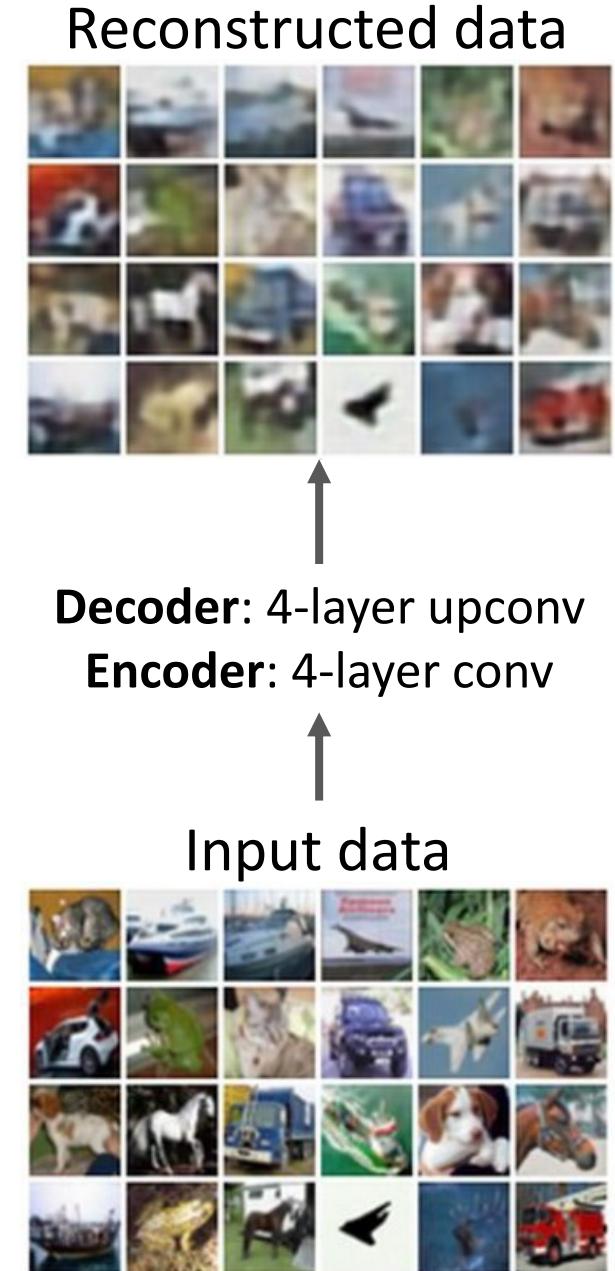
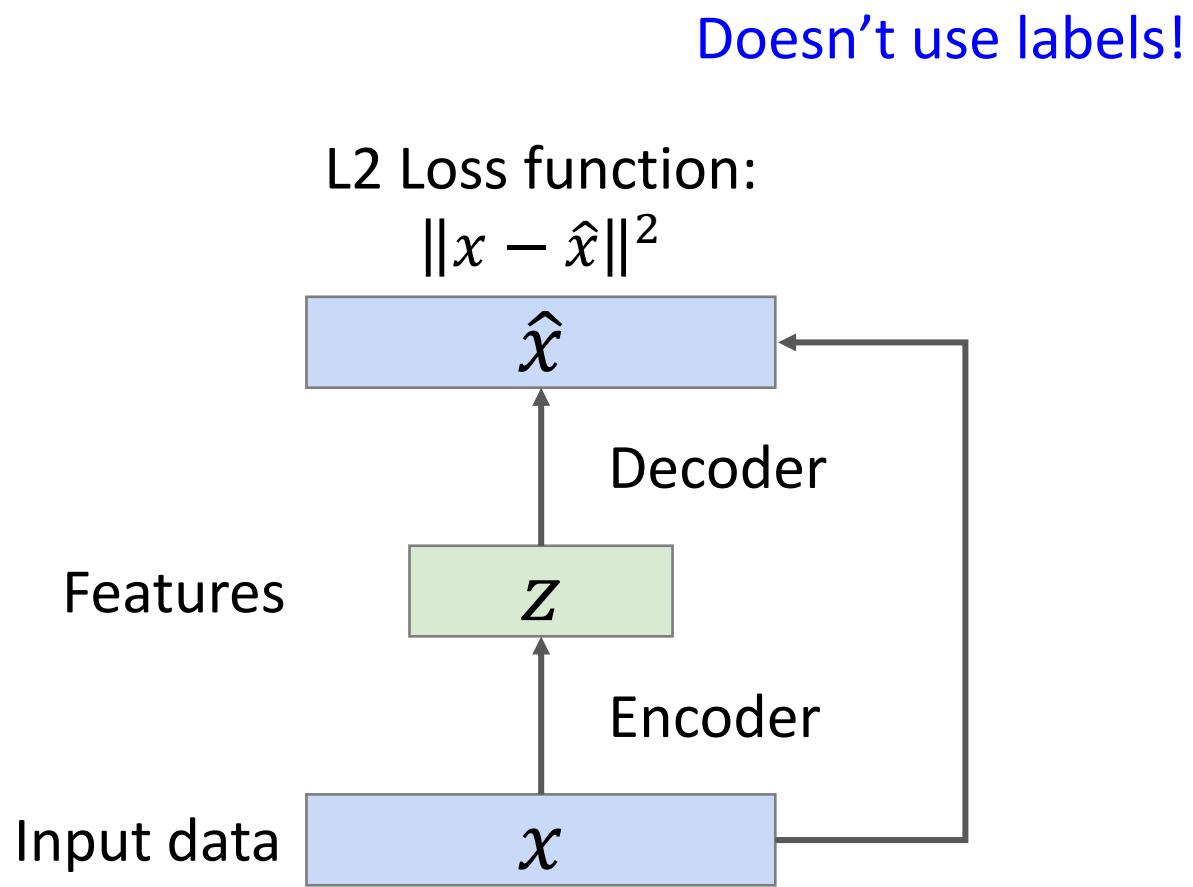
How to learn this
feature representation?

Train such that features
can be used to
reconstruct original
data “Autoencoding”
encoding input itself



AutoEncoder (AE)

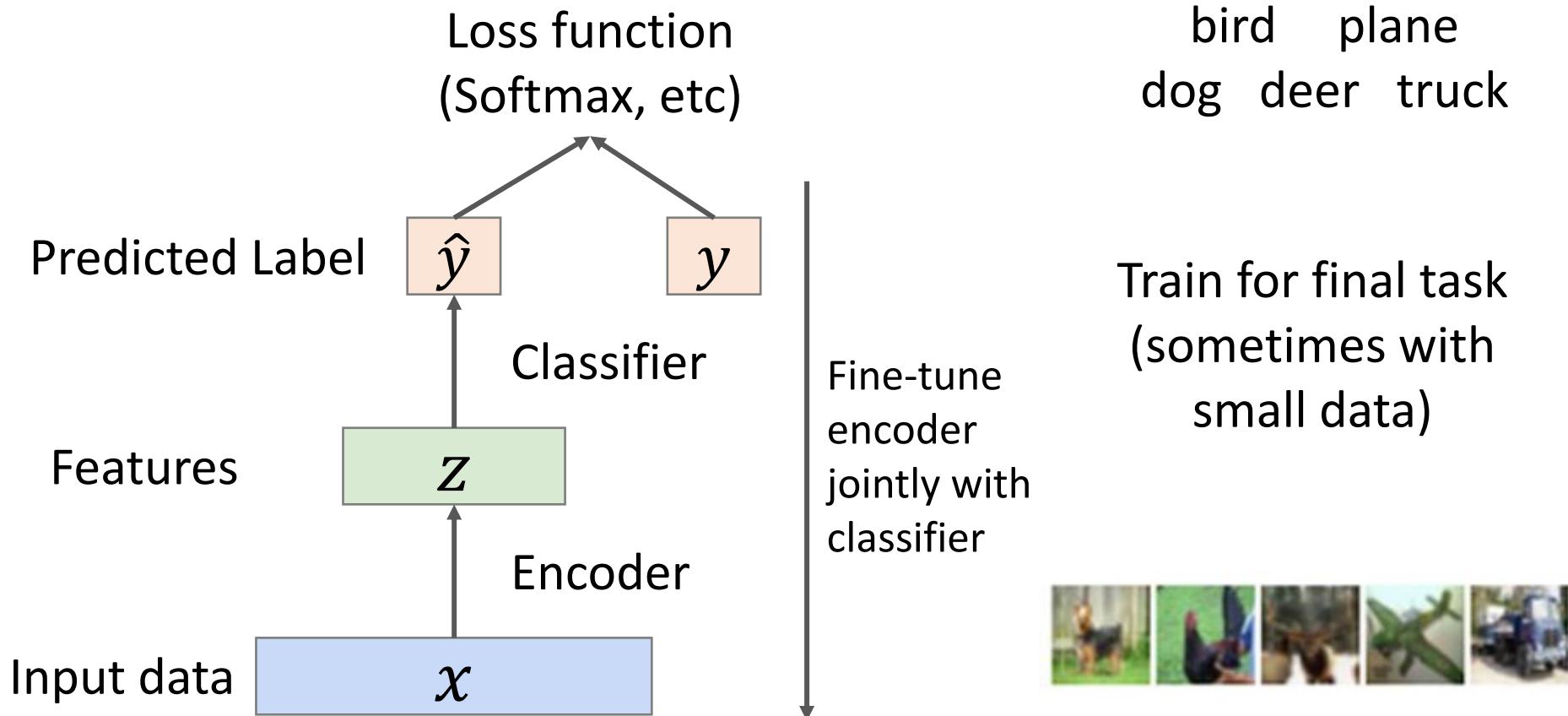
Train such that features
can be used to
reconstruct original data



AutoEncoder (AE)

Transfer from large,
unlabeled dataset to
small, labeled dataset.

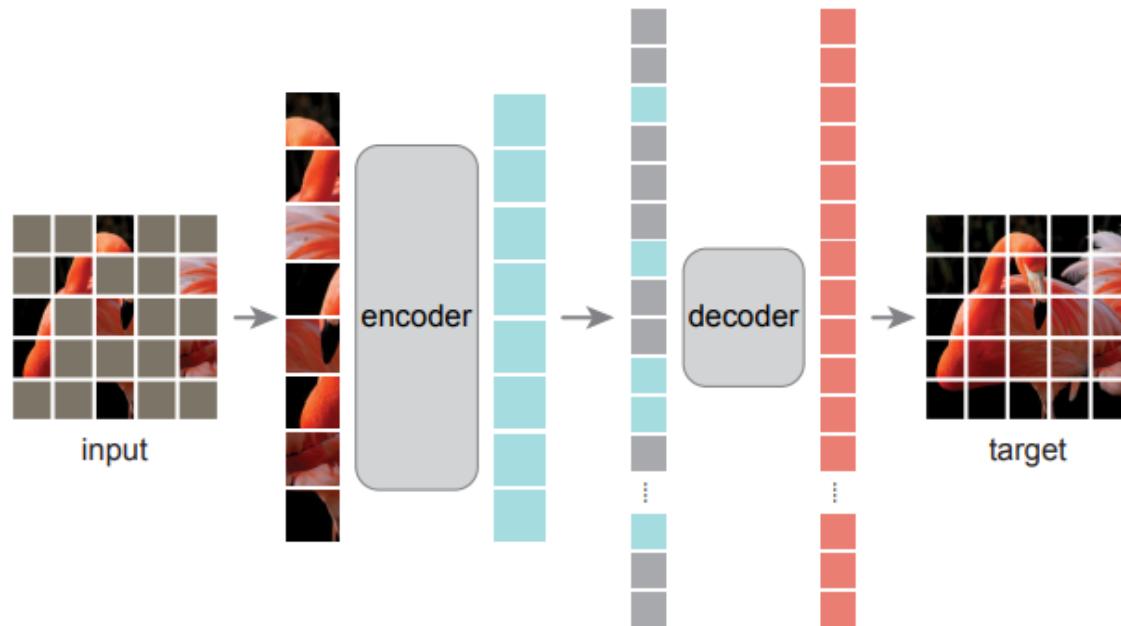
Encoder can be
used to initialize a
supervised model



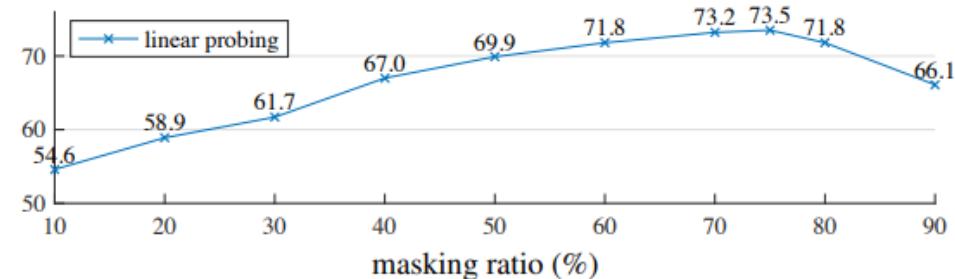
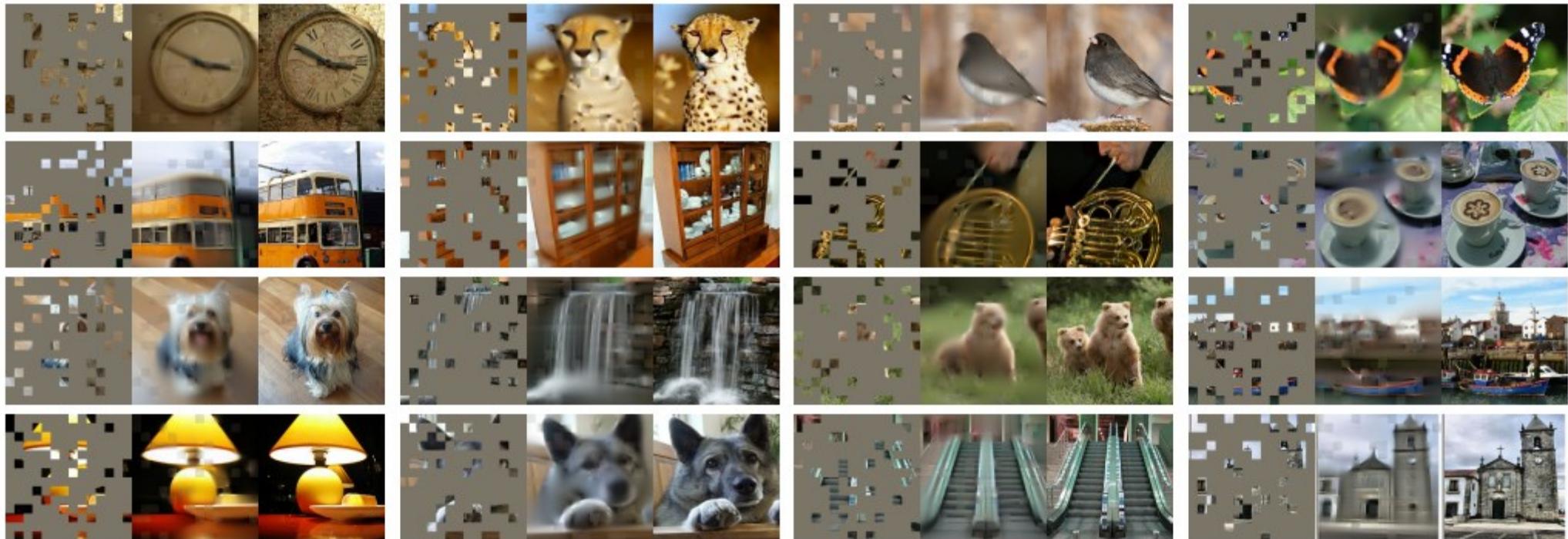
Masked Autoencoder (MAE)

- Pretext tasks
 - Mask random patches of the input image and reconstruct the missing pixels.

MAE encoder & decoder:
Vision Transformer



Masked Autoencoder (MAE)

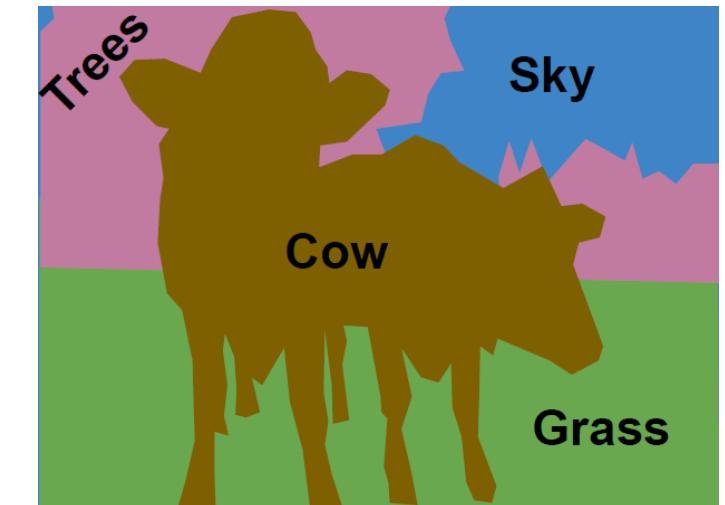


Semantic Segmentation

Label each pixel in the image with a category label



Don't differentiate instances, only care about pixels



Computer Vision Tasks

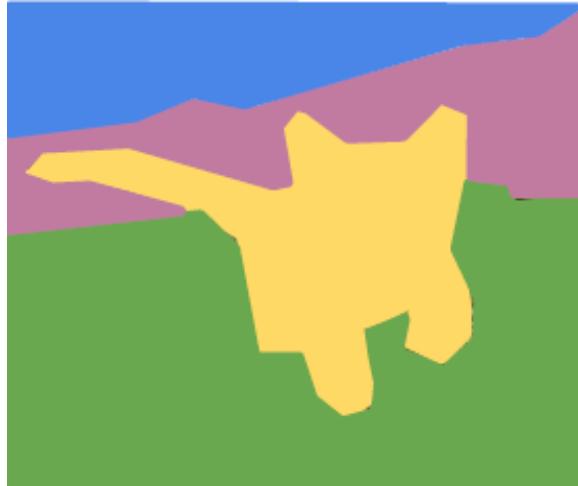
Classification



CAT

No spatial extent

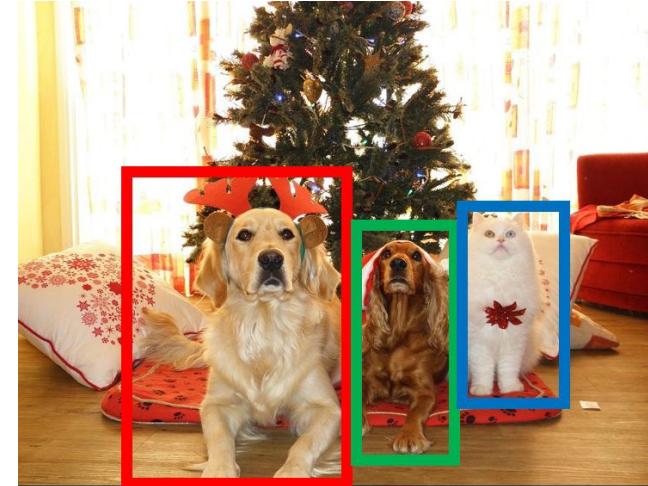
Semantic Segmentation



GRASS, CAT
TREE, SKY

No objects, just pixels

Object Detection



DOG, DOG, CAT

Multiple Object

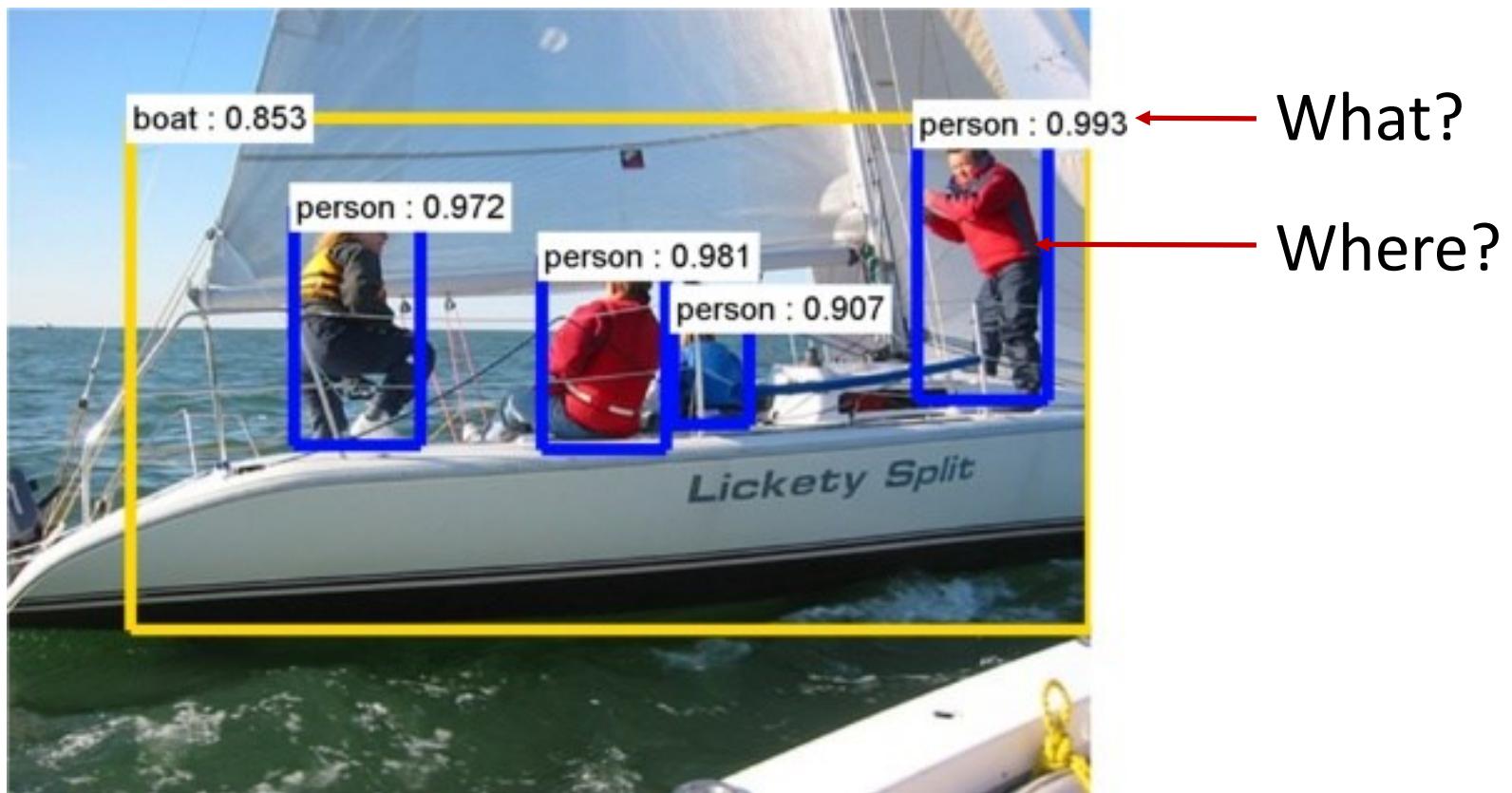
Instance Segmentation



DOG, DOG, CAT

Object Detection

What objects are in an image and where are they?



Steady Progress on Boxes and Masks

- R-CNN [Girshick et al. 2014]
- SPP-net [He et al. 2014]
- Fast R-CNN [Girshick. 2015]
- Faster R-CNN [Ren et al. 2015]
- R-FCN [Dai et al. 2016]
- Feature Pyramid Networks + Faster R-CNN [Lin et al. 2017]
- Mask R-CNN [He et al. 2017]
- Training with Large Minibatches (MegDet) [Peng, Xiao, Li, et al. 2017]
- Cascade R-CNN [Cai & Vasconcelos 2018]
- ...

Steady Progress on Boxes and Masks

2014

Rich feature hierarchies for accurate object detection and semantic segmentation
Tech report (v5)

Ross Girshick Jeff Donahue Trevor Darrell Jitendra Malik
UC Berkeley

2015

Fast R-CNN

Ross Girshick
Microsoft Research

2015

Faster R-CNN: Towards Real-Time Object Detection with Region Proposal Networks

Shaoqing Ren, Kaiming He, Ross Girshick, and Jian Sun

2017

Mask R-CNN

Kaiming He Georgia Gkioxari Piotr Dollár Ross Girshick
Facebook AI Research (FAIR)

2021

Masked Autoencoders Are Scalable Vision Learners

Kaiming He^{*,†} Xinlei Chen^{*} Saining Xie Yanghao Li Piotr Dollár Ross Girshick

^{*}equal technical contribution [†]project lead

Facebook AI Research (FAIR)

Alexander Kirillov^{1,2,4} Eric Mintun² Nikhila Ravi^{1,2} Hanzi Mao² Chloe Rolland³ Laura Gustafson³
Tete Xiao³ Spencer Whitehead Alexander C. Berg Wan-Yen Lo Piotr Dollár⁴ Ross Girshick⁴

¹project lead

²joint first author

³equal contribution

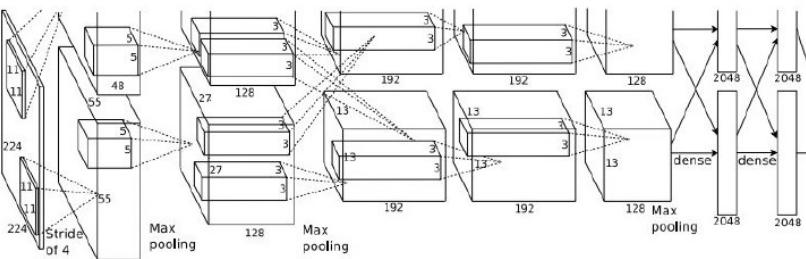
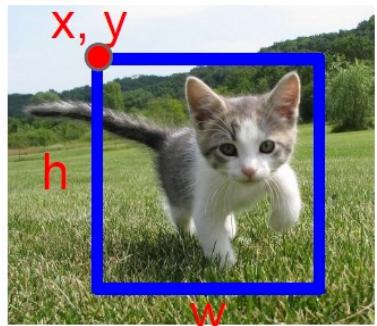
⁴directional lead

Meta AI Research, FAIR

2023

Segment Anything

Object Detection: Single Object



Treat localization as a
regression problem!

Fully
Connected:
4096 to 1000

Vector: Fully
Connected:
4096 to 4

Class Scores
Cat: 0.9
Dog: 0.05
Car: 0.01

Multitask Loss

Box
Coordinates
(x, y, w, h)

Correct label:
Cat

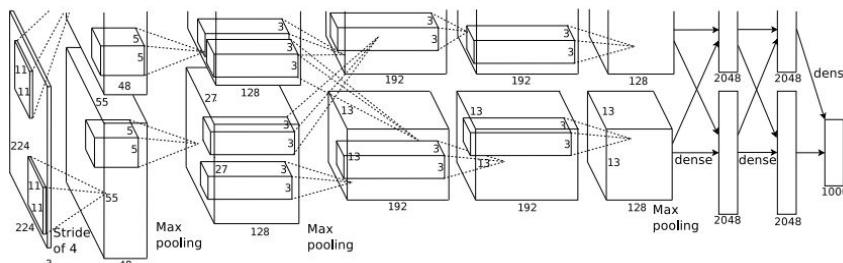
CrossEntropy
Loss

+ → Loss

→ L2 Loss

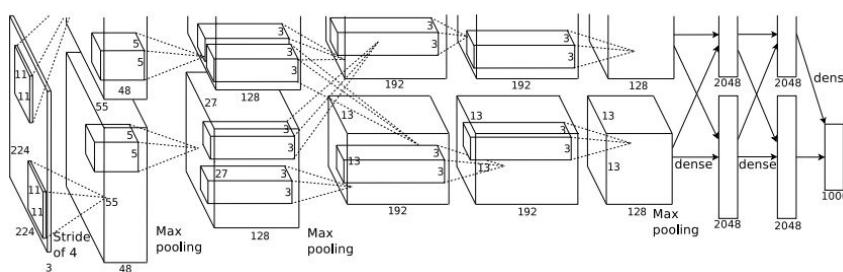
Correct box:
(x', y', w', h')

Object Detection: Multiple Objects



CAT: (x, y, w, h)

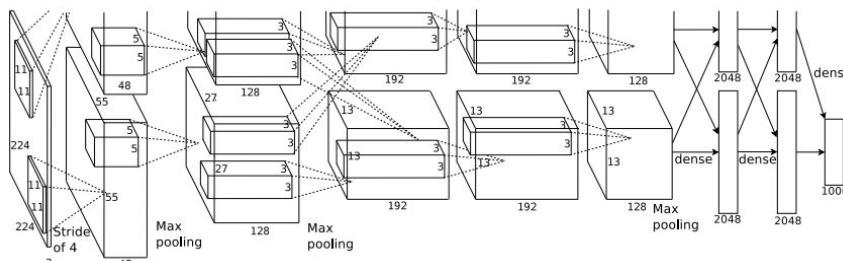
Each image needs a
different number of outputs!



DOG: (x, y, w, h)

DOG: (x, y, w, h)

CAT: (x, y, w, h)



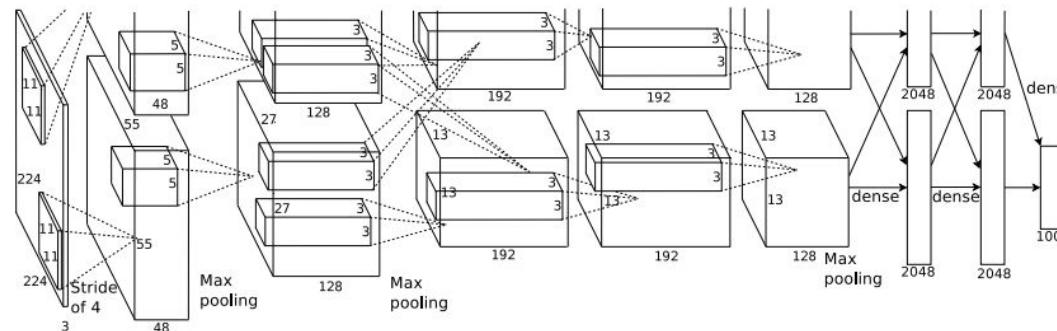
DUCK: (x, y, w, h)

DUCK: (x, y, w, h)

...

Object Detection: Multiple Objects

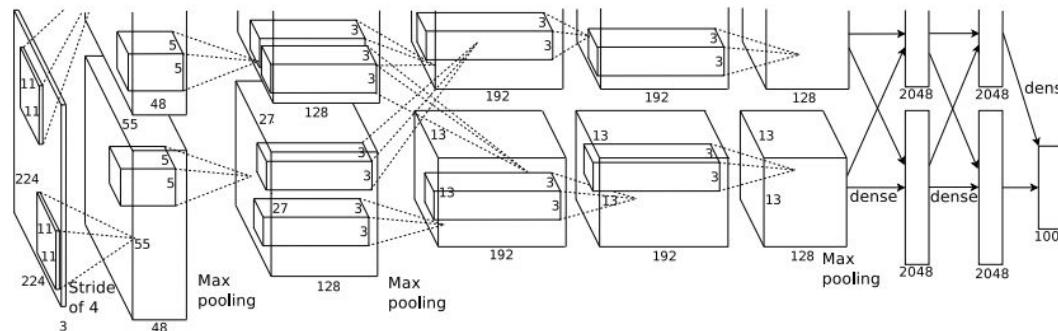
Apply a CNN to many different crops of the image, CNN classifies each crop as object or background



Dog? No
Cat? No
Background? YES

Object Detection: Multiple Objects

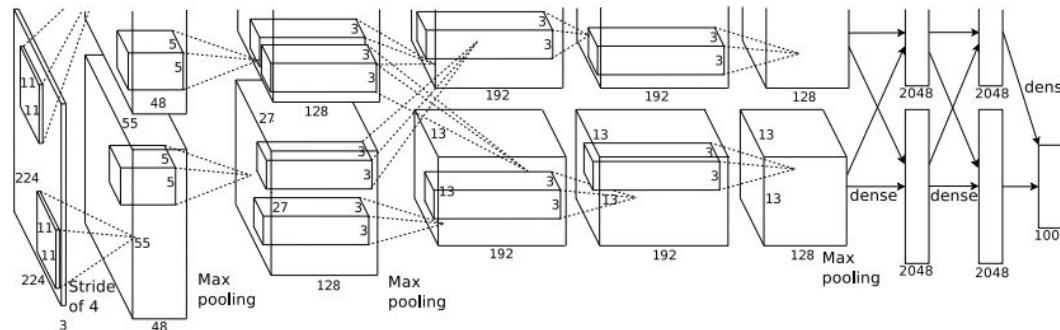
Apply a CNN to many different crops of the image, CNN classifies each crop as object or background



Dog? Yes
Cat? No
Background? No

Object Detection: Multiple Objects

Apply a CNN to many different crops of the image, CNN classifies each crop as object or background

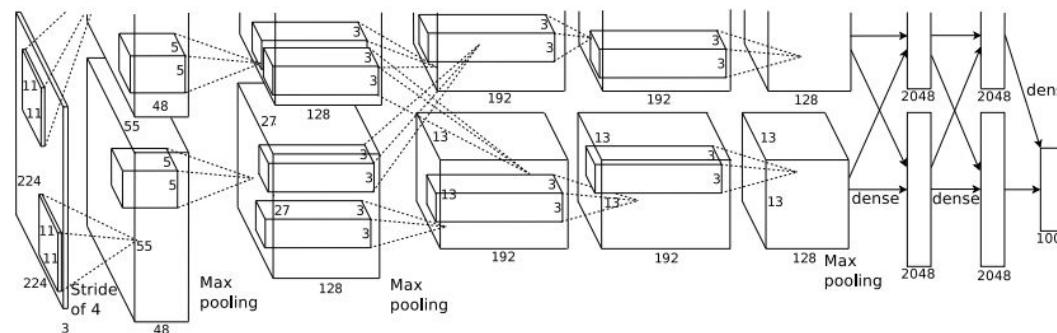


Dog? Yes
Cat? No
Background? No

Object Detection: Multiple Objects



Apply a CNN to many different crops of the image, CNN classifies each crop as object or background

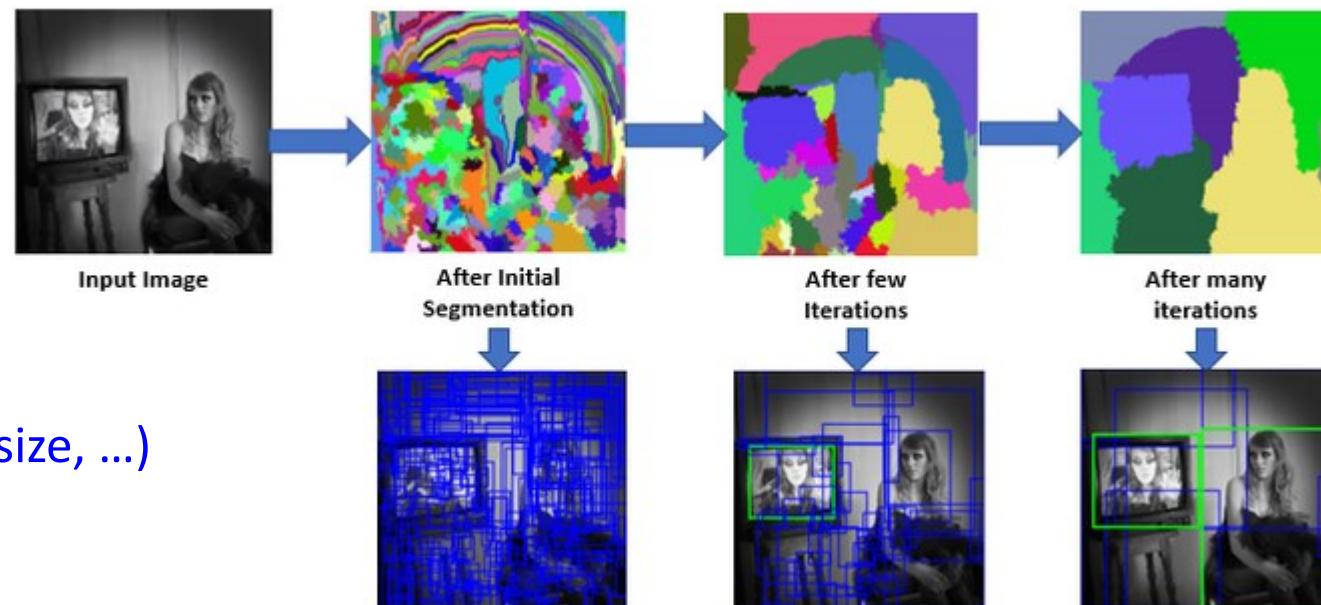


Dog? No
Cat? Yes
Background? No

Problem: Need to apply CNN to huge number of locations, scales, and aspect ratios, very computationally expensive!

Region Proposals: Selective Search

- Find image regions that are likely to contain objects
- Relatively fast to run; e.g. Selective Search gives 2000 region proposals in a few seconds on CPU



R-CNN (Region-based Convolutional Neural Net)

Per-image computation



I:

Per-region computation

R-CNN (Region-based Convolutional Neural Net)

Per-image computation

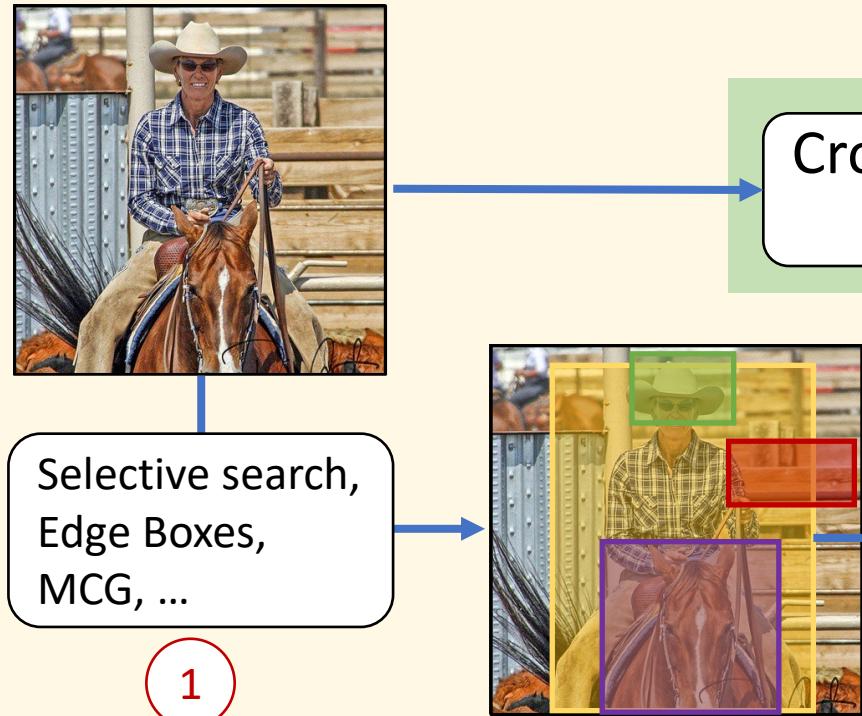


Per-region computation

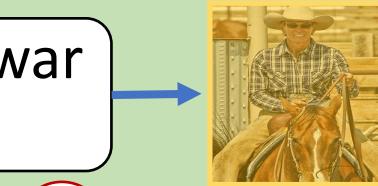
1. Use an off-the-shelf *Region of Interest* (RoI) proposal algorithm (~2k proposals per image)

R-CNN

Per-image computation



Per-region computation for each $r_i \in r(I)$

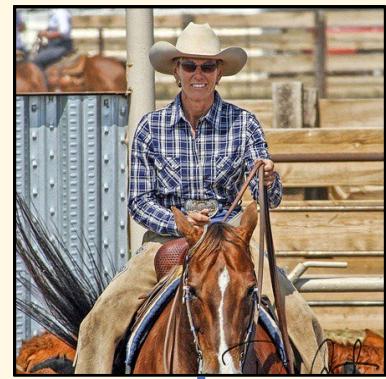


2

2. Crop and warp each proposal image window to obtain a fixed-size network input (224x224)

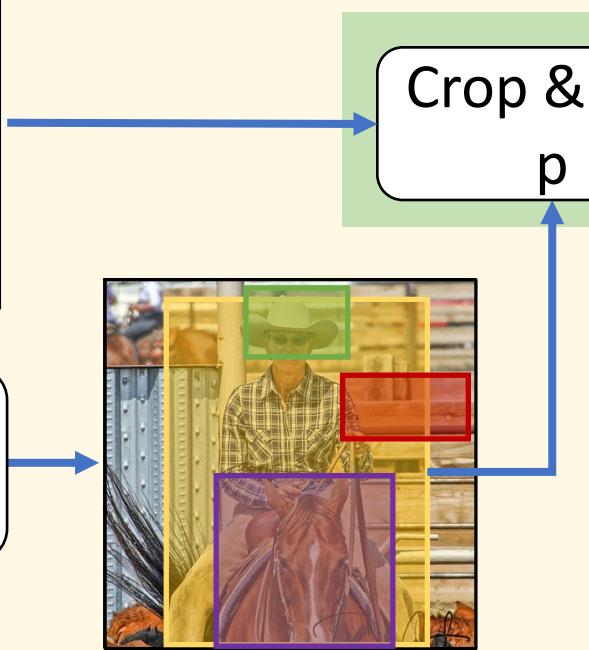
R-CNN

Per-image computation



Selective search,
Edge Boxes,
MCG, ...

1



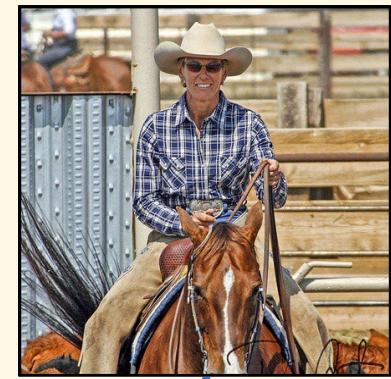
Per-region computation for each $r_i \in r(I)$



3. Forward propagate the fixed-size network input to get a feature representation (ImageNet-pretrained)

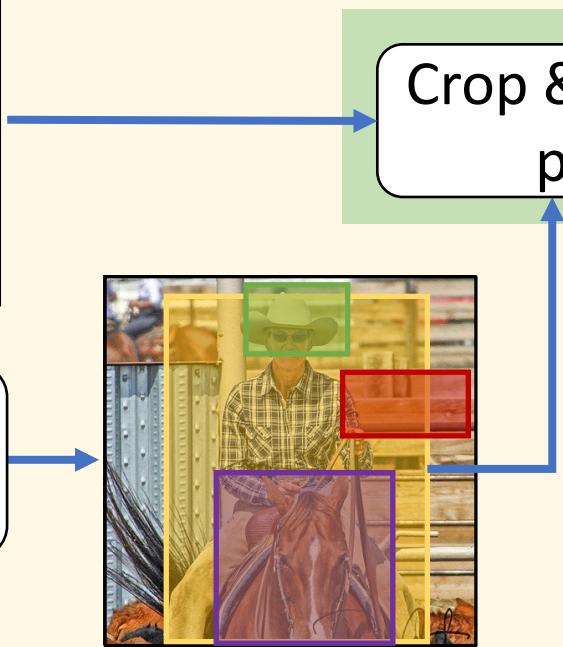
R-CNN

Per-image computation



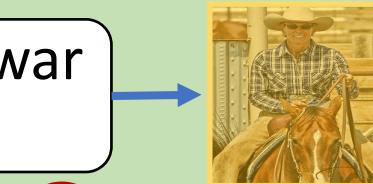
Selective search,
Edge Boxes,
MCG, ...

1



Crop & war
 p

2



ConvNet(r_i)

3



Linear classifier

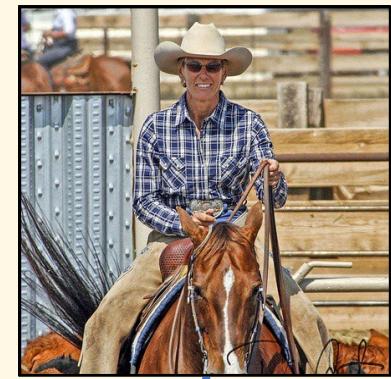
4

4. Object classification

Per-region computation for each $r_i \in r(I)$

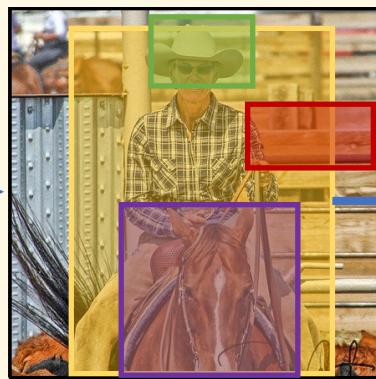
R-CNN

Per-image computation



Selective search,
Edge Boxes,
MCG, ...

1



Crop & war
 p

2



Per-region computation for each $r_i \in r(I)$

ConvNet(r_i)

3

Linear classifier

4

Box regressor

5

5. Refine proposal localization
with bounding-box regression

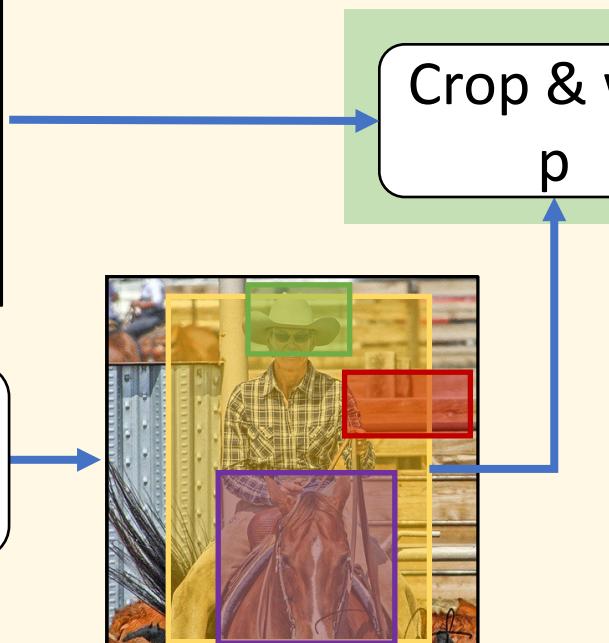
Slow R-CNN

Per-image computation

Problem: Very slow!
Need to do $\sim 2k$ independent forward passes for each image!

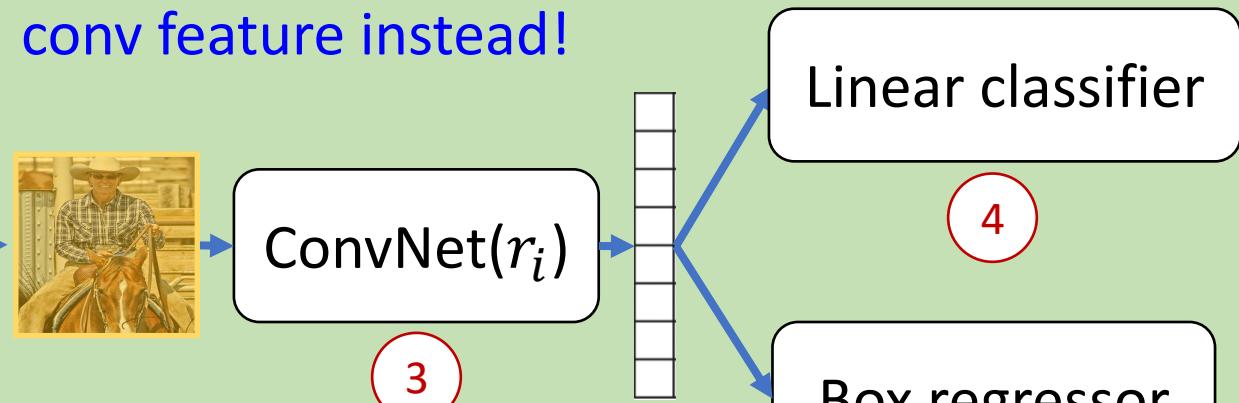


Selective search,
Edge Boxes,
MCG, ...



Per-region computation for each $r_i \in r(I)$

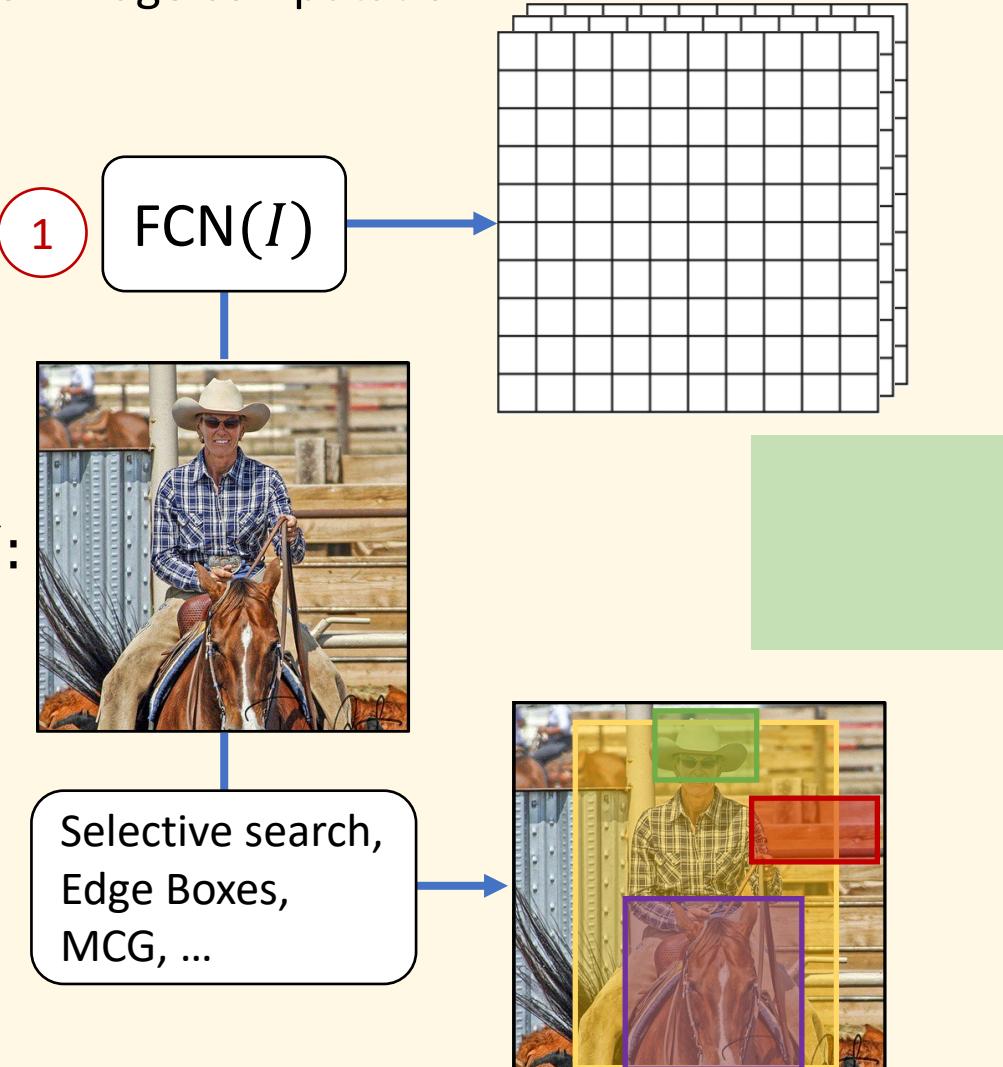
Idea: Pass the image through convnet before cropping! Crop the conv feature instead!



5. Refine proposal localization with bounding-box regression

Fast R-CNN

Per-image computation

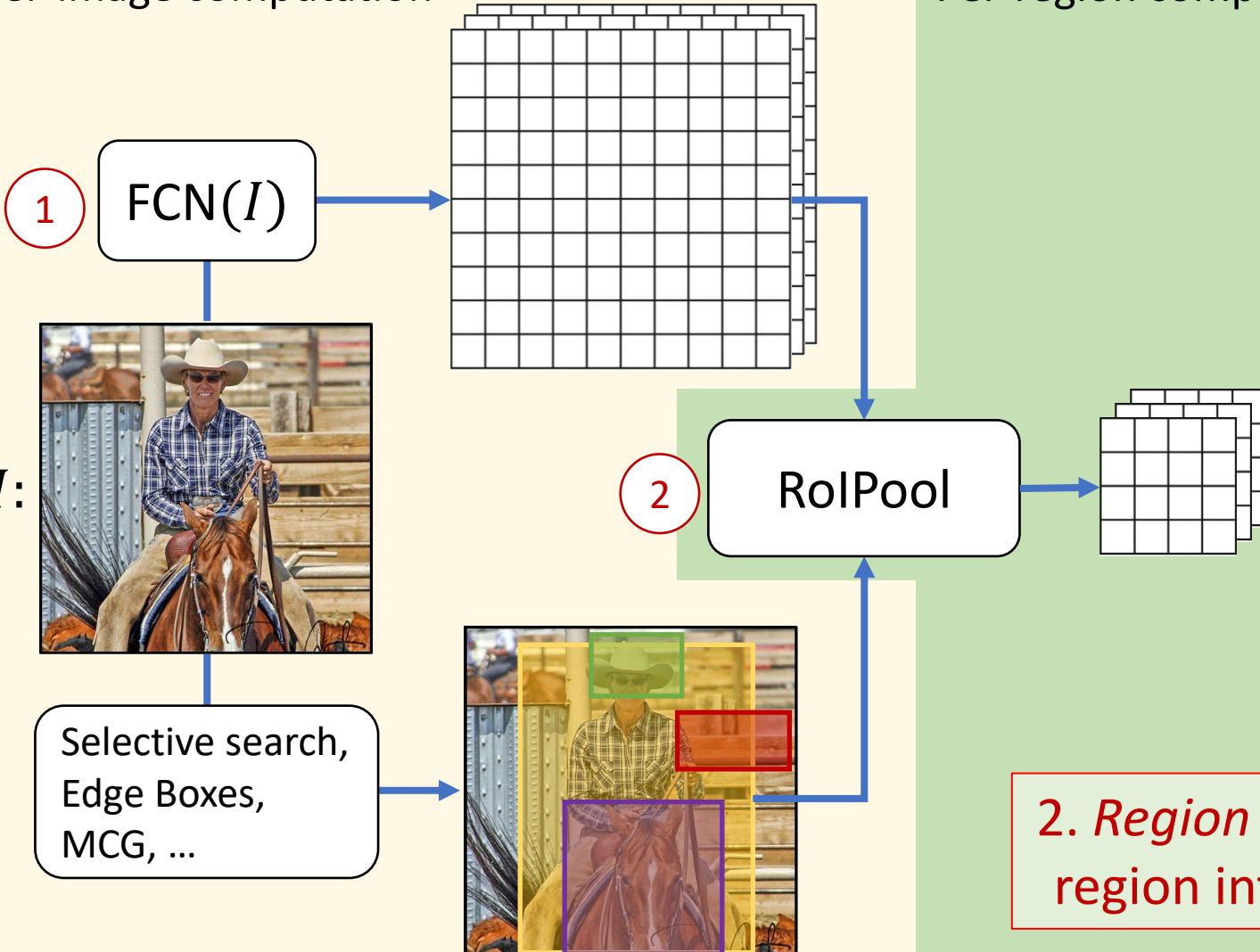


Per-region computation for each $r_i \in r(I)$

1. Fully convolutional network (FCN) maps the image to a lower resolution spatial feature map

Fast R-CNN

Per-image computation

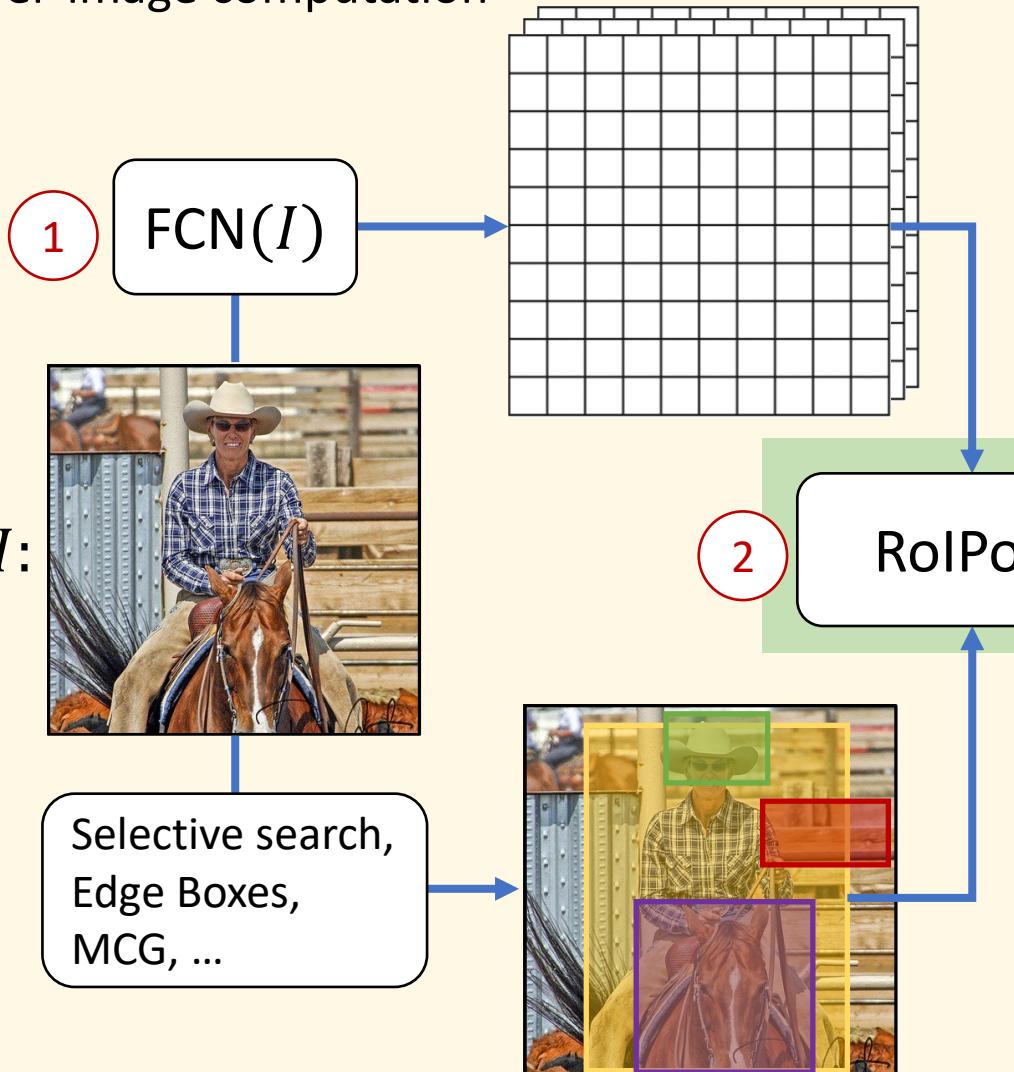


Per-region computation for each $r_i \in r(I)$

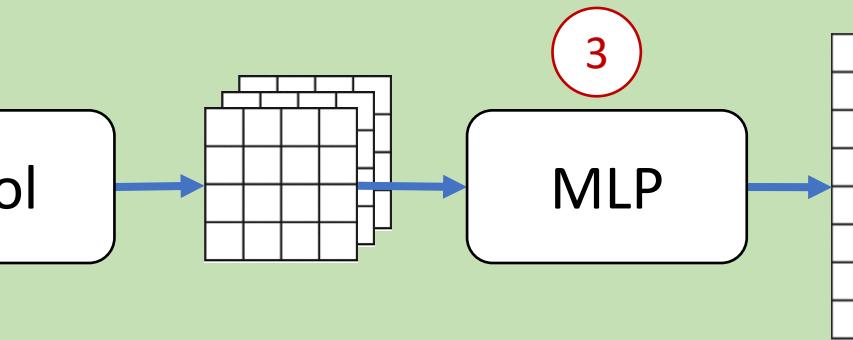
2. *Region of interest (RoI) pooling converts each region into a fixed dimensional representation*

Fast R-CNN

Per-image computation



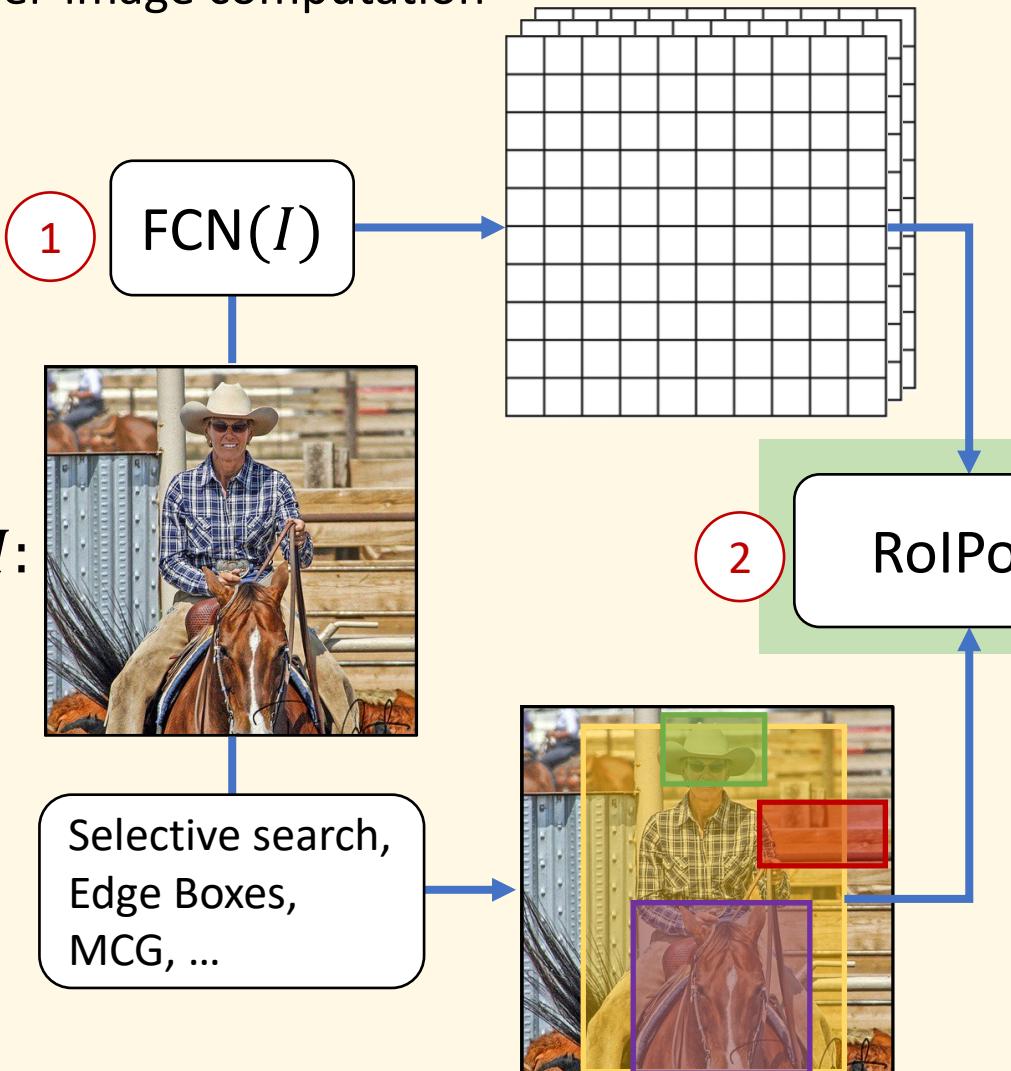
Per-region computation for each $r_i \in r(I)$



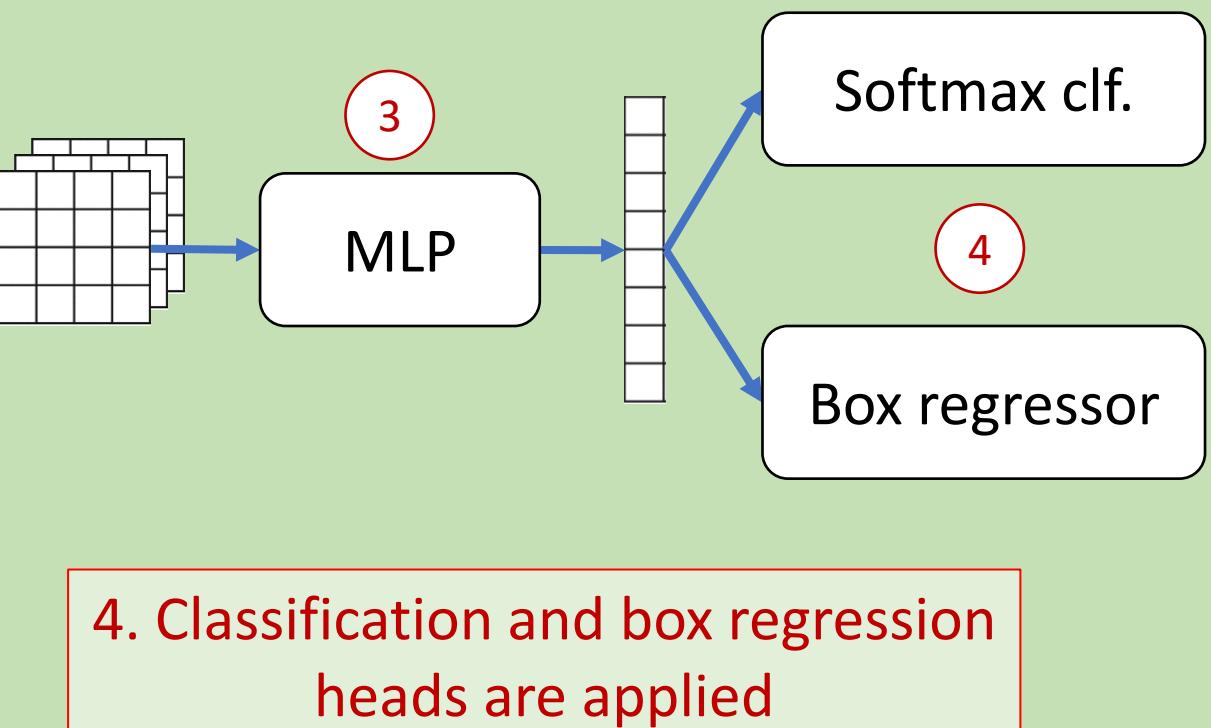
3. A lightweight MLP processes each region feature map
(Alternatively, the MLP could be a small ConvNet, etc.)

Fast R-CNN

Per-image computation

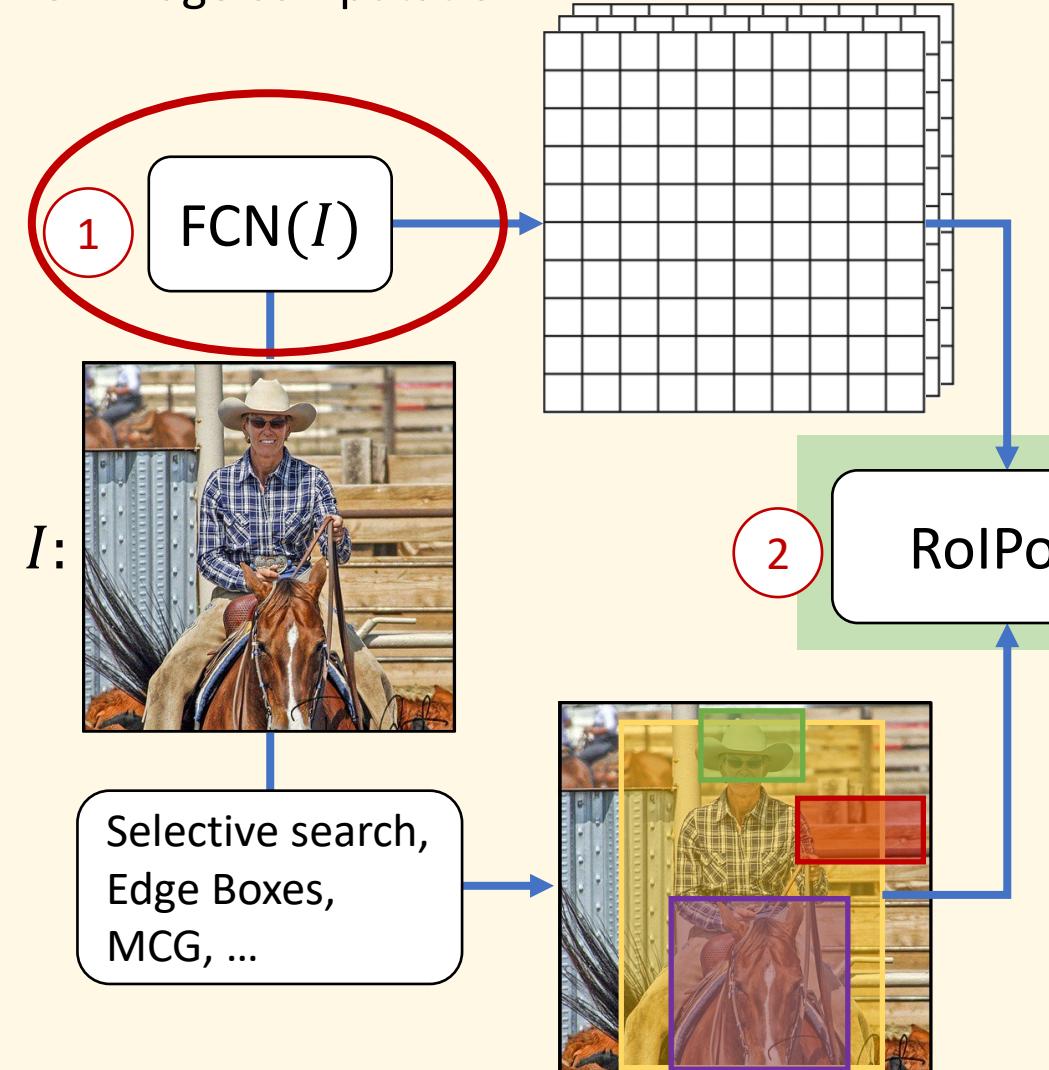


Per-region computation for each $r_i \in r(I)$

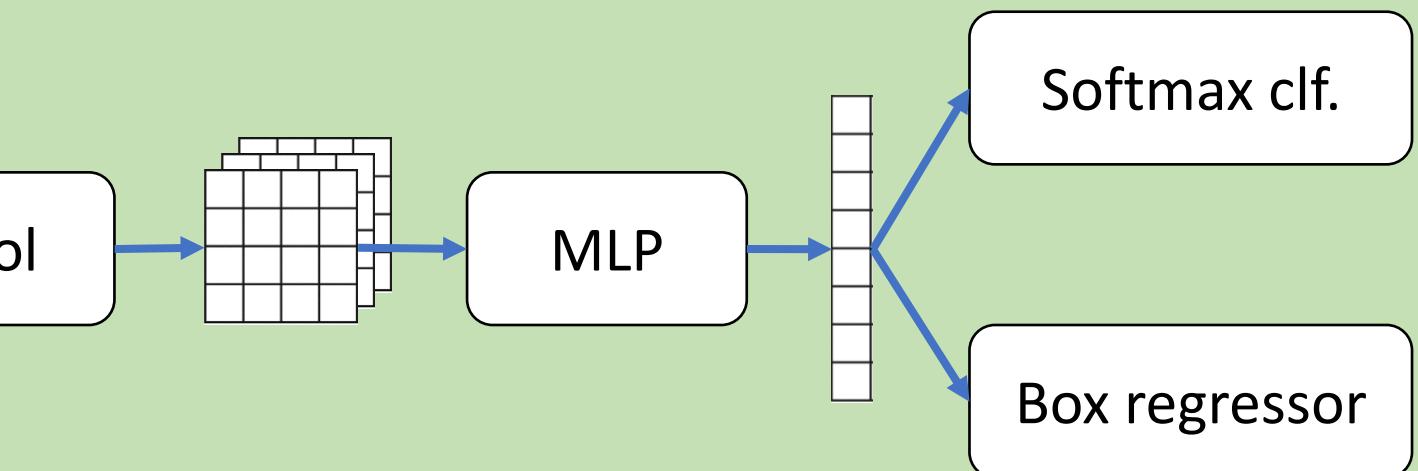


Fast R-CNN

Per-image computation



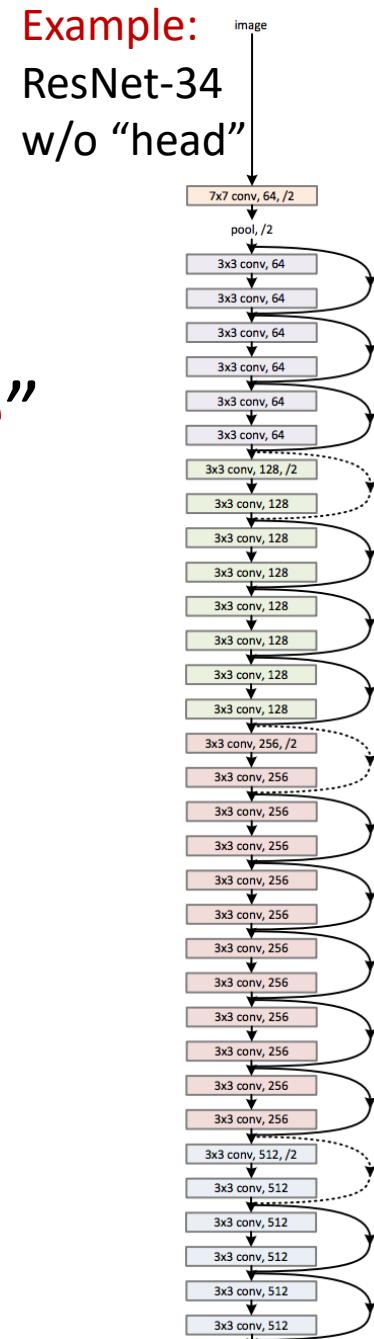
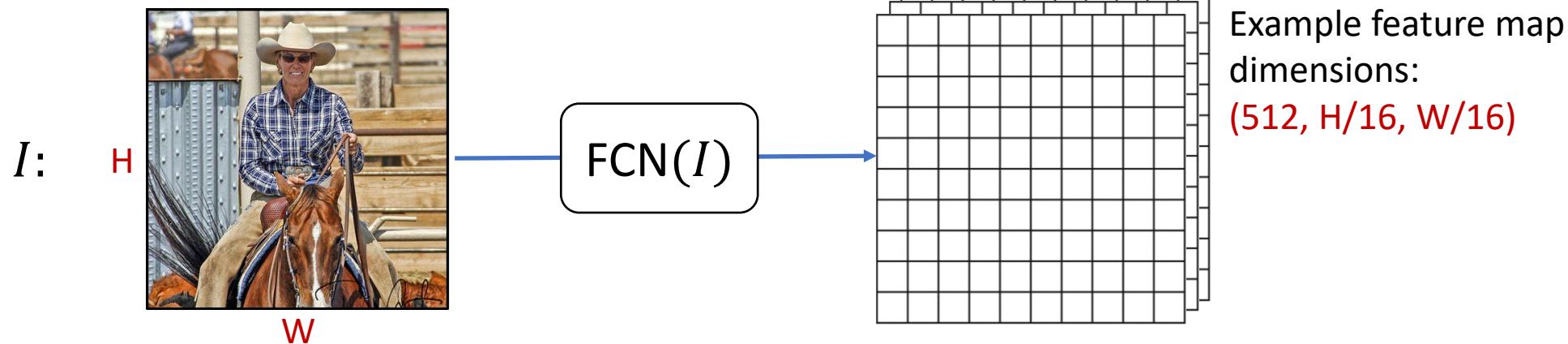
Per-region computation for each $r_i \in r(I)$



Whole-image, Fully Conv. Network (FCN)

Use any standard ConvNet as the “backbone architecture”

- AlexNet, VGG, ResNet, Inception, Inception-ResNet, ResNeXt, DenseNet, NAS*, ...
- Remove global pooling / FC
- Output spatial dims are proportional to input spatial dims



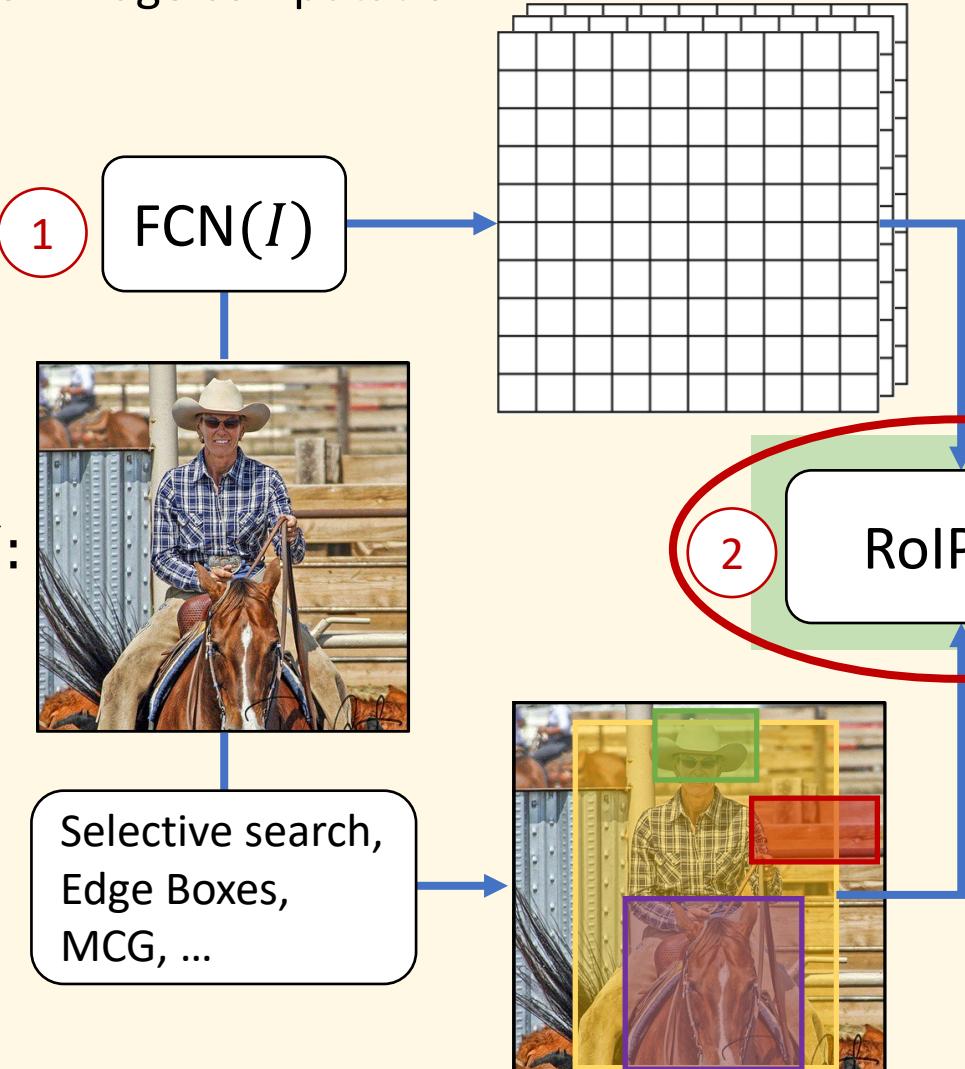
The Engine of Recognition

A good network lift all boats (“features matter”)

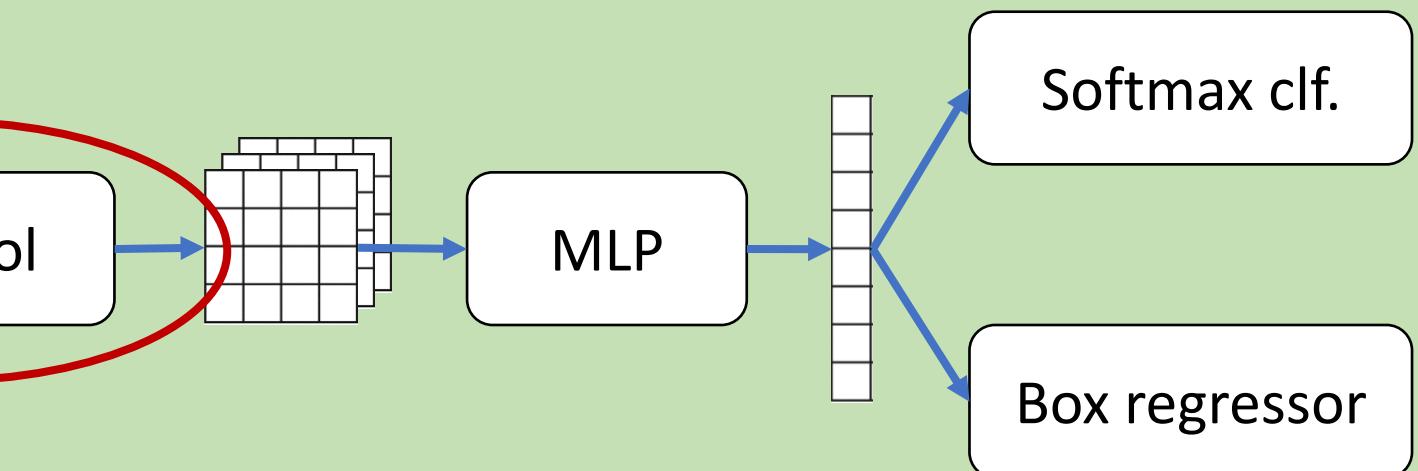
- AlexNet
- VGG
- GoogleNet / Inception
- ResNet
- ResNeXt
- NAS-derived networks
- ...

Fast R-CNN

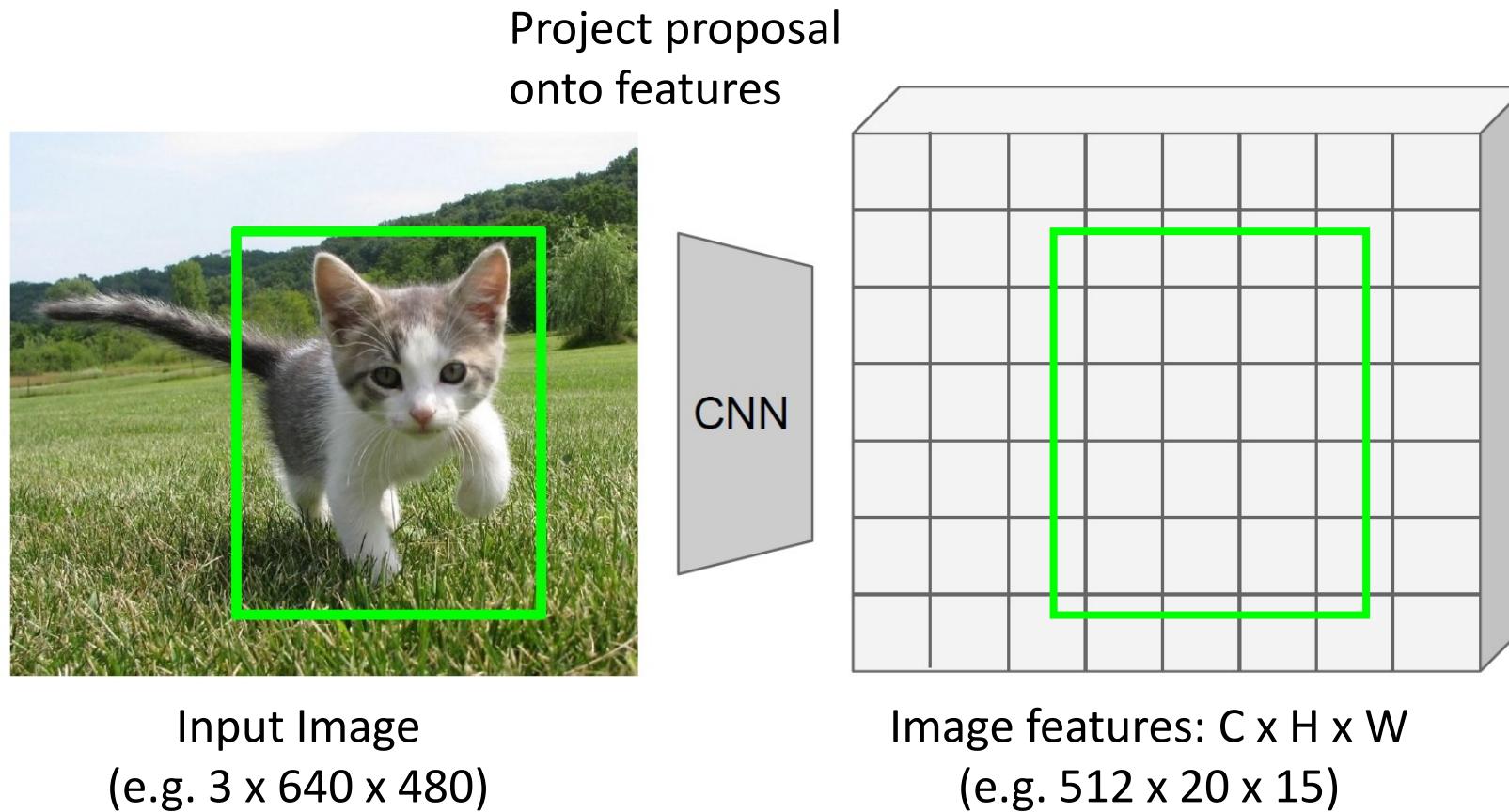
Per-image computation



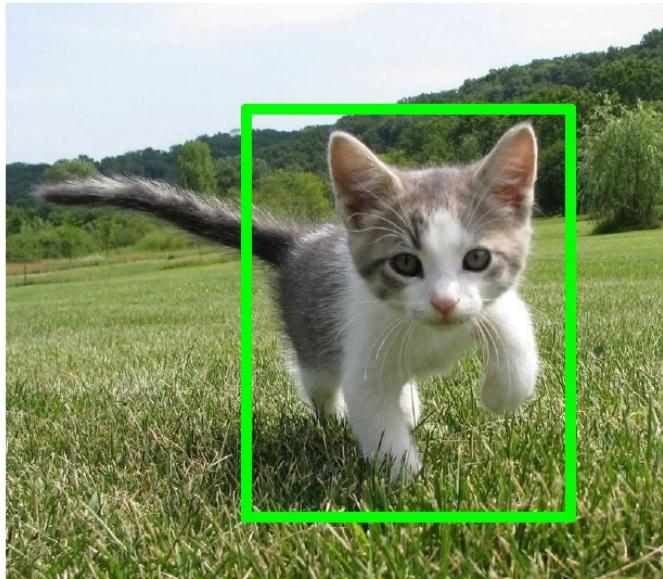
Per-region computation for each $r_i \in r(I)$



Cropping Features: RoIPool



Cropping Features: RoIPool



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

Project proposal
onto features



“Snap” to grid cells

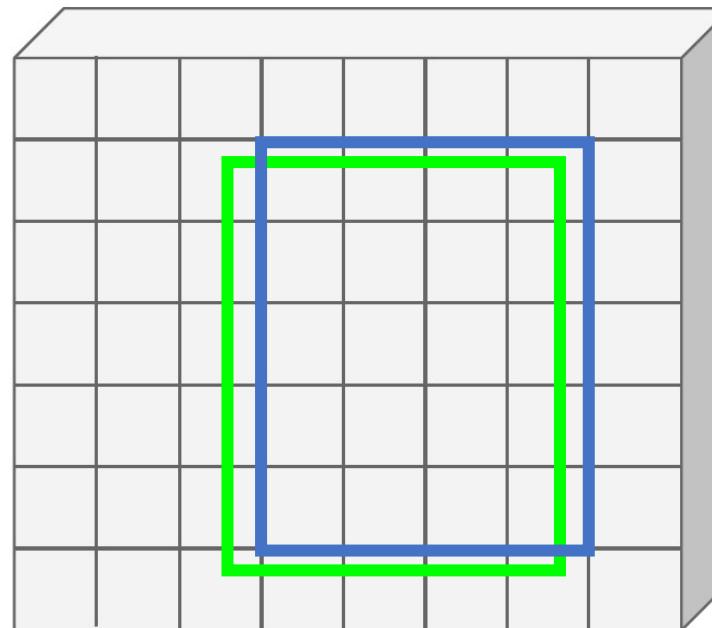
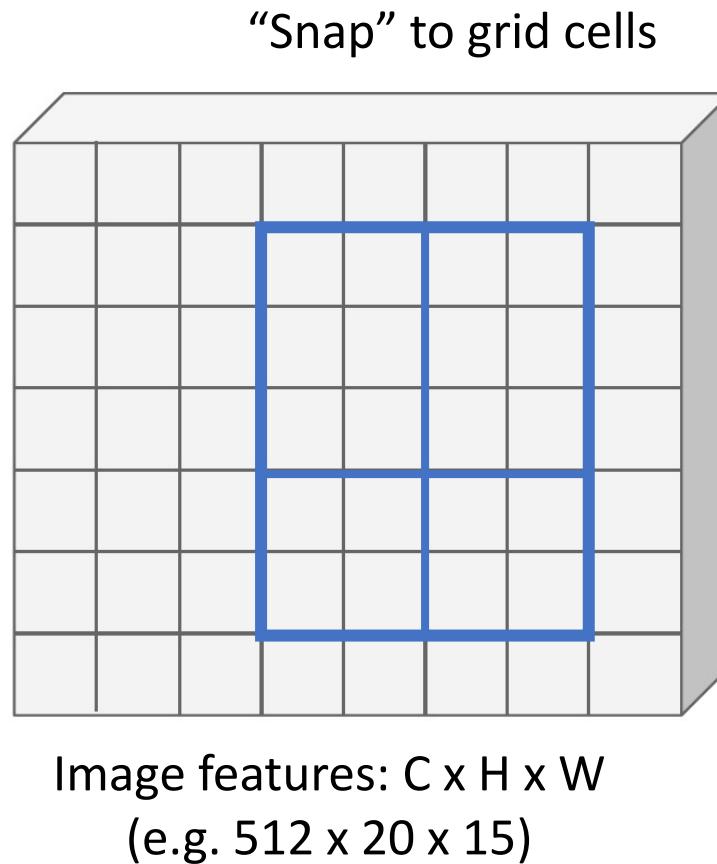
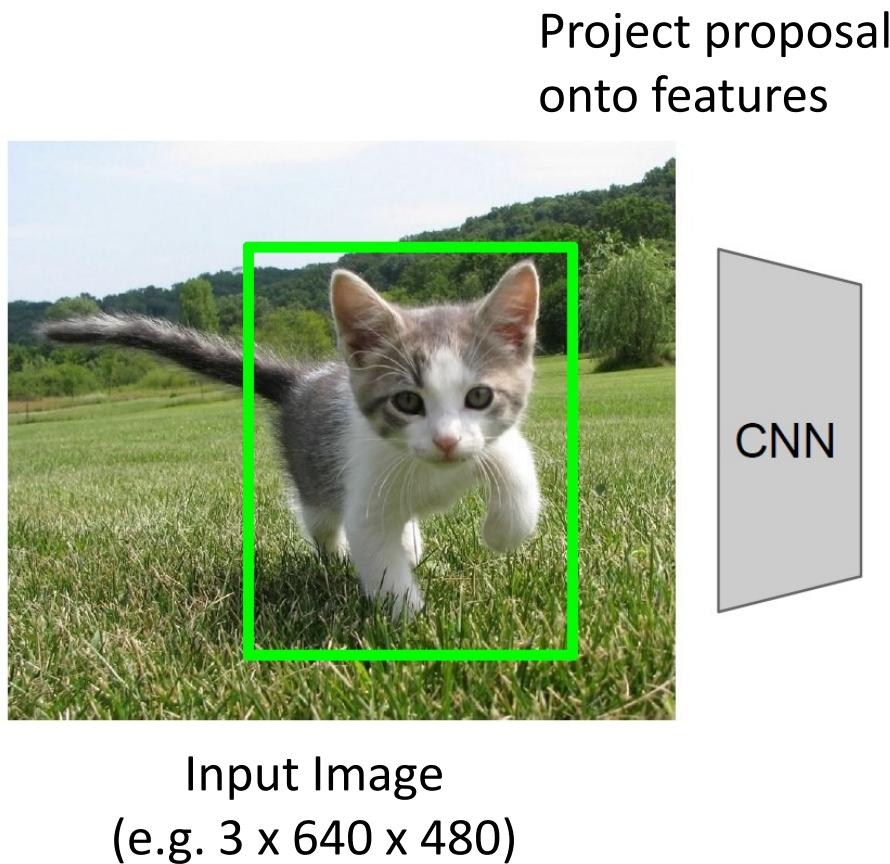


Image features: $C \times H \times W$
(e.g. $512 \times 20 \times 15$)

Divide into 2×2
grid of (roughly)
equal subregions

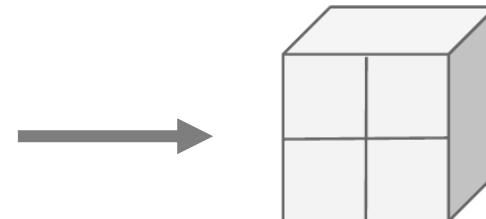
Q: how do we resize the
 $512 \times 5 \times 4$ region to,
e.g., a $512 \times 2 \times 2$ tensor?

Cropping Features: RoIPool



Divide into 2×2 grid of (roughly) equal subregions

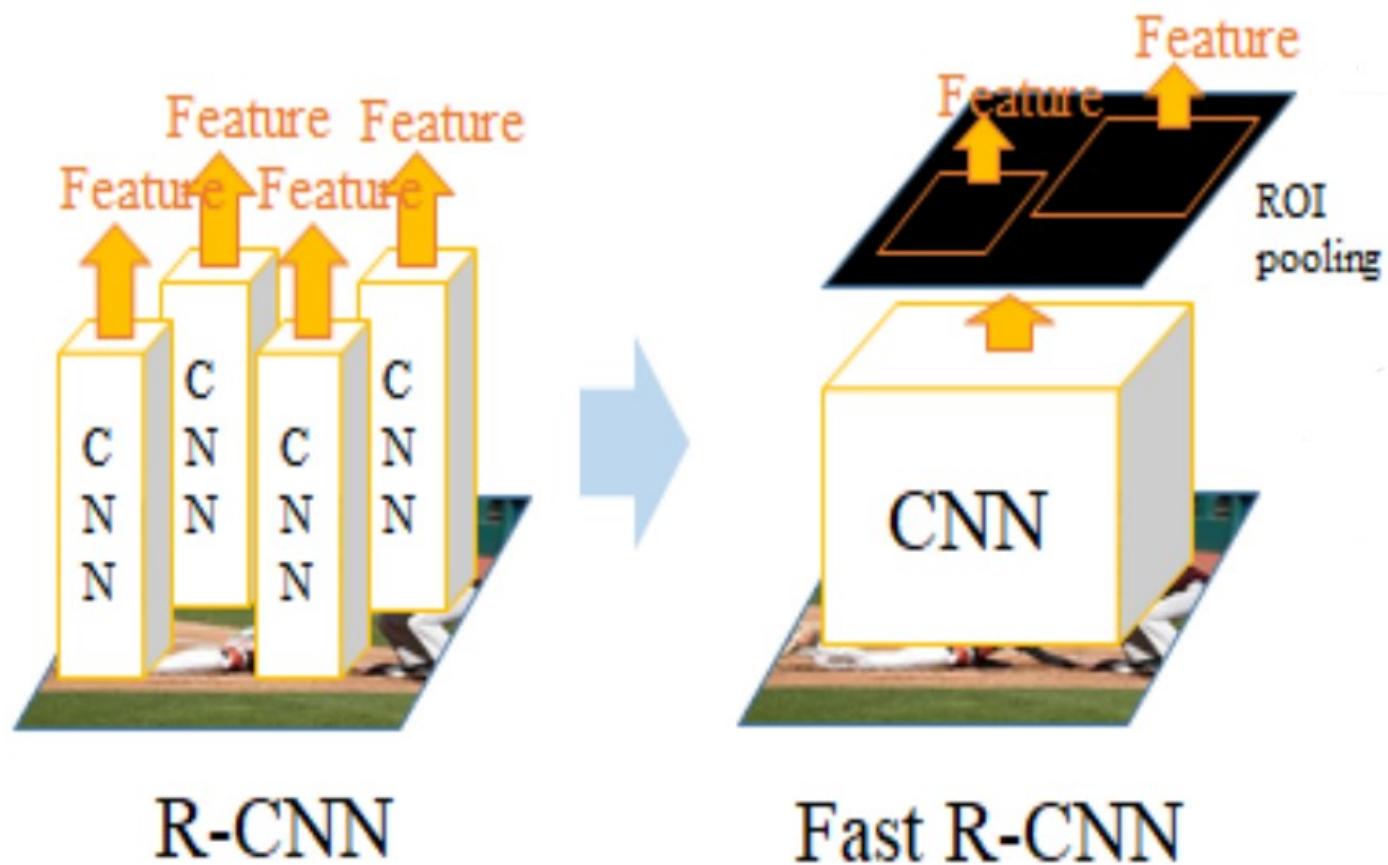
Max-pool within each subregion



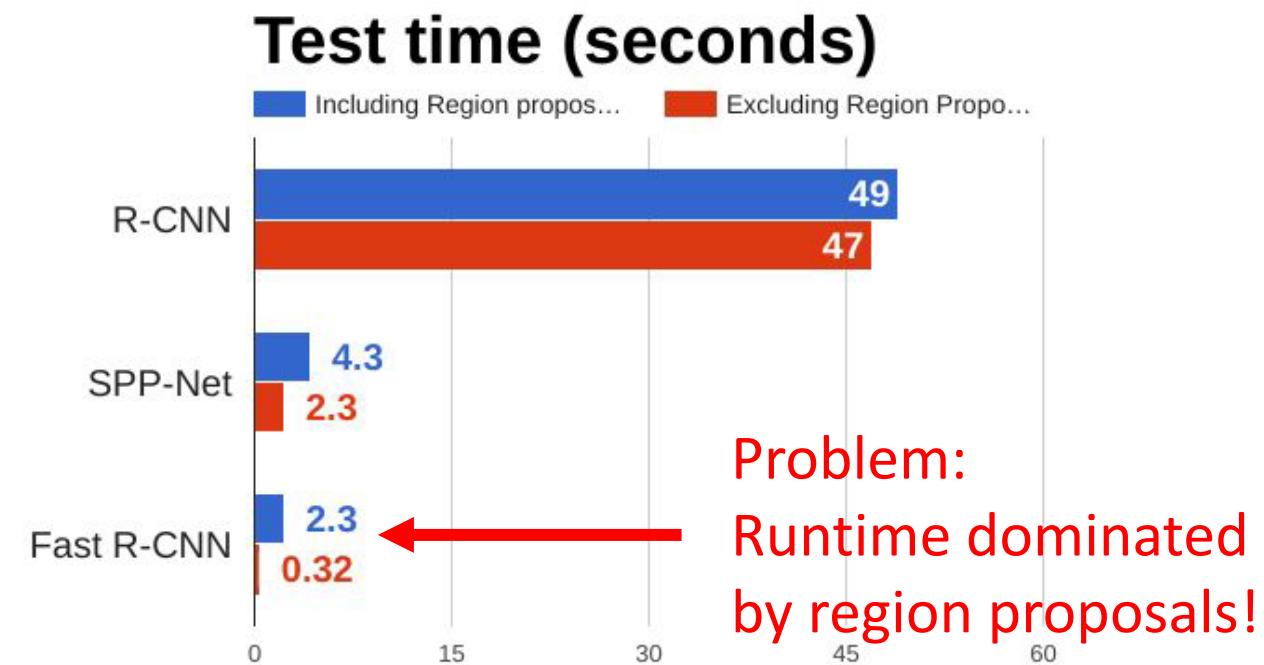
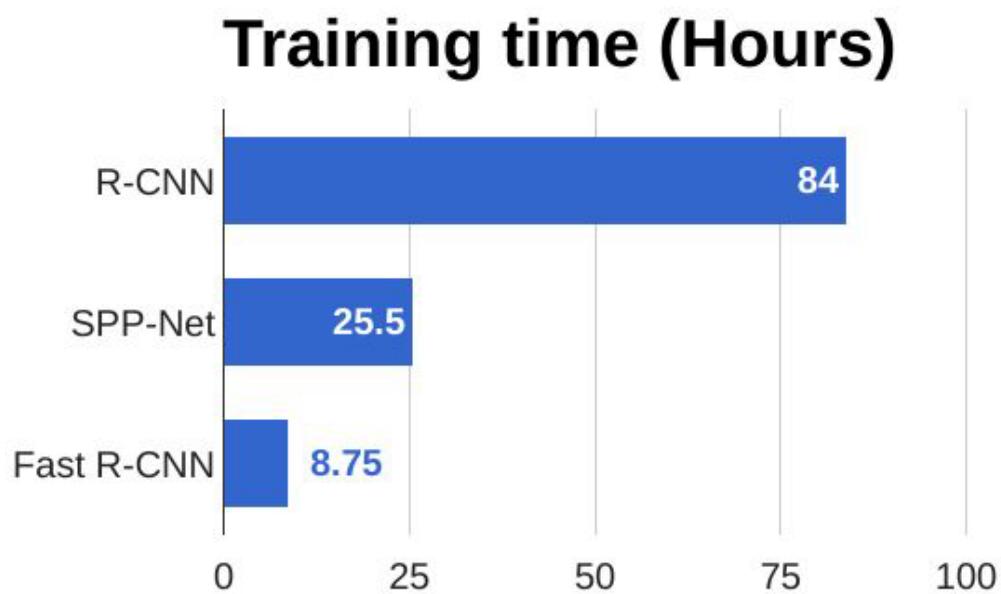
Region features
(here $512 \times 2 \times 2$;
In practice e.g $512 \times 7 \times 7$)

Region features always the same size even if input regions have different sizes!

R-CNN vs Fast R-CNN



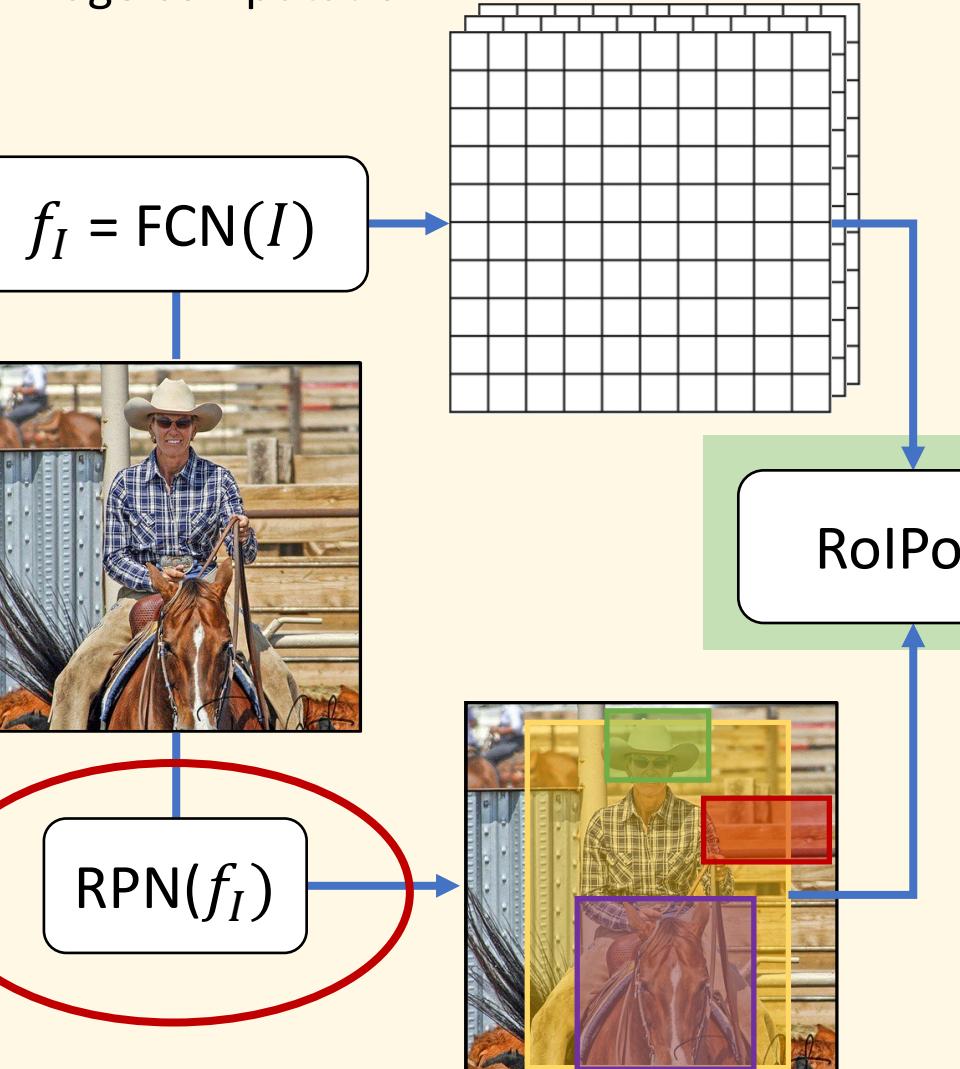
R-CNN vs Fast R-CNN



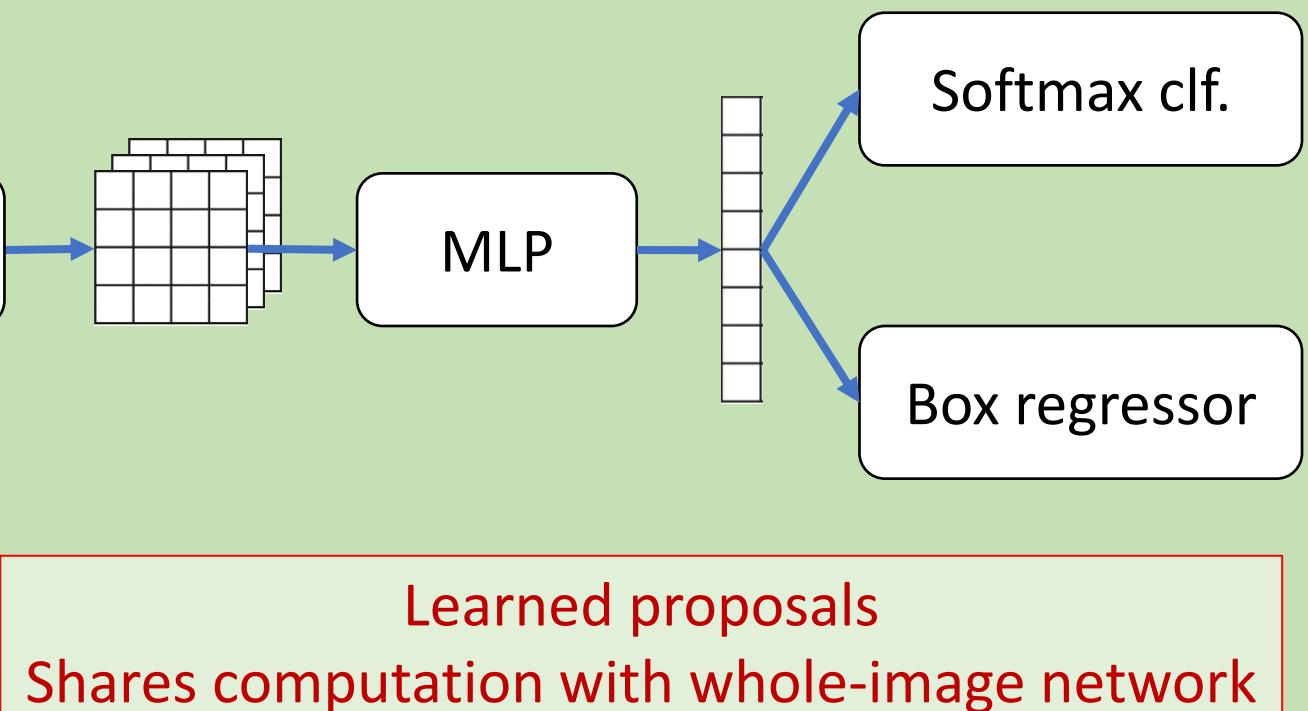
Girshick et al, "Rich feature hierarchies for accurate object detection and semantic segmentation", CVPR 2014.
He et al, "Spatial pyramid pooling in deep convolutional networks for visual recognition", ECCV 2014
Girshick, "Fast R-CNN", ICCV 2015

Faster R-CNN

Per-image computation

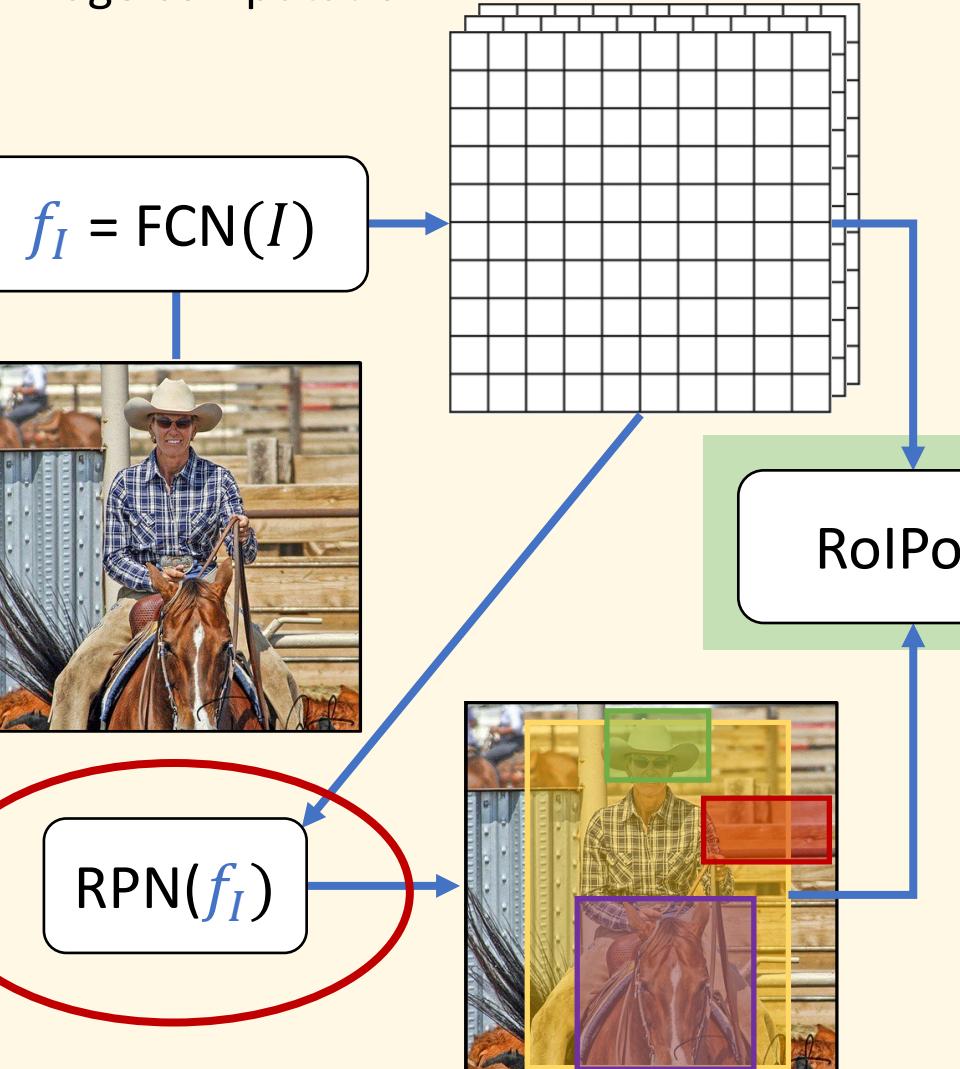


Per-region computation for each $r_i \in r(I)$

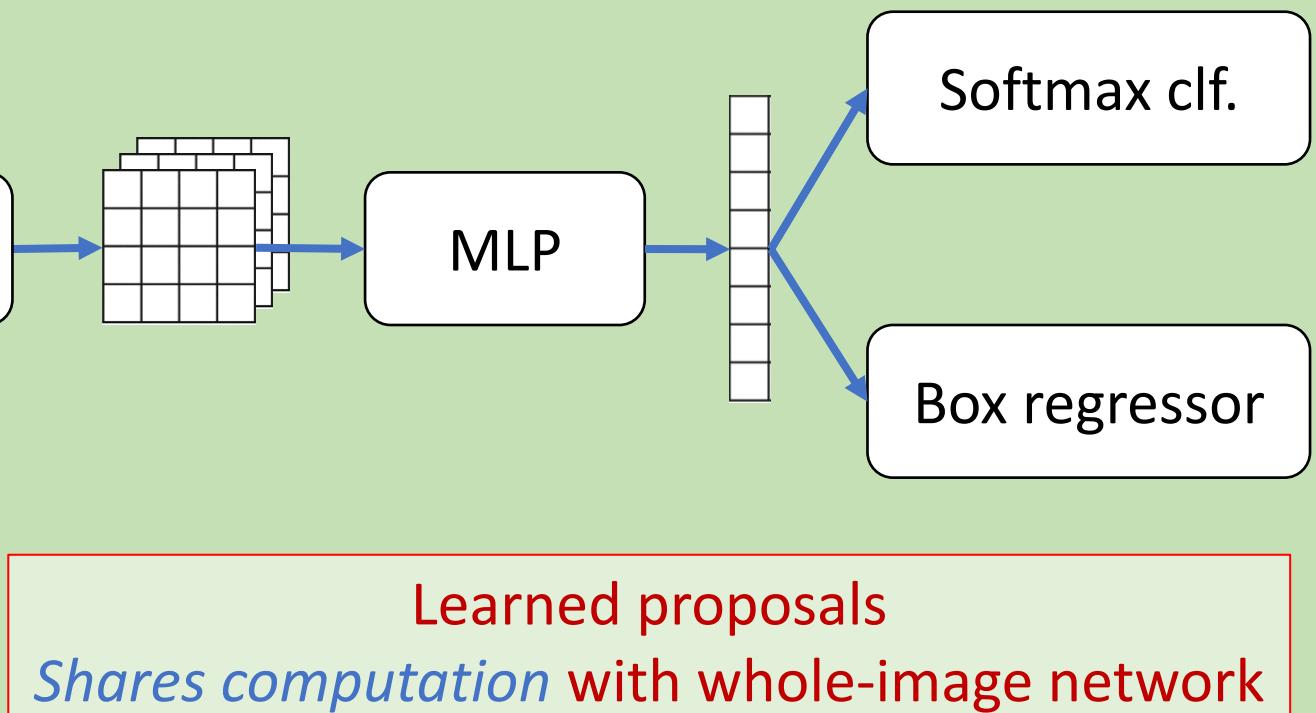


Faster R-CNN

Per-image computation



Per-region computation for each $r_i \in r(I)$



RPN: Region Proposal Network



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

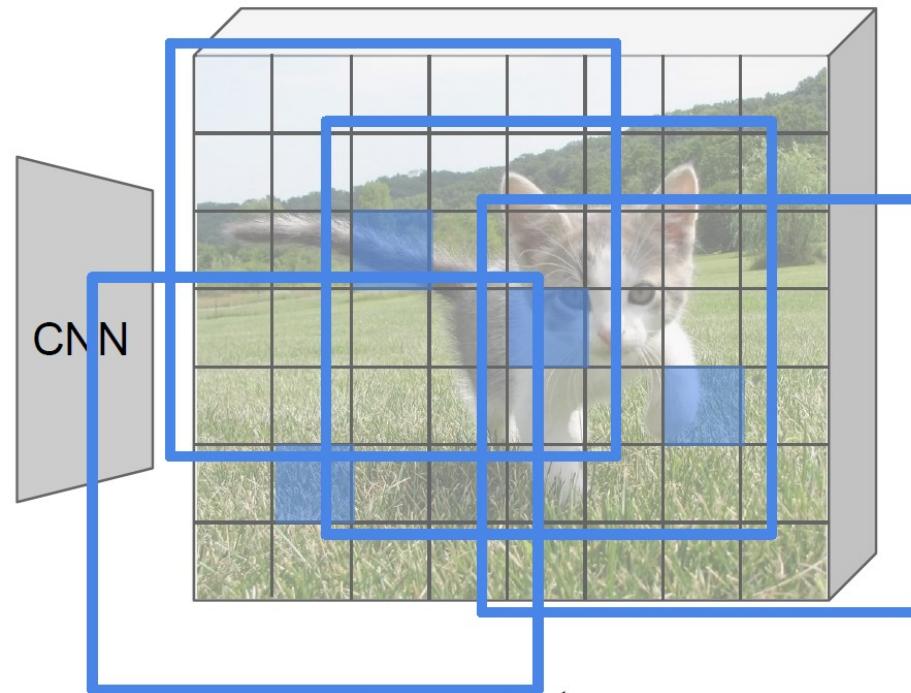
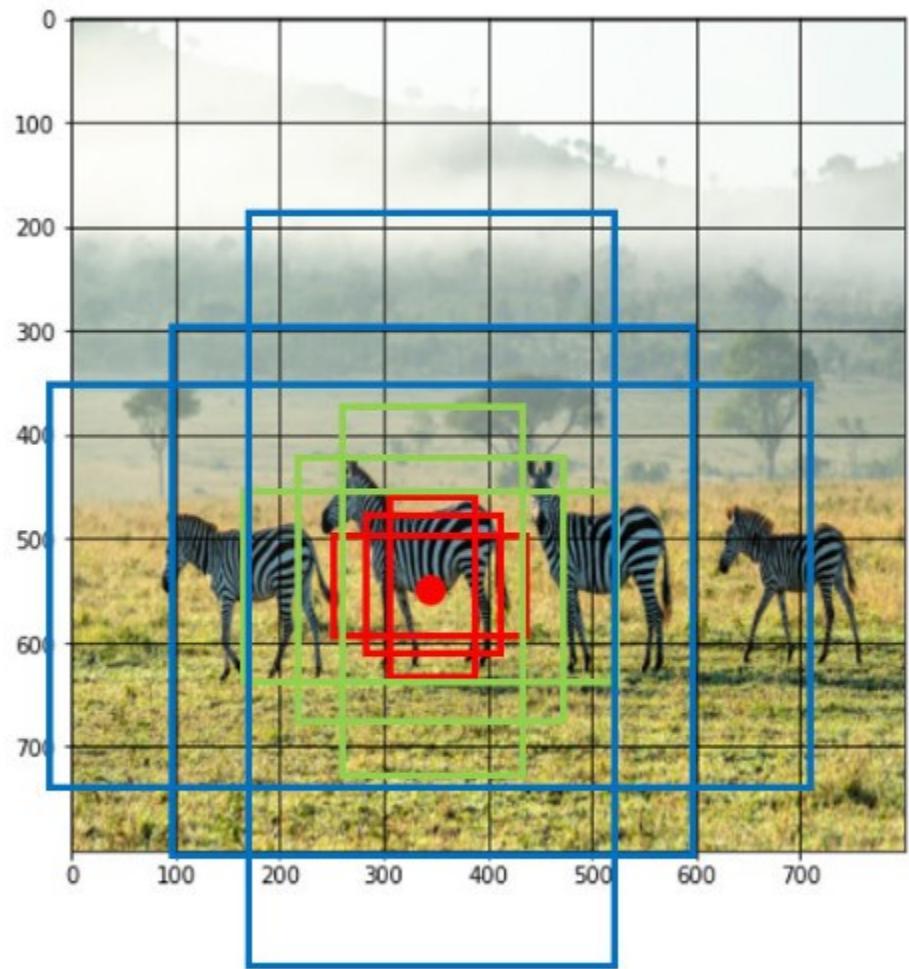
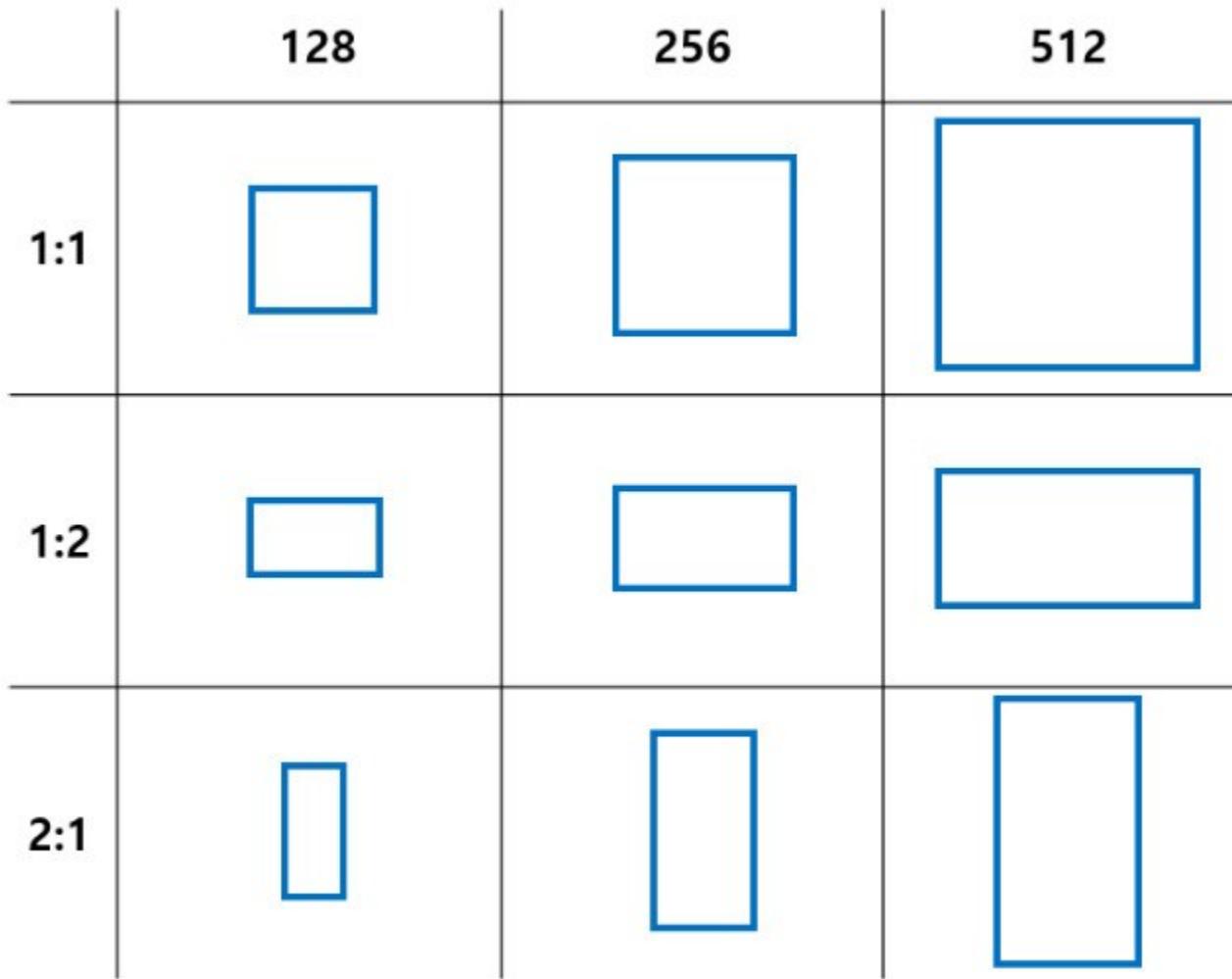


Image features: $C \times H \times W$
(e.g. $512 \times 20 \times 15$)

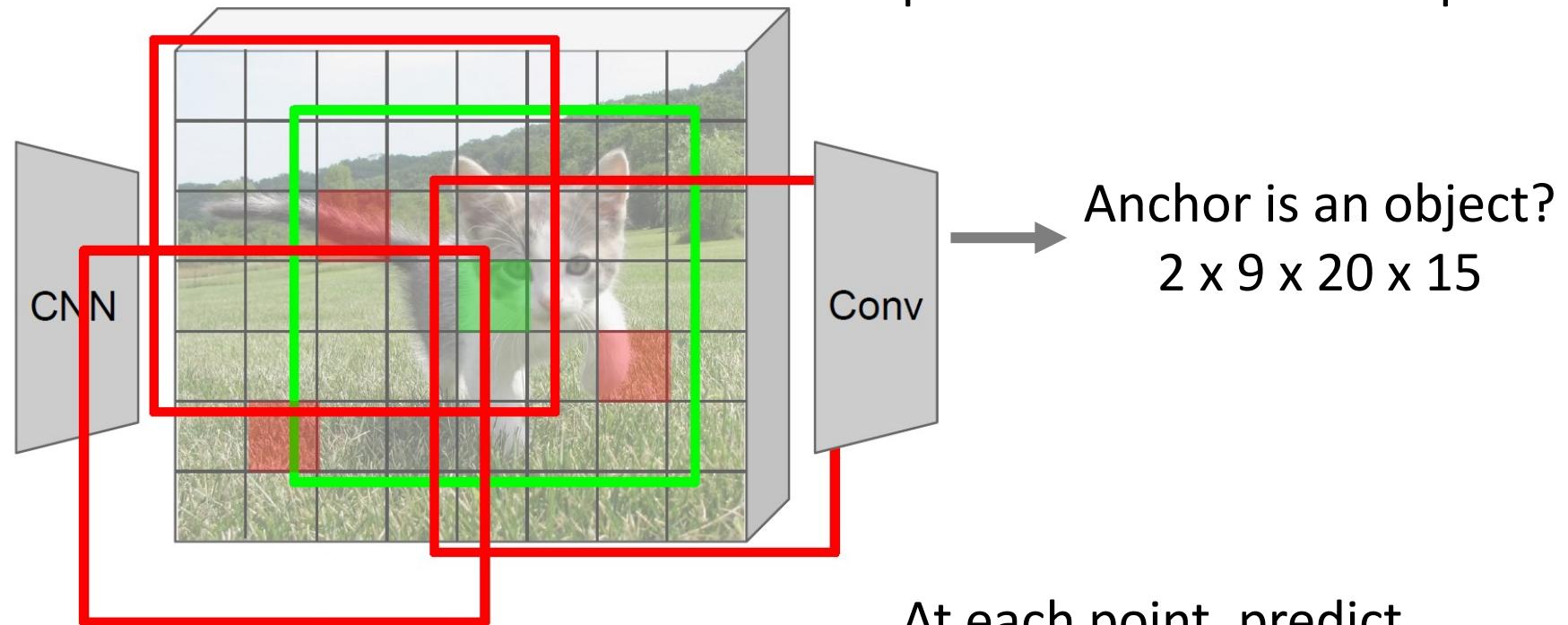
RPN: Region Proposal Network



RPN: Region Proposal Network



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

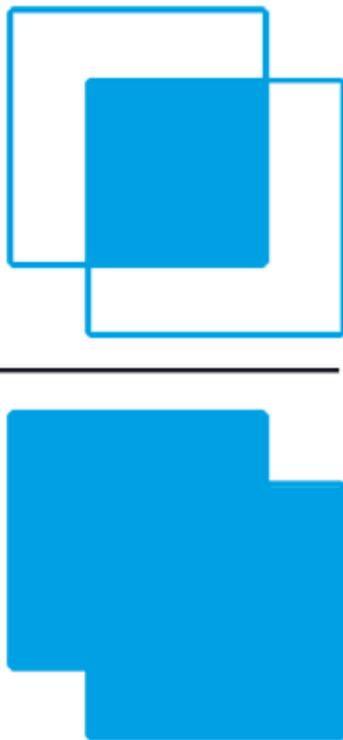


At each point, predict whether the corresponding anchor contains an object (binary classification)

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

Intersection over Union

$$\text{IoU} = \frac{\text{Area of Overlap}}{\text{Area of Union}}$$



IoU of ground truth box and anchor box

Highest IoU

Positive box

$\text{IoU} > 0.7$

Positive box

$\text{IoU} < 0.3$

Negative box

Others

do not use in training

RPN: Region Proposal Network



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

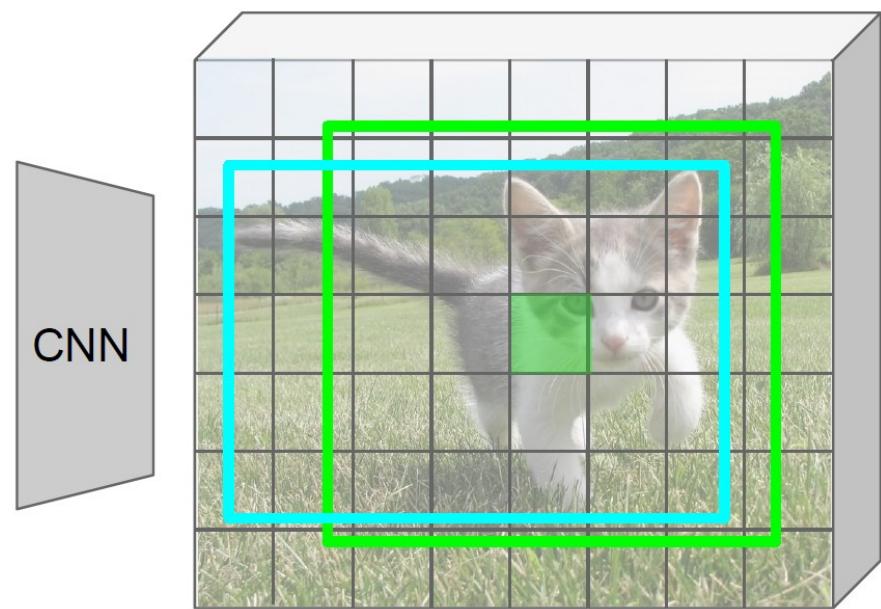


Image features: C x H x W
(e.g. $512 \times 20 \times 15$)

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

Anchor is an object?
 $2 \times 9 \times 20 \times 15$

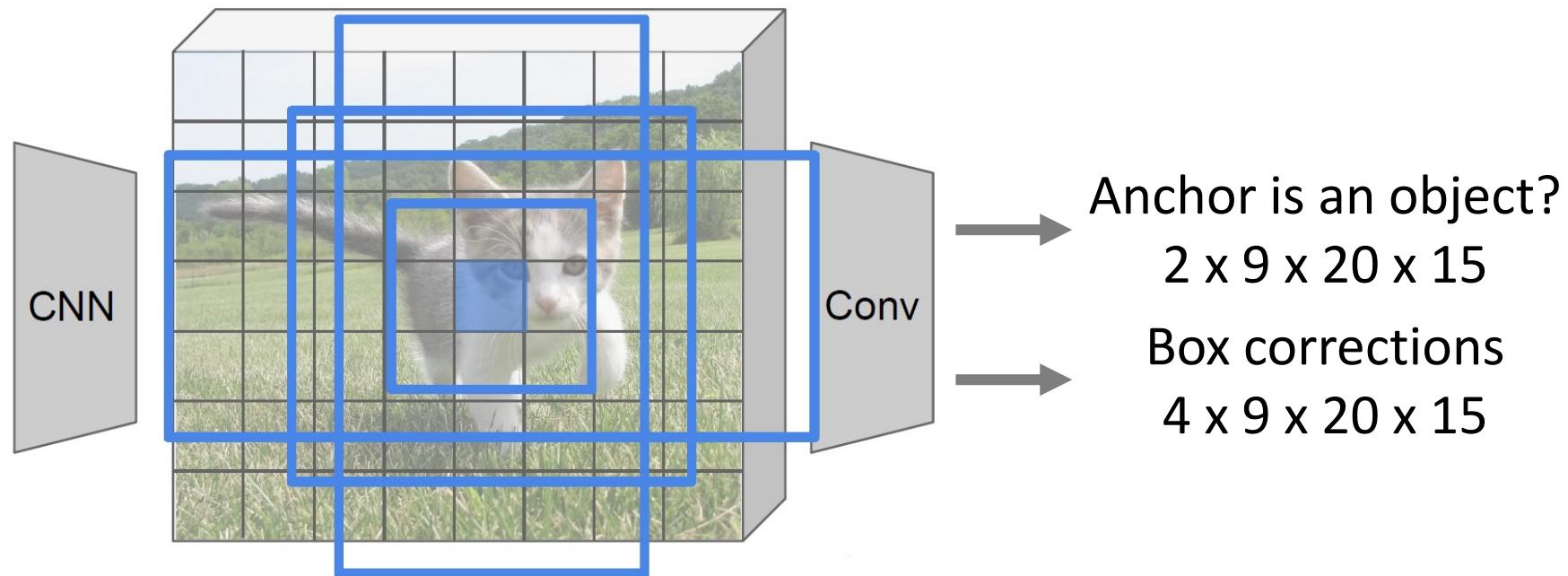
Box corrections
 $4 \times 9 \times 20 \times 15$

For positive boxes, also predict
a corrections from the anchor to
the ground-truth box (regress 4
numbers per pixel)

RPN: Region Proposal Network

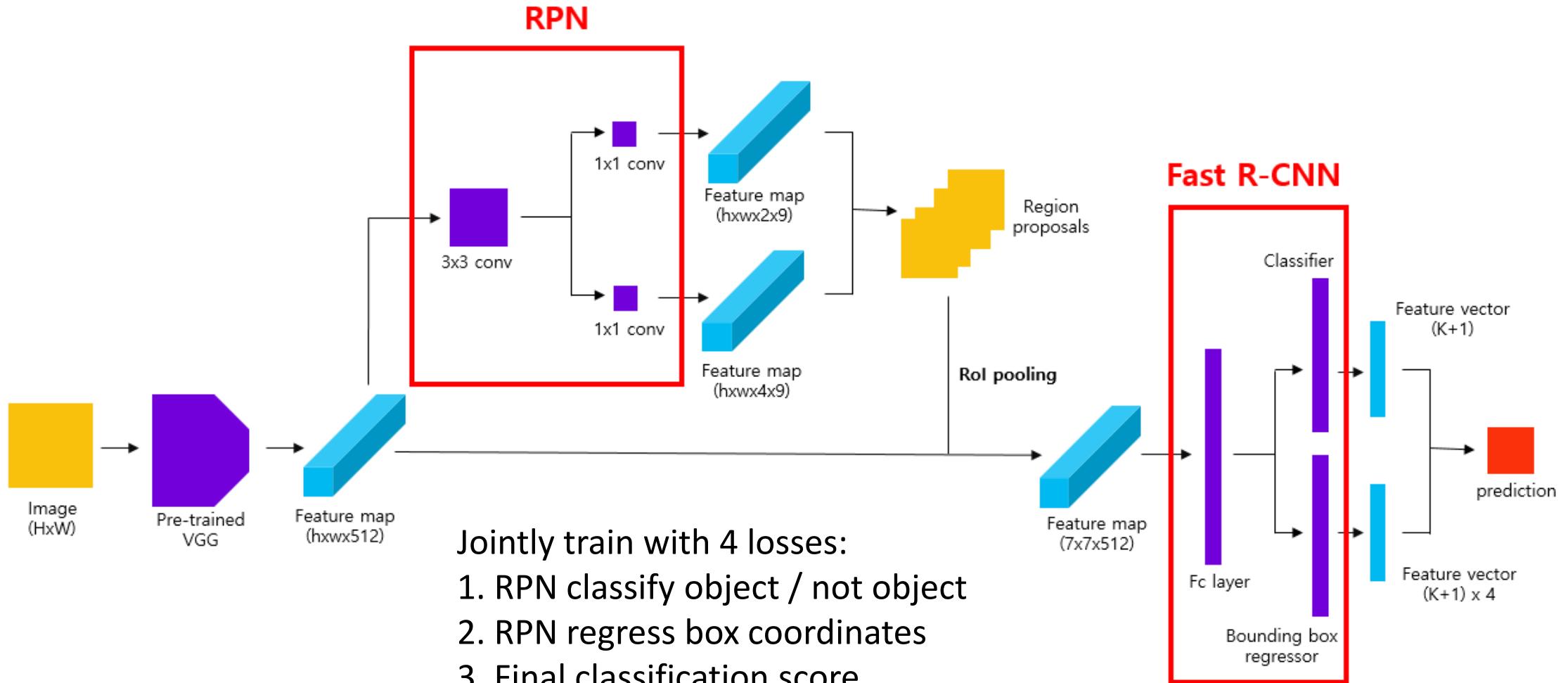


Input Image
(e.g. $3 \times 640 \times 480$)

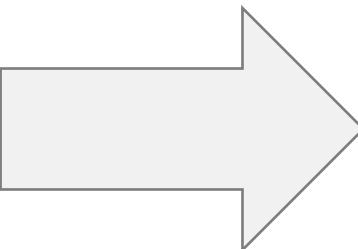
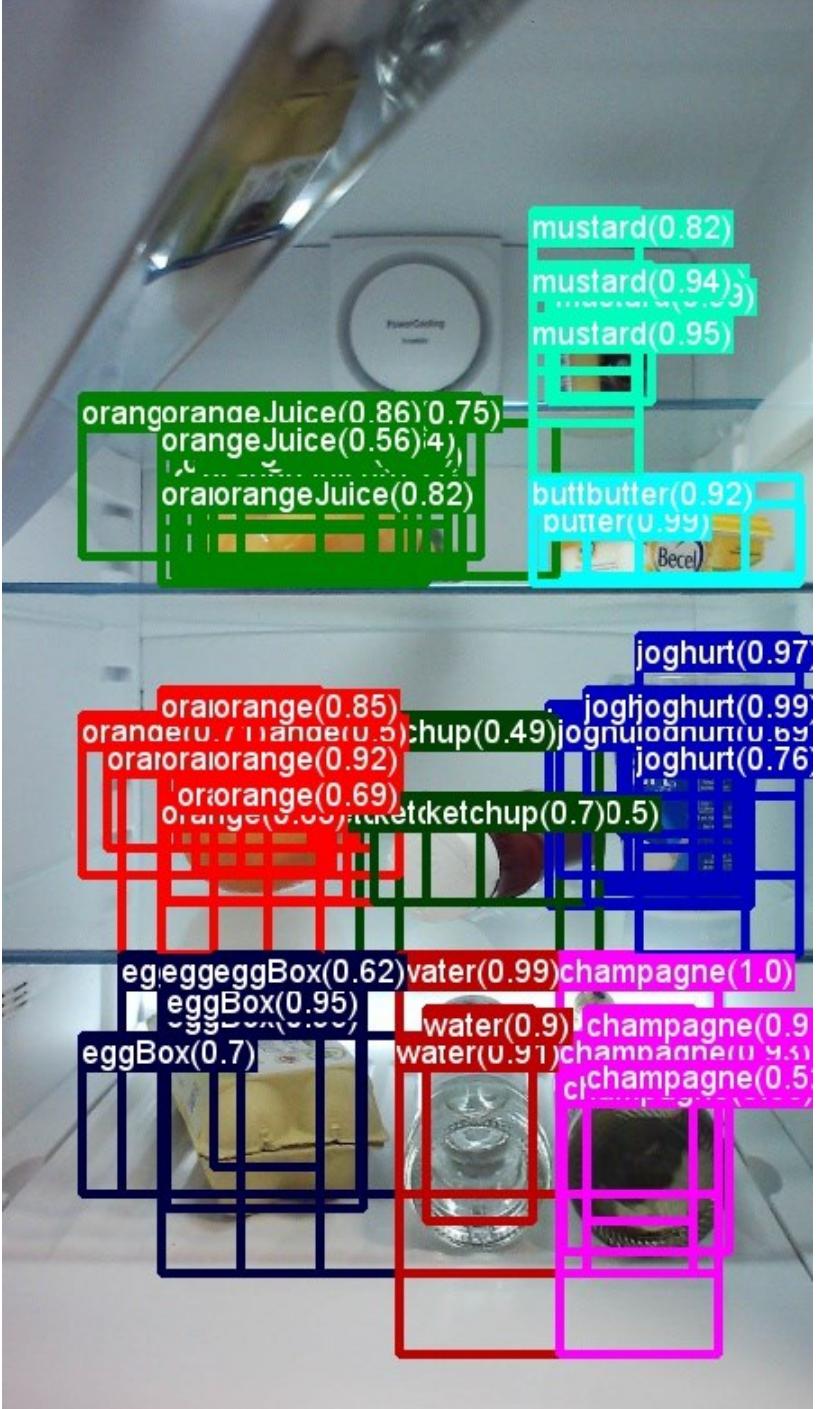


Sort the $9 \times 20 \times 15$ boxes by their “objectness” score, take top ~ 300 as our proposals

RPN: Region Proposal Network



Non-Maximum Suppression

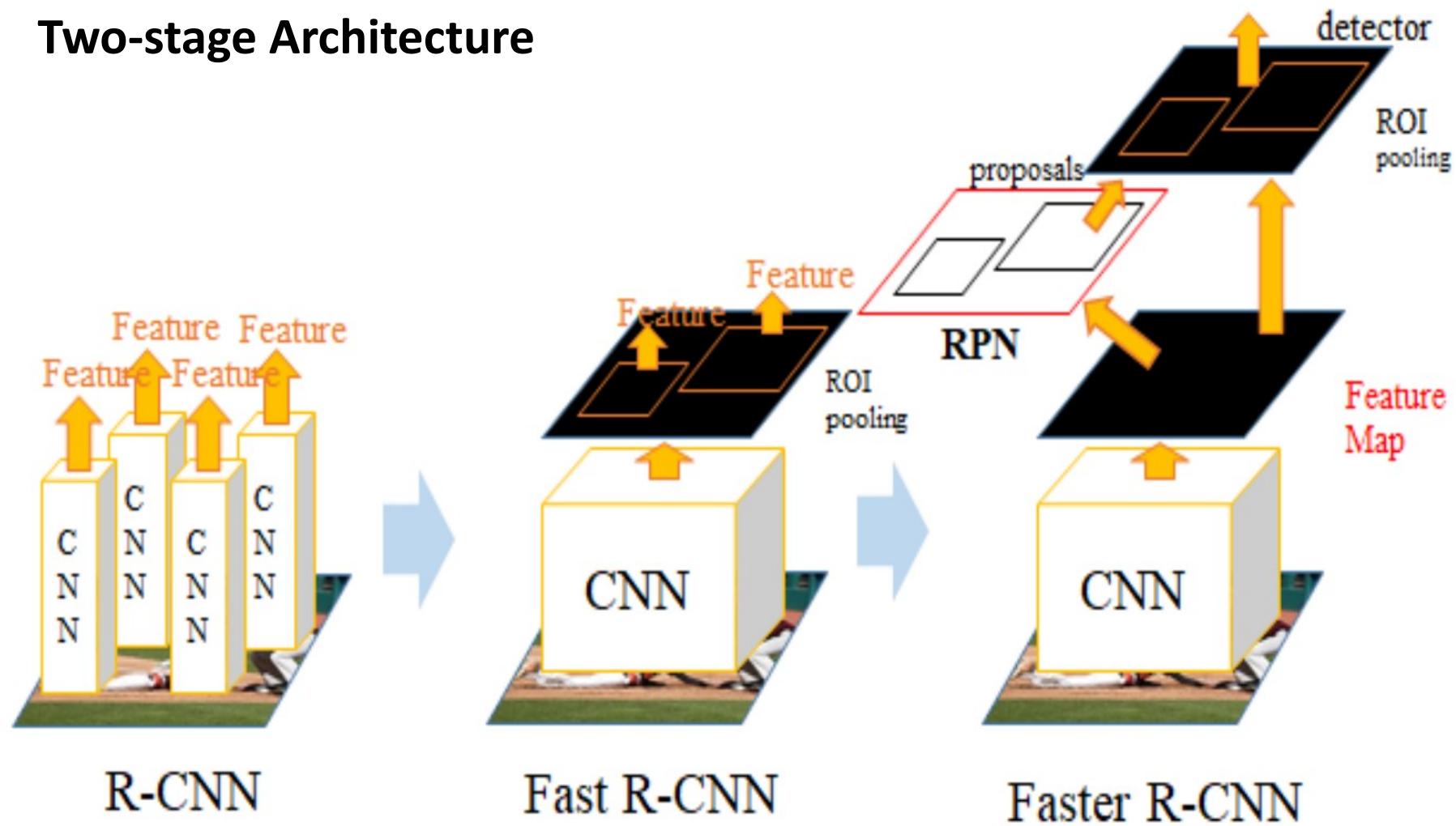


If $\text{IoU} > 0.7$, remove the box with lower confidence



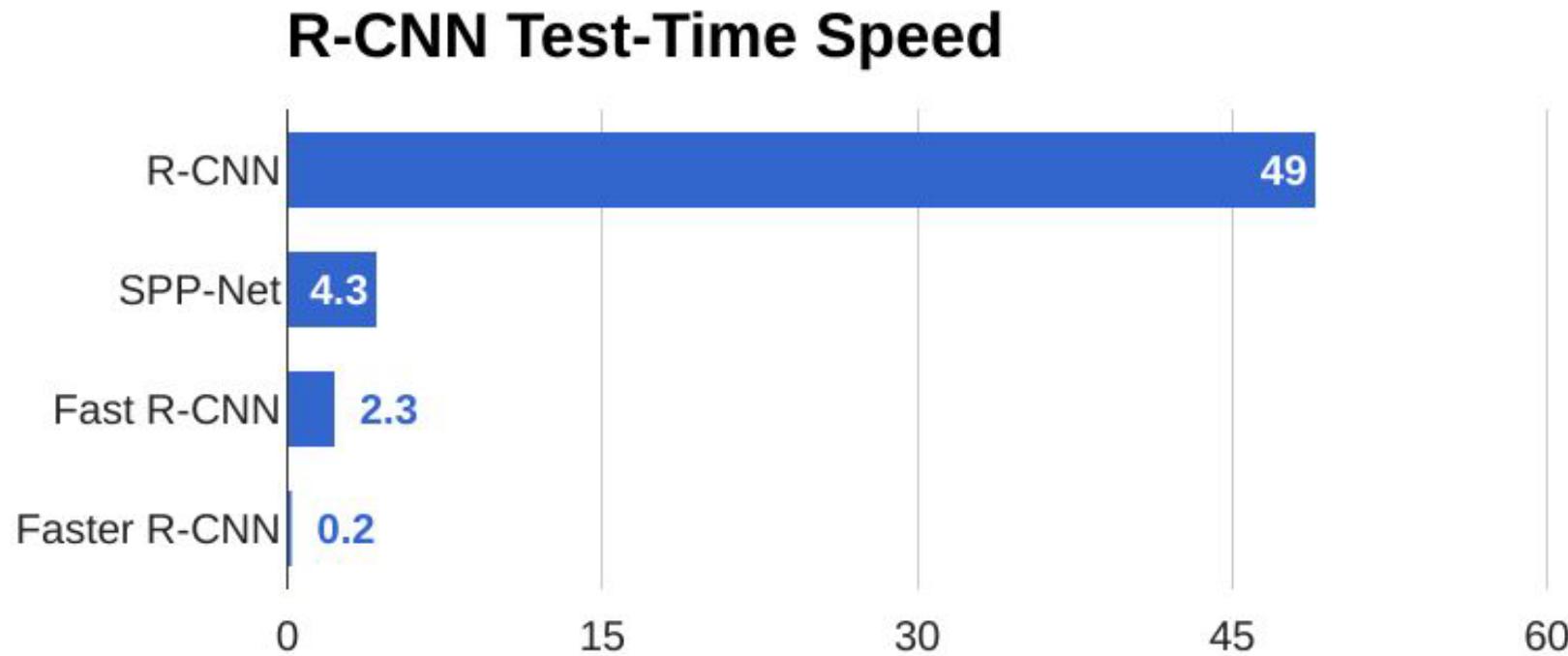
R-CNN vs Fast R-CNN vs Faster R-CNN

Two-stage Architecture



Faster R-CNN

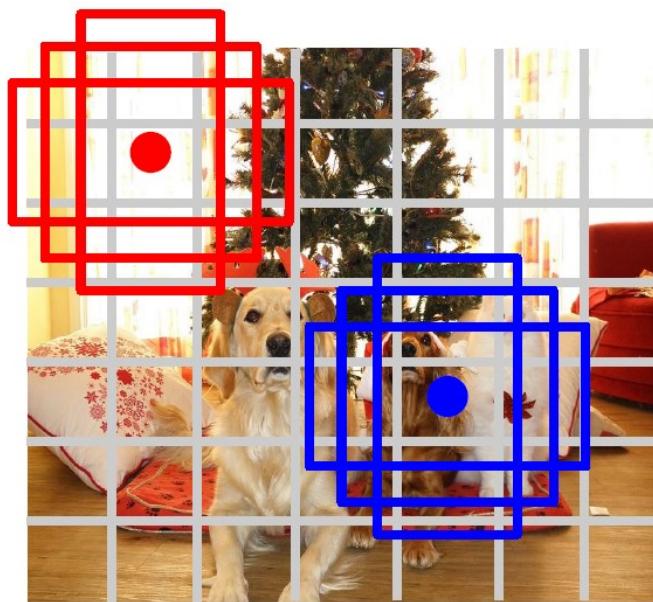
- Make CNN do proposal!



Single-Stage Object Detectors: YOLO



Input image
 $3 \times H \times W$



Divide image into grid
 $512 \times 7 \times 7$

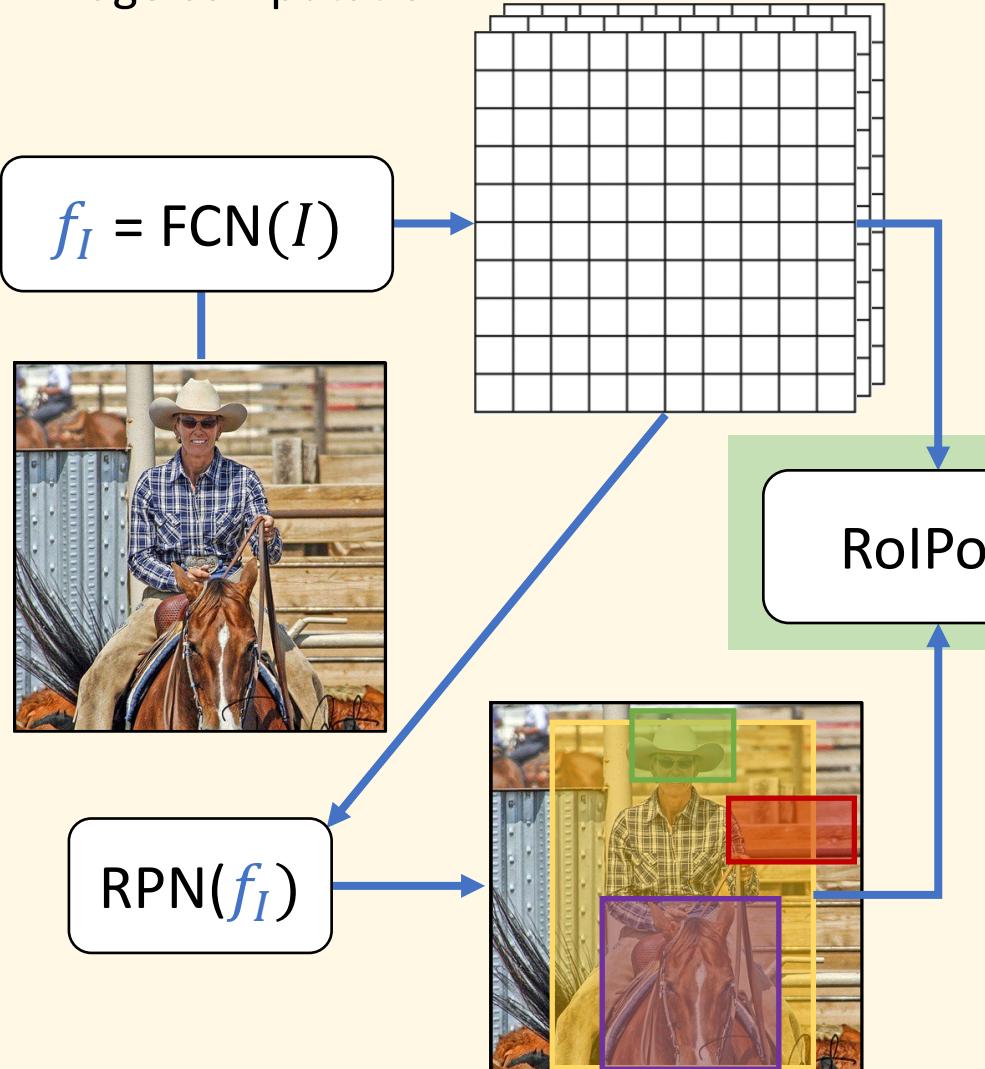
Image a set of base boxes
centered at each grid cell
Here $B = 3$

- 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)
 - Looks a lot like RPN, but category-specific!

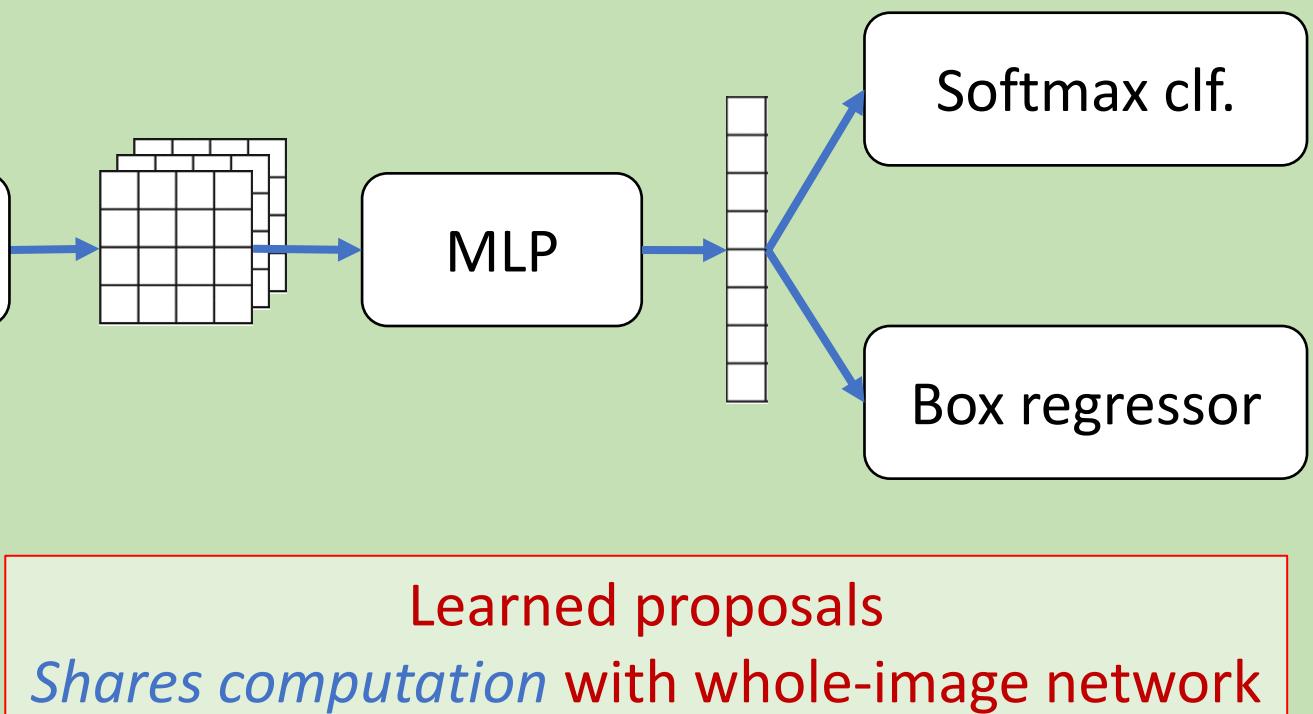
Output:
 $7 \times 7 \times (5 * B + C)$

Faster R-CNN

Per-image computation

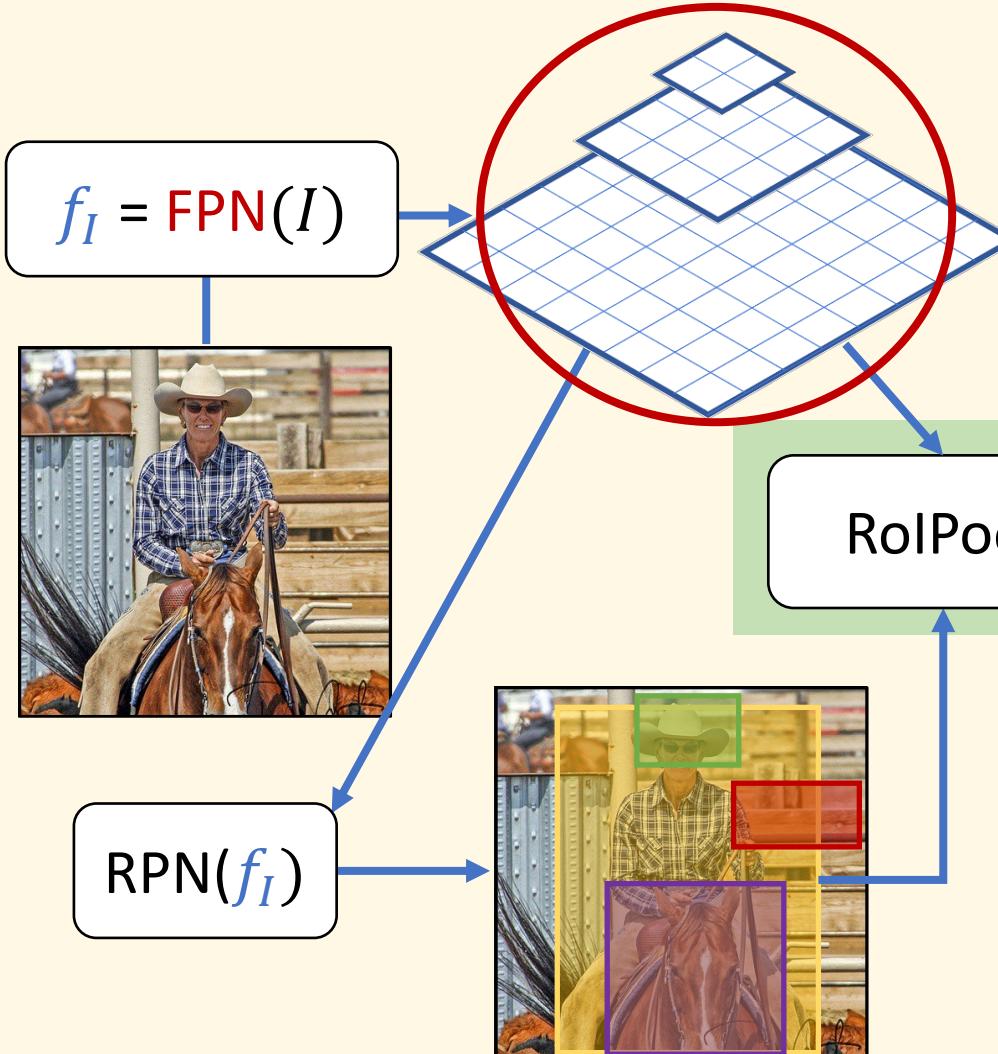


Per-region computation for each $r_i \in r(I)$

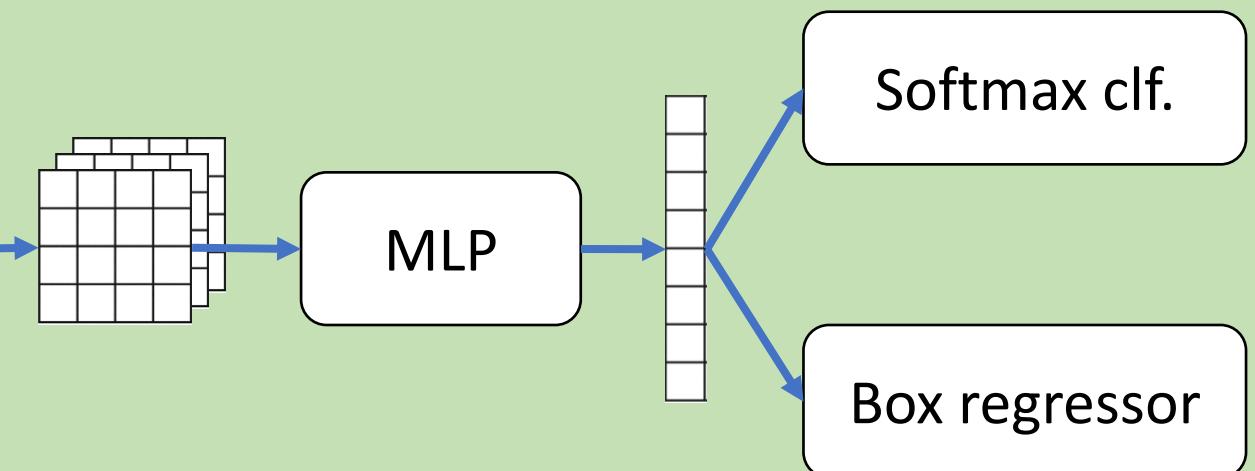


Faster R-CNN with a Feature Pyramid Network

Per-image computation

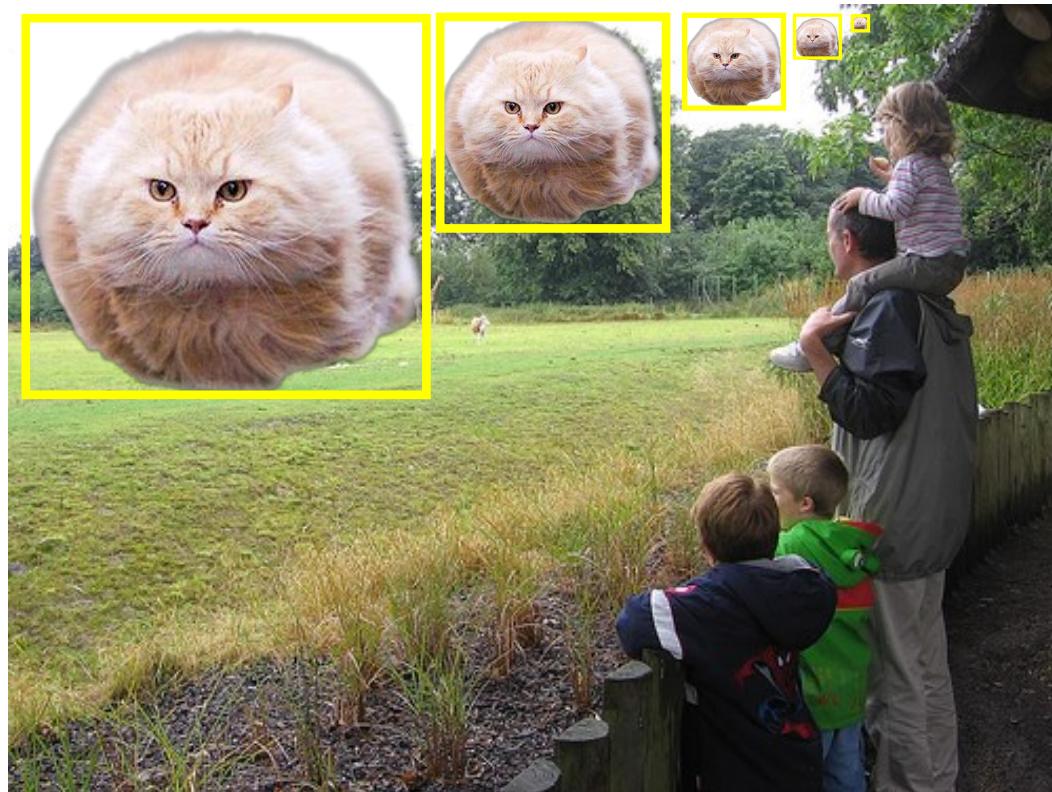


Per-region computation for each $r_i \in r(I)$



The whole-image feature representation can be improved by making it *multi-scale*

FPN: Improving Scale Invariance and Equivariance



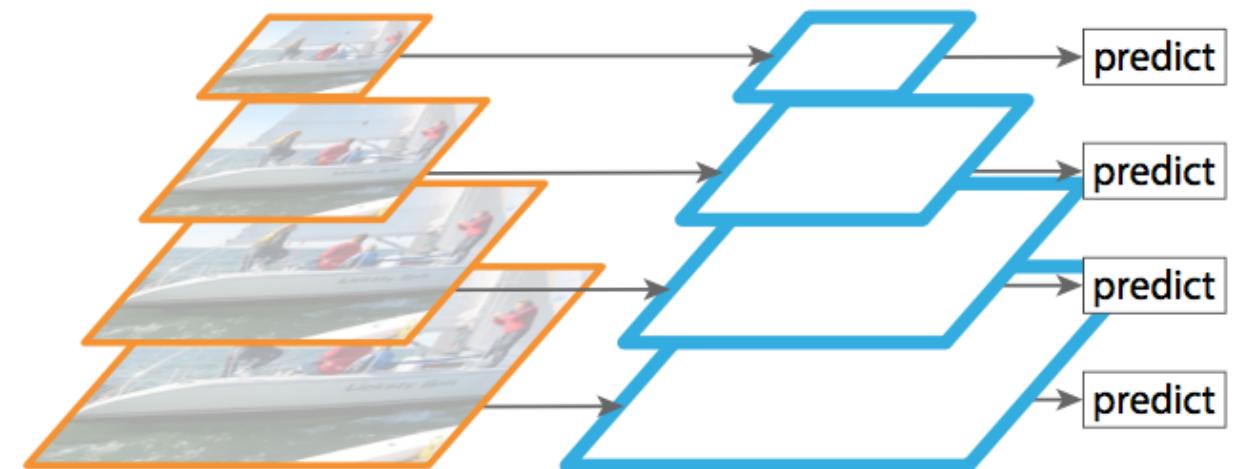
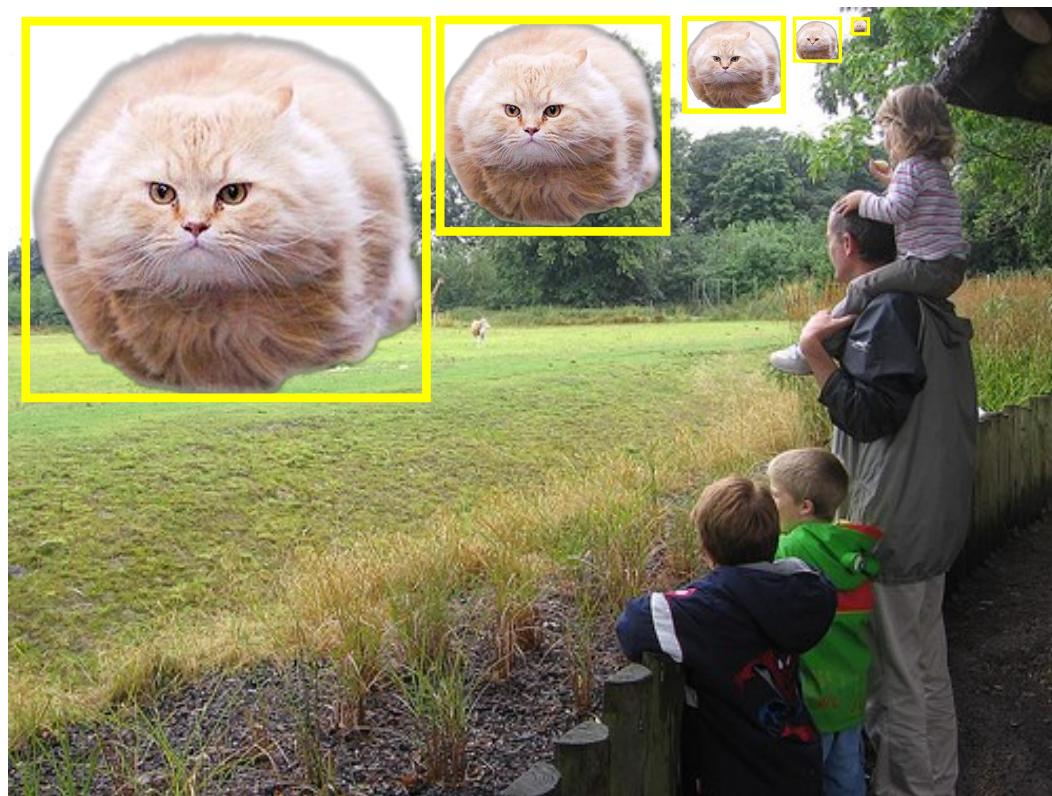
Detectors need to

1. Classify (invariance) and
2. Localize (equivariance)

objects over a **wide range of scales**

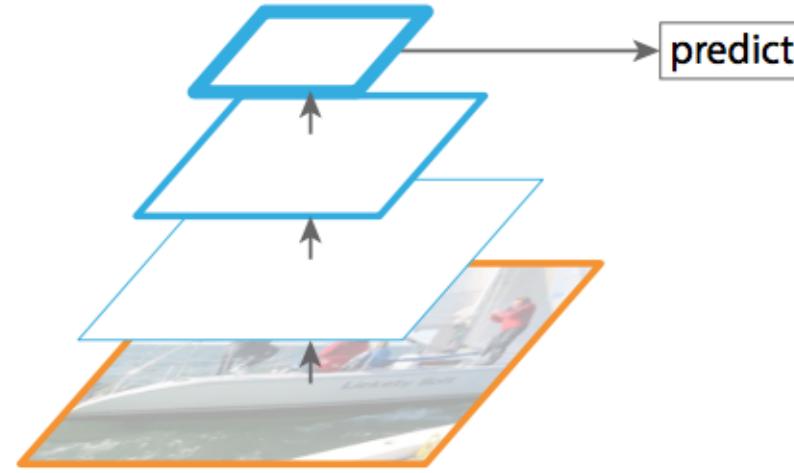
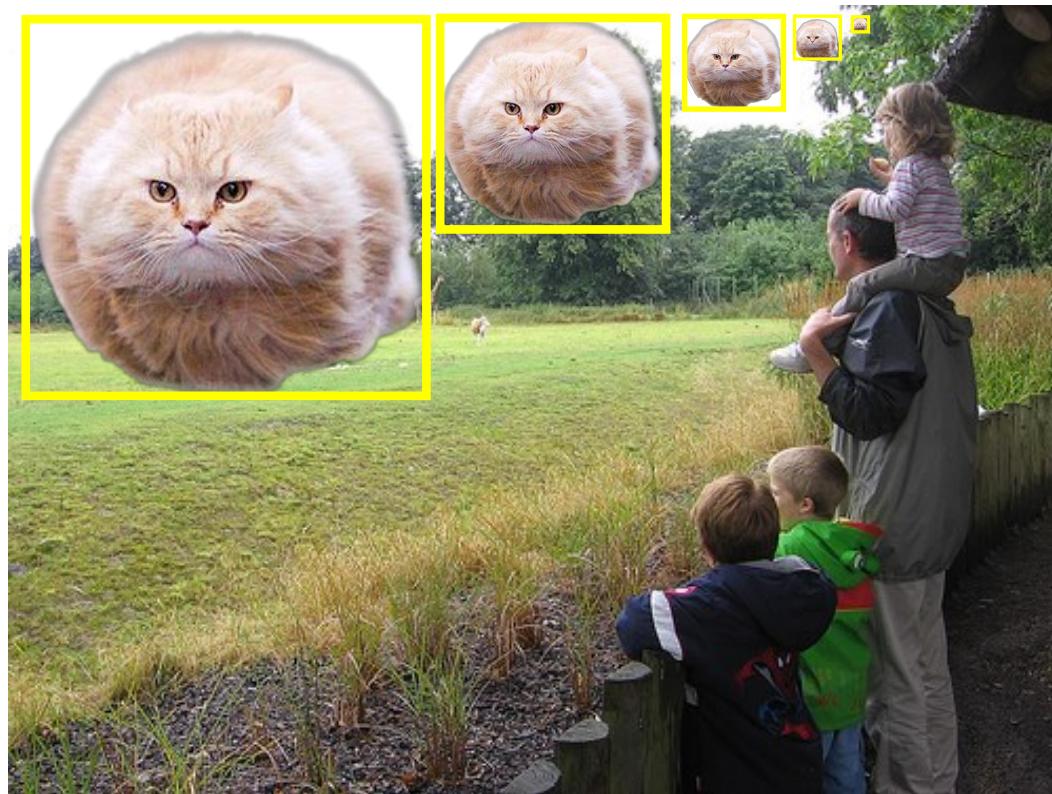
FPN improves this ability

Strategy 1: Image Pyramid



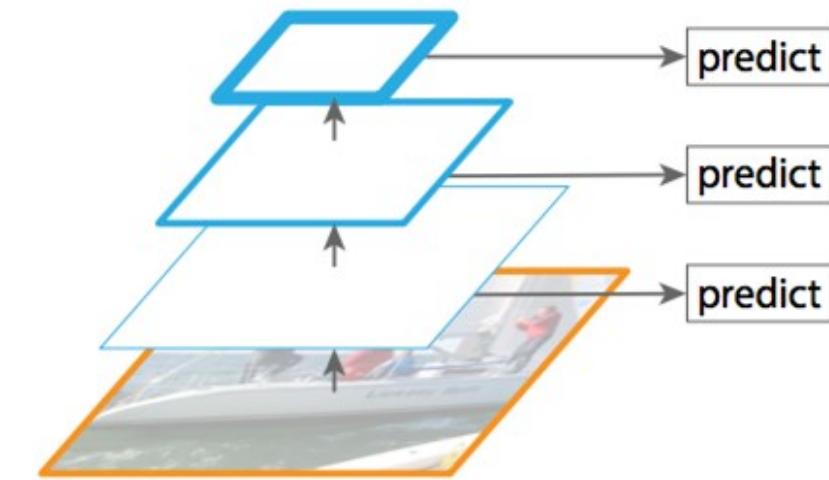
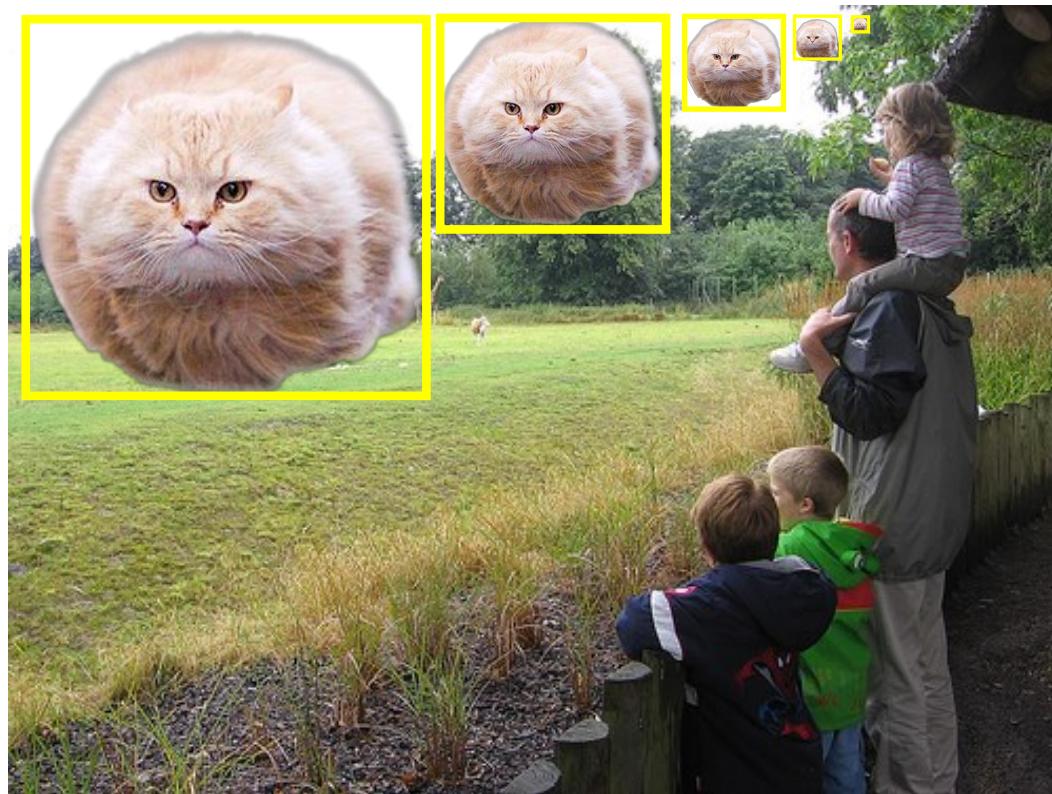
(a) Featurized image pyramid
Standard solution – *slow!*
(E.g., Viola & Jones, HOG, DPM, SPP-net,
multi-scale Fast R-CNN, ...)

Strategy 2: Multi-scale Features (Single-scale Map)



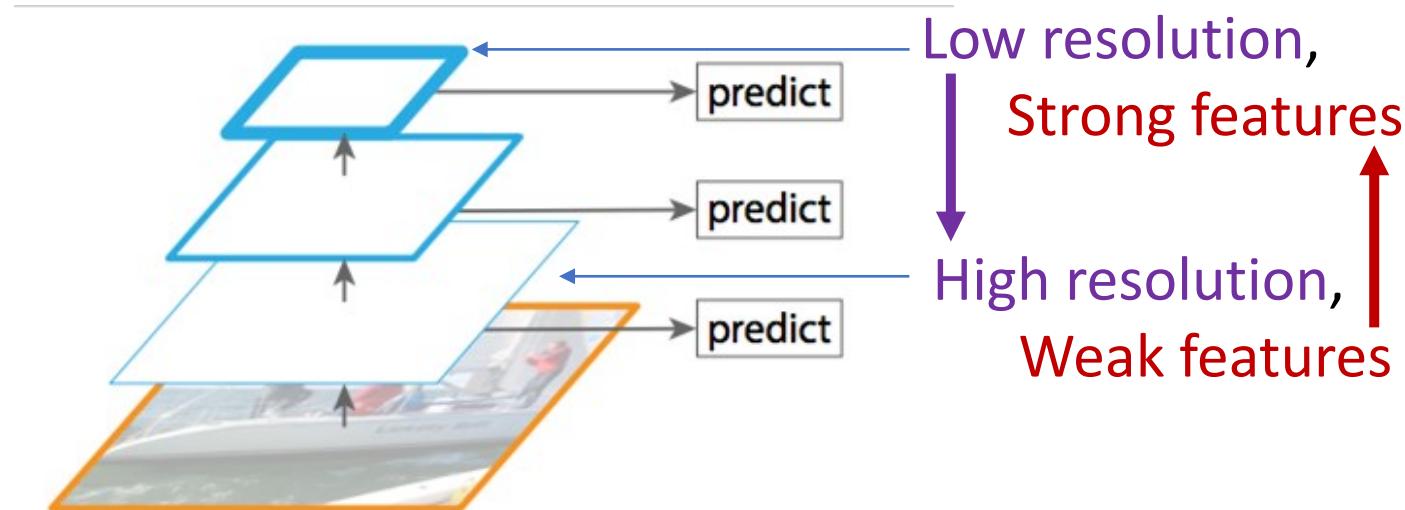
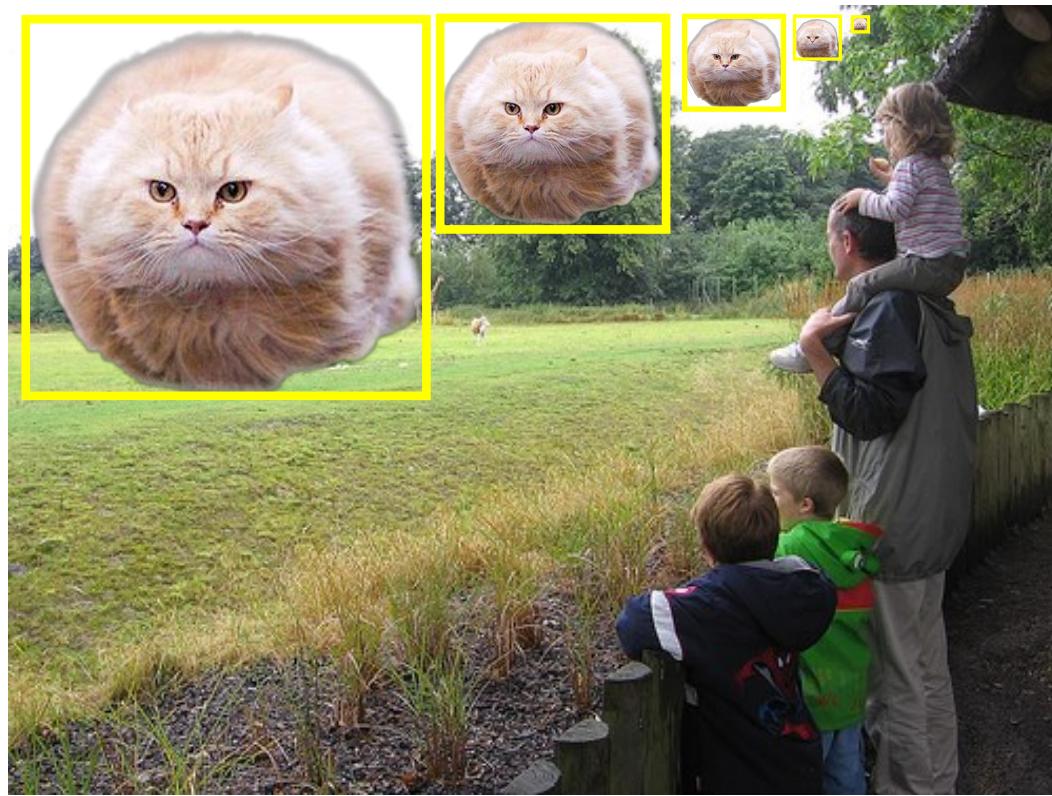
(b) Single feature map
Leave it all to the features – fast, suboptimal
(E.g., Fast/er R-CNN, YOLO, ...)

Strategy 3: Naïve In-network Pyramid



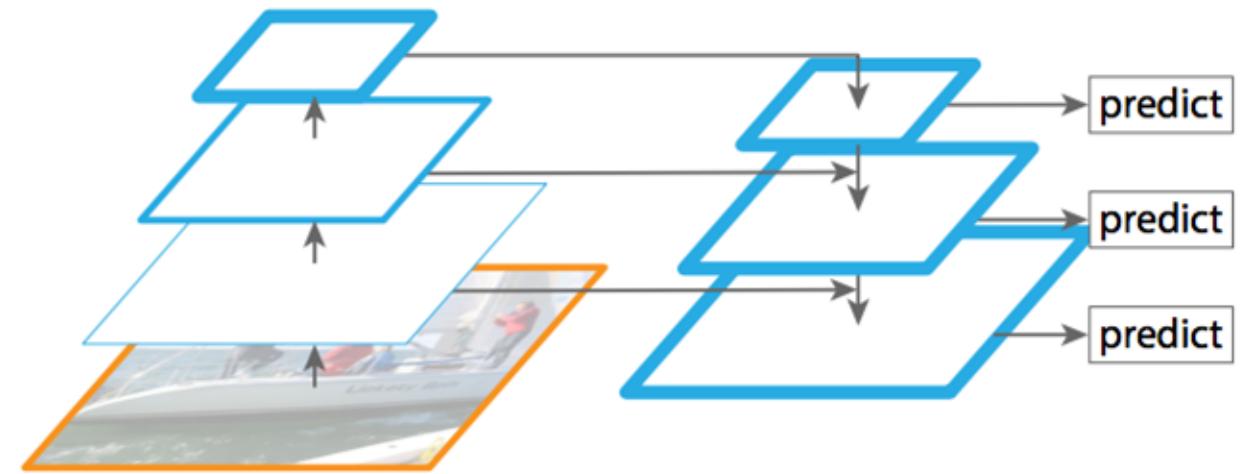
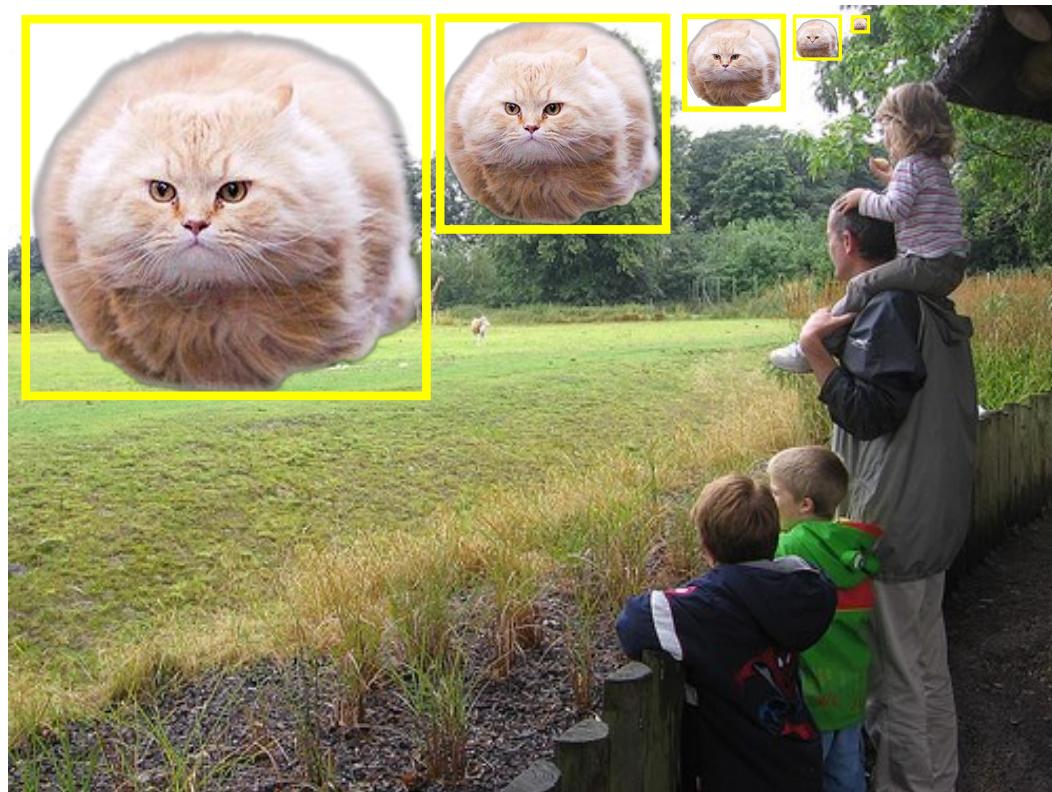
(c) Pyramidal feature hierarchy
Use the internal pyramid – *fast, suboptimal*
(E.g., \approx SSD, ...)

Strategy 3: Naïve In-network Pyramid



(c) Pyramidal feature hierarchy
Use the internal pyramid – *fast, suboptimal*
(E.g., \approx SSD, ...)

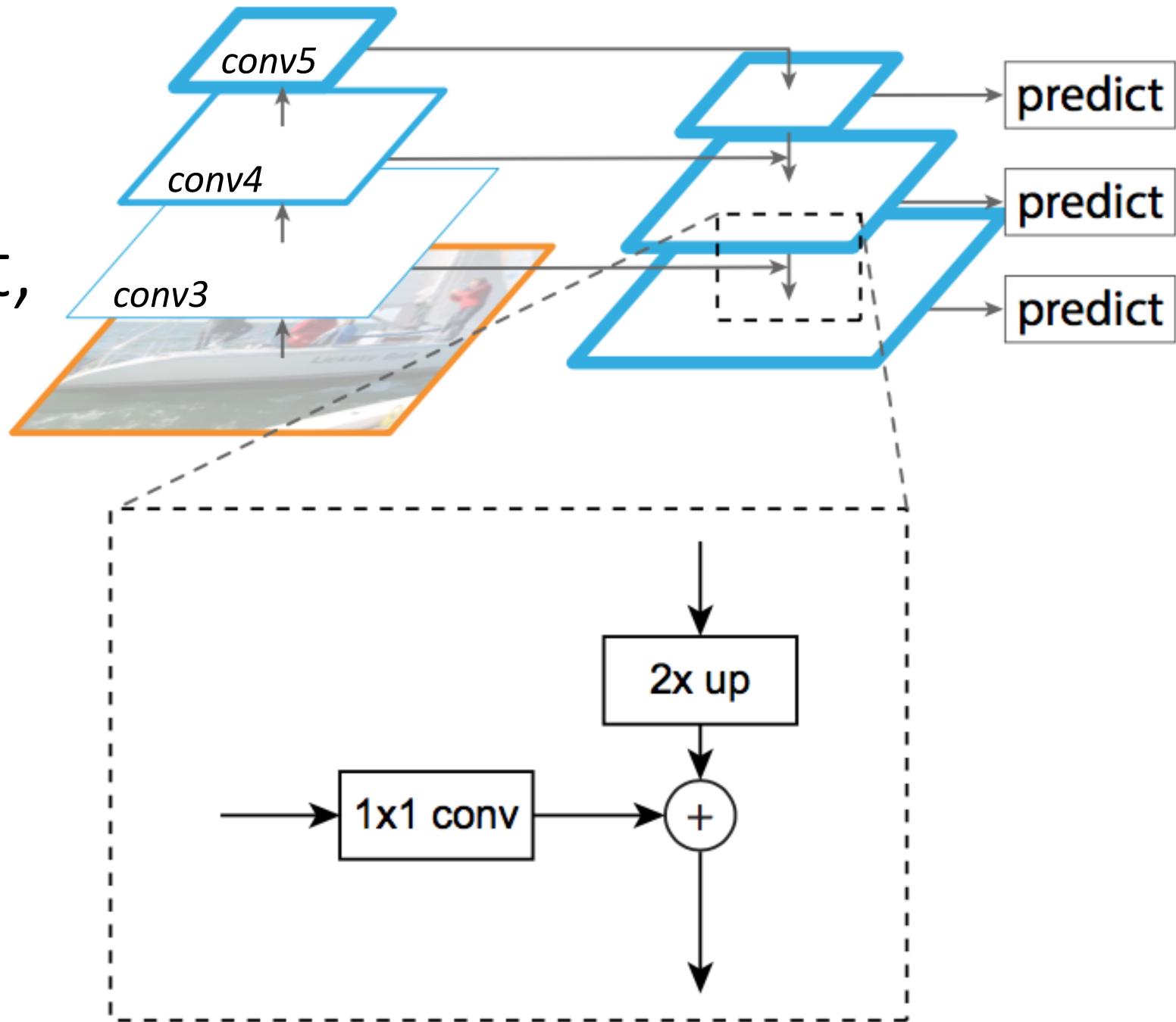
Strategy 4: Feature Pyramid Network



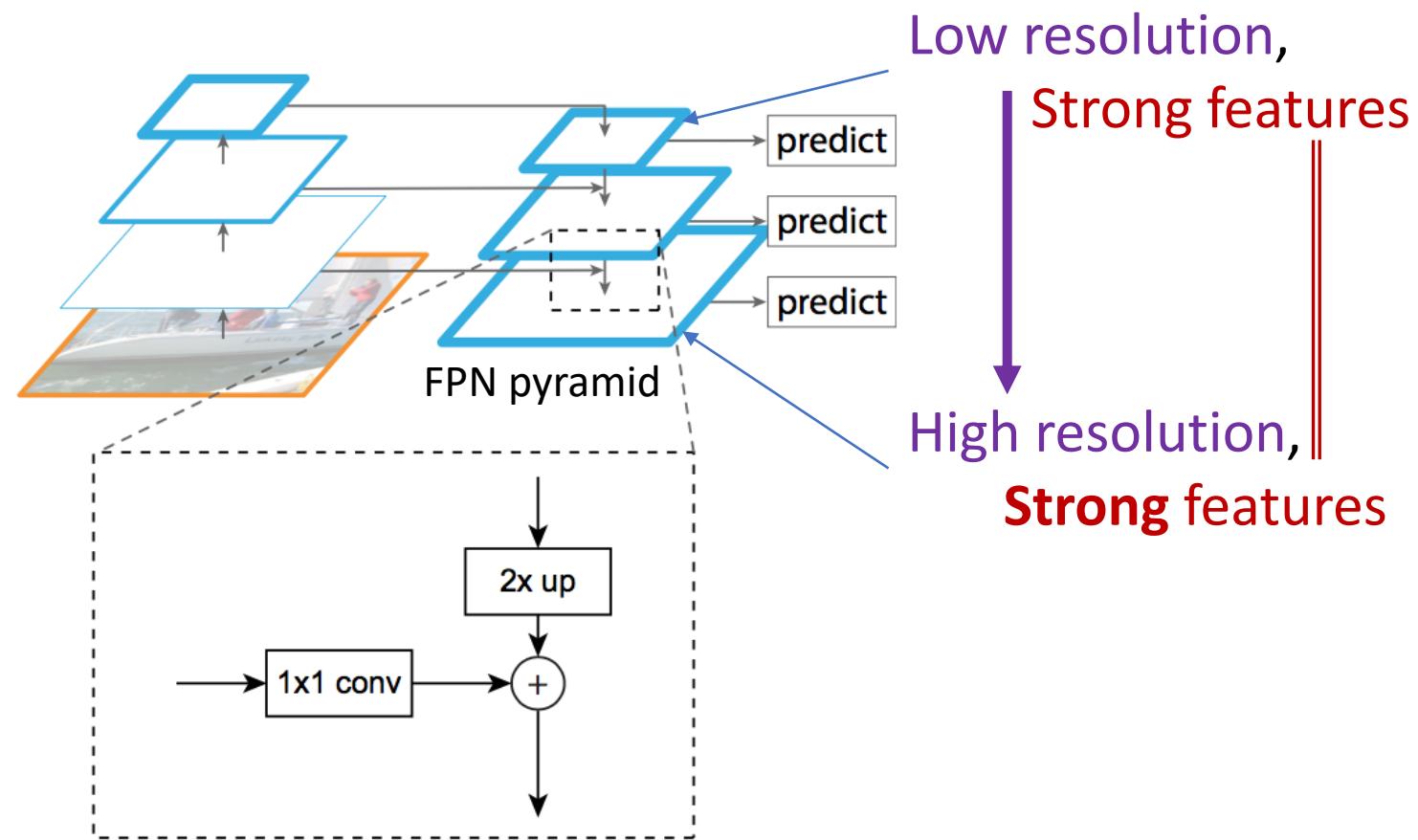
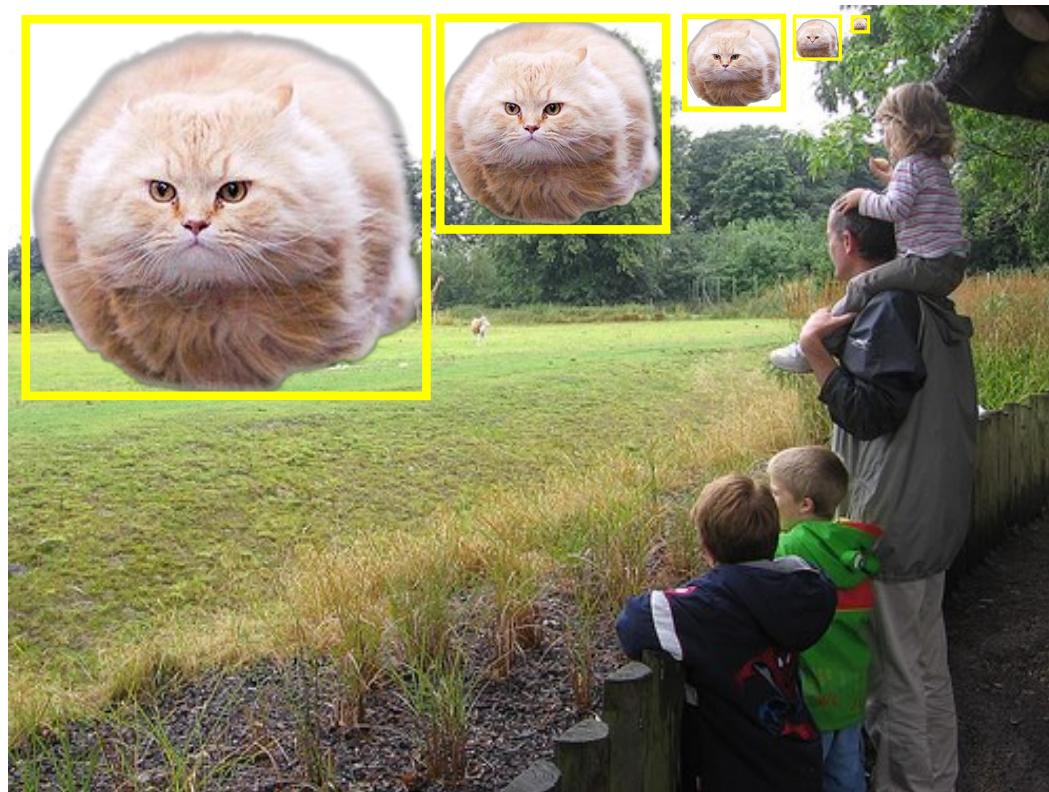
(d) Feature Pyramid Network

Top-down enrichment of high-res features –
fast, less suboptimal

FPN:
Light-weight,
Top-down
Refinement
Module



No Compromise on Feature Quality, Still Fast



FPN – A Generic Backbone Modification

Generates a feature pyramid—*useful in many applications!*

- RPN
- Fast/er R-CNN
- Mask R-CNN
- RetinaNet
- ...

Instance Segmentation

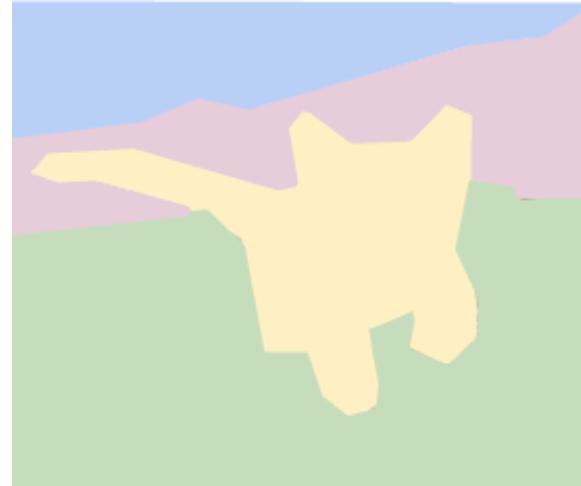
Classification



CAT

No spatial extent

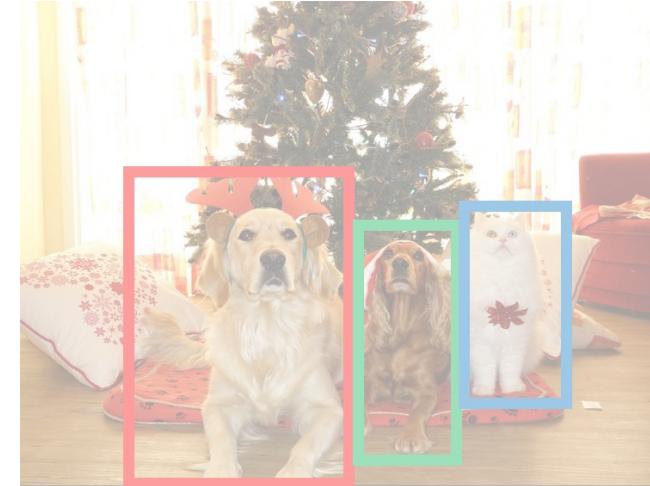
Semantic
Segmentation



GRASS, CAT
TREE, SKY

No objects, just pixels

Object
Detection



DOG, DOG, CAT

Multiple Object

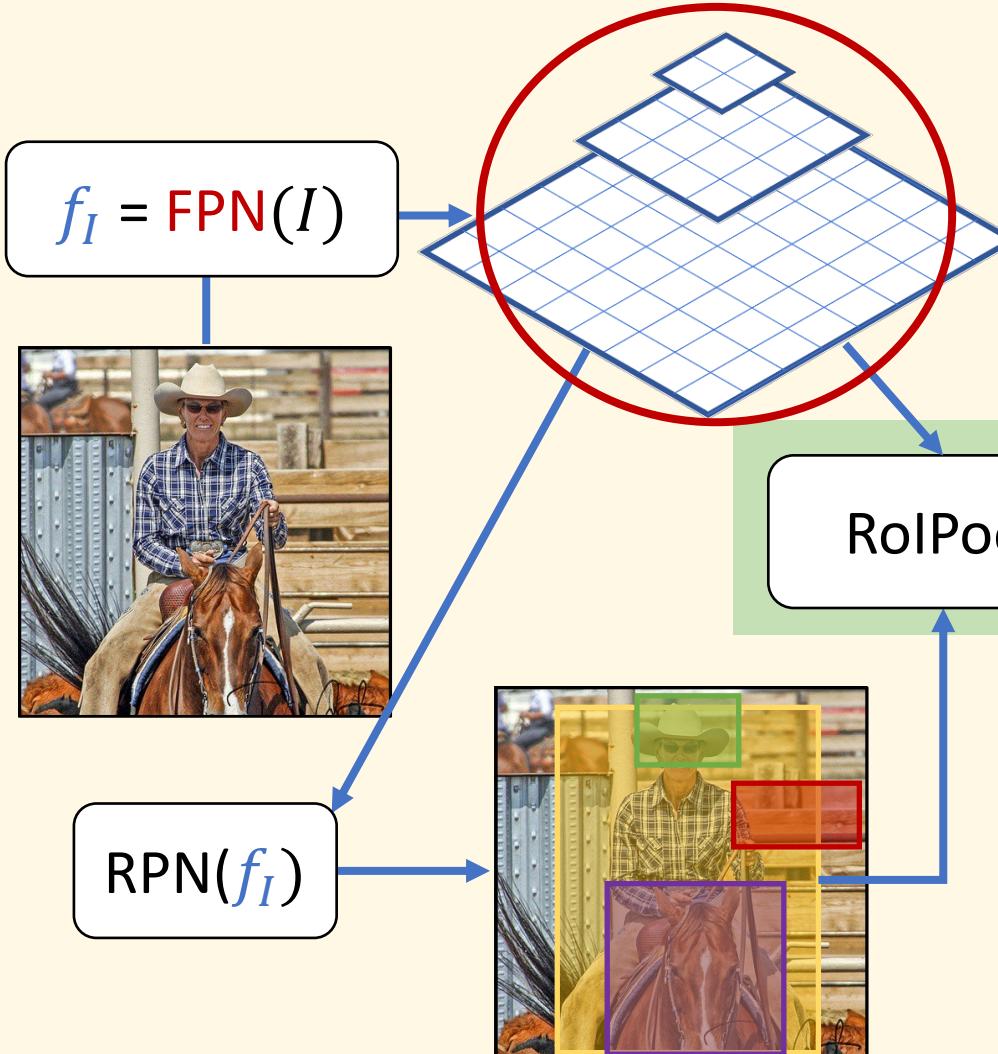
Instance
Segmentation



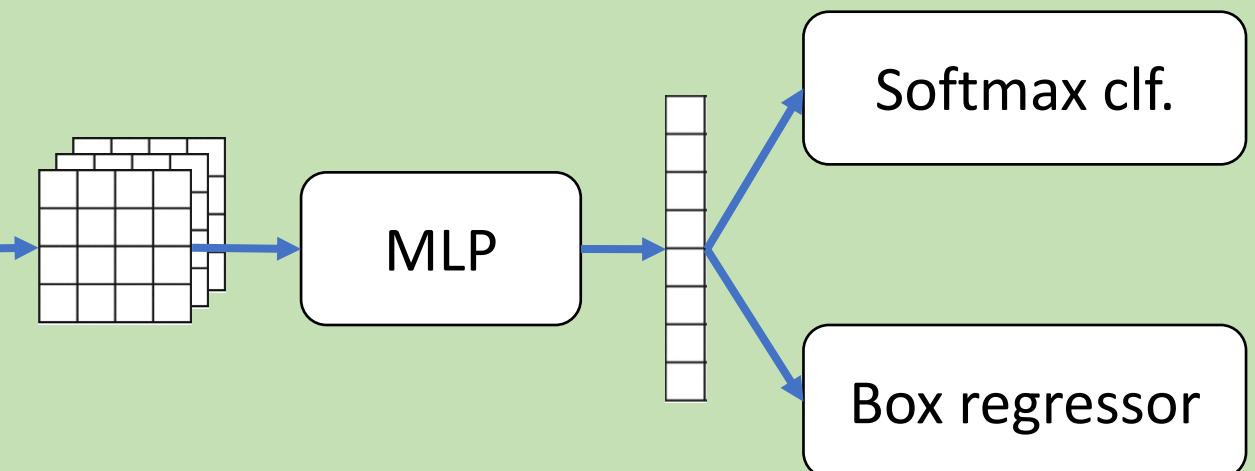
DOG, DOG, CAT

Faster R-CNN with a Feature Pyramid Network

Per-image computation



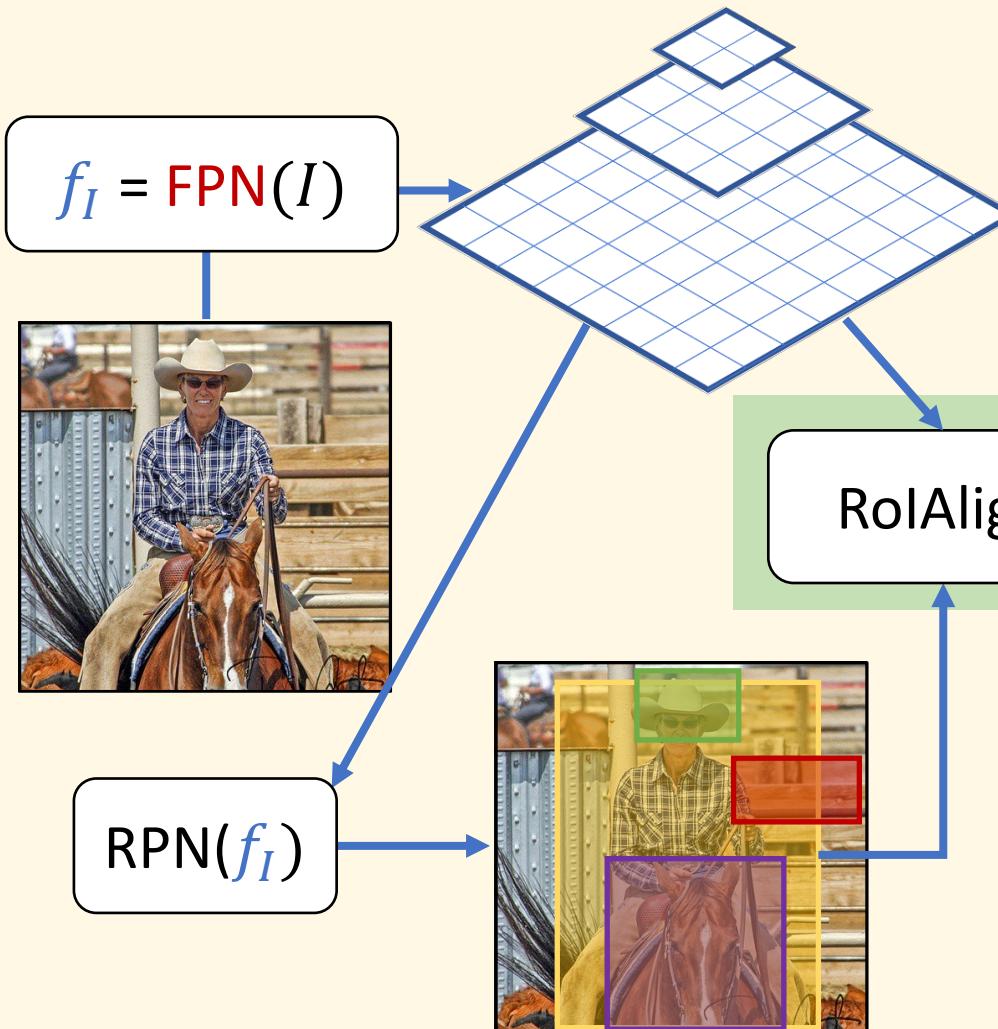
Per-region computation for each $r_i \in r(I)$



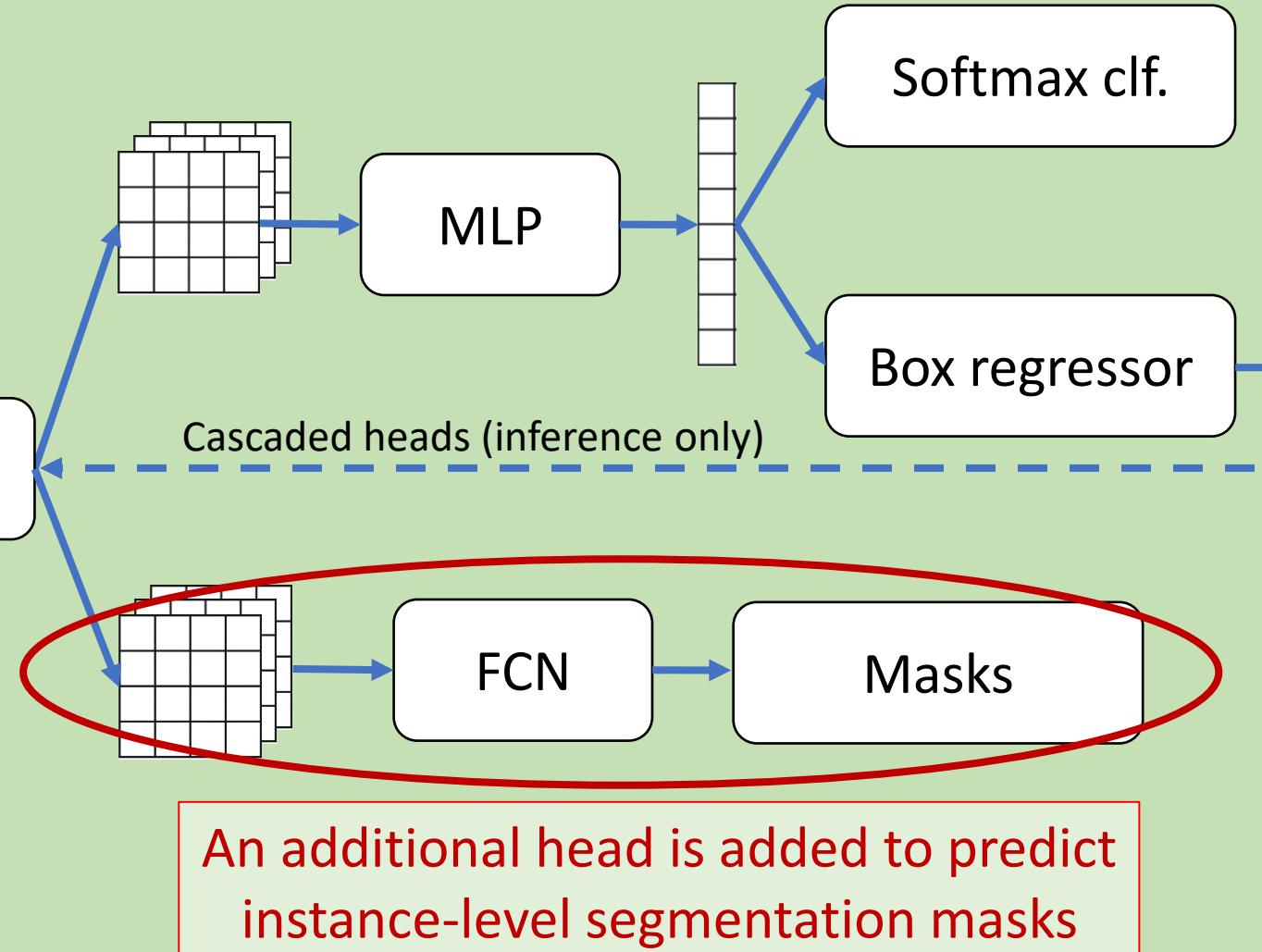
The whole-image feature representation can be improved by making it *multi-scale*

Mask R-CNN

Per-image computation

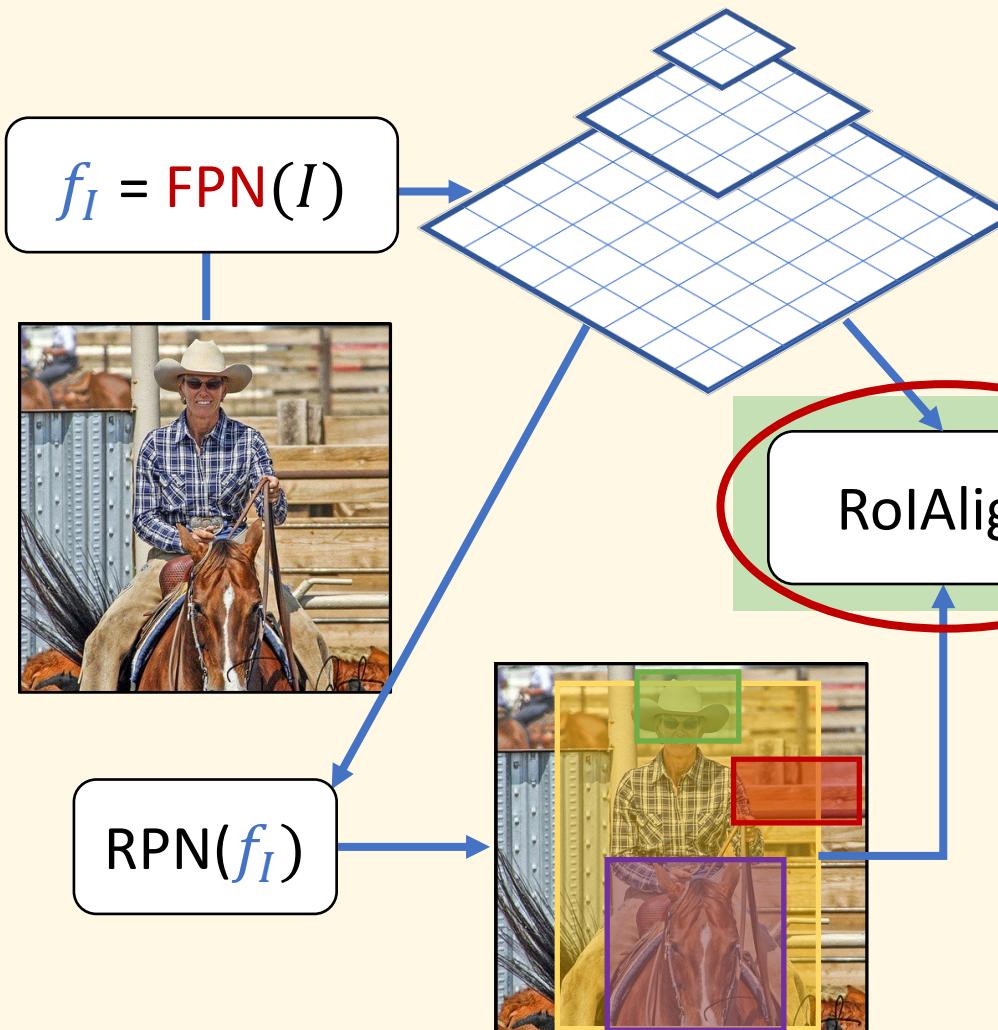


Per-region computation for each $r_i \in r(I)$

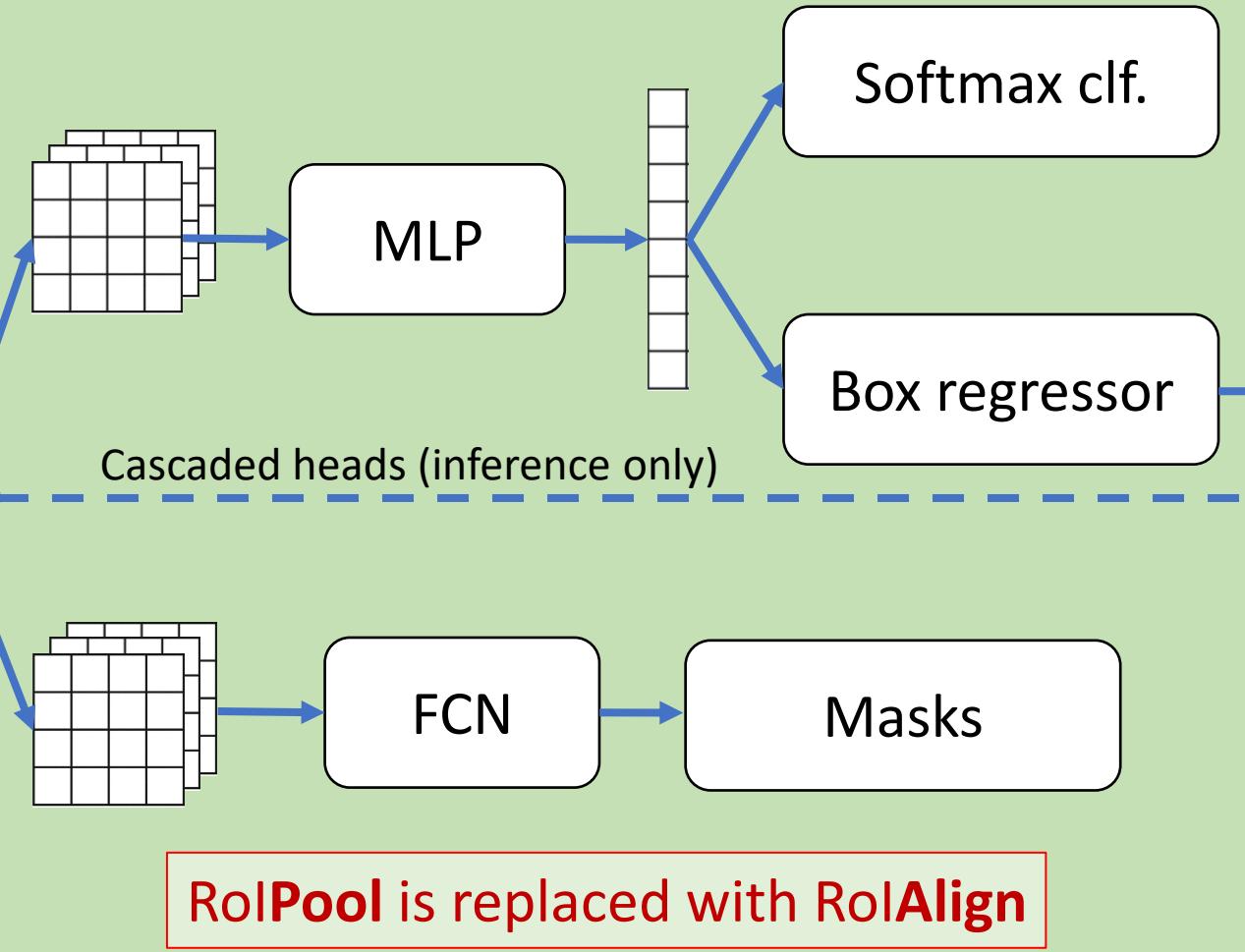


Mask R-CNN

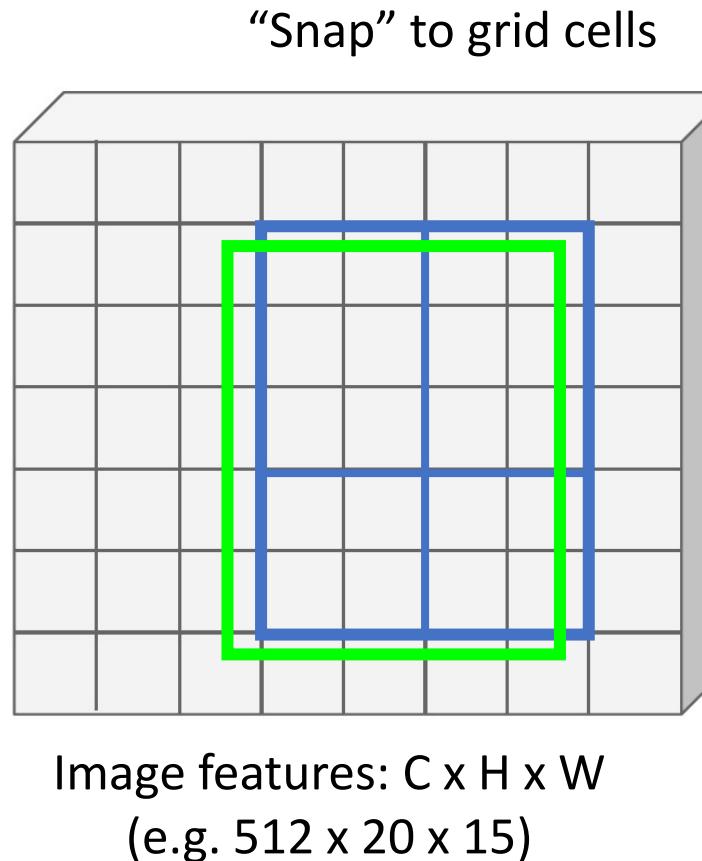
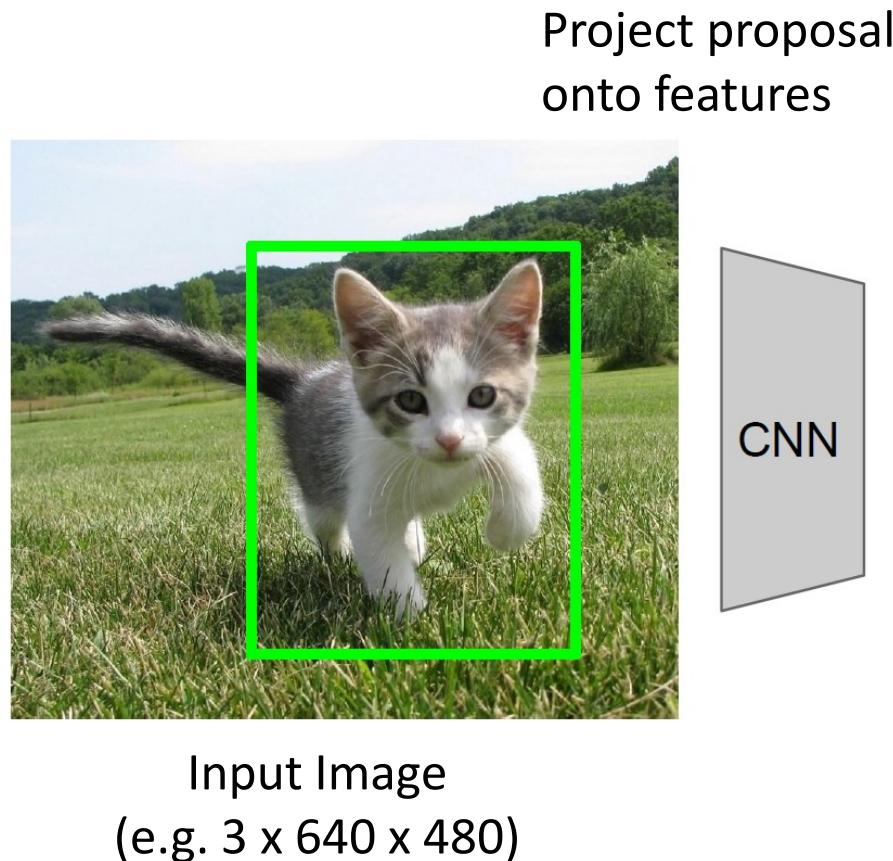
Per-image computation



Per-region computation for each $r_i \in r(I)$



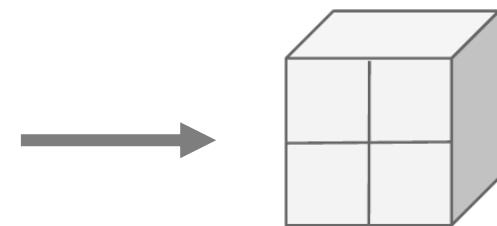
Cropping Features: RoIPool



Problem: Region features slightly misaligned

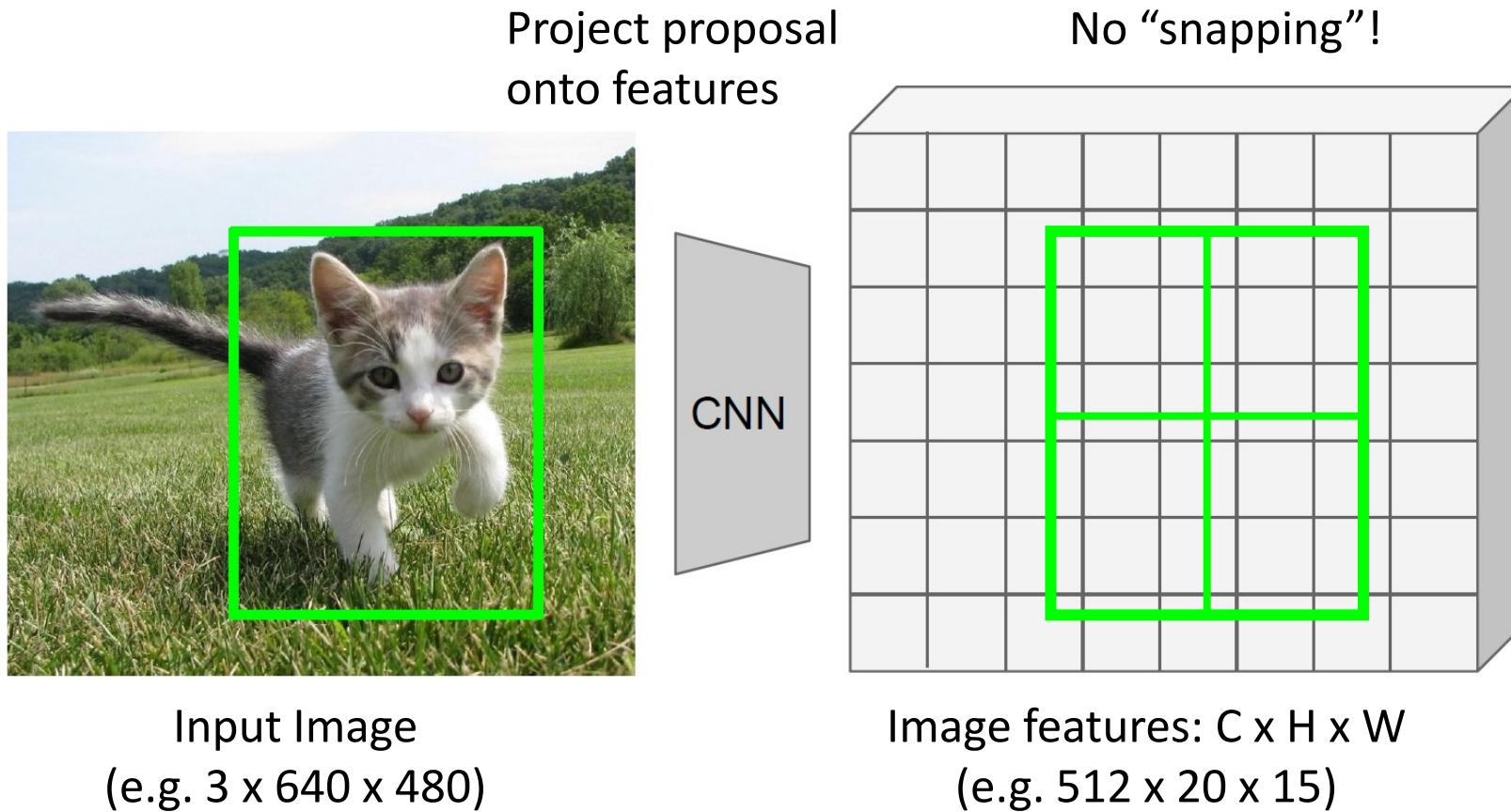
Divide into 2×2 grid of (roughly) equal subregions

Max-pool within each subregion

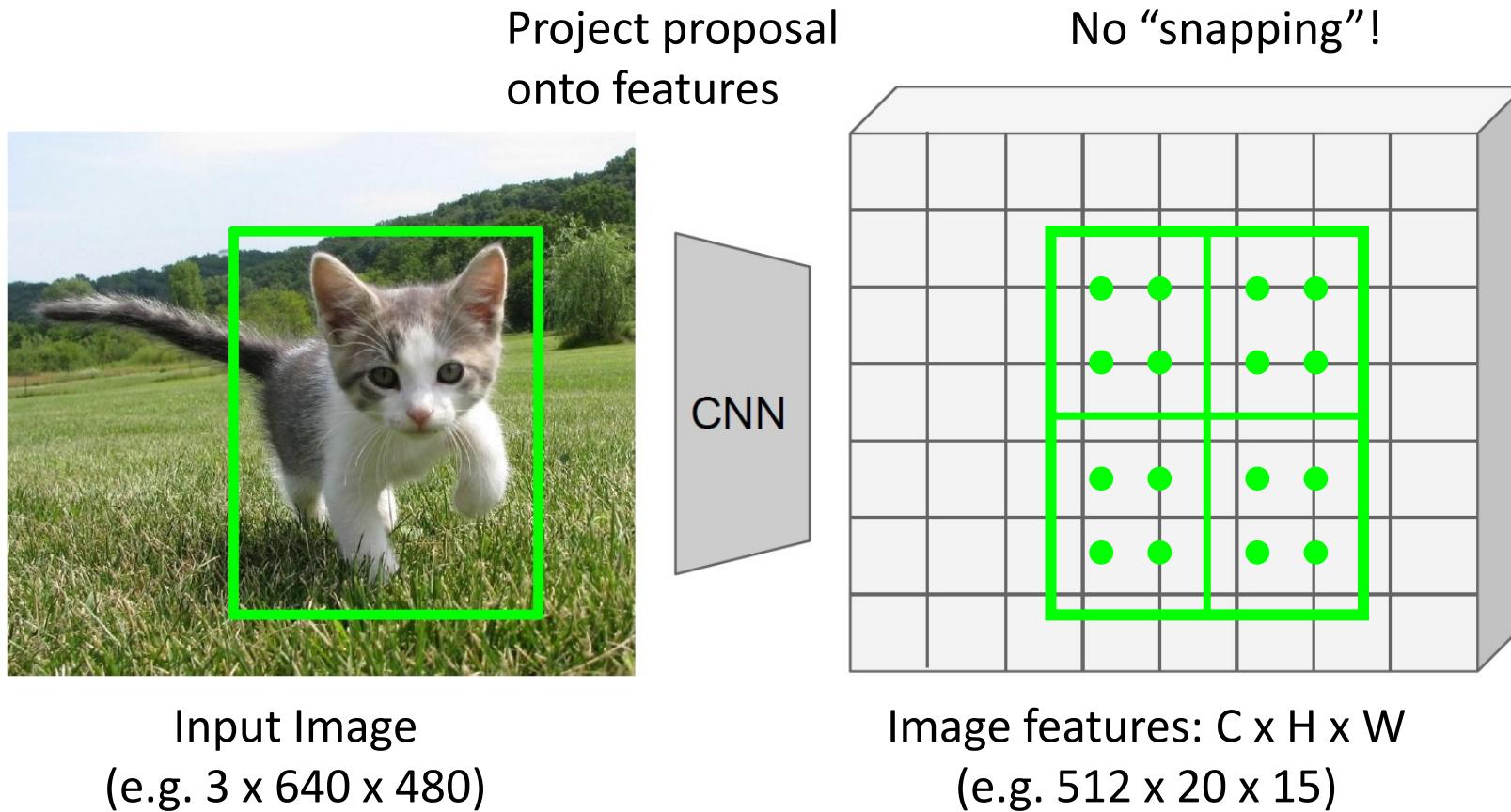


Region features always the same size even if input regions have different sizes!

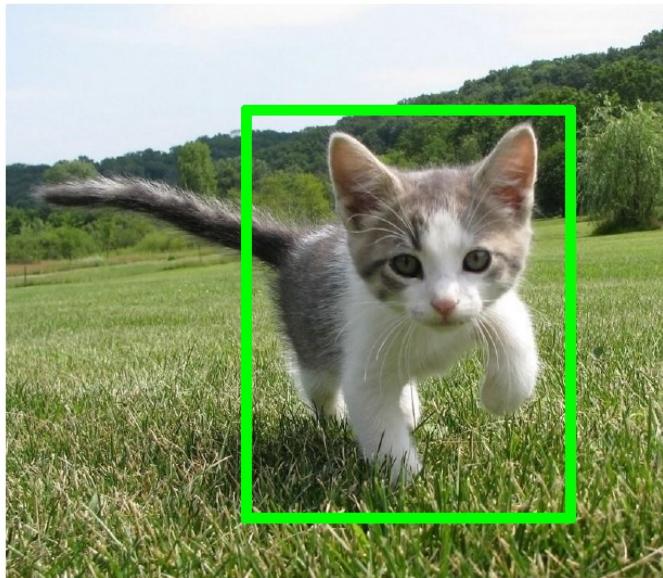
Cropping Features: RoIAlign



Cropping Features: RoIAlign



Cropping Features: RoIAlign



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

Project proposal
onto features



No “snapping”!

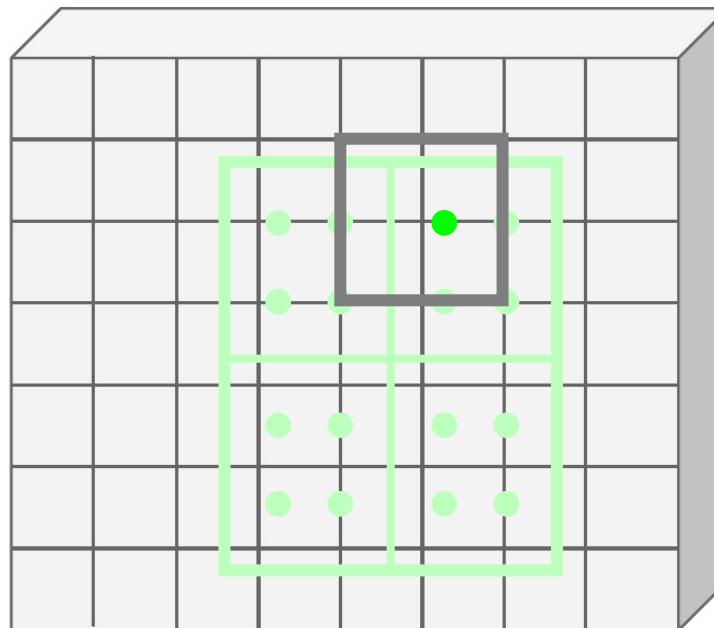
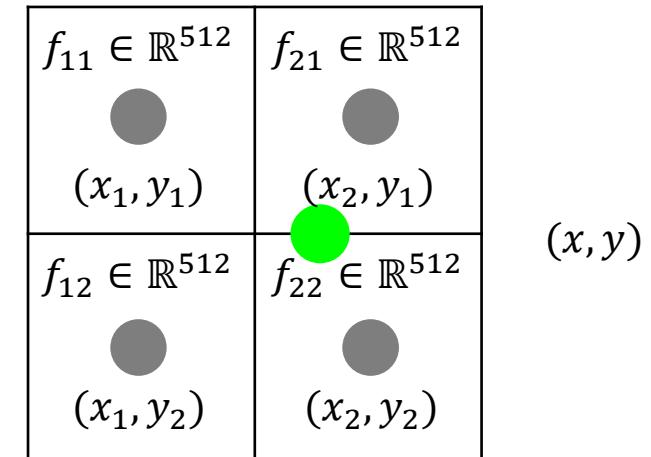


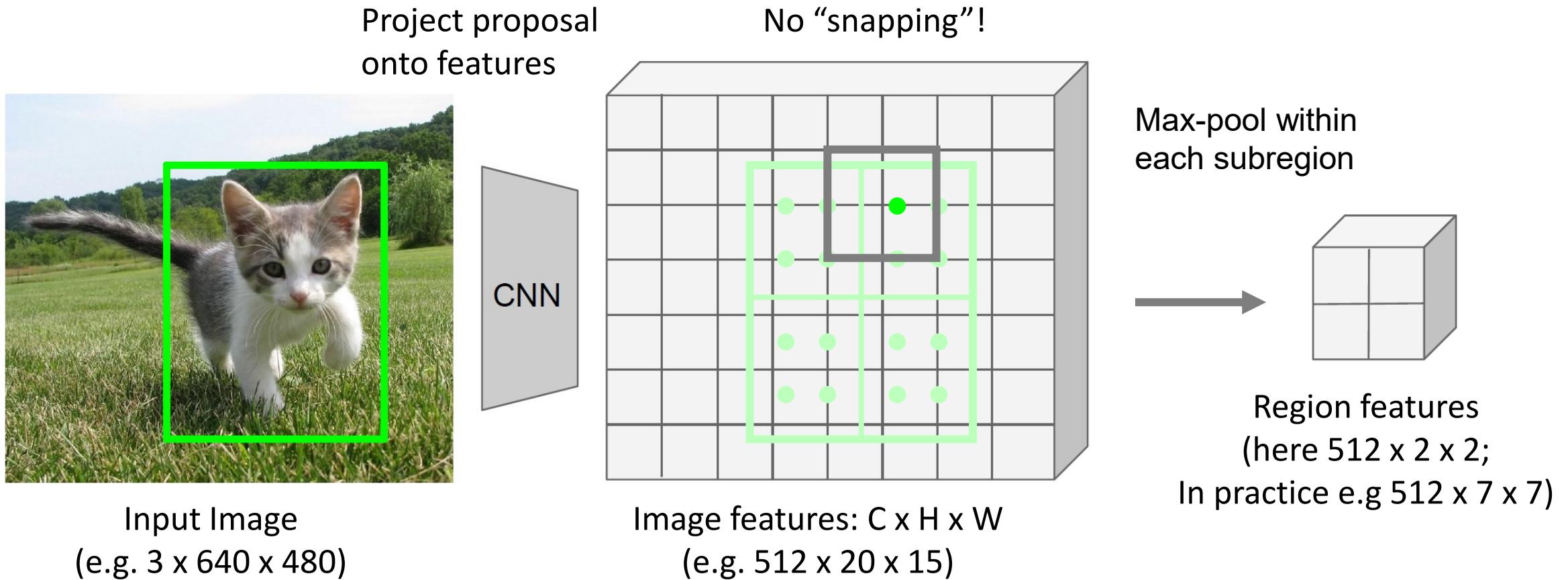
Image features: $C \times H \times W$
(e.g. $512 \times 20 \times 15$)



Feature f_{xy} for point (x, y)
is a linear combination of
features at its four
neighboring grid cells:

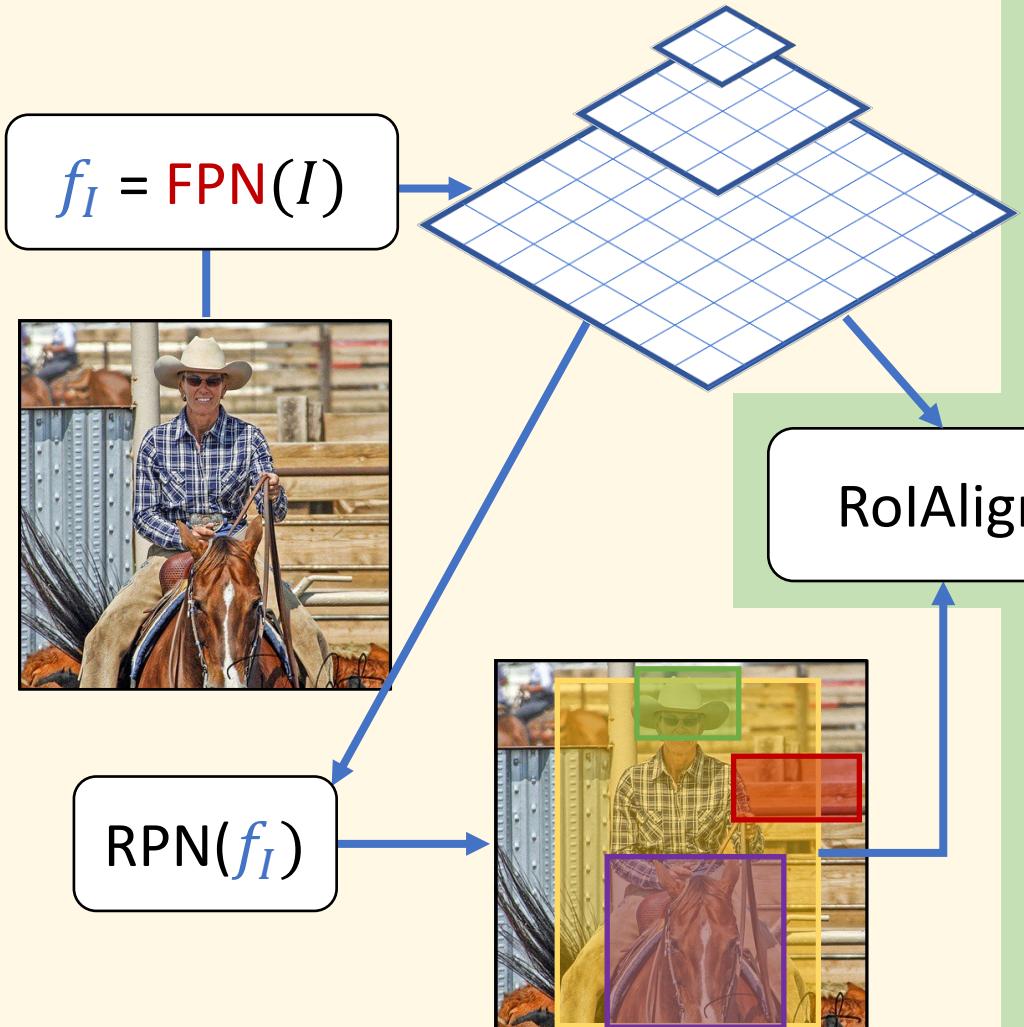
$$f_{xy} = \sum_{i,j=1}^2 f_{i,j} \max(0, 1 - |x - x_i|) \max(0, 1 - |y - y_j|)$$

Cropping Features: RoIAlign

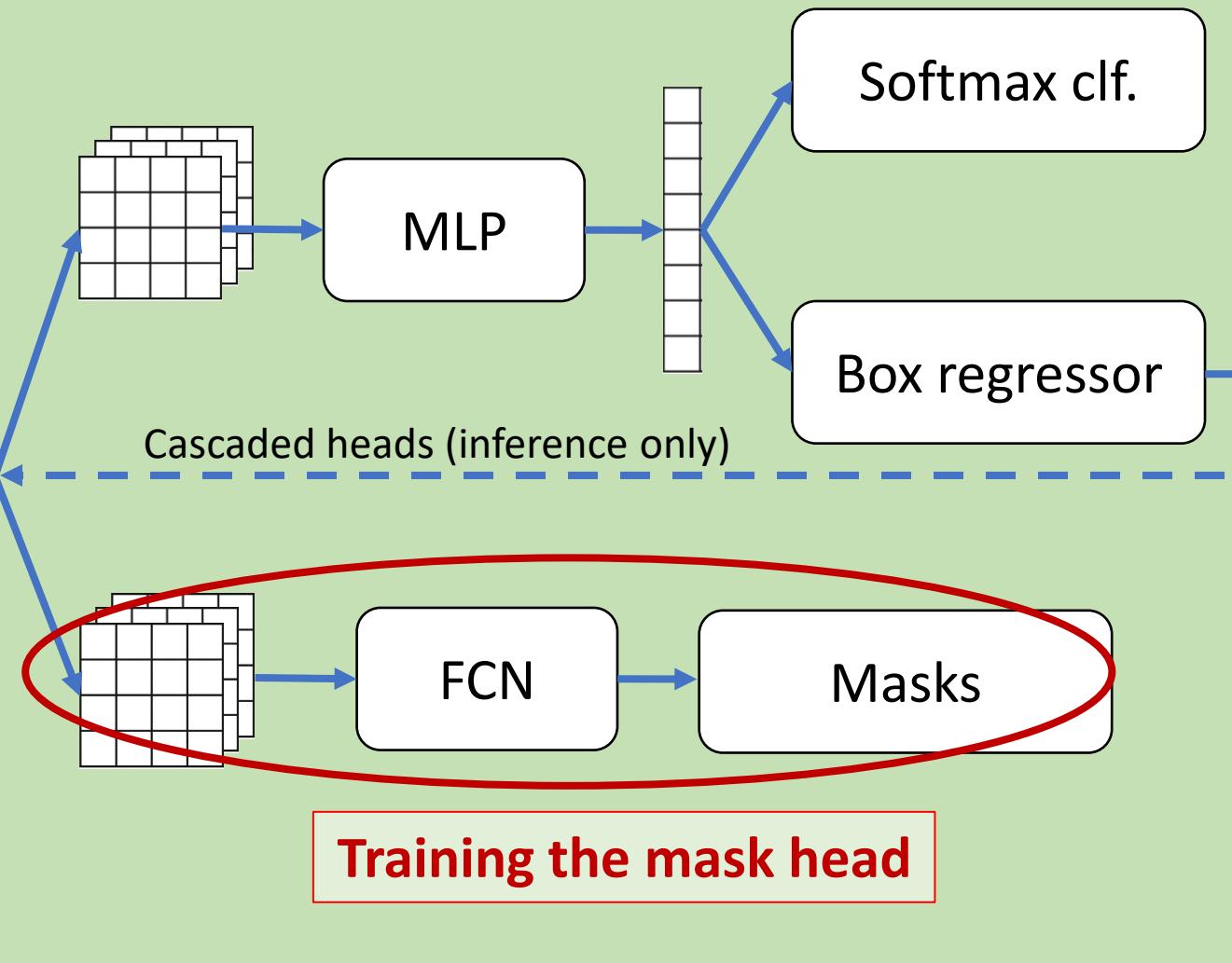


Mask R-CNN

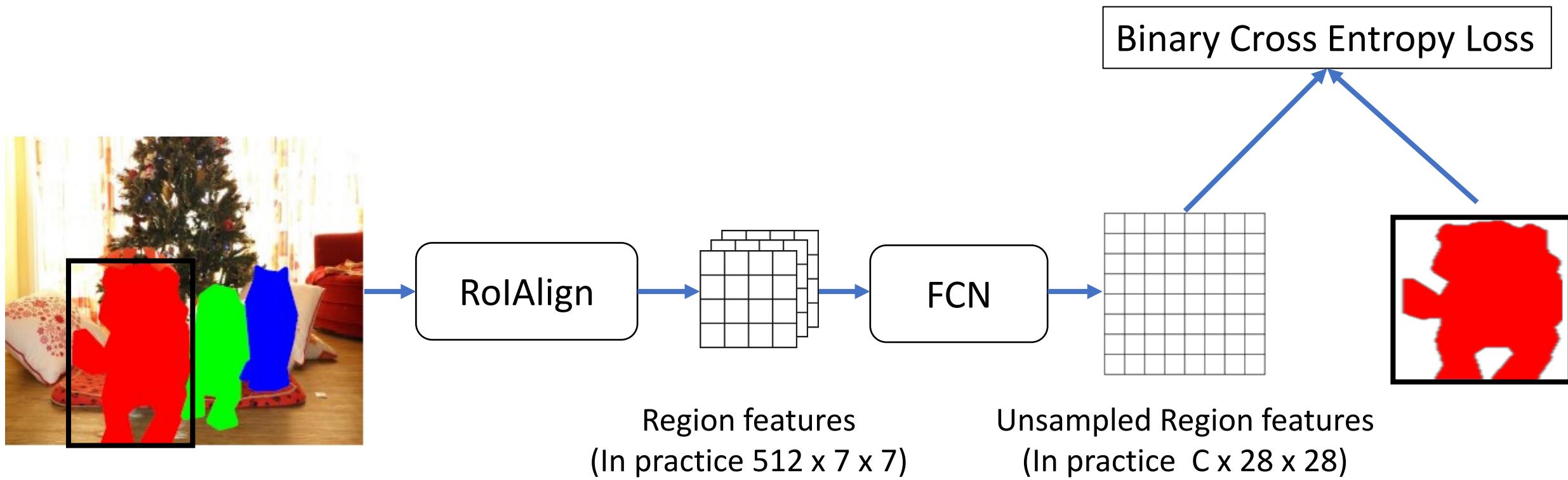
Per-image computation



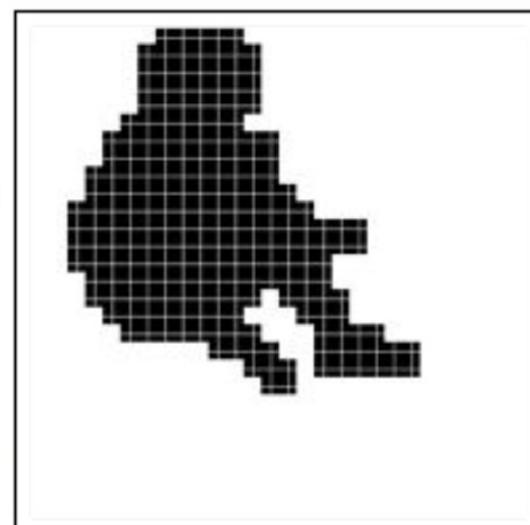
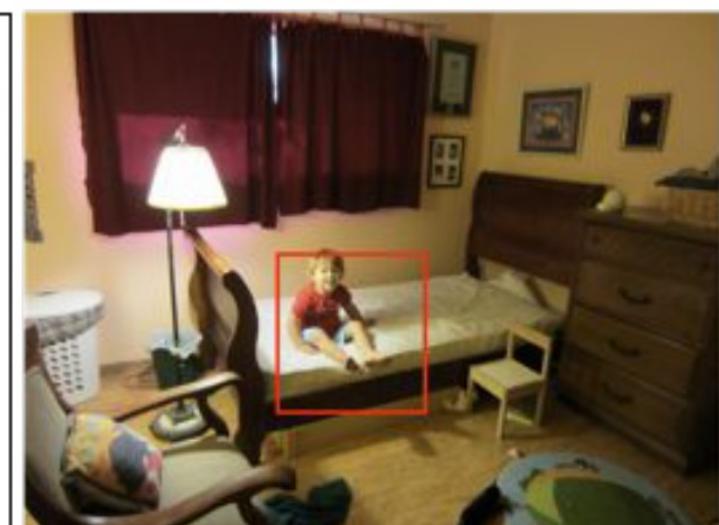
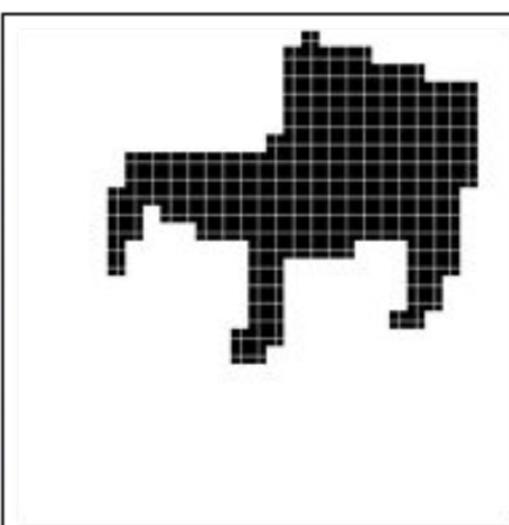
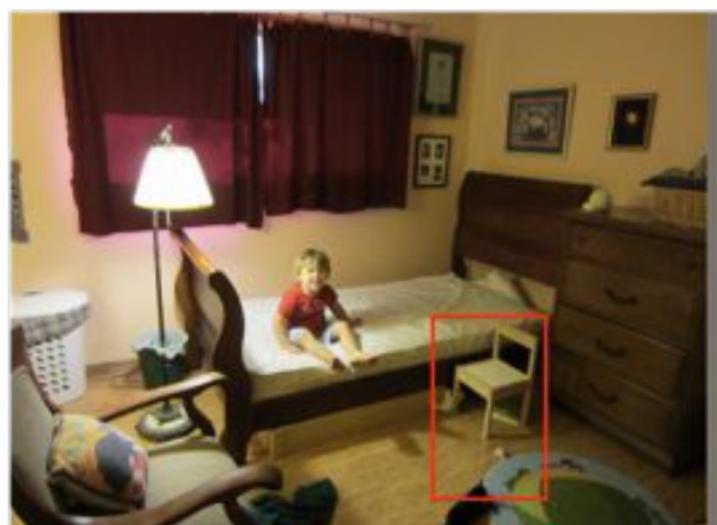
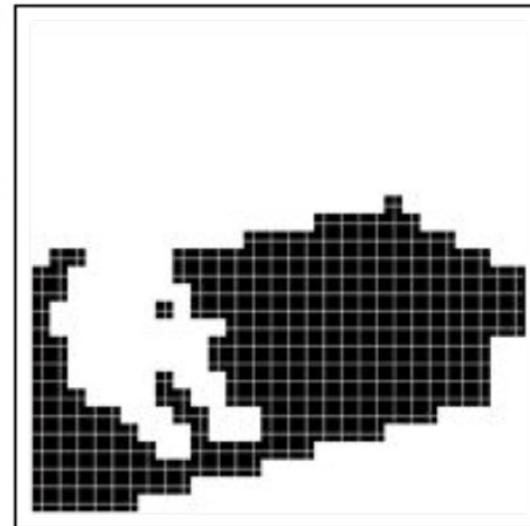
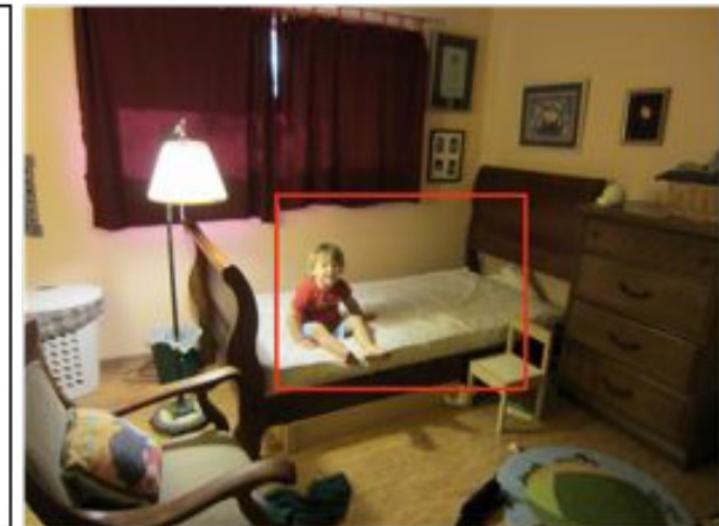
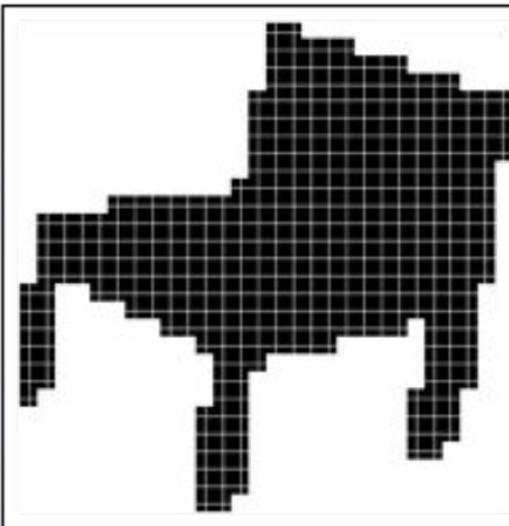
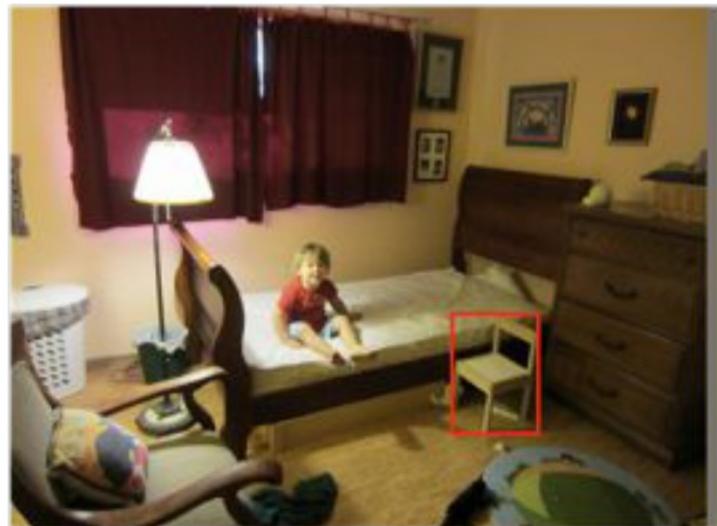
Per-region computation for each $r_i \in r(I)$



Training the mask head



Mask R-CNN: Example Mask Training Targets

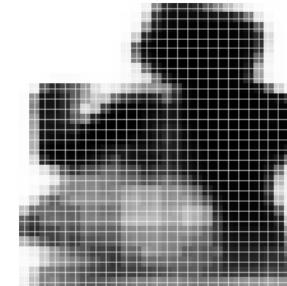


Mask Prediction



Validation image with box detection shown in red

28x28 soft prediction



Resized soft prediction



Final mask







