

Image Detection: Past, Present and Future

Presented by: Wei Ji

Overview

Topics to cover

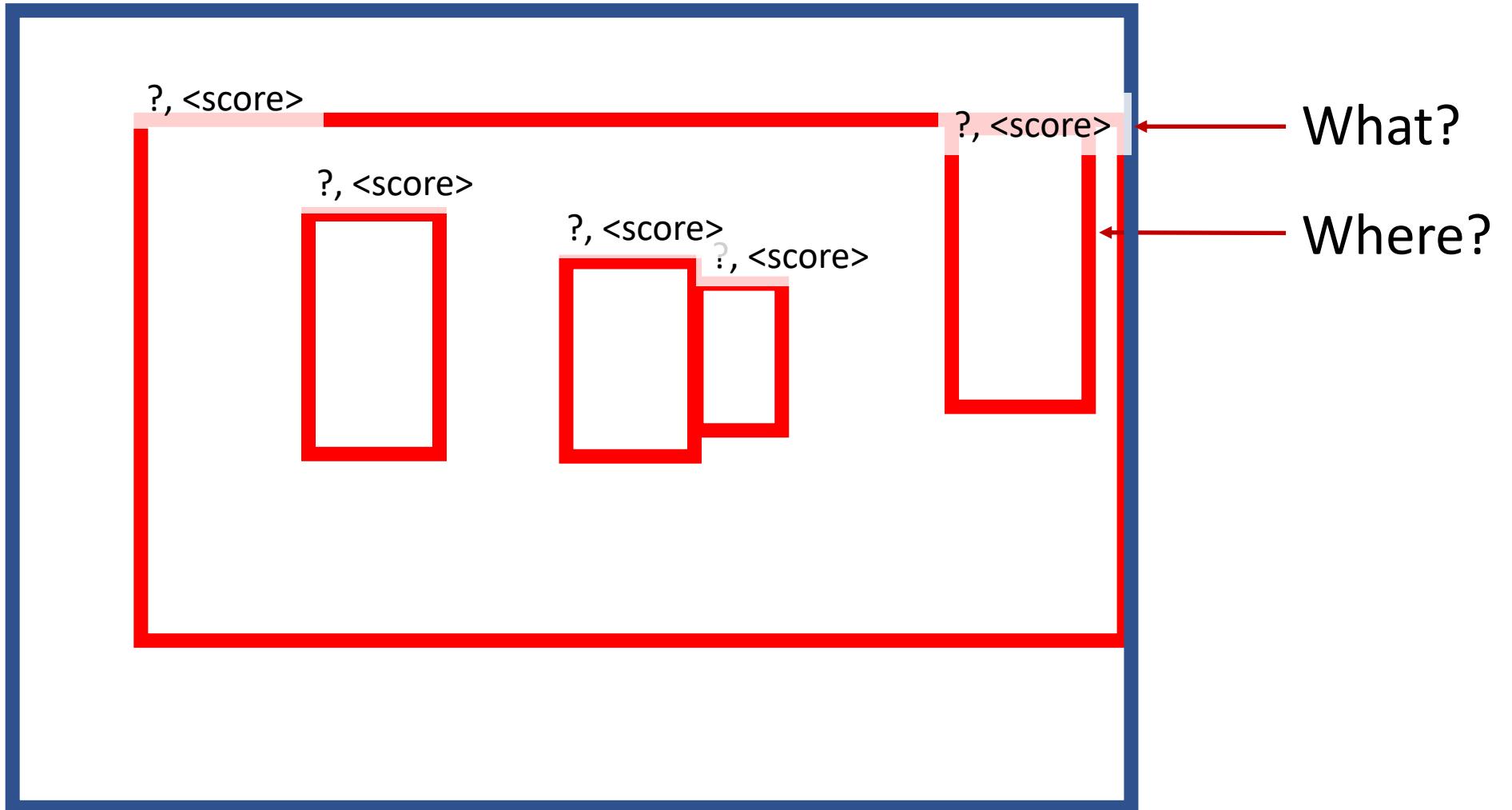
- Object detection intro
- The Generalized R-CNN framework
 - **R-CNN**
 - Fast R-CNN
 - Faster R-CNN
 - **Mask R-CNN**
- Future works + Discussion

Object Detection with Bounding Boxes

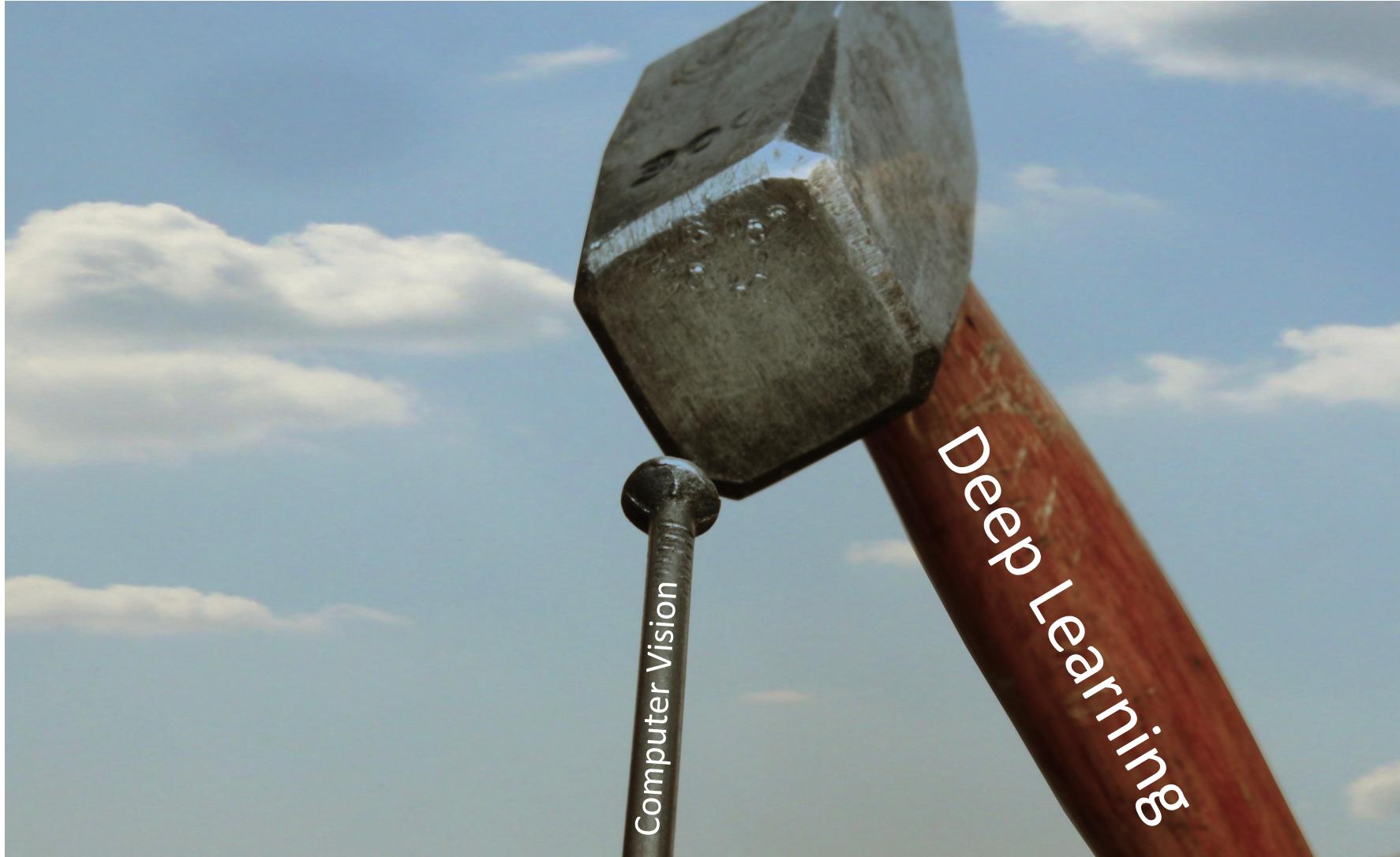


“Object detection”

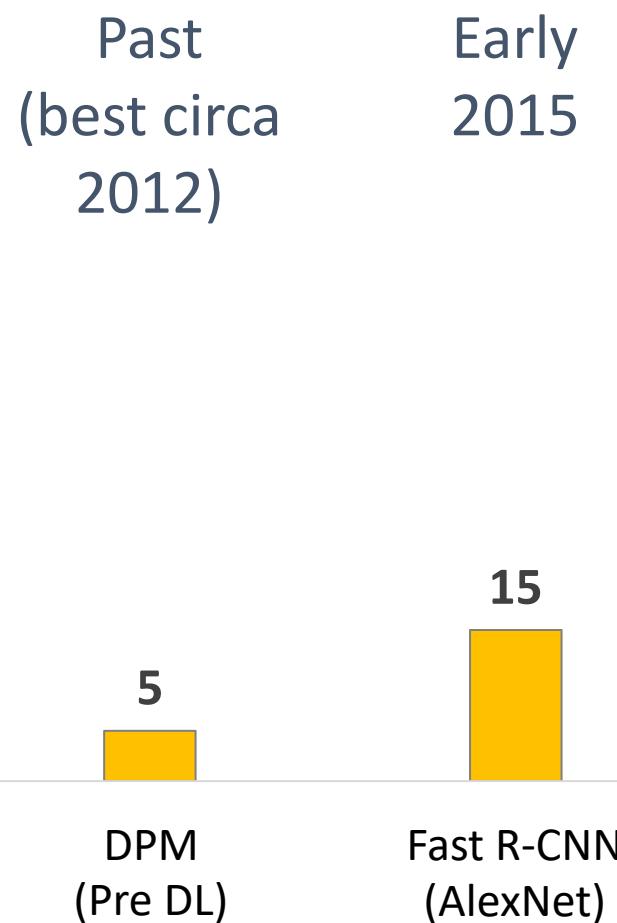
Object Detection



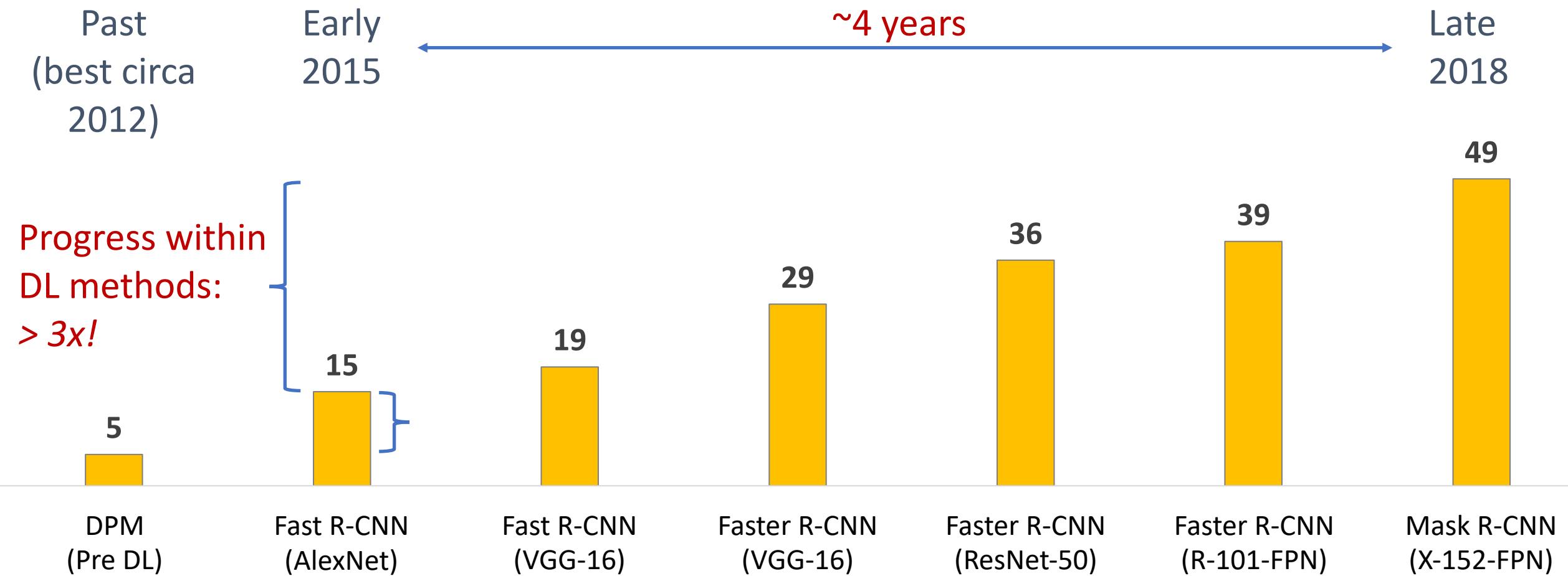
Modern Object Detection: Is this Picture Correct?



COCO Object Detection Average Precision (%)

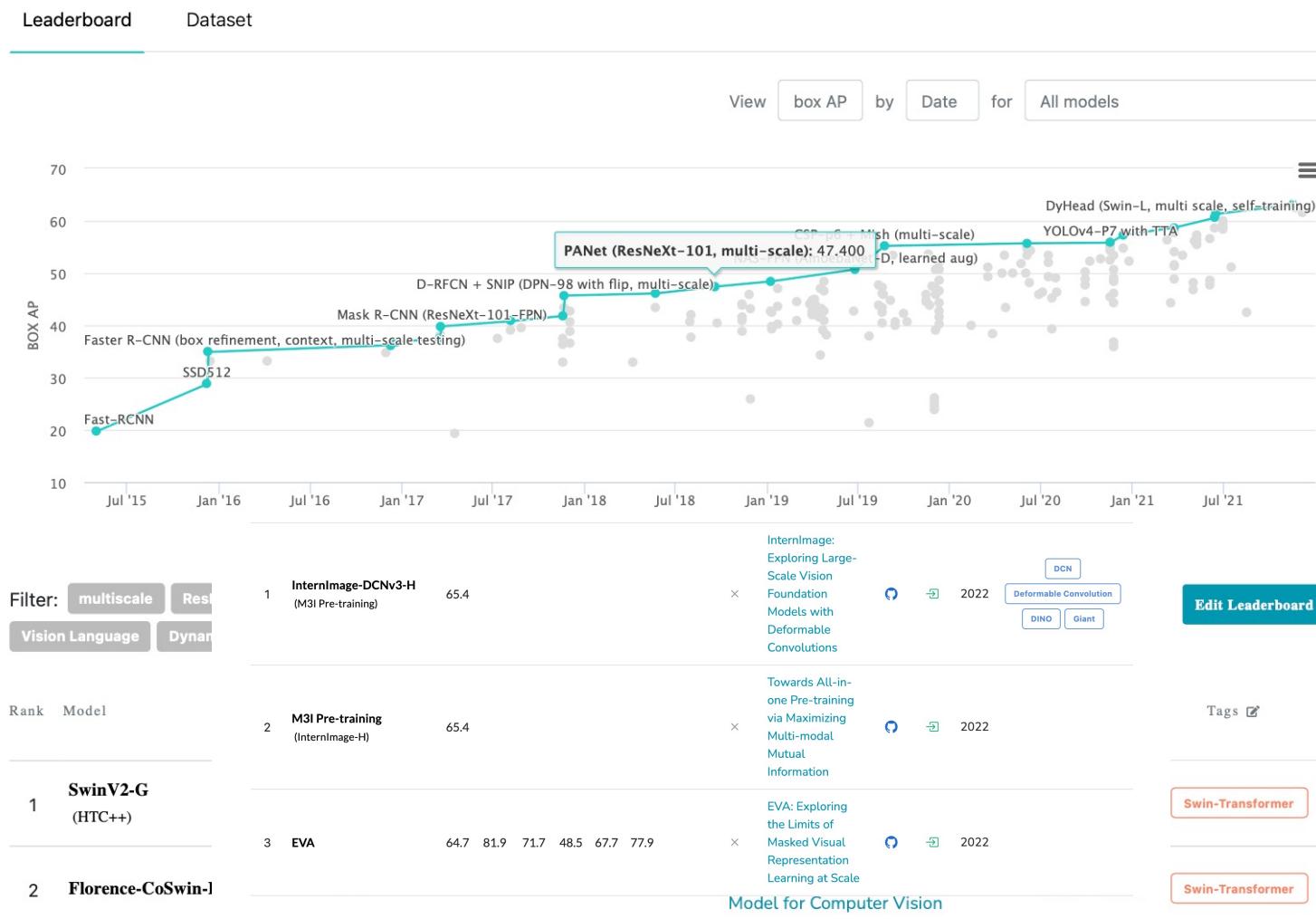


COCO Object Detection Average Precision (%)



COCO Object Detection Average Precision (%)

Object Detection on COCO test-dev



Progress within latest DL methods: *only 65.4%*



Past → Present

Detection Average Precision (%)

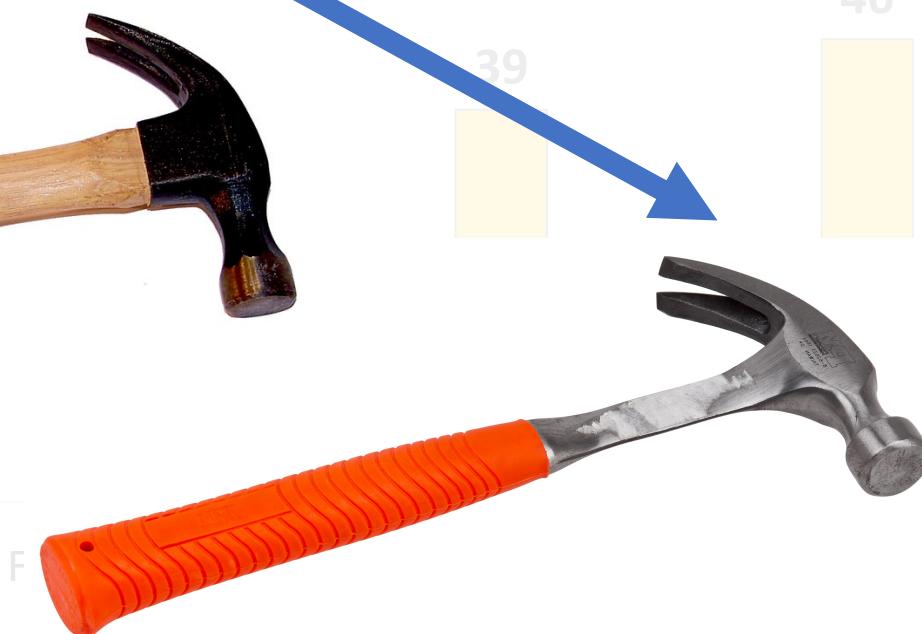
2.5 years

Building
a better
hammer

Late
2017

46

39



Progress within
DL methods:

Also 3x!

5

15

DPM
(Pre DL)

Fast R-CNN
(AlexNet)

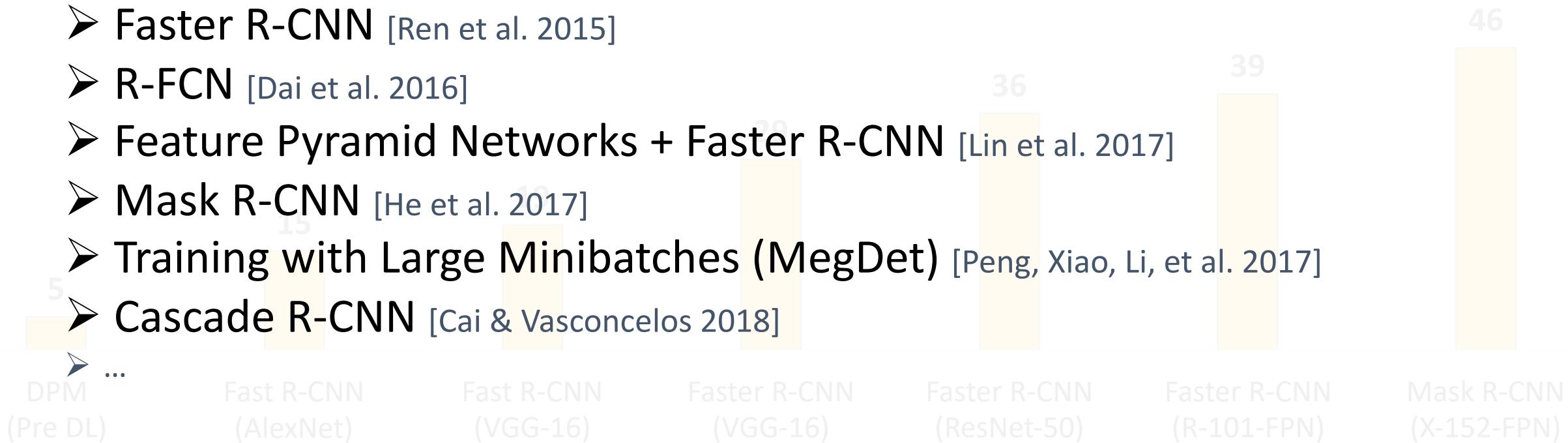
Fast R-CNN
(VGG-16)

Faster R-CNN
(VGG-16)

IN
N

Steady Progress on Boxes and Masks

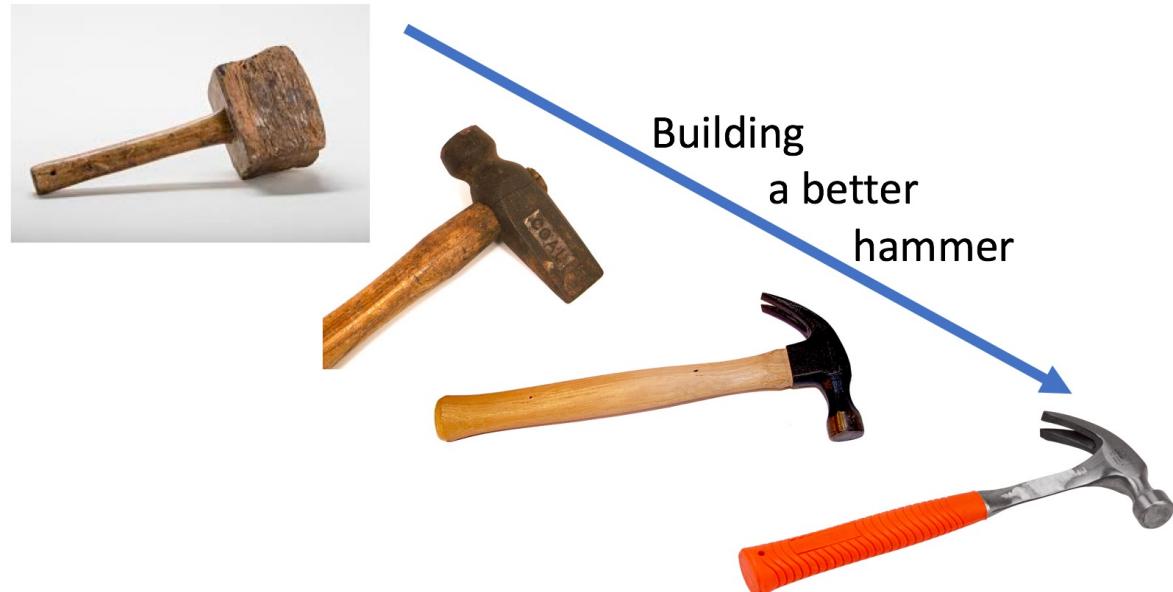
- R-CNN [Girshick et al. 2014]
- SPP-net [He et al. 2014]
- Fast R-CNN [Girshick. 2015]
- Faster R-CNN [Ren et al. 2015]
- R-FCN [Dai et al. 2016]
- Feature Pyramid Networks + Faster R-CNN [Lin et al. 2017]
- Mask R-CNN [He et al. 2017]
- Training with Large Minibatches (MegDet) [Peng, Xiao, Li, et al.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



Overview

Topics to cover

- Object detection intro
- The Generalized R-CNN framework
 - **R-CNN**
 - Fast R-CNN
 - Faster R-CNN
 - **Mask R-CNN**
- Future works
 - **SwinV2**



Overview

Topics to cover

- Object detection intro
- The Generalized R-CNN framework
 - **R-CNN**
 - Fast R-CNN
 - Faster R-CNN
 - **Mask R-CNN**
- Future works

Rolling Pin



Vanilla Solution – Translate to Image Recognition

28 x 28



Hundreds of windows per image

https://github.com/layumi/2015_Face_Detection

Face Detection – Sliding Window is all you need



We return locations that the face classifier predicts high confidence.

Overview

Topics to cover

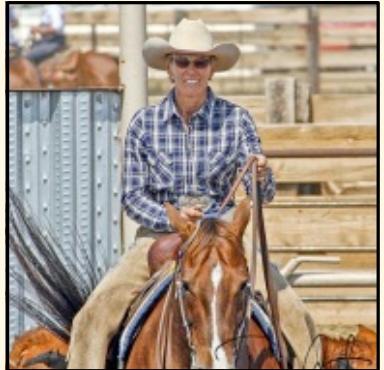
- Object detection intro
- The Generalized R-CNN framework
 - **R-CNN**
 - Fast R-CNN
 - Faster R-CNN
 - **Mask R-CNN**
- Future works

Old-school Hammer



R-CNN (Region-based Convolutional Neural Net)

Per-image computation



I:

Per-region computation

R-CNN (Region-based Convolutional Neural Net)

Per-image computation

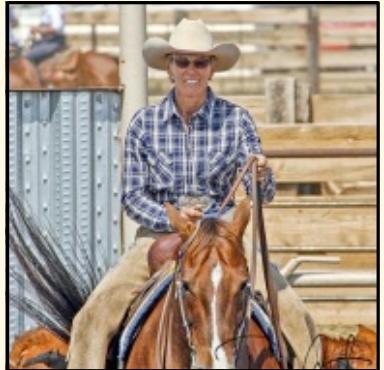


I:

Per-region computation

R-CNN (Region-based Convolutional Neural Net)

Per-image computation

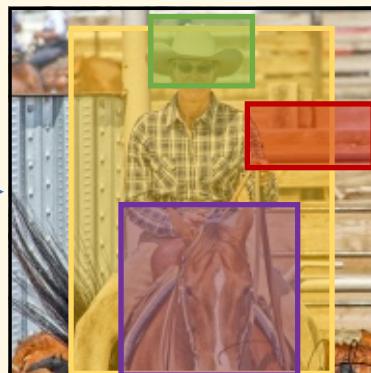
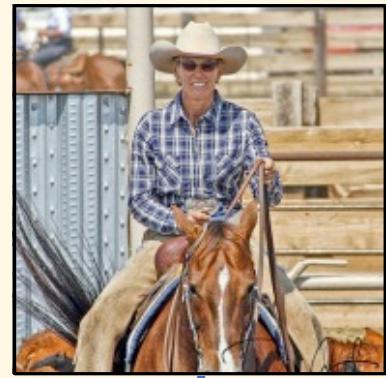


I:

Per-region computation

R-CNN (Region-based Convolutional Neural Net)

Per-image computation



Selective search,
Edge Boxes,
MCG, ...

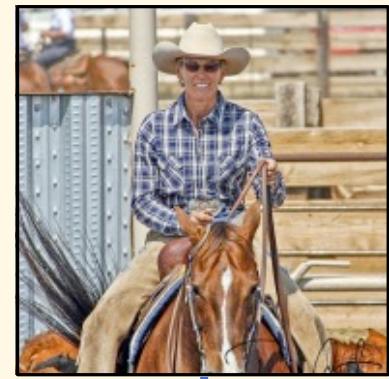
1

Per-region computation

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

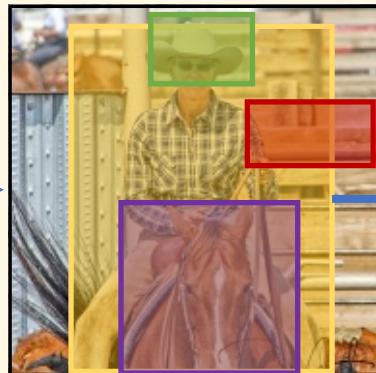
R-CNN

Per-image computation



Selective search,
Edge Boxes,
MCG, ...

1



Crop &
warp

2

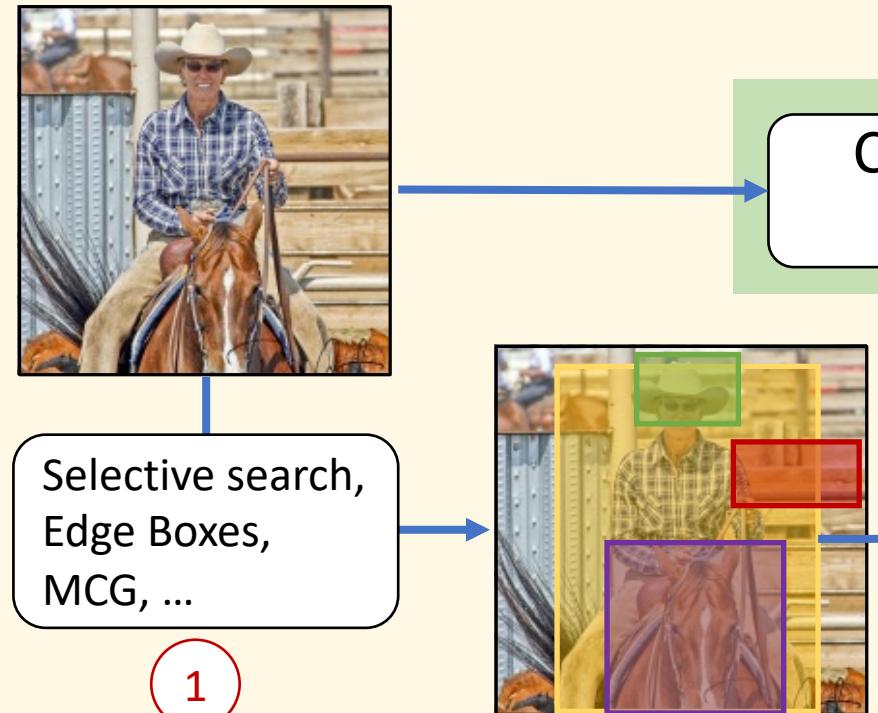


Per-region computation for each $r_i \in r(I)$

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

R-CNN

Per-image computation



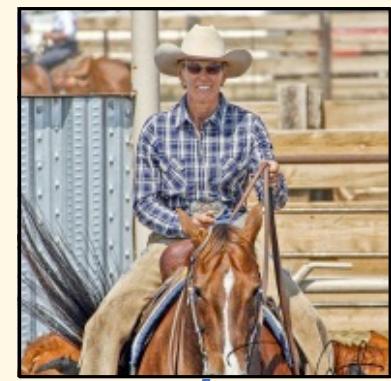
Per-region computation for each $r_i \in r(I)$



3. Forward propagate the fixed-size network input to get a feature representation

R-CNN

Per-image computation



Selective search,
Edge Boxes,
MCG, ...

1

Crop &
warp

2



ConvNet(r_i)

3



Linear
classifier

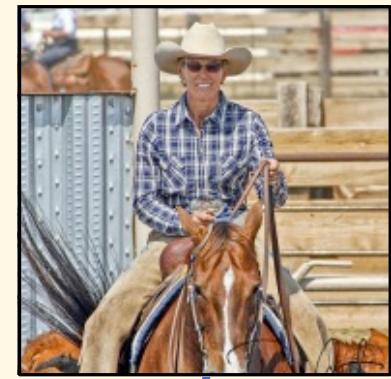
4

4. Object classification

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

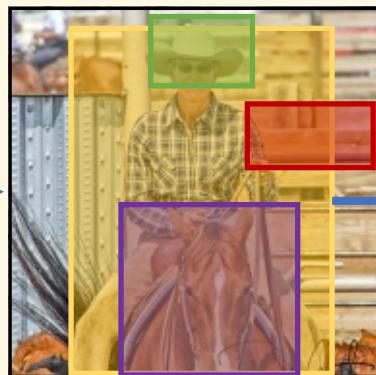
R-CNN

Per-image computation



Selective search,
Edge Boxes,
MCG, ...

1



Crop &
warp

2



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

ConvNet(r_i)

3

Linear
classifier

4

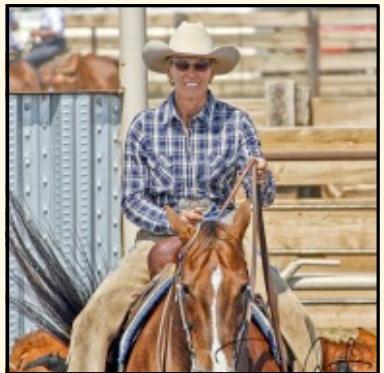
Box regressor

5

5. Refine proposal localization
with bounding-box regression

Generalized R-CNN Framework

Per-image computation



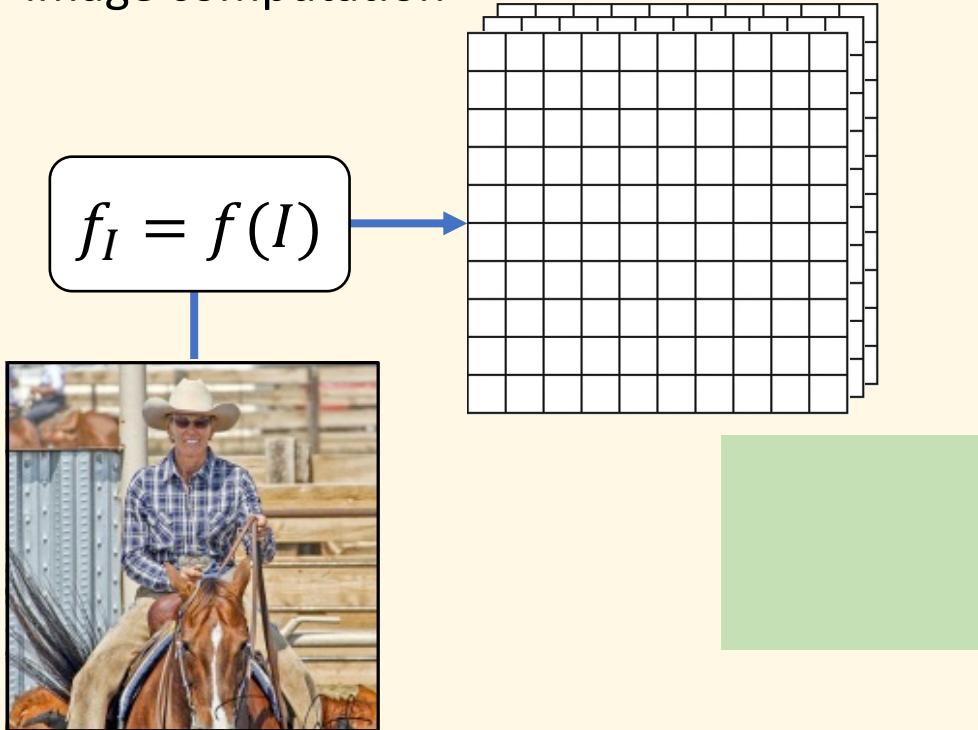
I :

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

Input image
per-image operations | per-region operations

Generalized R-CNN Framework

Per-image computation

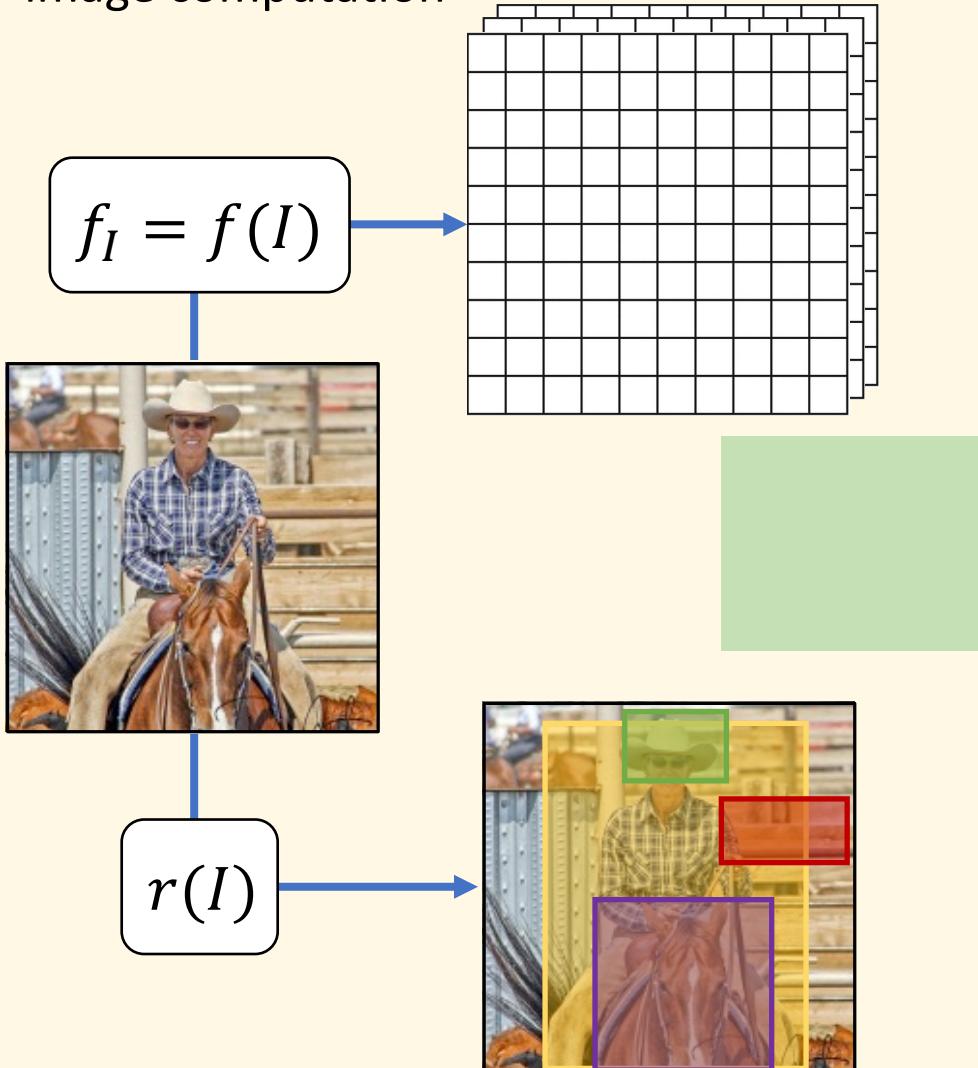


Per-region computation for each $r_i \in r(I)$

Transformation of the input image
into a featurized representation

Generalized R-CNN Framework

Per-image computation

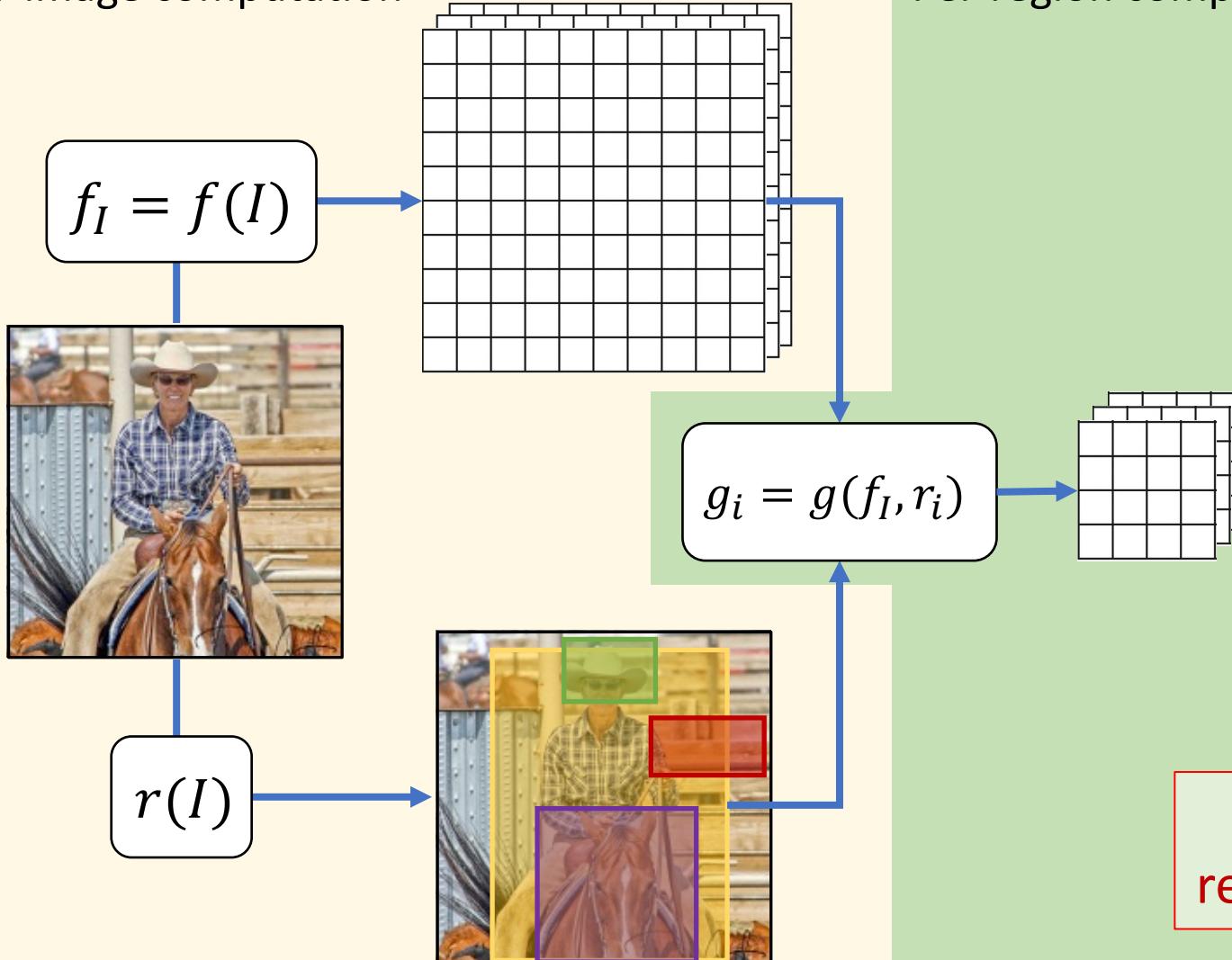


Per-region computation for each $r_i \in r(I)$

Region of Interest proposals
computed for the image

Generalized R-CNN Framework

Per-image computation

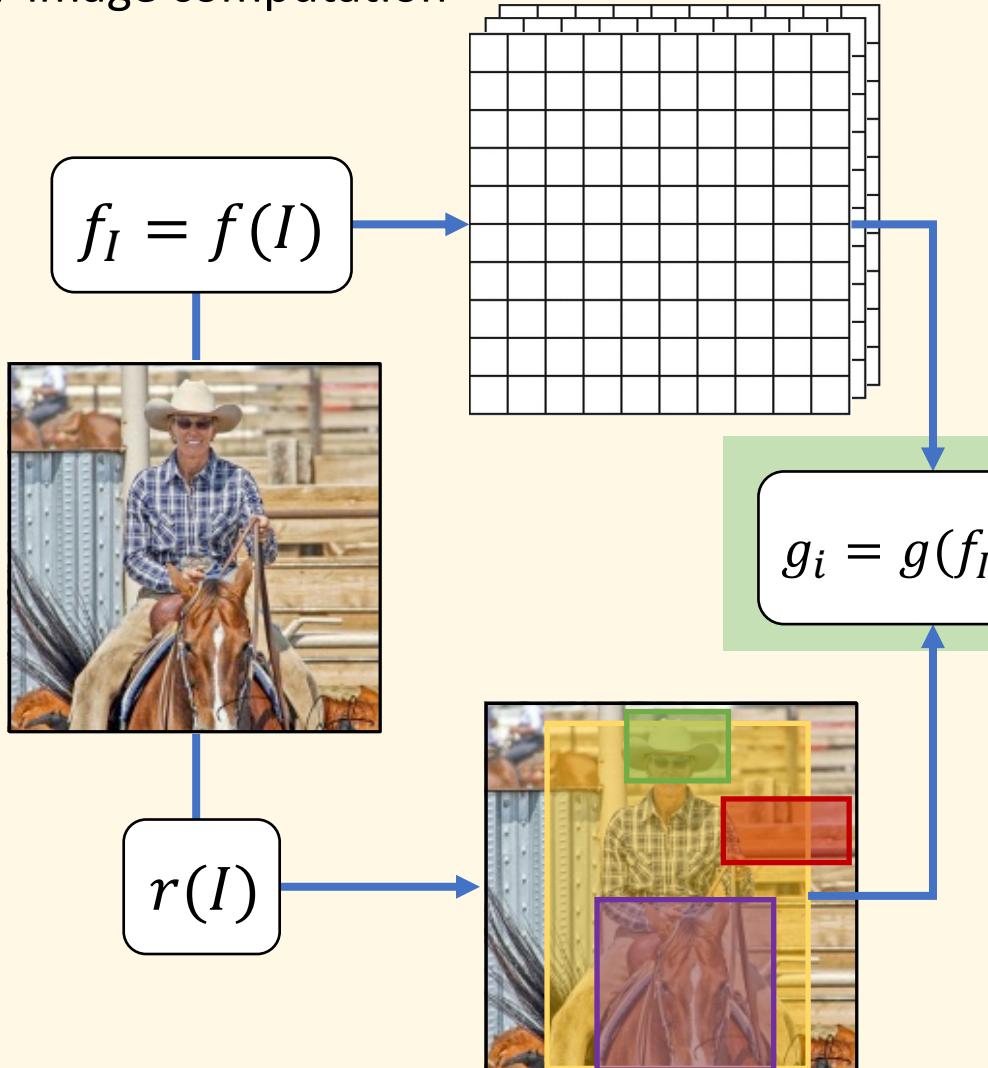


Per-region computation for each $r_i \in r(I)$

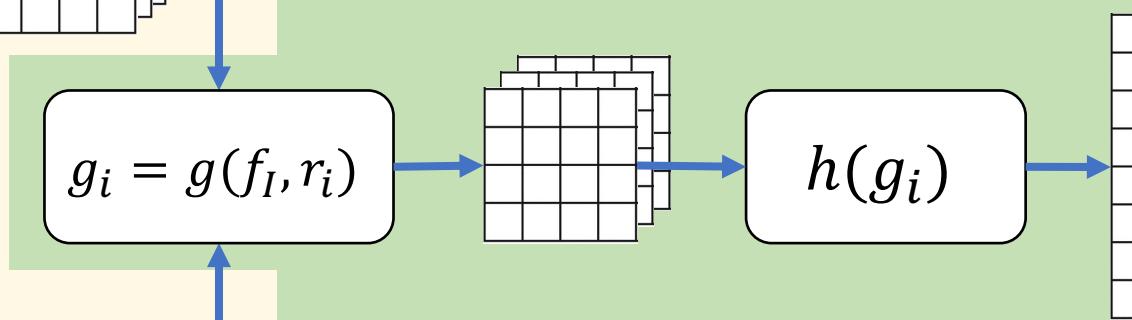
Compute a featurized representation of each proposal

Generalized R-CNN Framework

Per-image computation



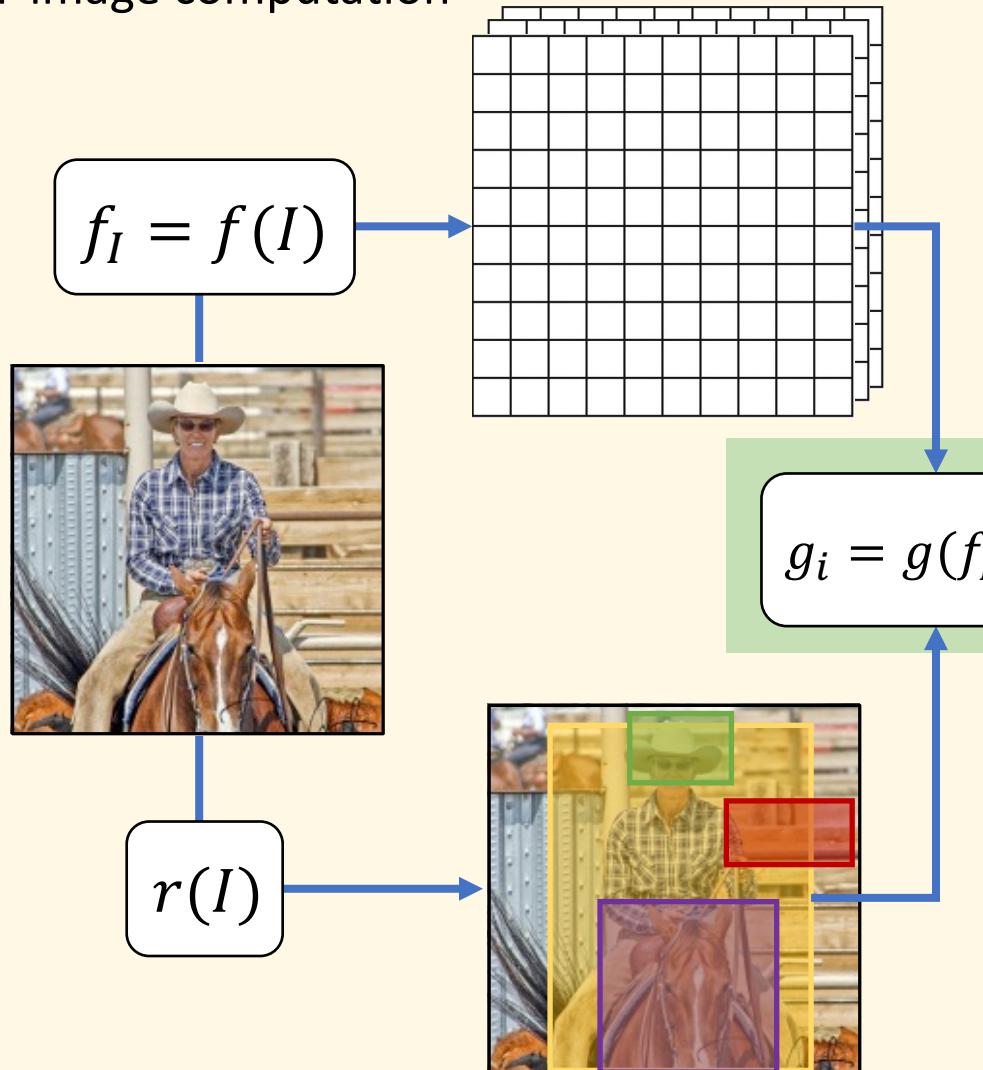
Per-region computation for each $r_i \in r(I)$



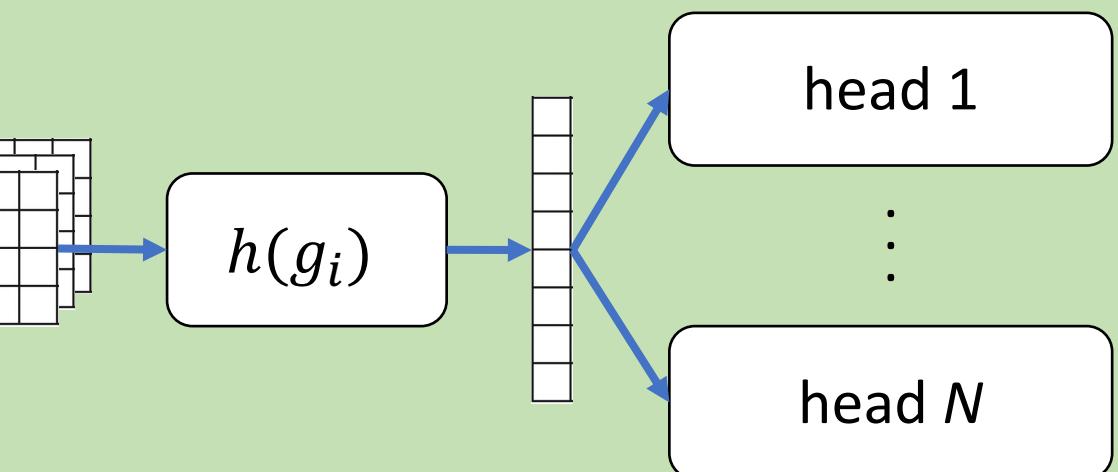
Apply additional processing
to each proposal feature

Generalized R-CNN Framework

Per-image computation



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

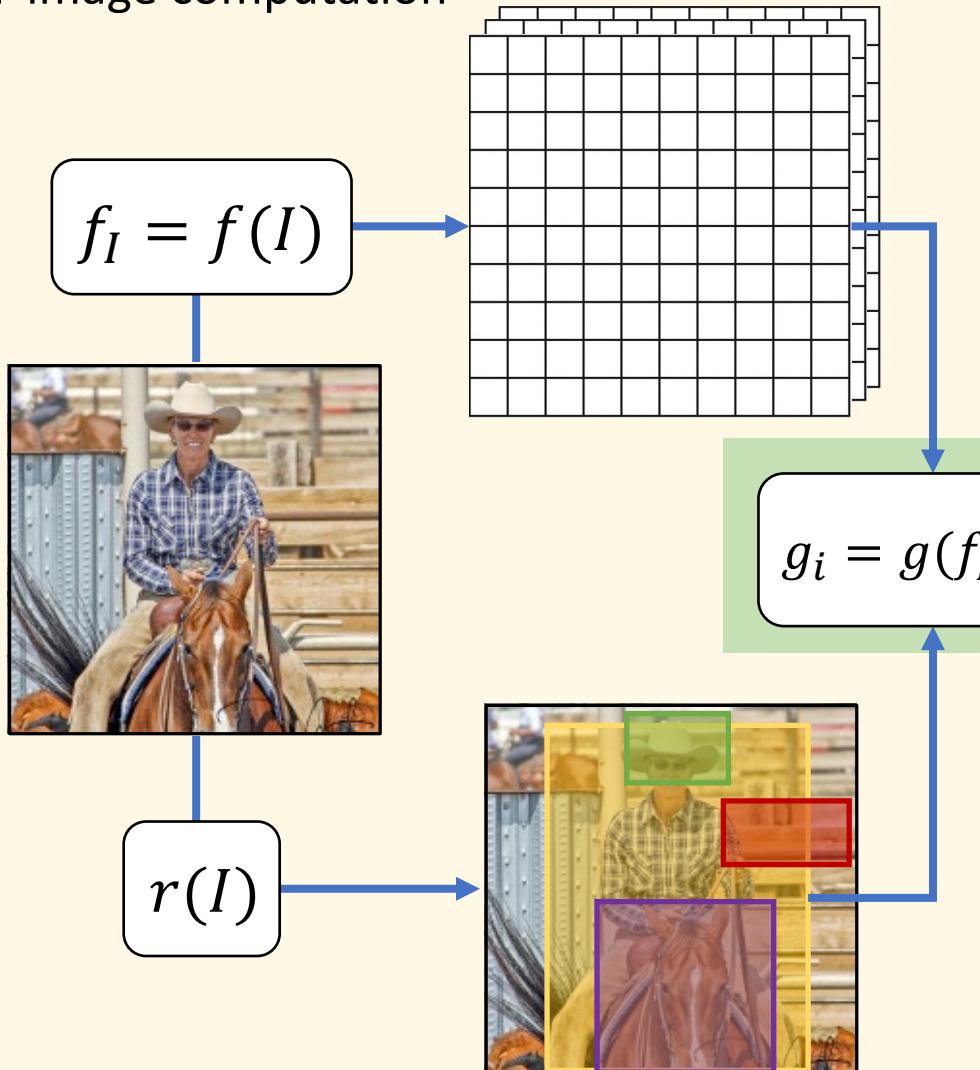


Apply multiple ‘heads’
to make task-specific predictions

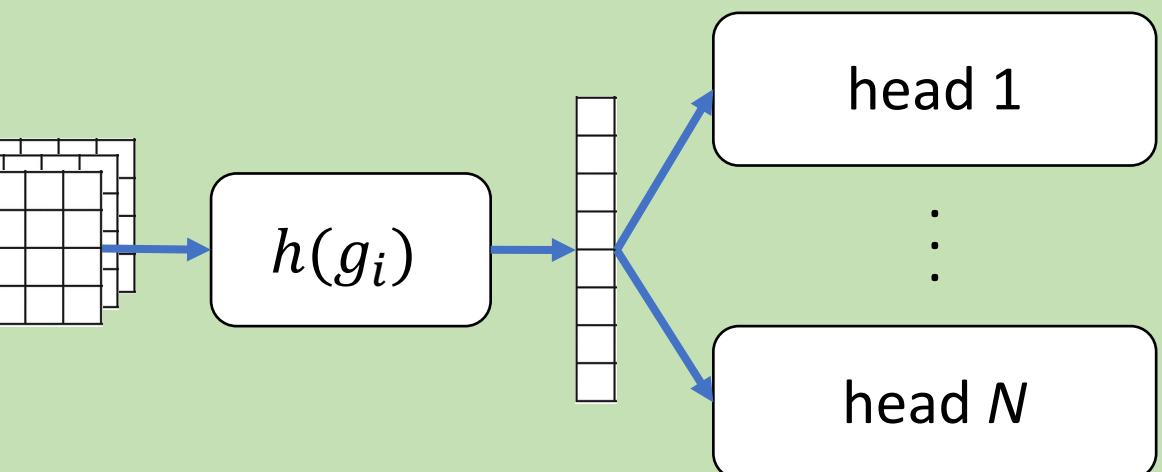
What Does R-CNN Look Like in the Generalized Framework?

R-CNN in the Generalized Framework

Per-image computation



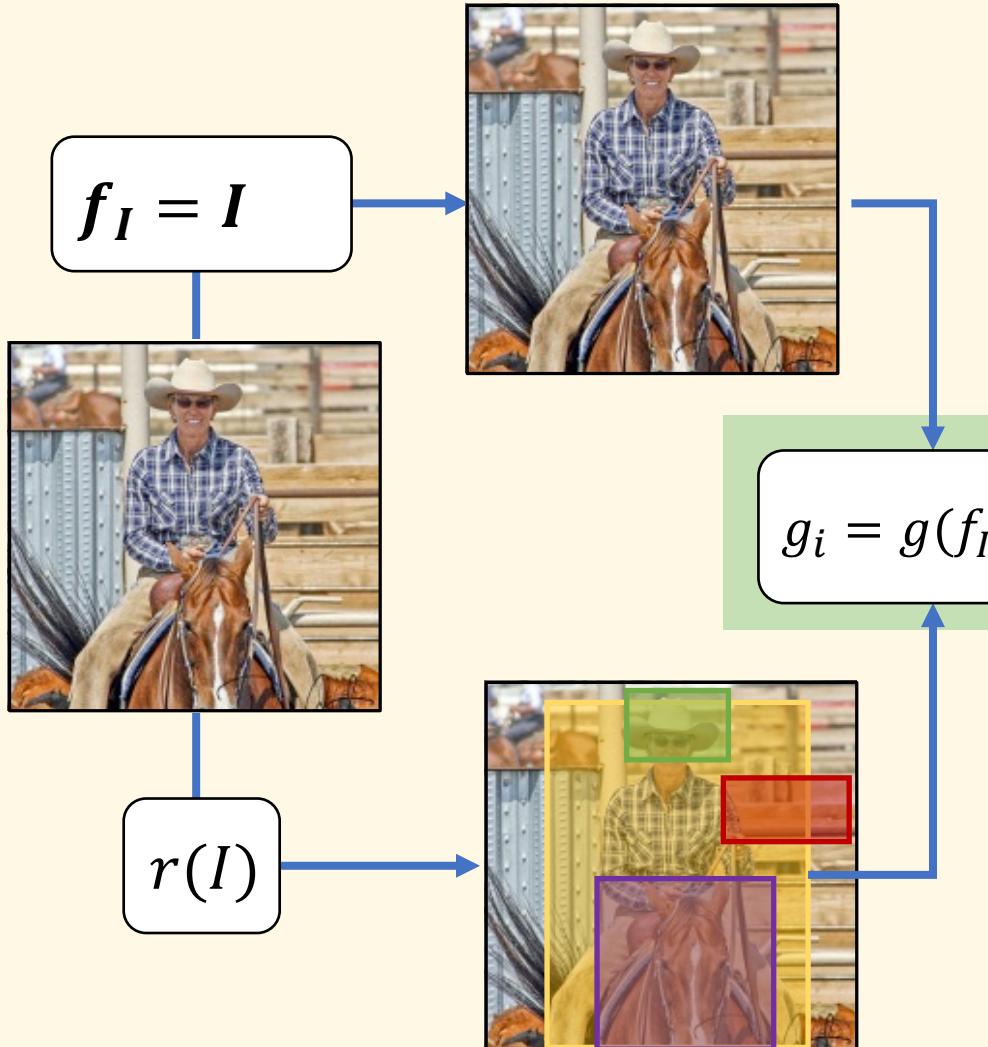
Per-region computation for each $r_i \in r(I)$



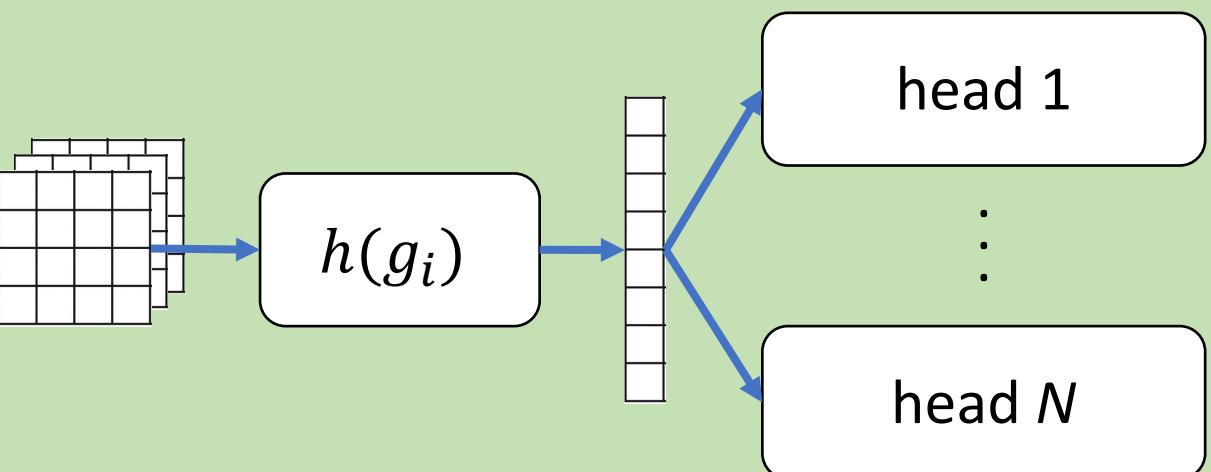
R-CNN in the generalized framework

R-CNN in the Generalized Framework

Per-image computation



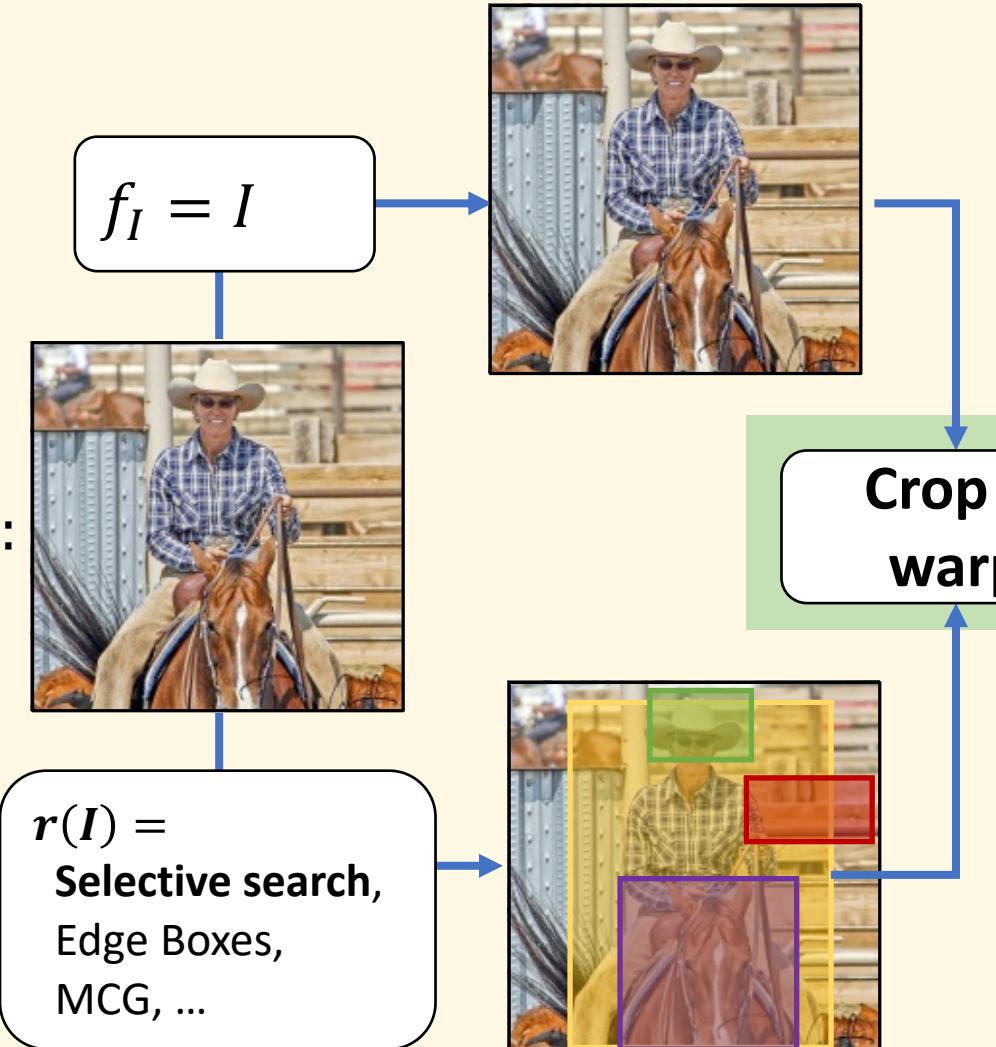
Per-region computation for each $r_i \in r(I)$



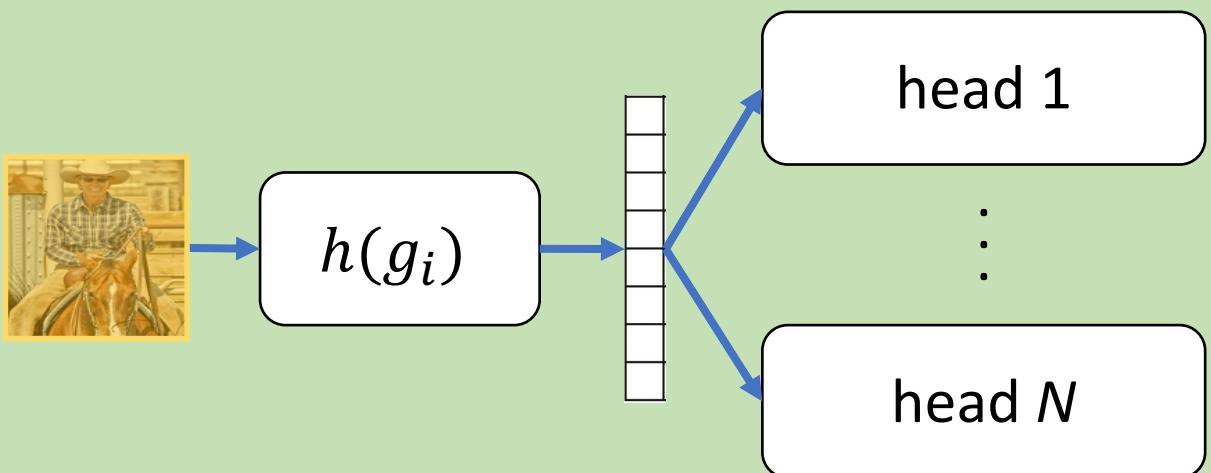
R-CNN in the generalized framework

R-CNN in the Generalized Framework

Per-image computation



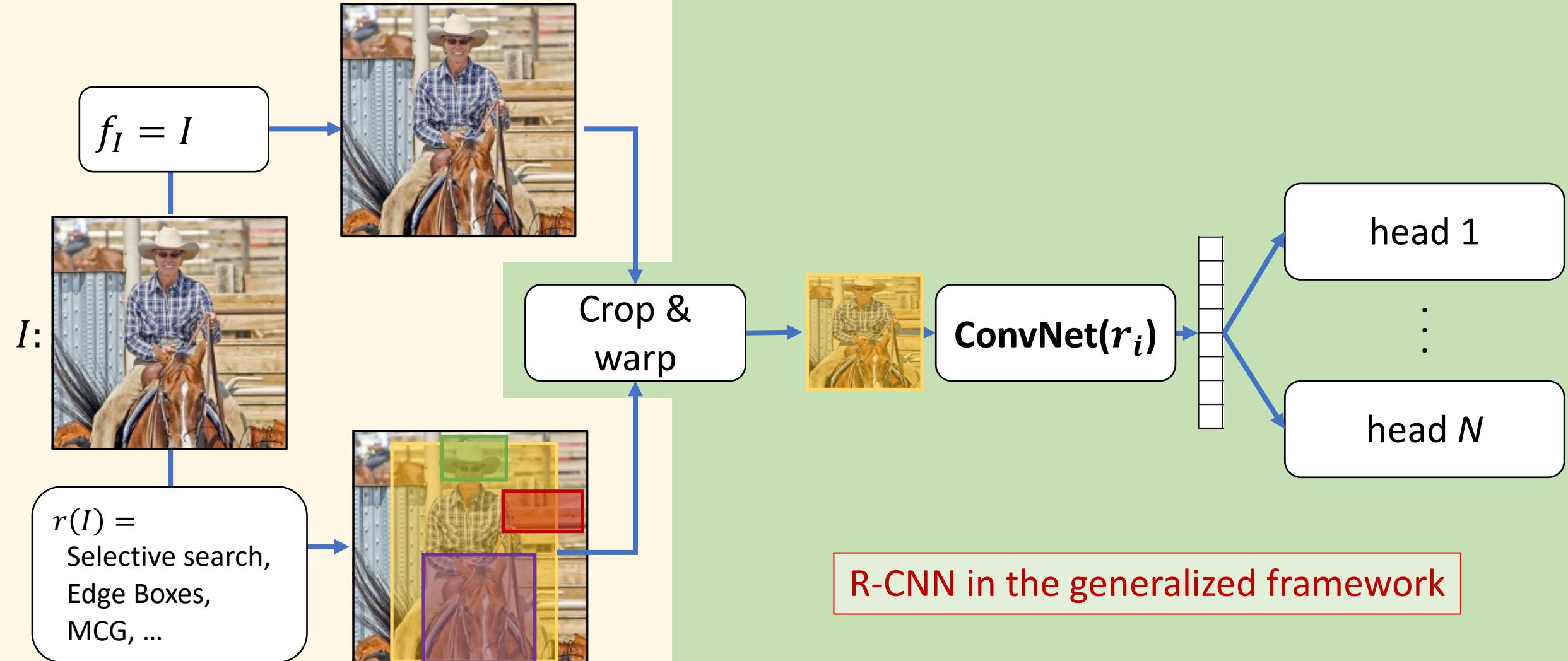
Per-region computation for each $r_i \in r(I)$



R-CNN in the generalized framework

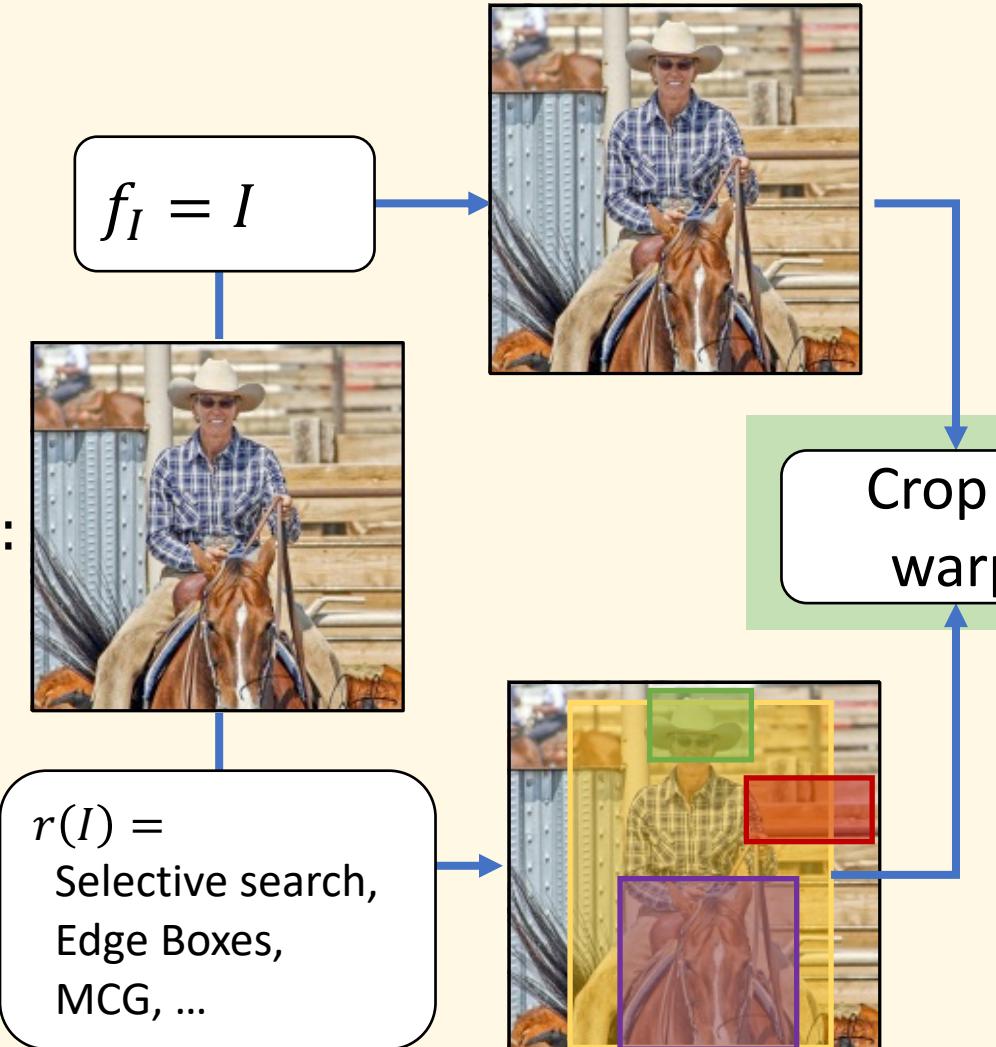
R-CNN in the Generalized Framework

Per-image computation

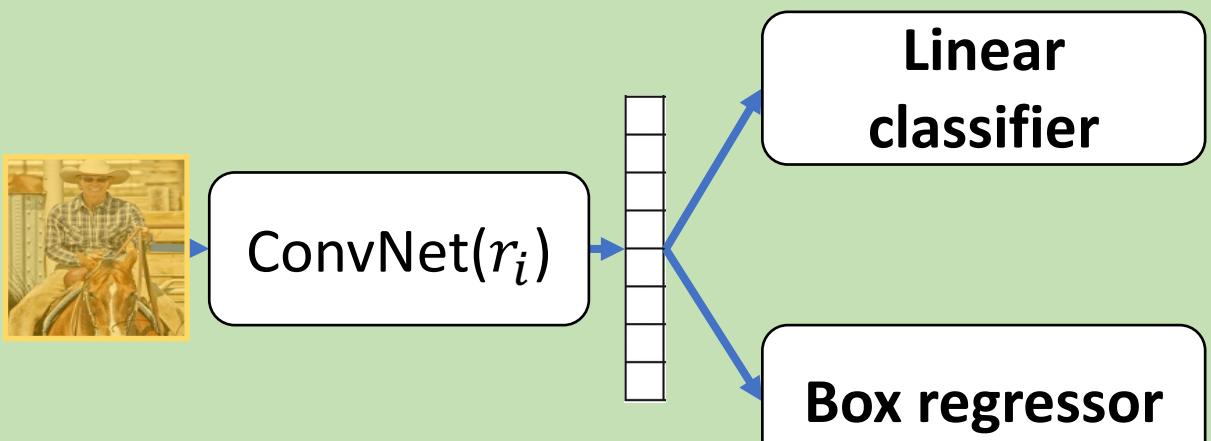


R-CNN in the Generalized Framework

Per-image computation



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



R-CNN in the generalized framework

The Problem with R-CNN

R-CNN



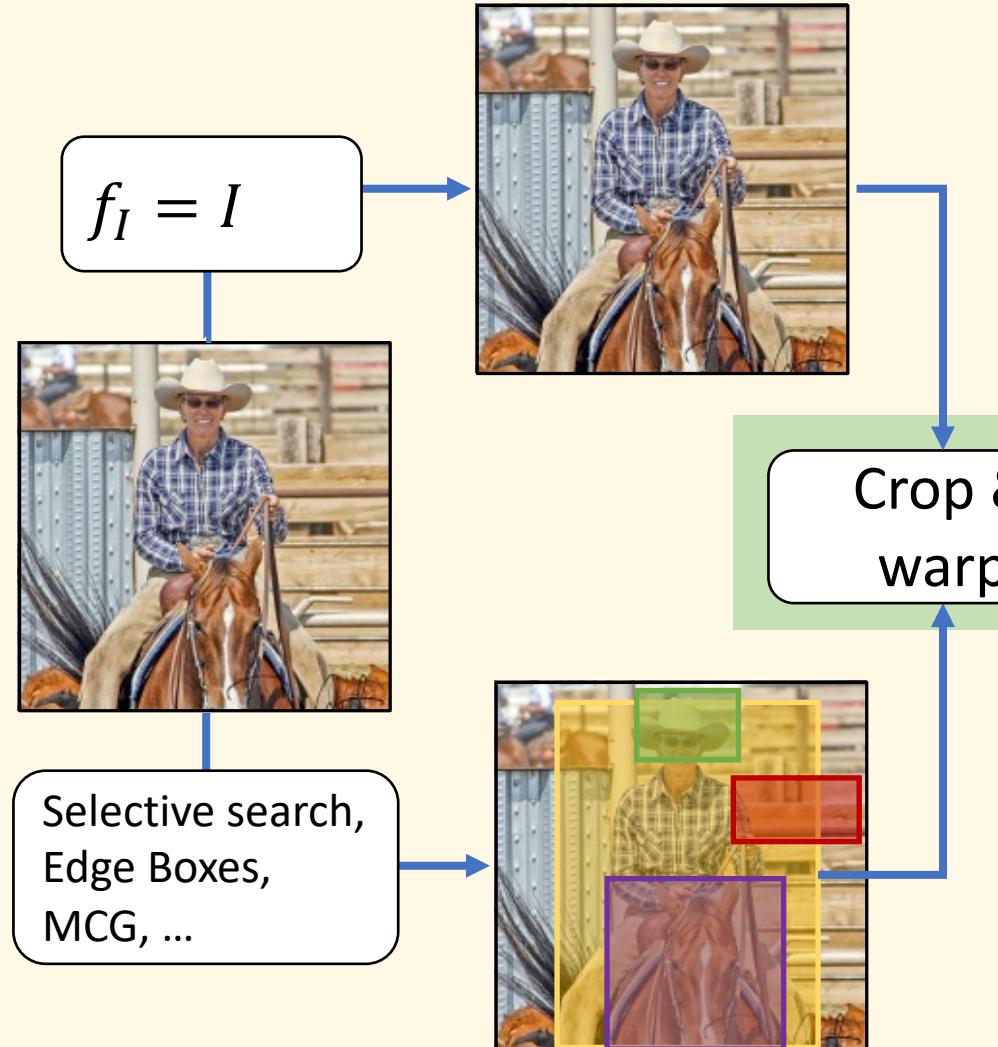
?



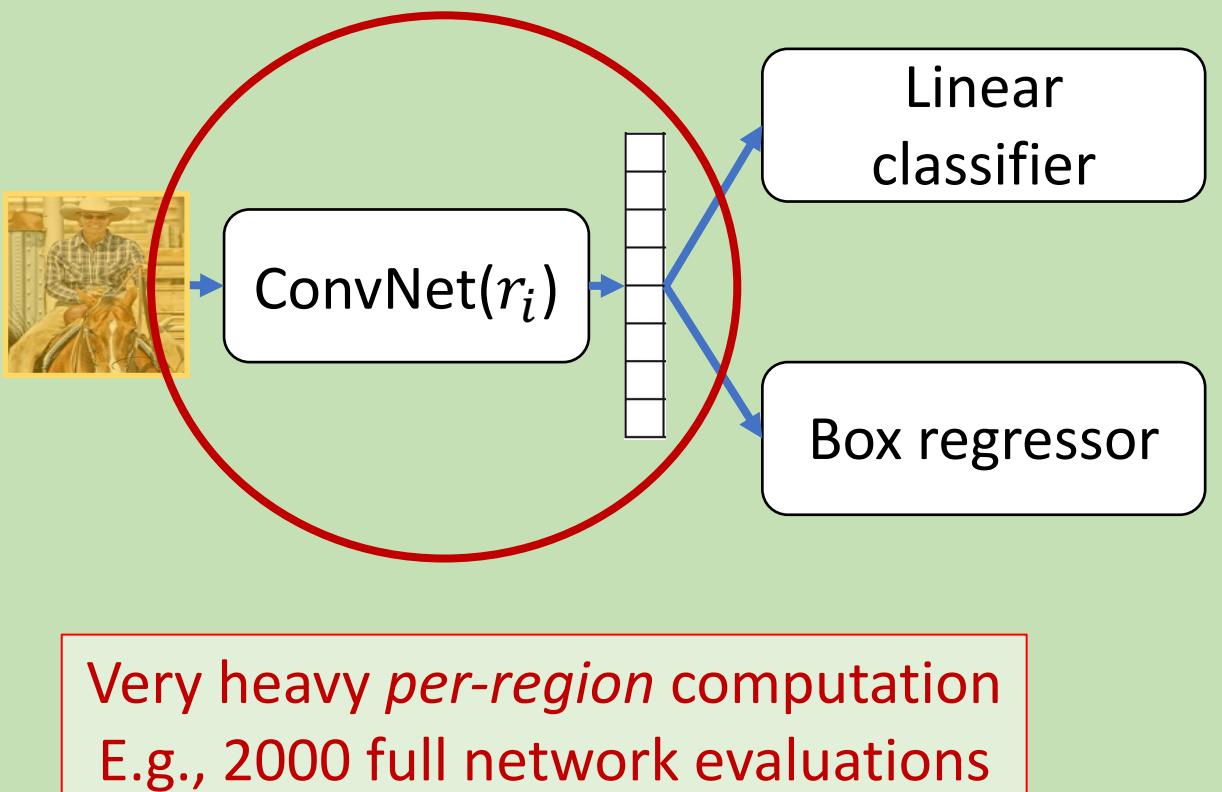
Building
a better
hammer

“Slow” R-CNN

Per-image computation



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



Fast R-CNN

References

- K. He, X. Zhang, S. Ren, and J. Sun. Spatial pyramid pooling in deep convolutional networks for visual recognition. In ECCV, 2014.
- R. Girshick. Fast R-CNN. In ICCV, 2015.

R-CNN



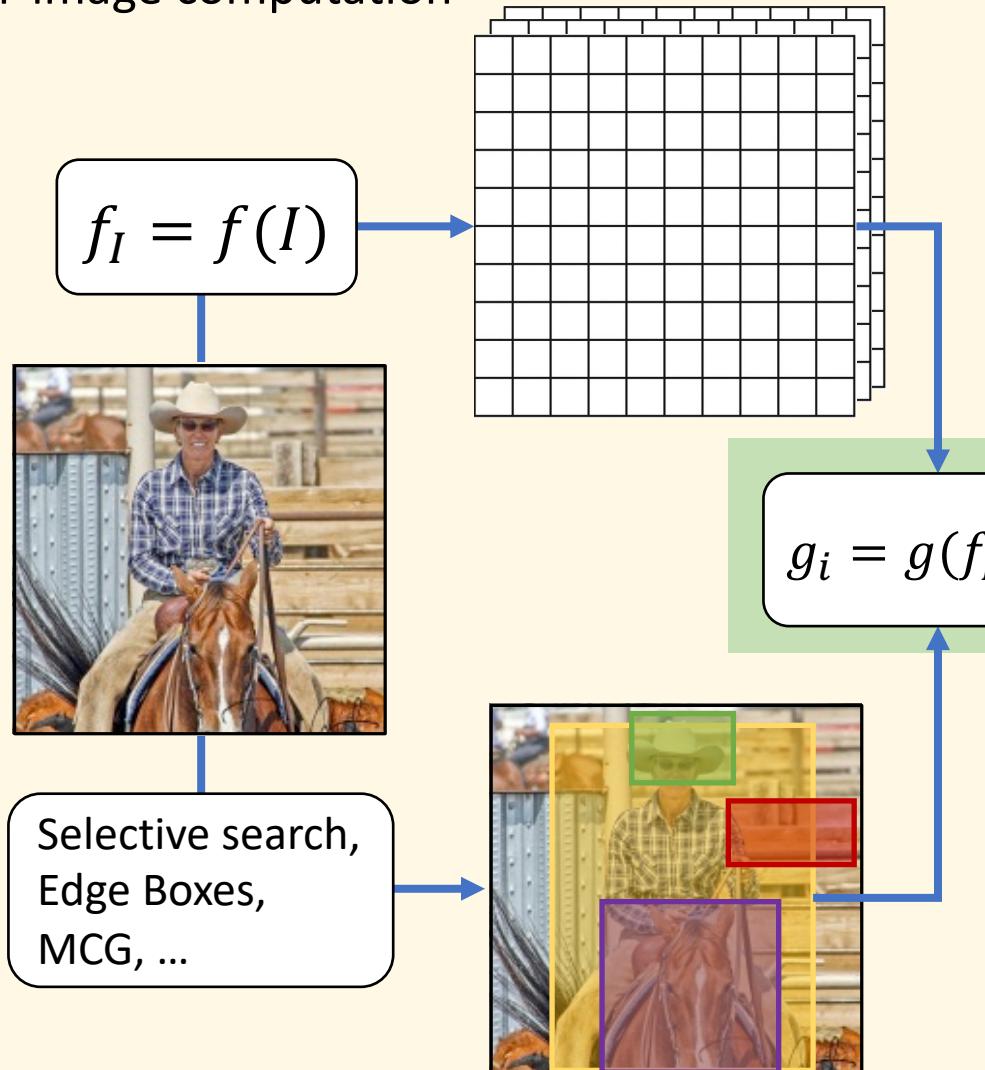
Fast R-CNN



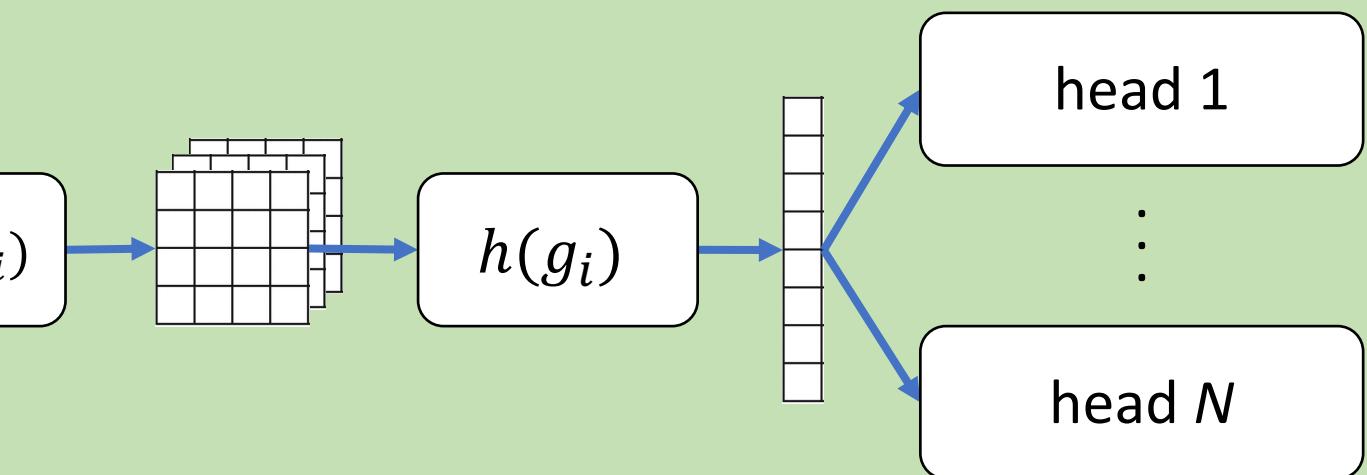
Building
a better
hammer

Generalized R-CNN → Fast R-CNN

Per-image computation



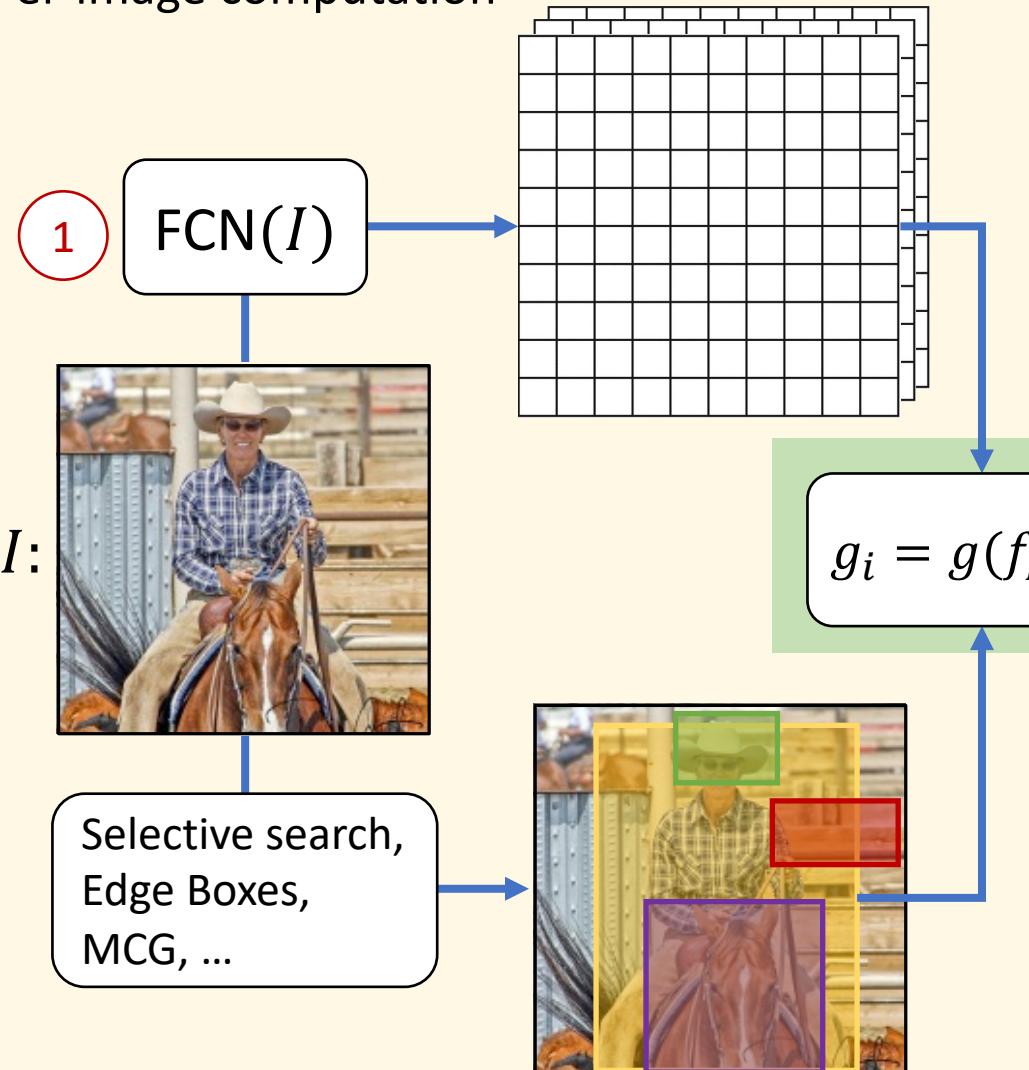
Per-region computation for each $r_i \in r(I)$



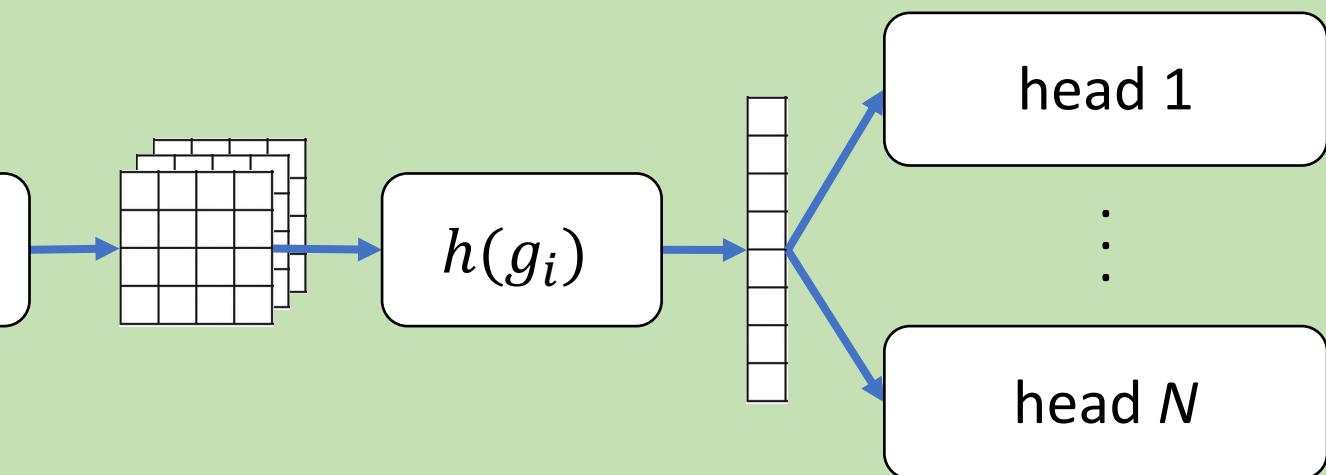
Light-weight per-region computation
End-to-end trainable version of SPP-net [He et al. 2014]

Generalized R-CNN → Fast R-CNN

Per-image computation



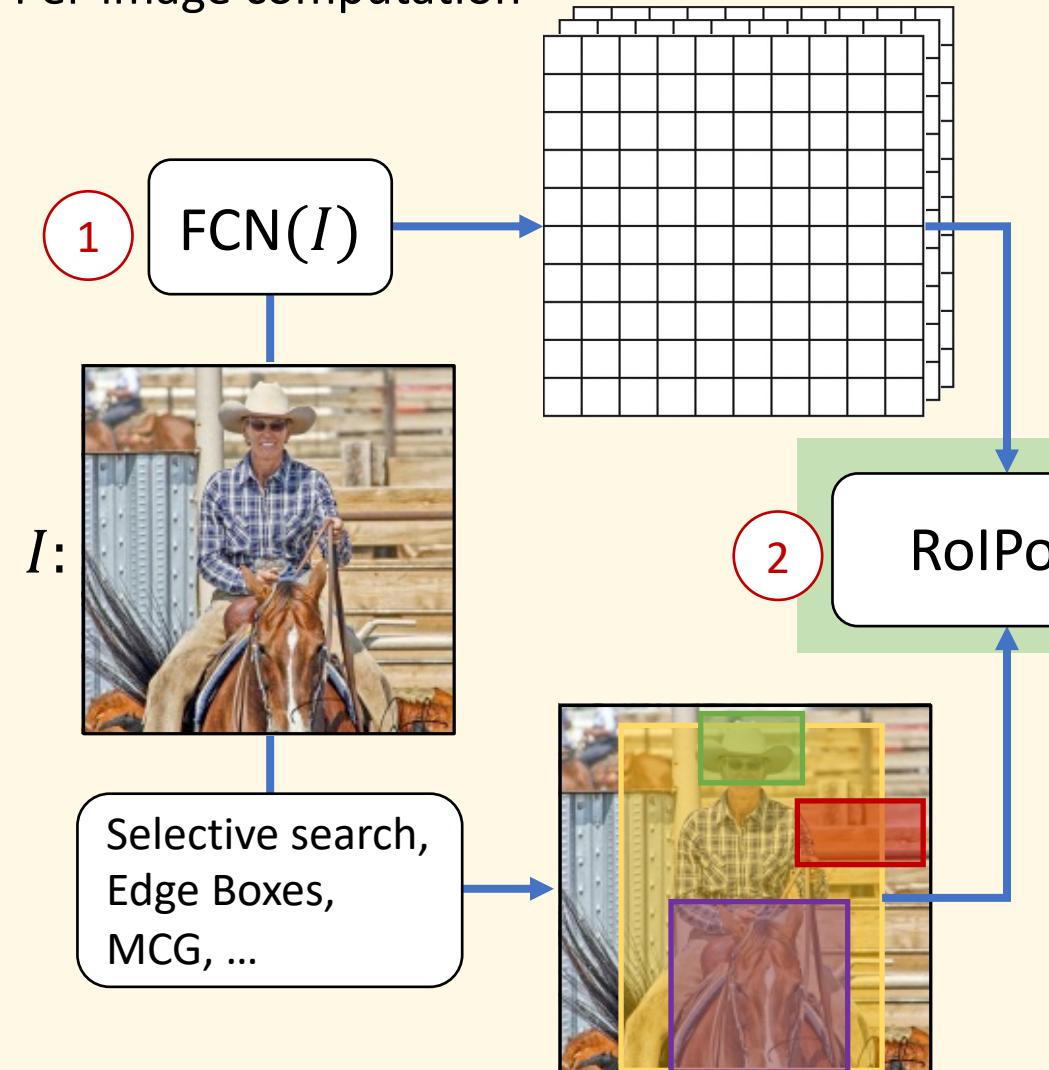
Per-region computation for each $r_i \in r(I)$



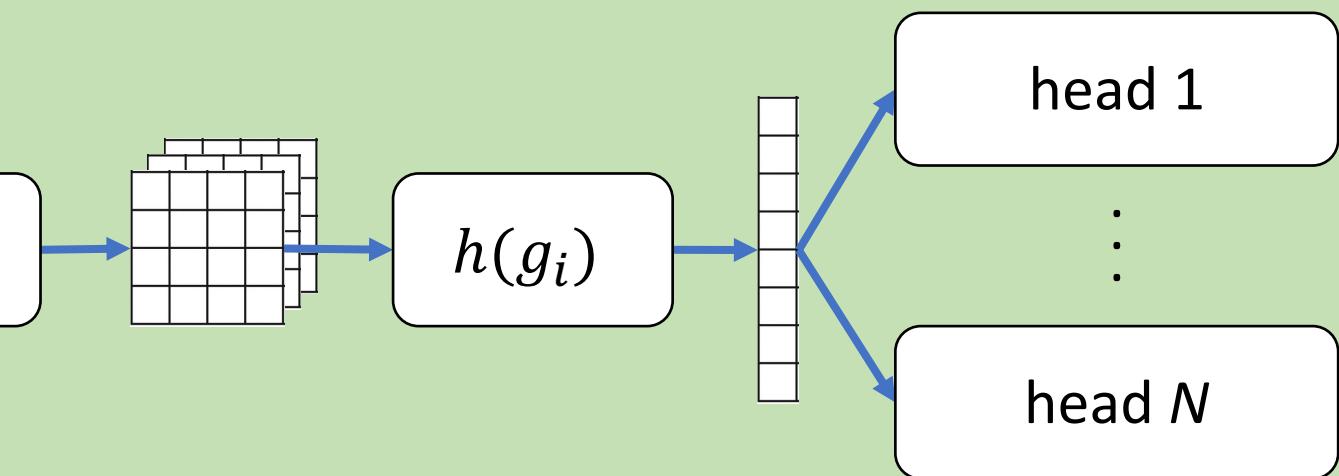
1. Fully convolutional network (FCN) maps the image to a lower resolution spatial feature map

Generalized R-CNN → Fast R-CNN

Per-image computation



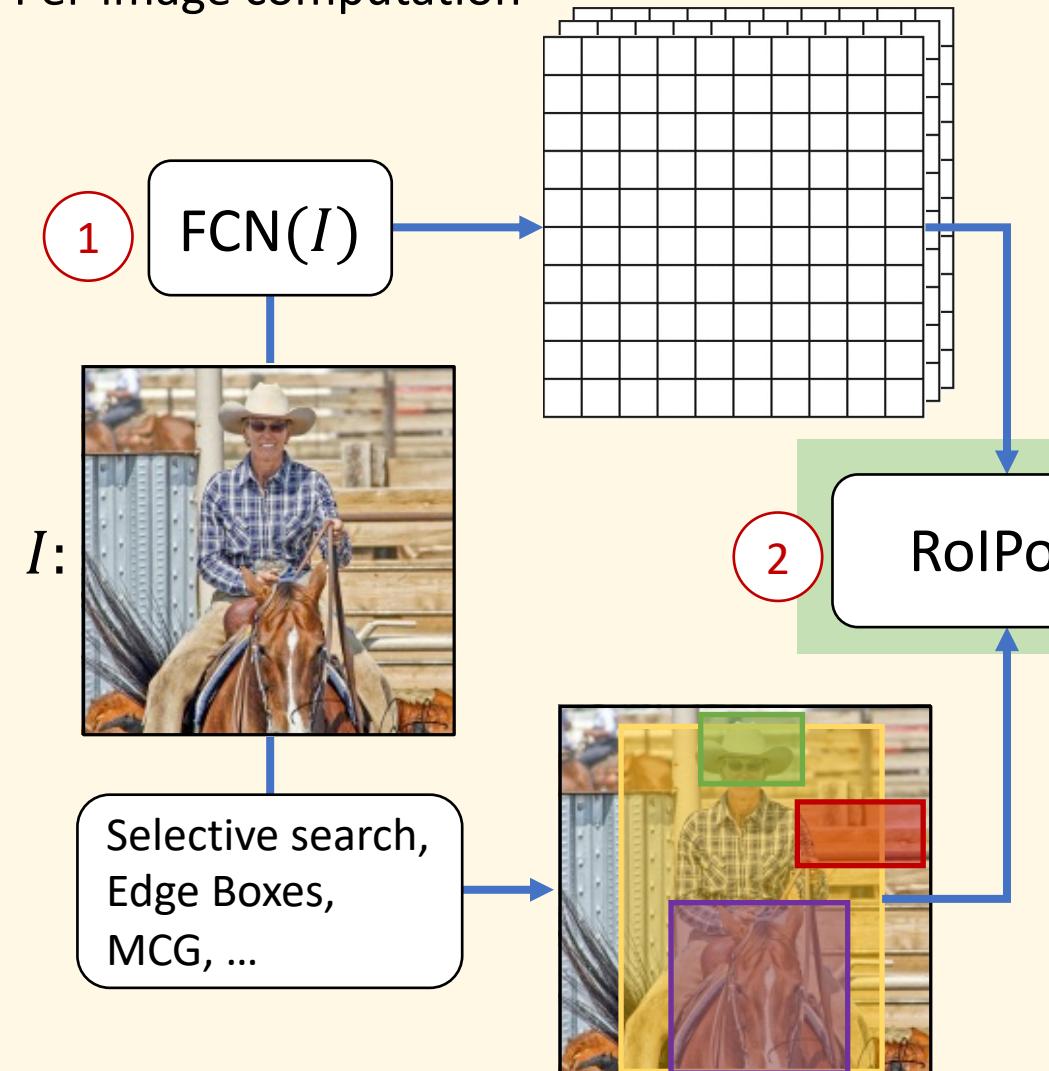
Per-region computation for each $r_i \in r(I)$



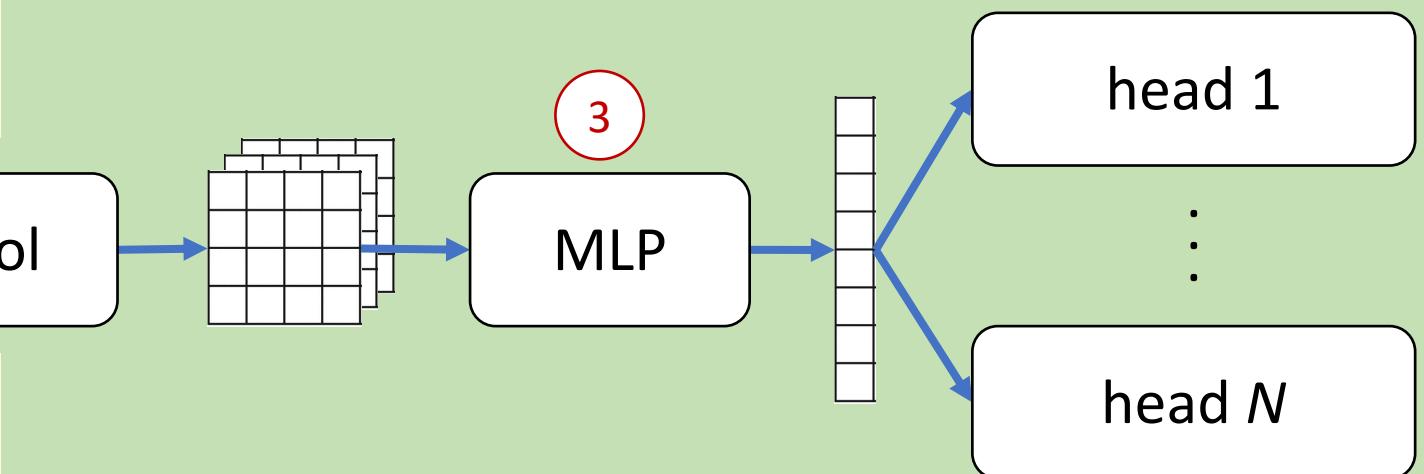
2. *Region of interest (RoI) pooling converts each region into a fixed dimensional representation*

Generalized R-CNN → Fast R-CNN

Per-image computation



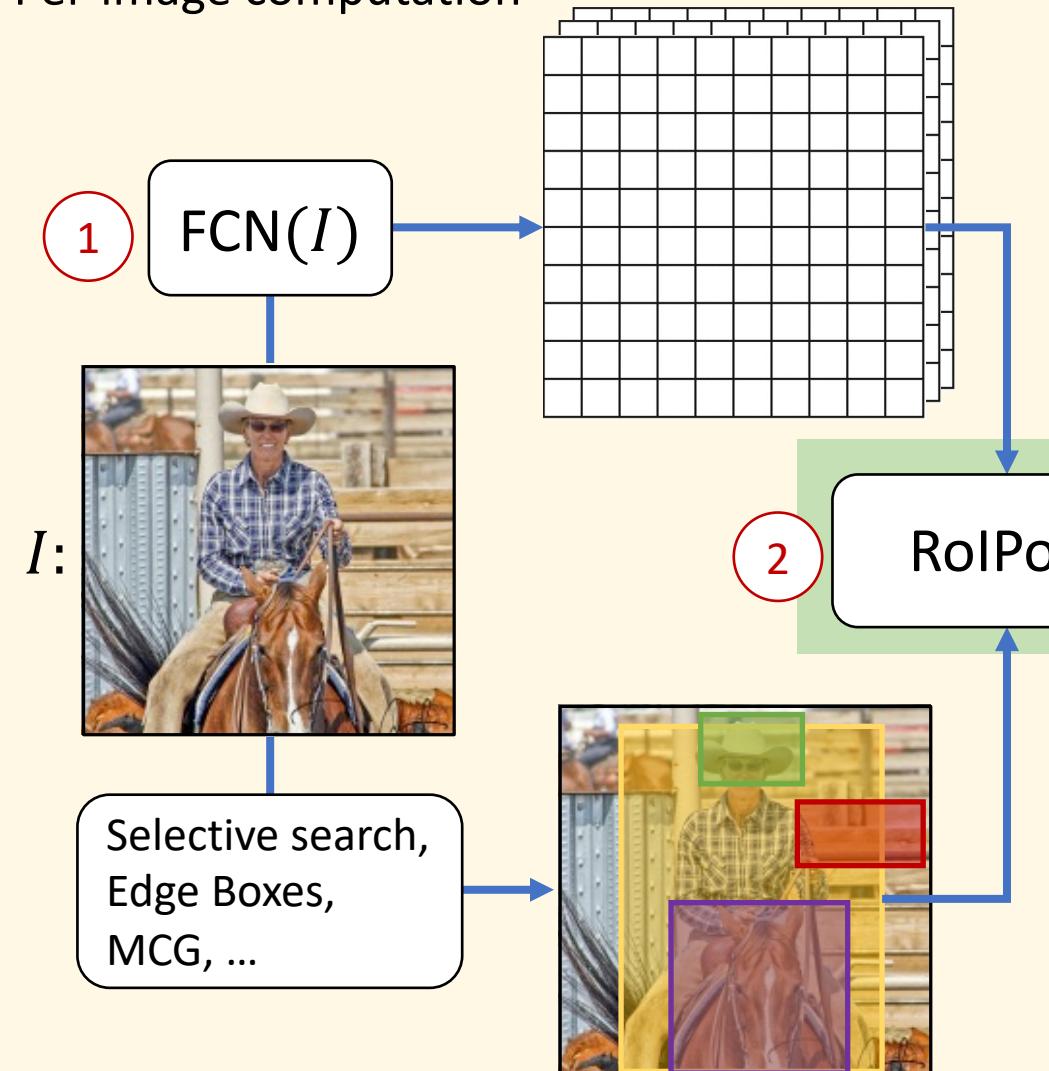
Per-region computation for each $r_i \in r(I)$



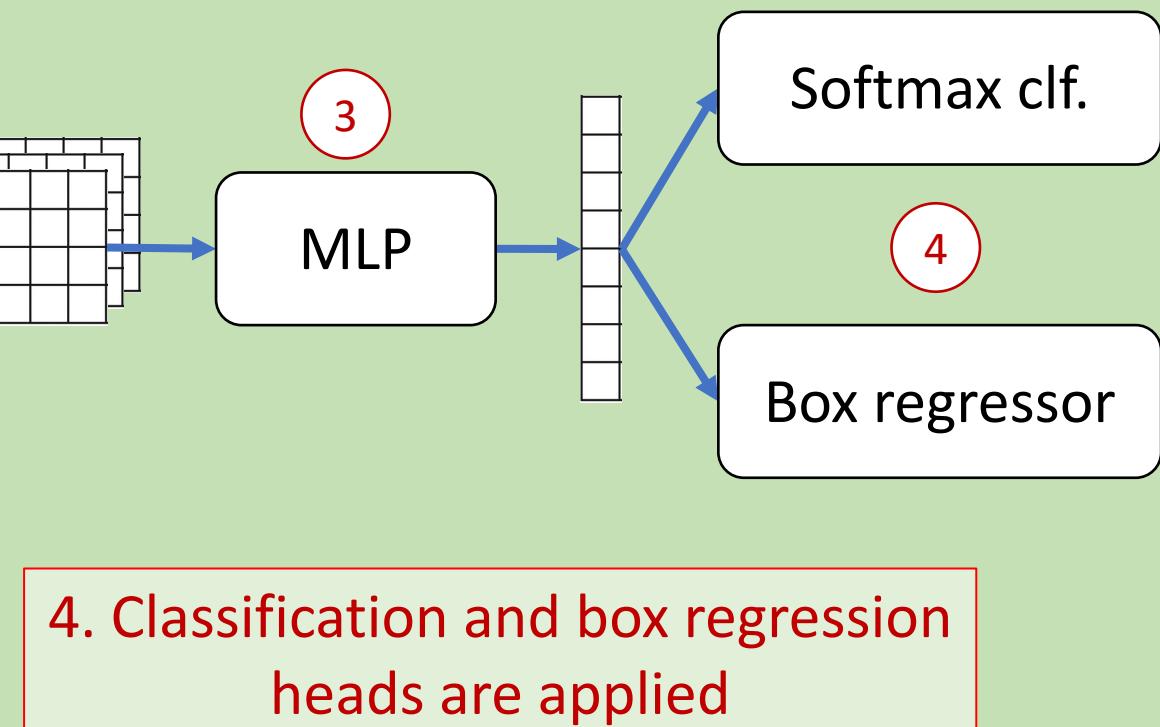
3. A lightweight MLP processes each region feature map
(Alternatively, the MLP could be a small ConvNet, etc.)

Generalized R-CNN → Fast R-CNN

Per-image computation

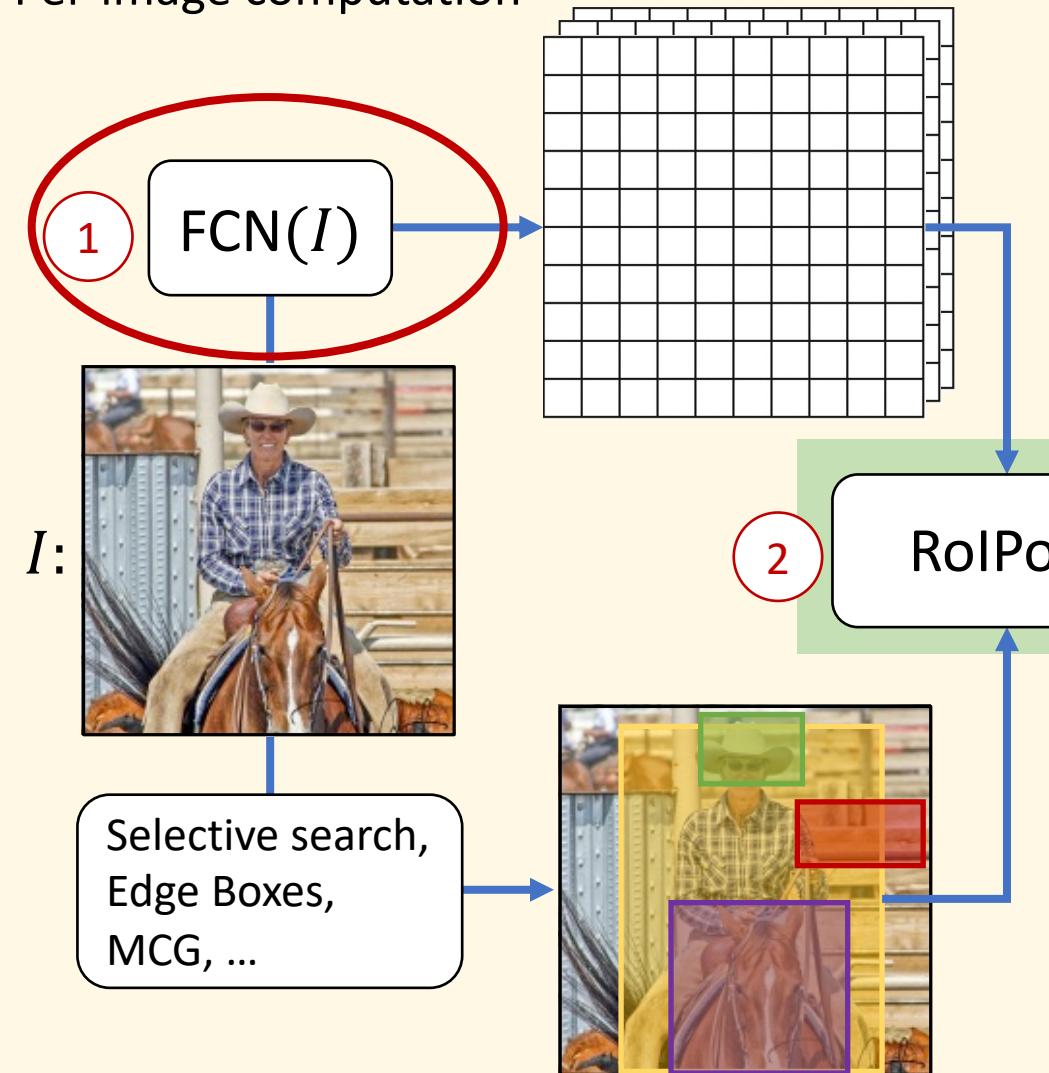


Per-region computation for each $r_i \in r(I)$

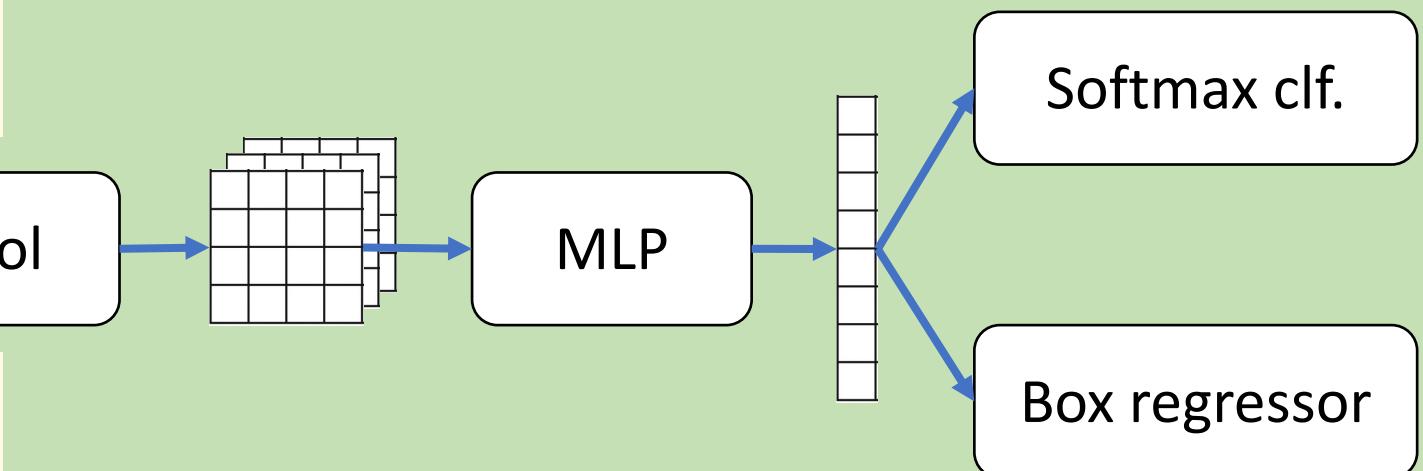


Fast R-CNN

Per-image computation



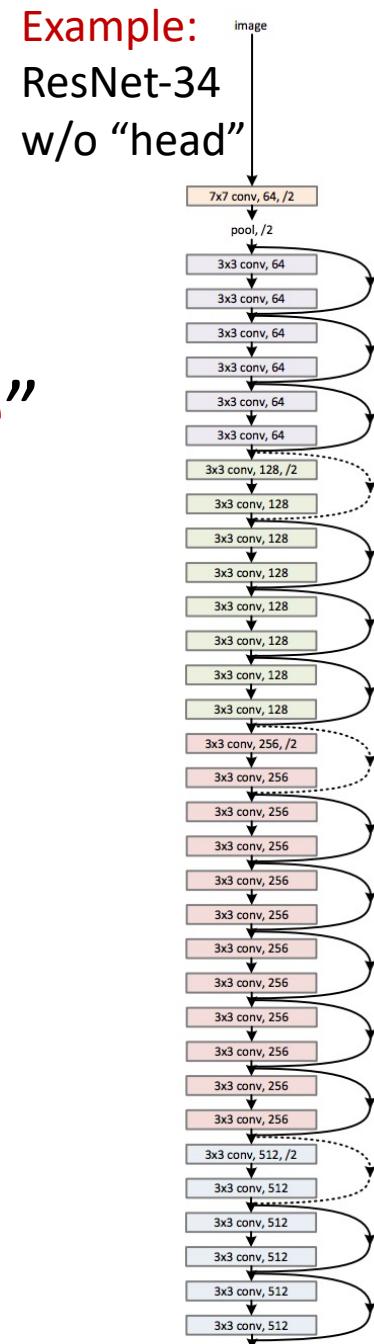
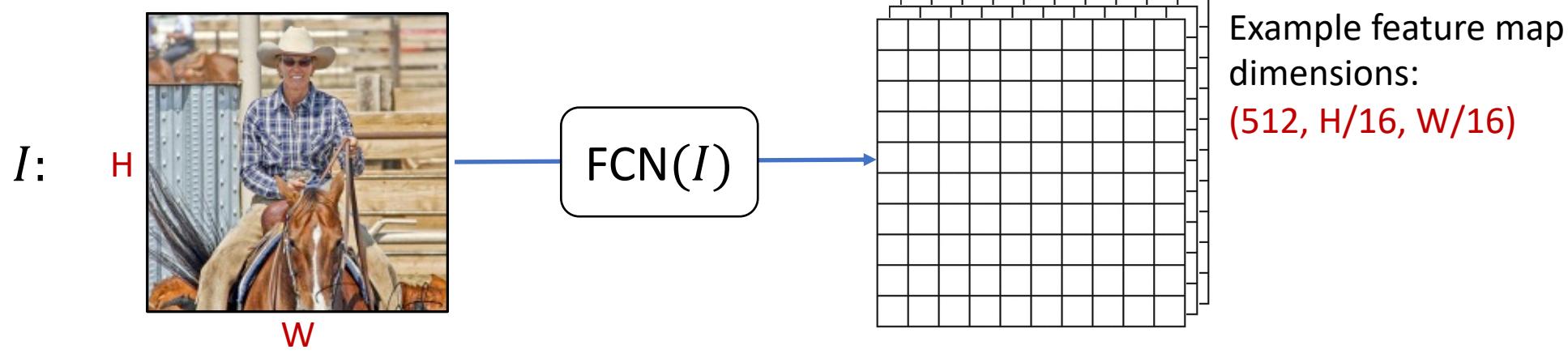
Per-region computation for each $r_i \in r(I)$



Whole-image, Fully Conv. Network (FCN)

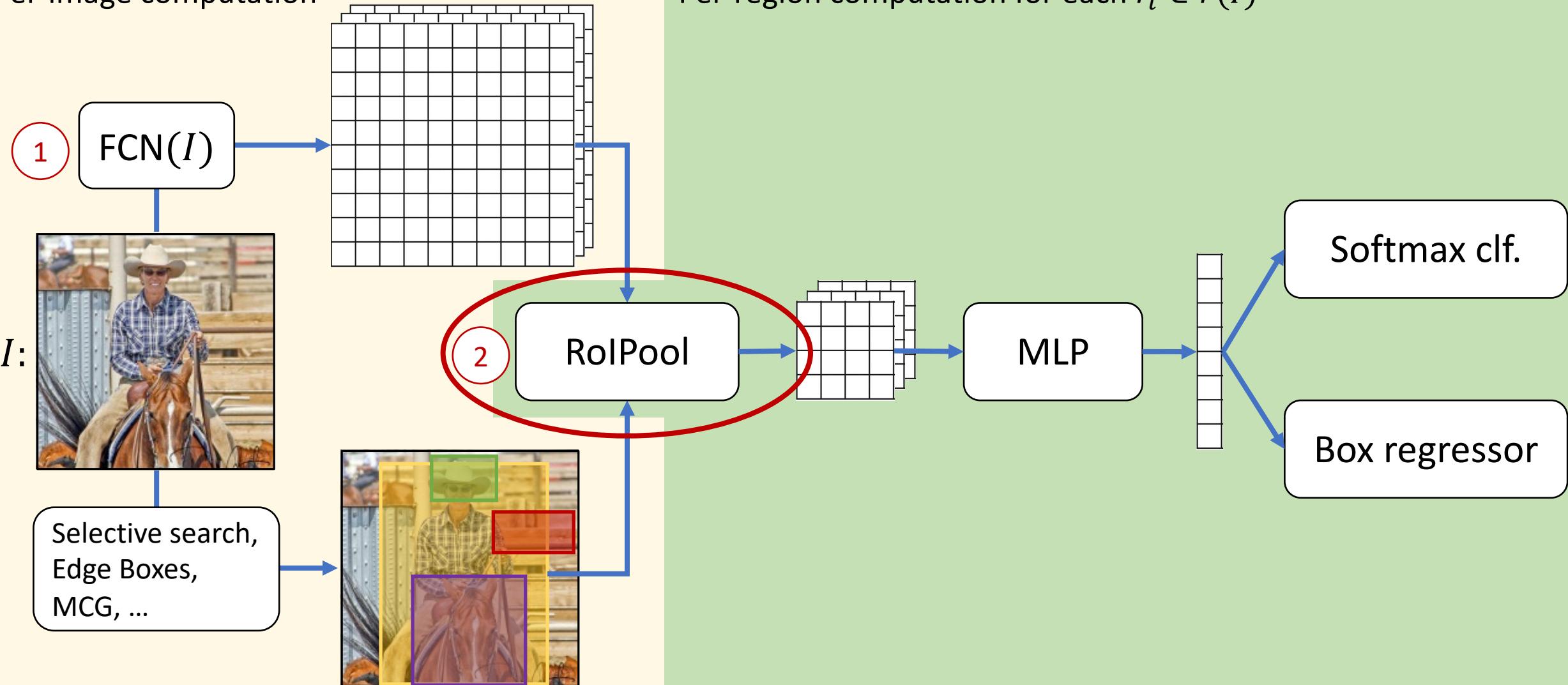
Use any standard ConvNet as the “backbone architecture”

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

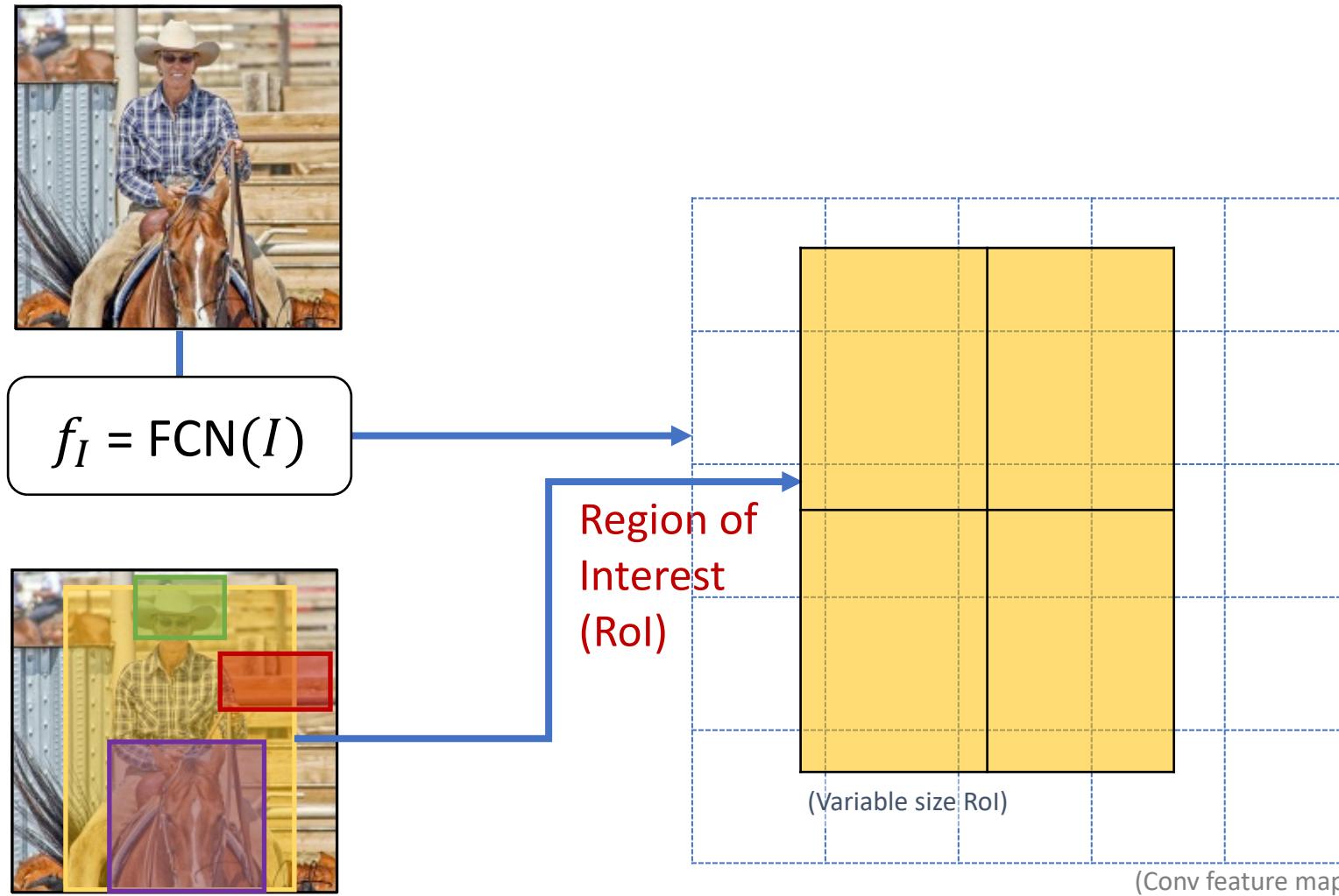


Fast R-CNN

Per-image computation

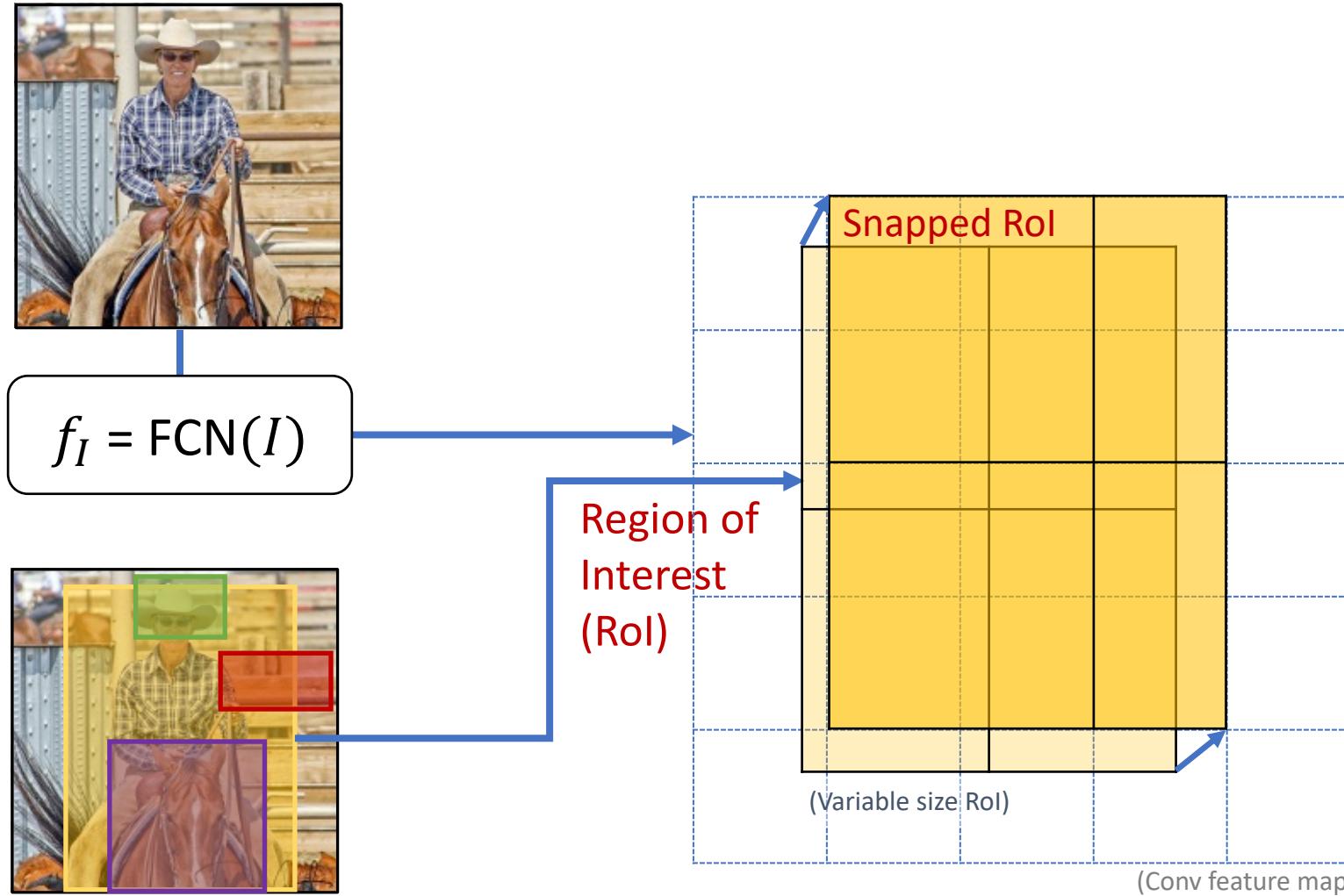


RoIPool Operation (on each Proposal)



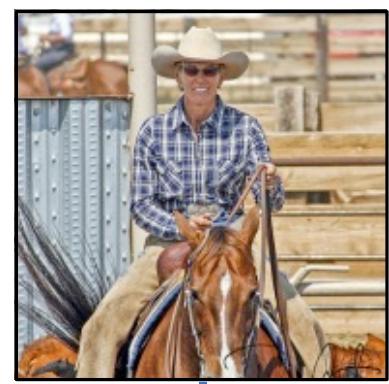
Key innovation in SPP-net
[He et al. 2014]

RoIPool Operation (on each Proposal)

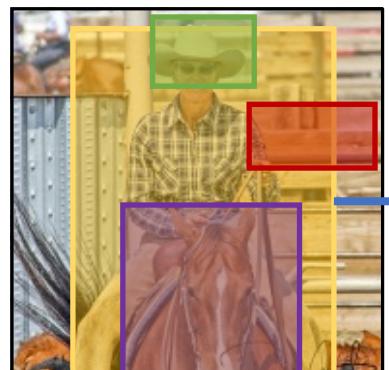


Key innovation in SPP-net
[He et al. 2014]

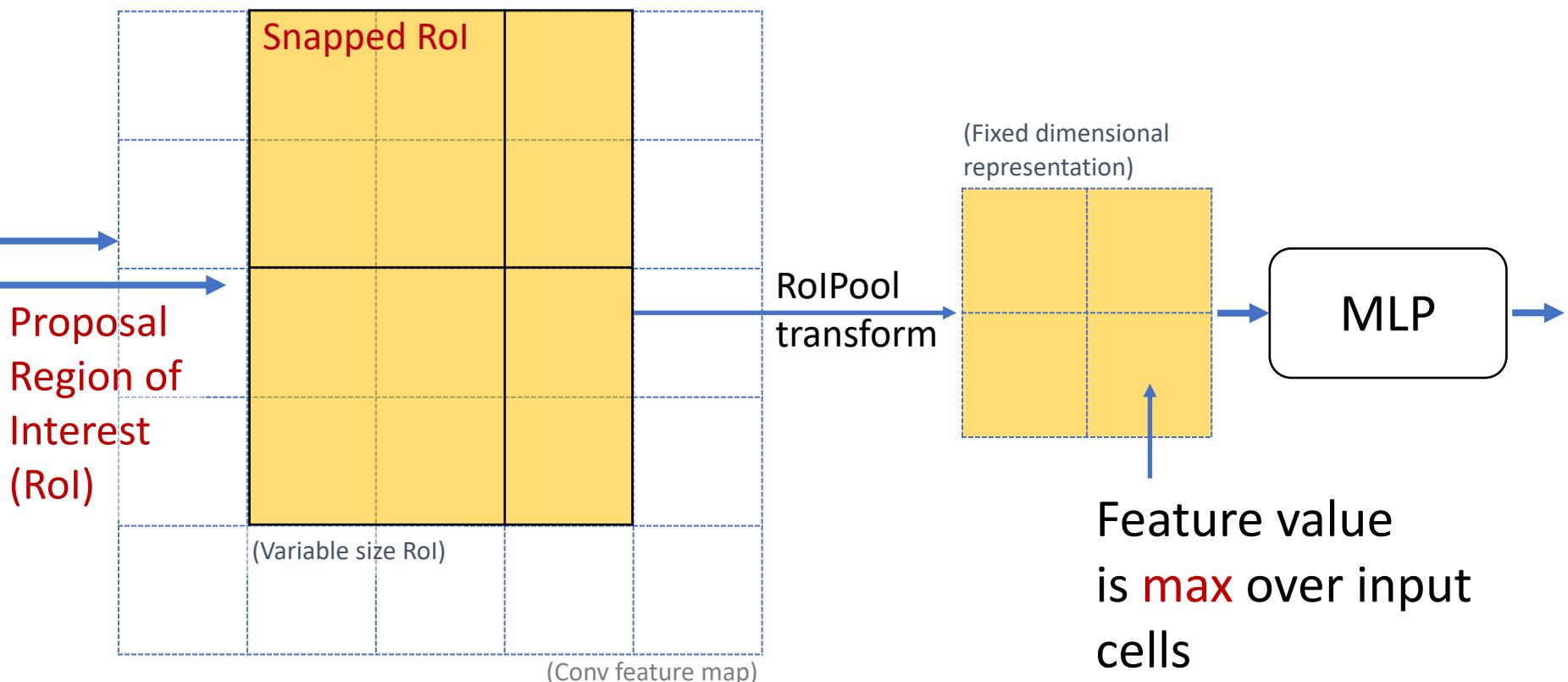
RoIPool Operation (on each Proposal)



$$f_I = \text{FCN}(I)$$



Transform **arbitrary size proposal** into a **fixed-dimensional representation** (e.g., 2x2)



Generalized R-CNN → Fast R-CNN

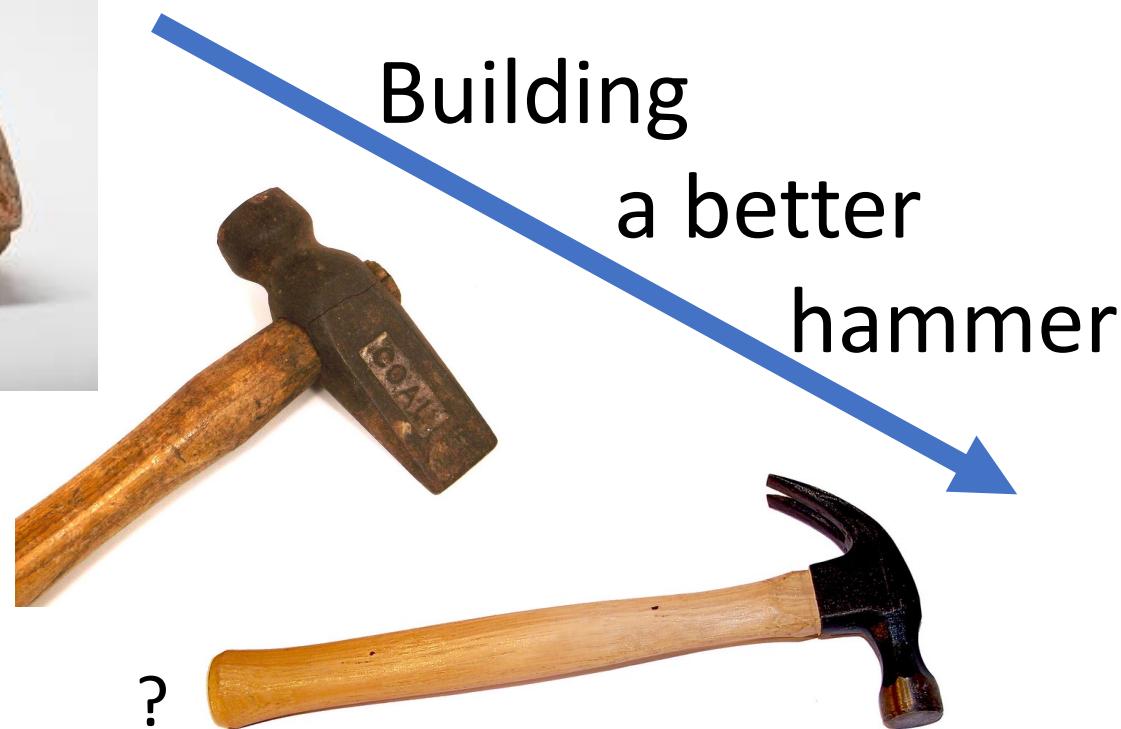
	Fast R-CNN			R-CNN			SPPnet
	S	M	L	S	M	L	$\dagger L$
train time (h)	1.2	2.0	9.5	22	28	84	25
train speedup	18.3×	14.0×	8.8×	1×	1×	1×	3.4×
test rate (s/im)	0.10	0.15	0.32	9.8	12.1	47.0	2.3
▷ with SVD	0.06	0.08	0.22	-	-	-	-
test speedup	98×	80×	146×	1×	1×	1×	20×
▷ with SVD	169×	150×	213×	-	-	-	-
VOC07 mAP	57.1	59.2	66.9	58.5	60.2	66.0	63.1
▷ with SVD	56.5	58.7	66.6	-	-	-	-

The Problem with Fast R-CNN

R-CNN

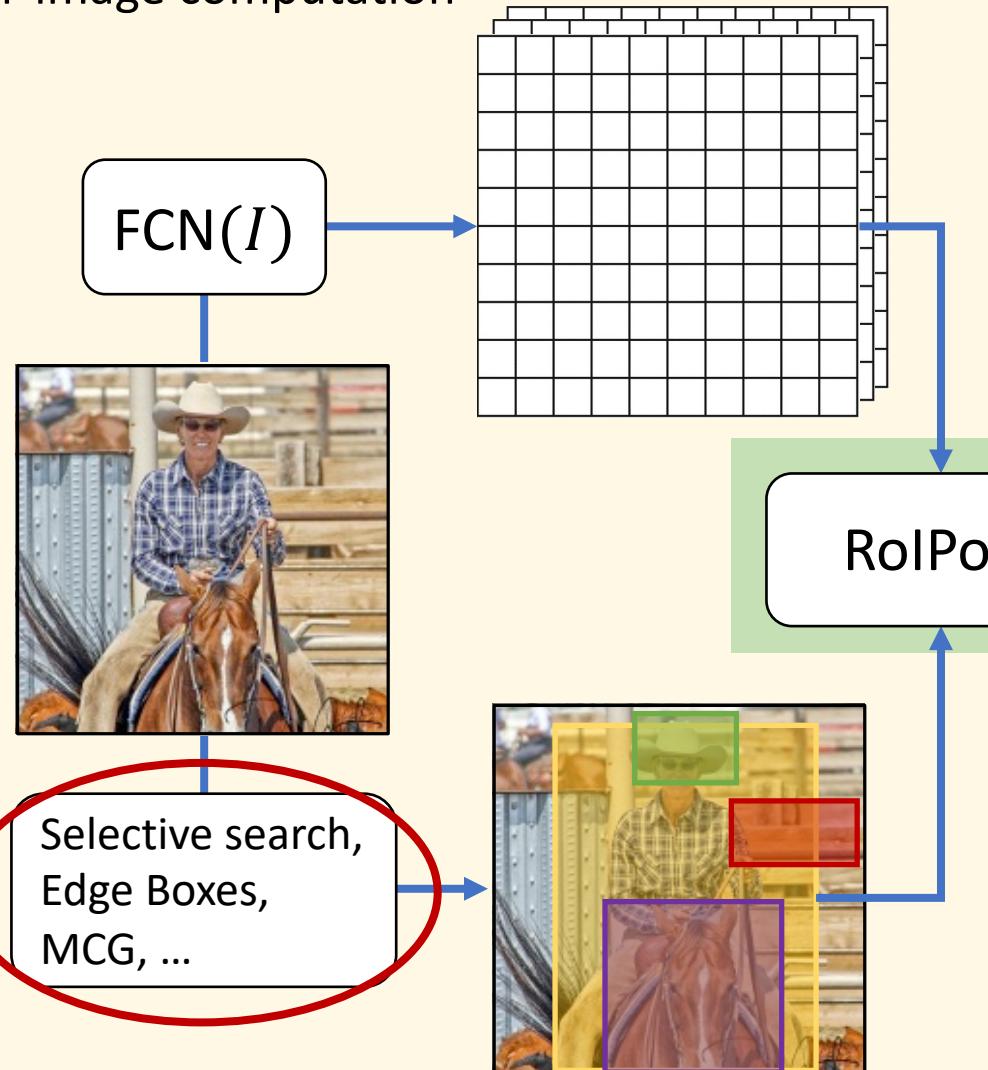


Fast R-CNN

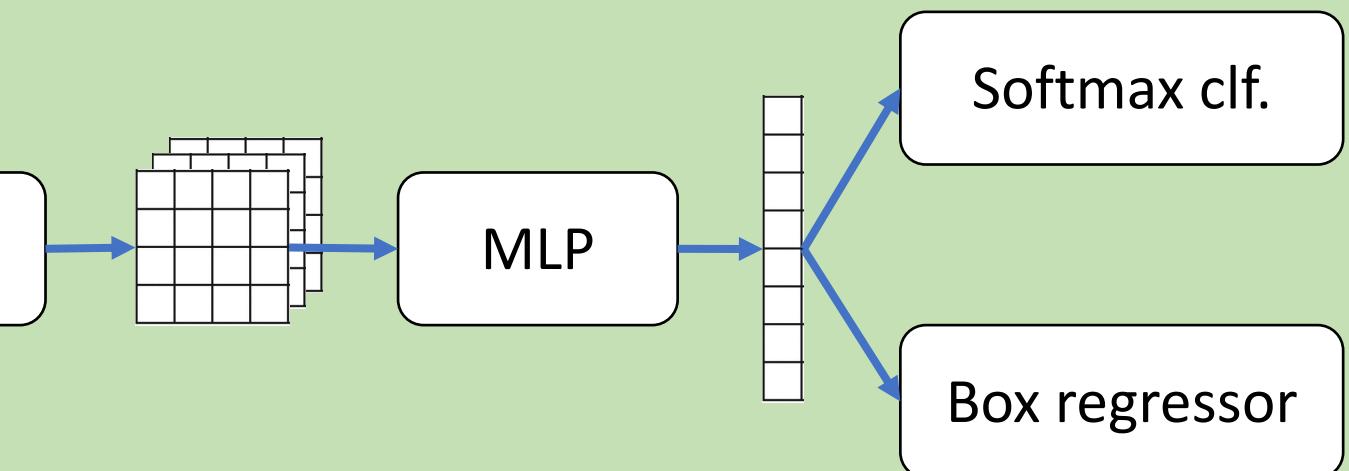


The Problem with Fast R-CNN

Per-image computation



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



Region proposals have very poor recall
(ok for PASCAL VOC, major bottleneck for COCO)
Also, they can be slow

Faster R-CNN

References

- D. Erhan, C. Szegedy, A. Toshev, and D. Anguelov. Scalable object detection using deep neural networks. In CVPR, 2014.
- P. O. Pinheiro, R. Collobert, and P. Dollar. Learning to segment object candidates. In NIPS, 2015.
- S. Ren, K. He, R. Girshick, and J. Sun. Faster R-CNN: Towards real-time object detection with region proposal networks. In NIPS, 2015.

R-CNN



Fast R-CNN



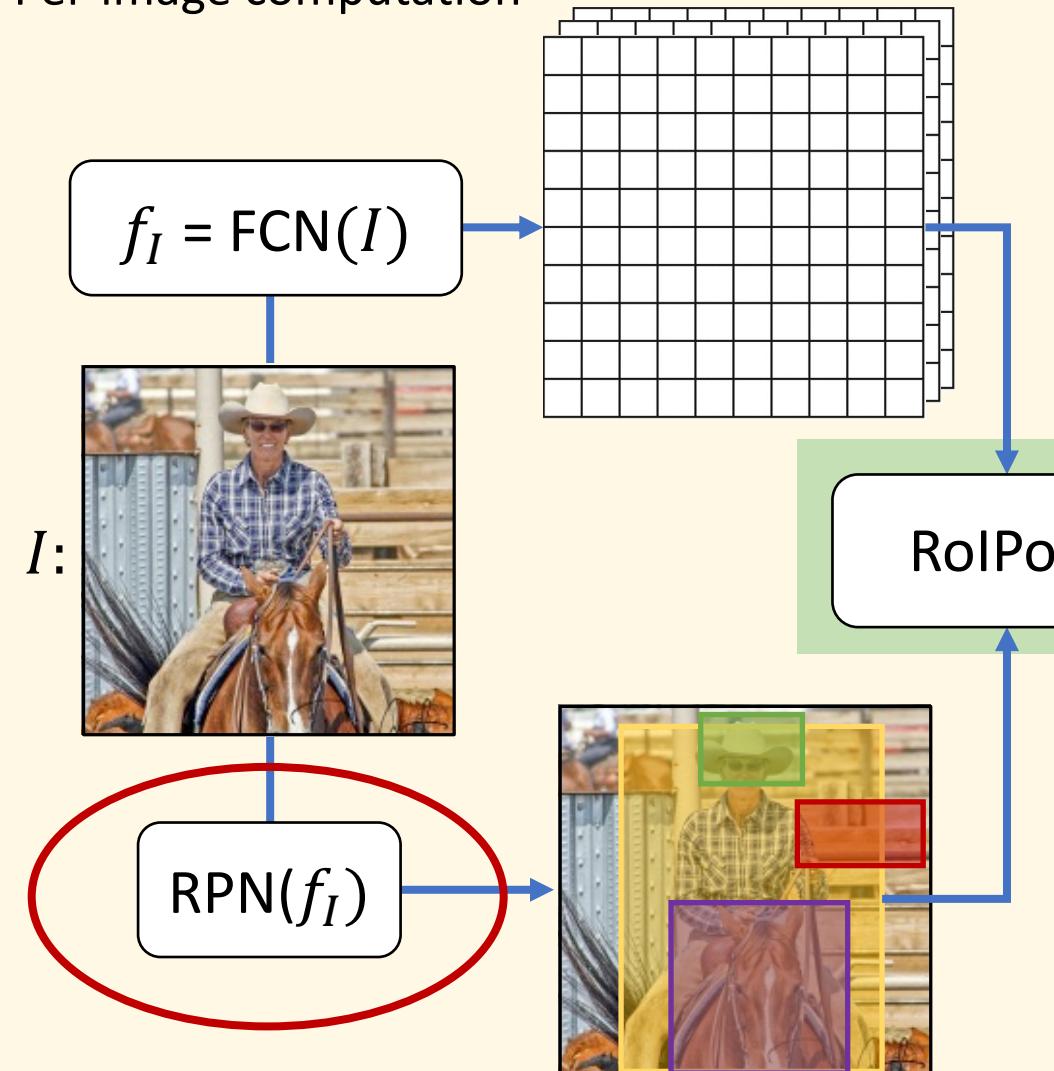
Faster R-CNN



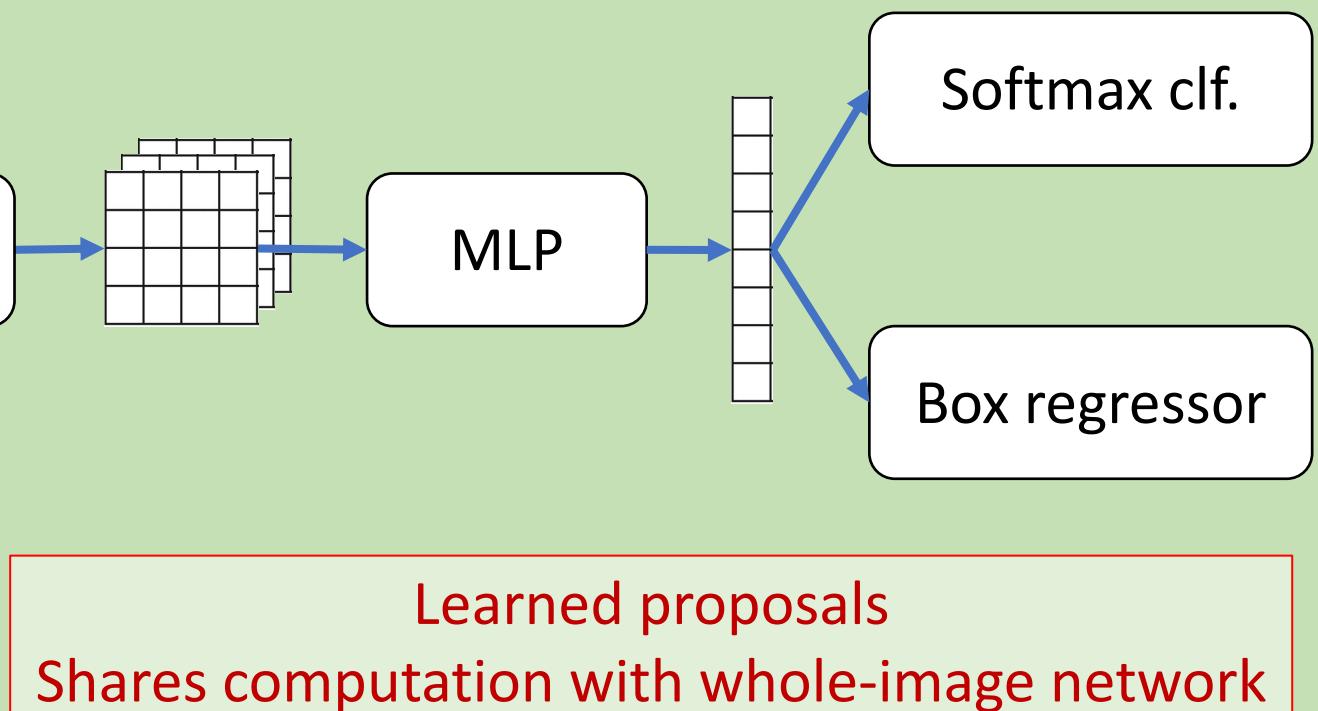
Building
a better
hammer

Faster R-CNN

Per-image computation

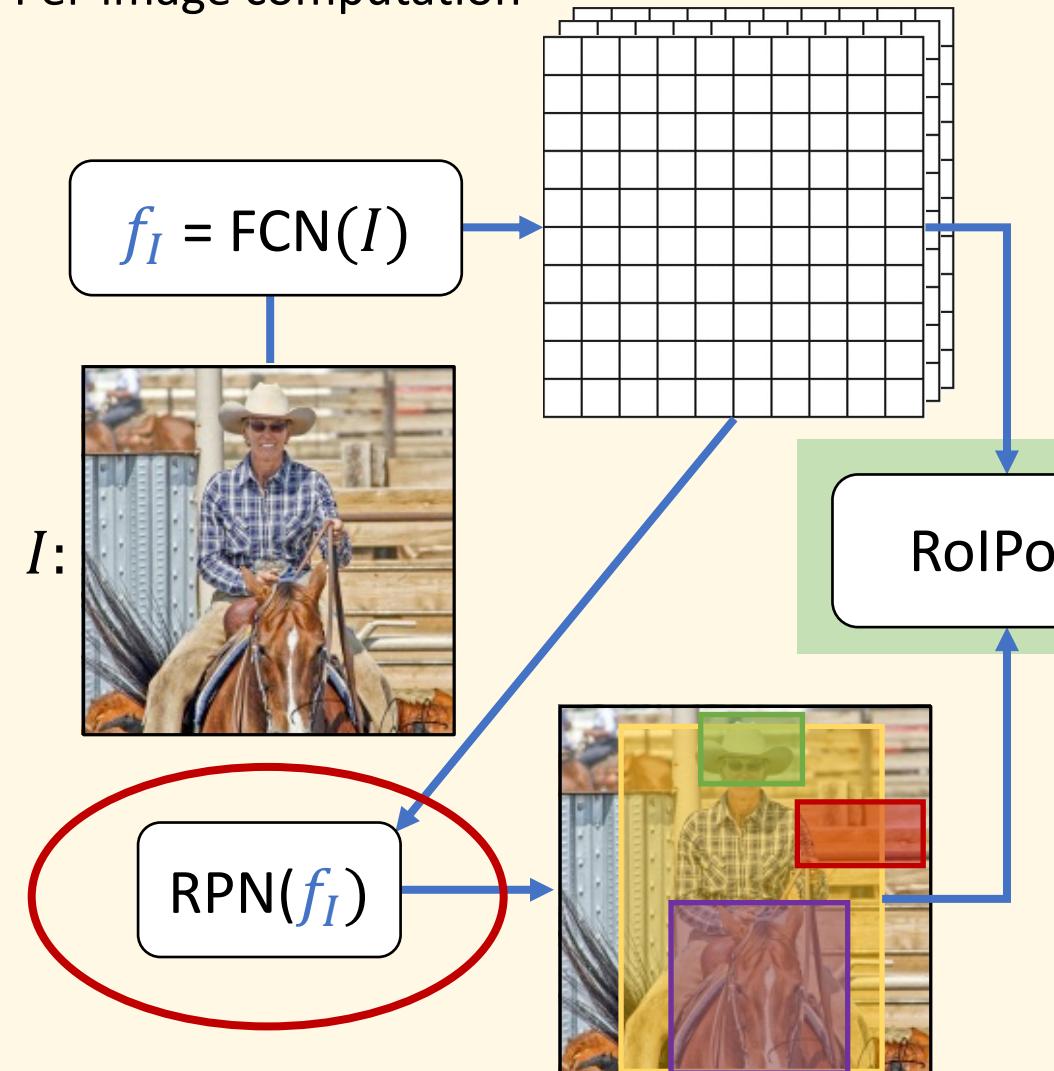


Per-region computation for each $r_i \in r(I)$

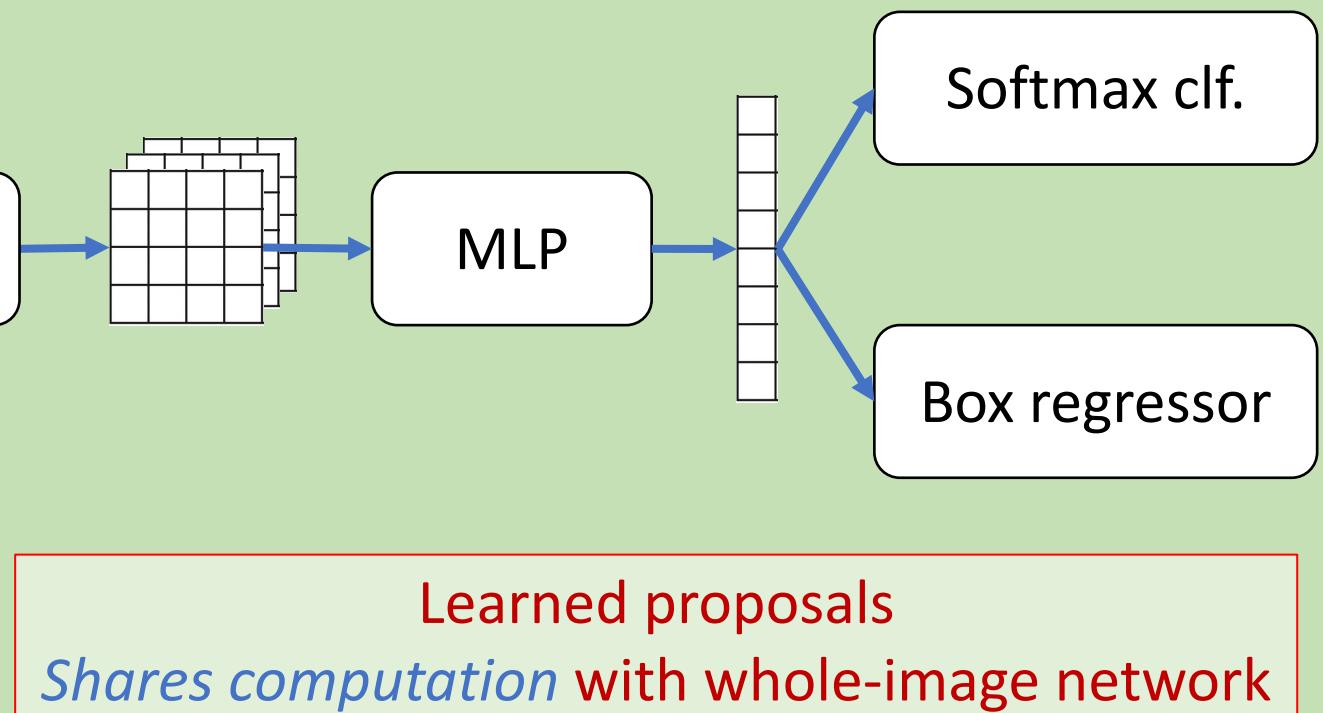


Faster R-CNN

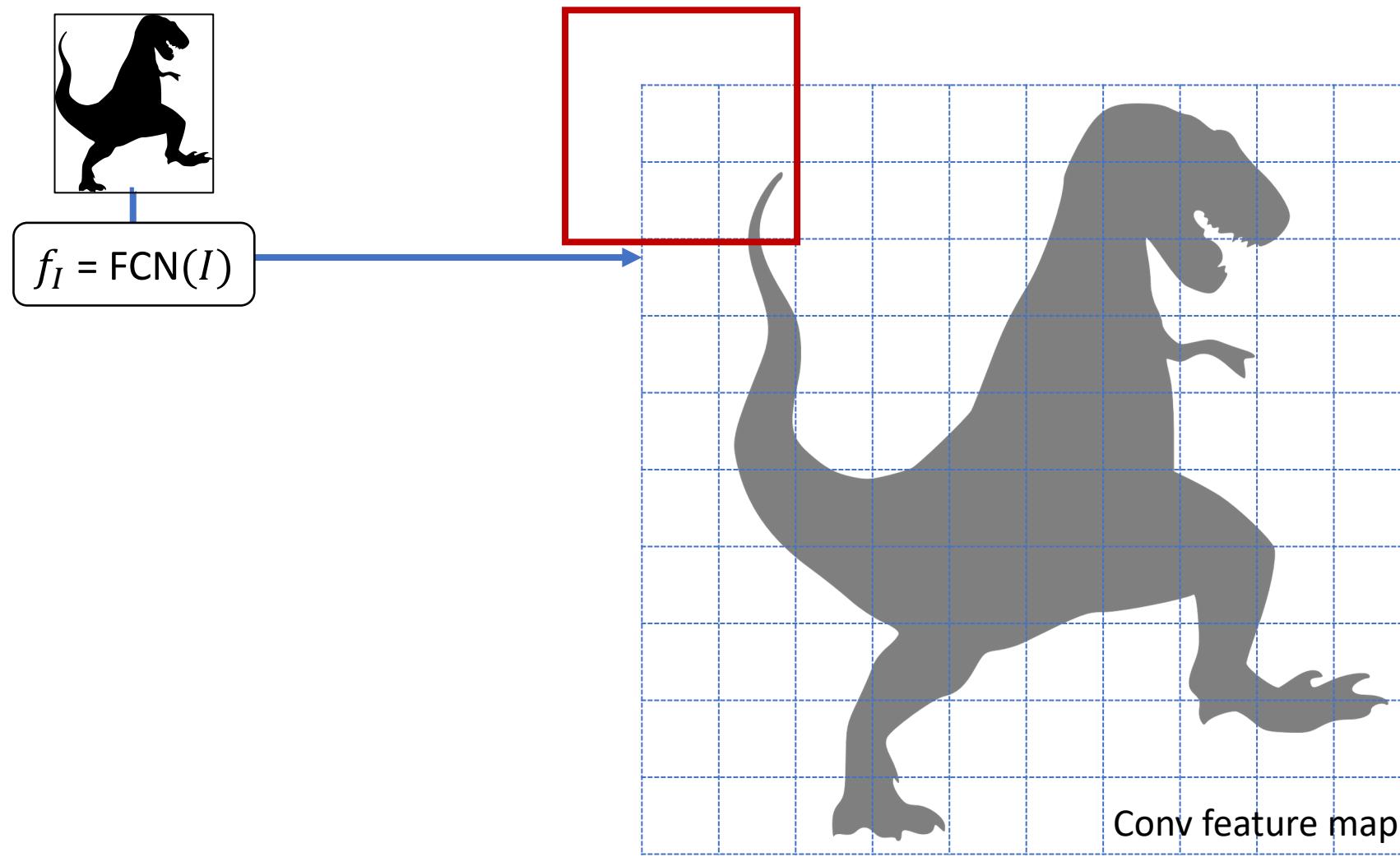
Per-image computation



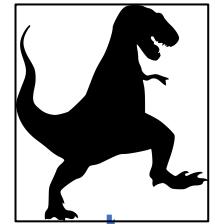
Per-region computation for each $r_i \in r(I)$



RPN: Region Proposal Network

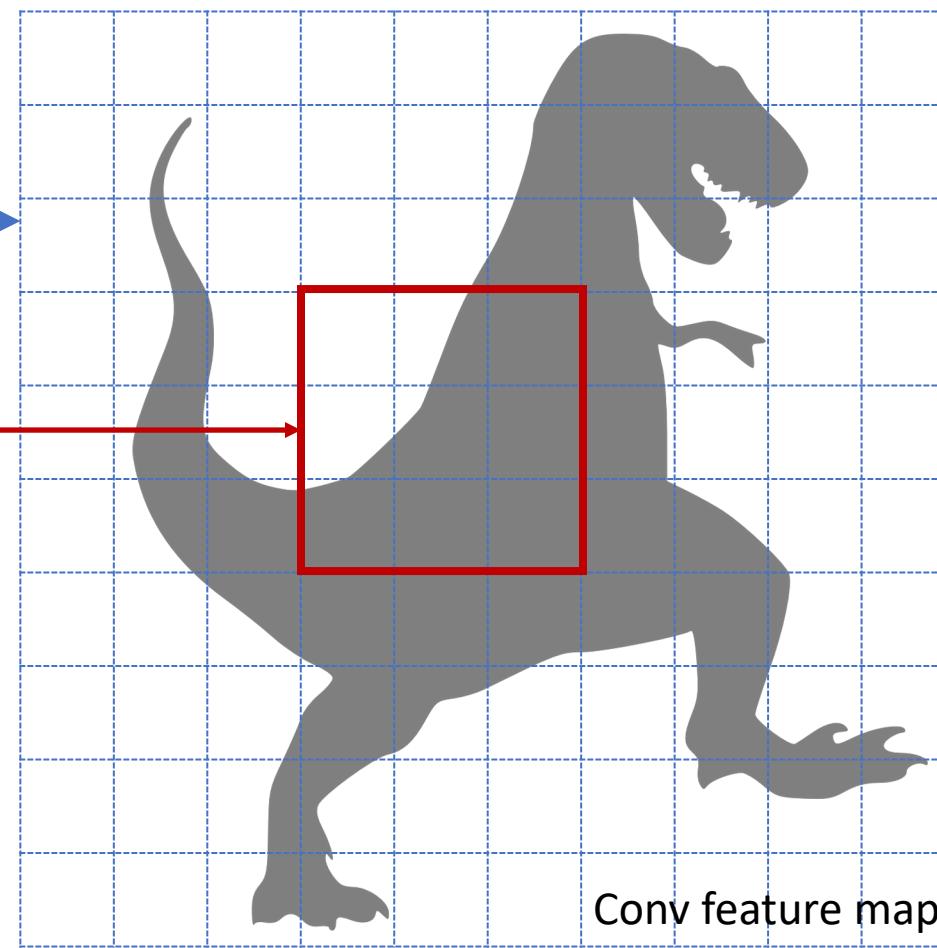


RPN: Region Proposal Network



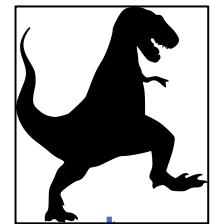
$$f_I = \text{FCN}(I)$$

3x3 “sliding window”
Scans the feature map
looking for objects



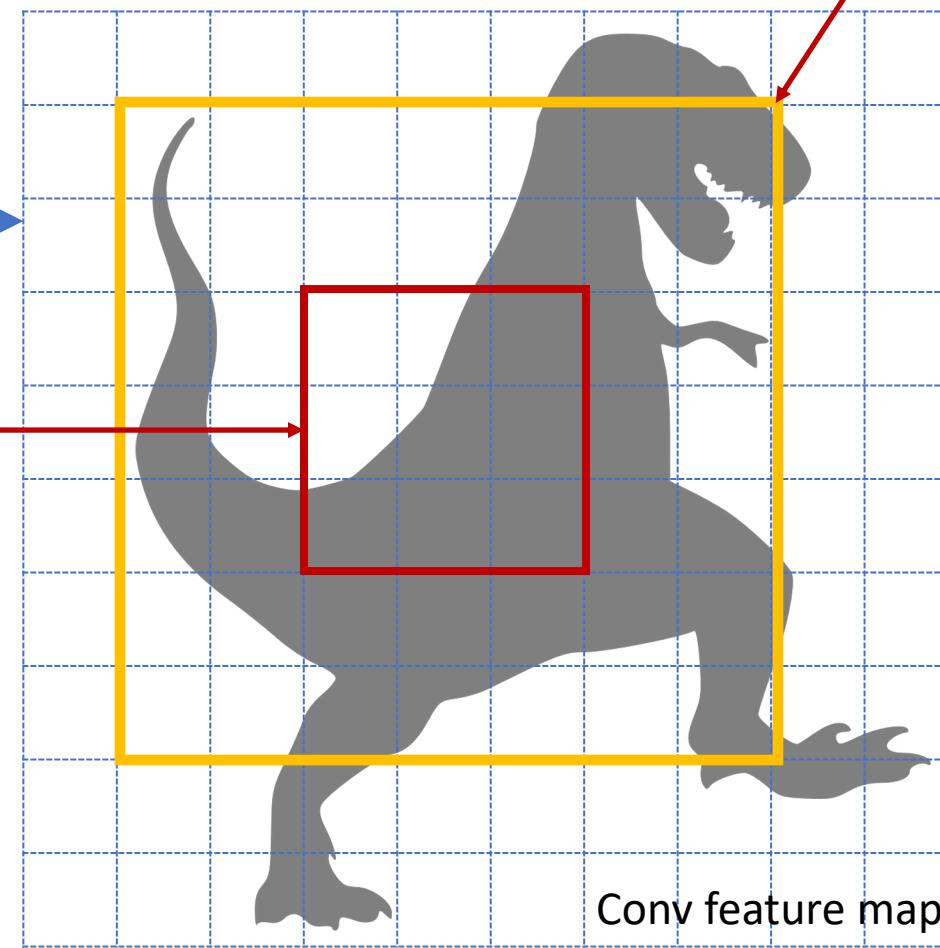
Conv feature map

RPN: Anchor Box



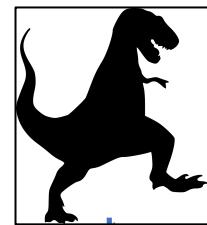
$$f_I = \text{FCN}(I)$$

3x3 “sliding window”
Scans the feature map
looking for objects

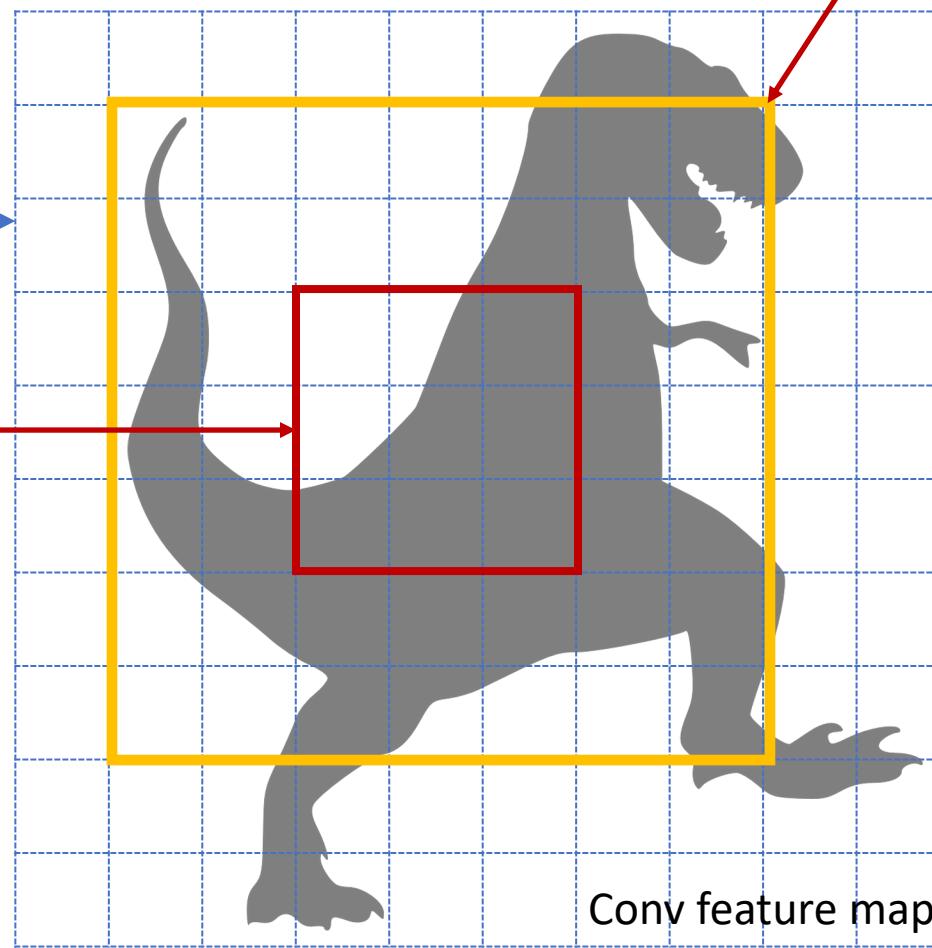


Anchor box: predictions are
w.r.t. this box, *not the 3x3
sliding window*

RPN: Anchor Box



$$f_I = \text{FCN}(I)$$



Anchor box: predictions are w.r.t. this box, *not the 3x3 sliding window*

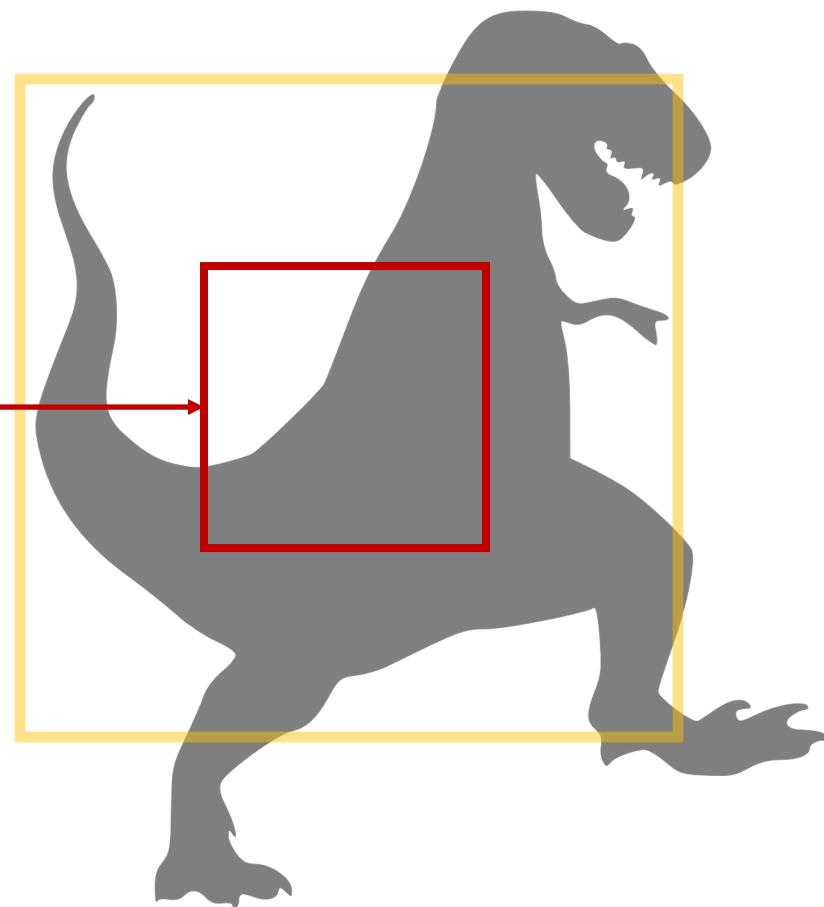
- 3x3 “sliding window”
- Objectness classifier [0, 1]
- Box regressor predicting (dx, dy, dh, dw)

RPN: Prediction (on object)

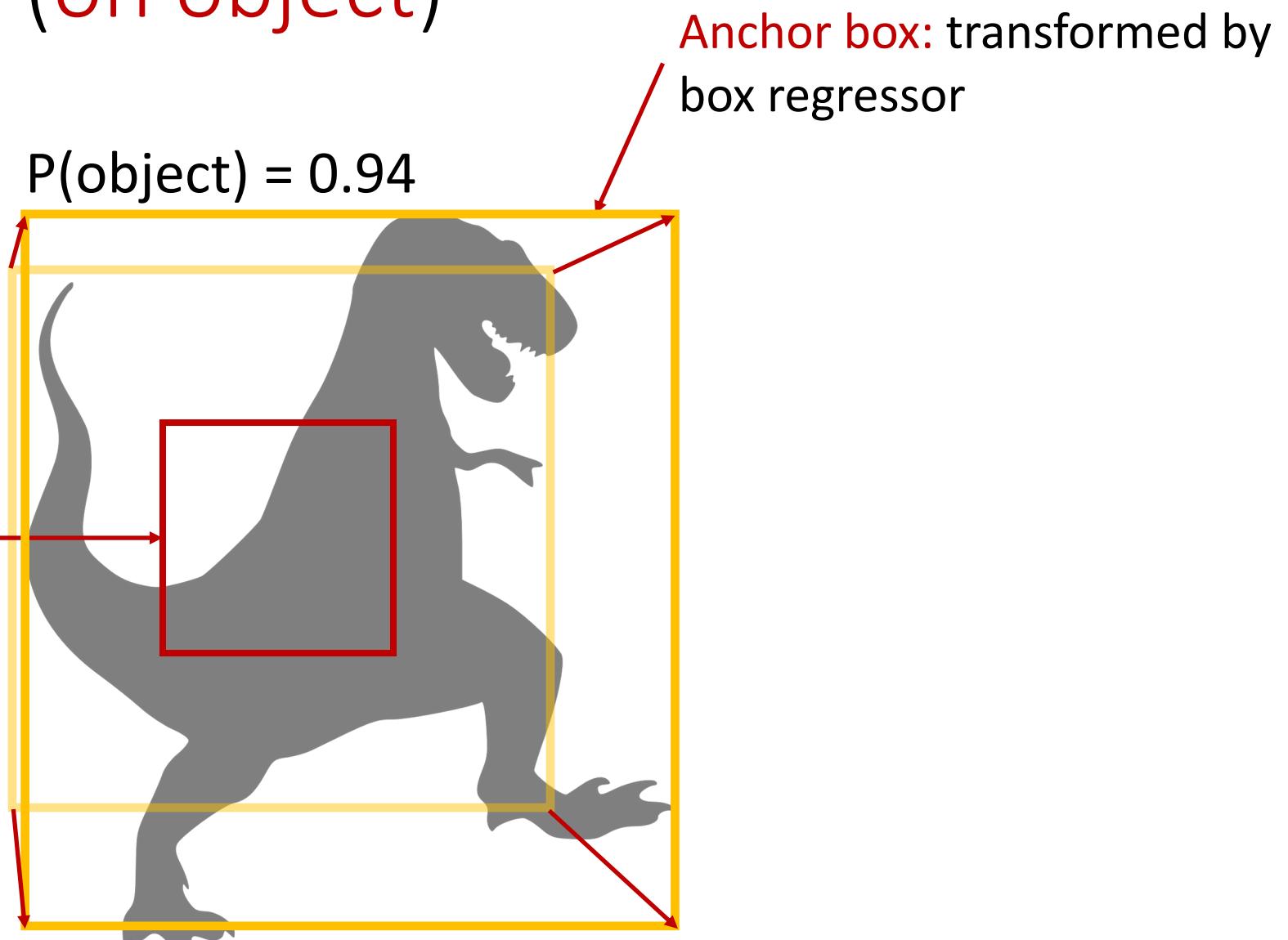
Objectness score

$$P(\text{object}) = 0.94$$

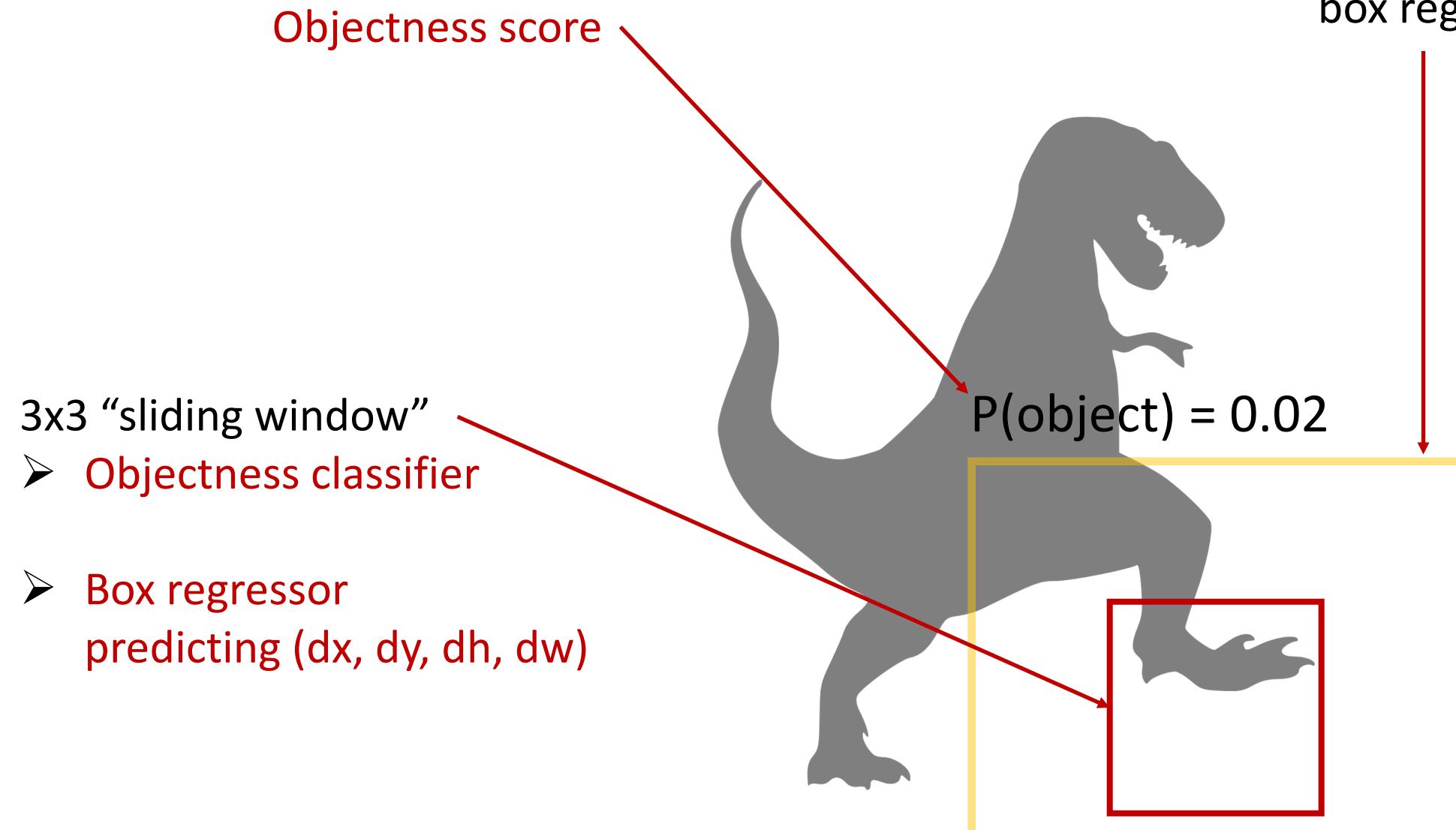
- 3x3 “sliding window”
- Objectness classifier [0, 1]
 - Box regressor
predicting (dx , dy , dh , dw)



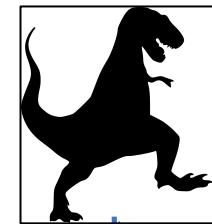
RPN: Prediction (on object)



RPN: Prediction (off object)

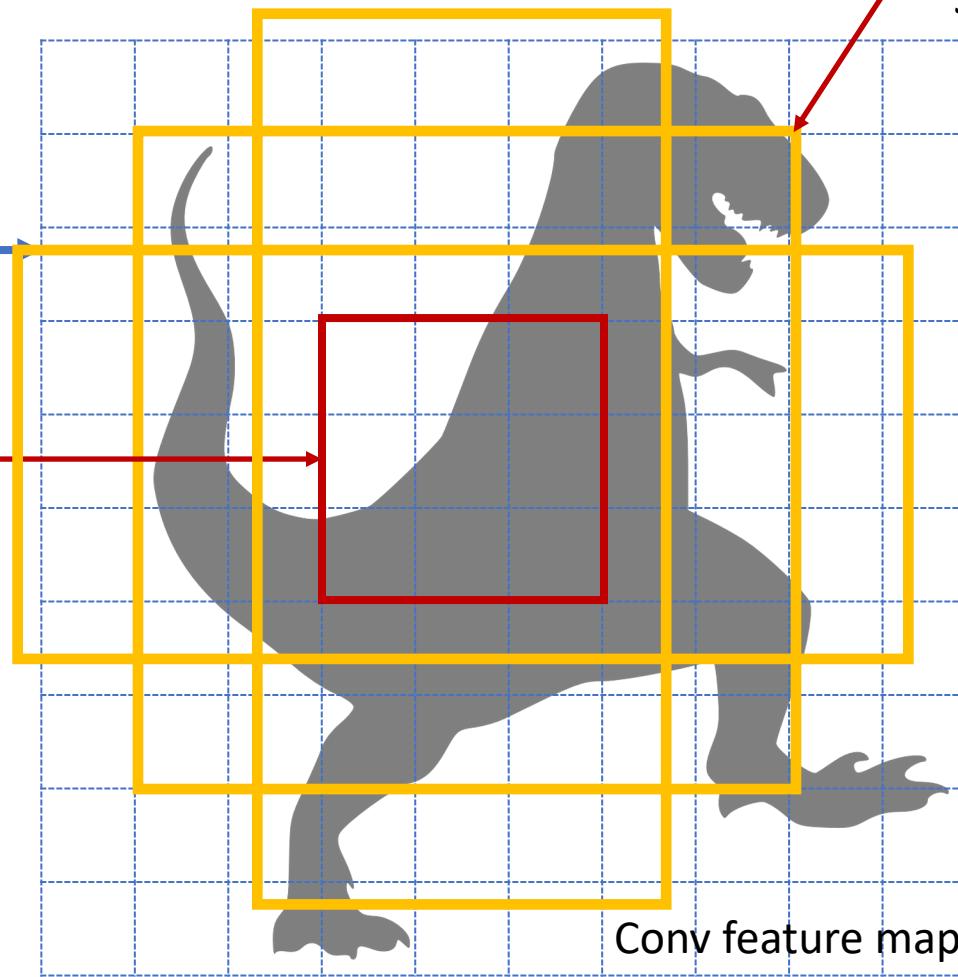


RPN: Multiple Anchors



$$f_I = \text{FCN}(I)$$

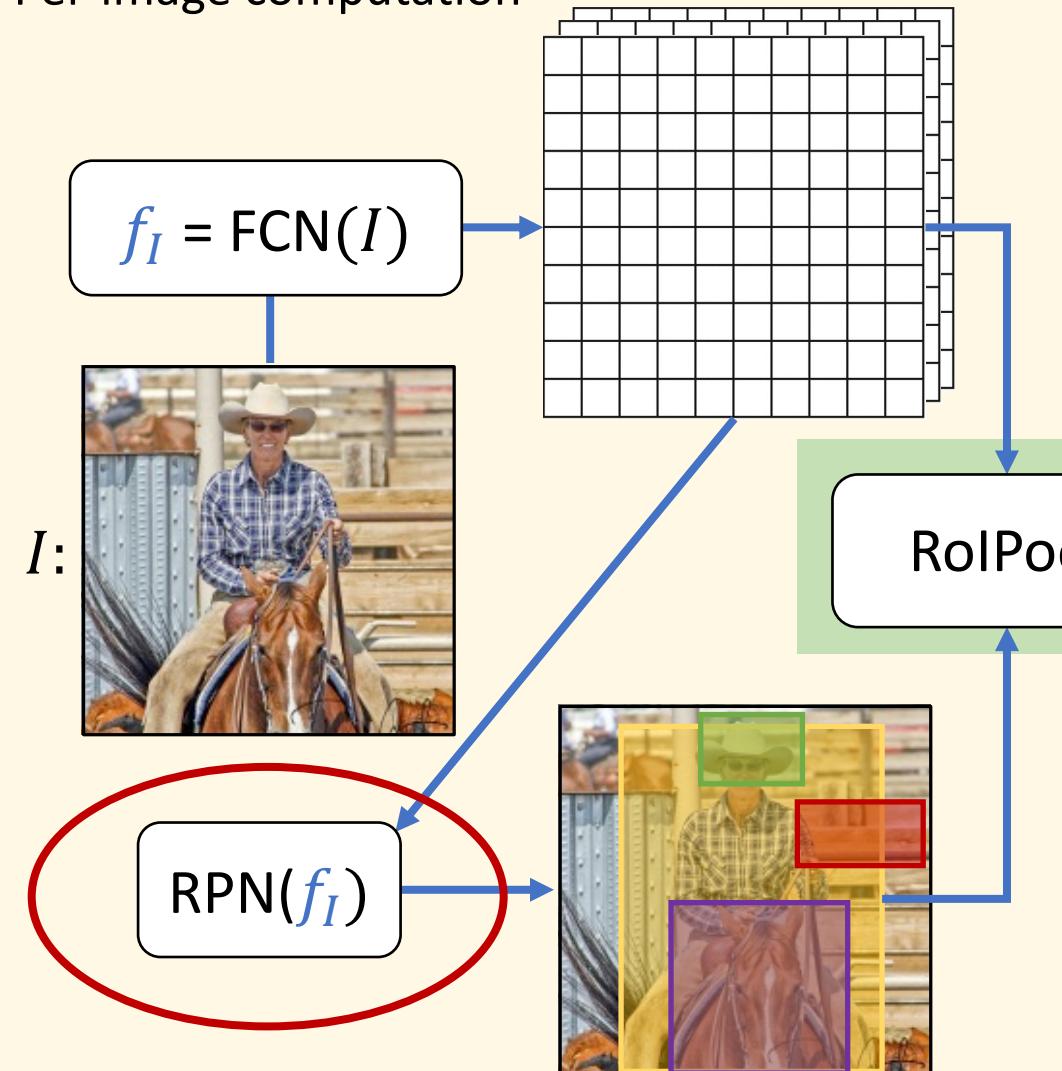
- 3x3 “sliding window”
- K objectness classifiers
 - K box regressors



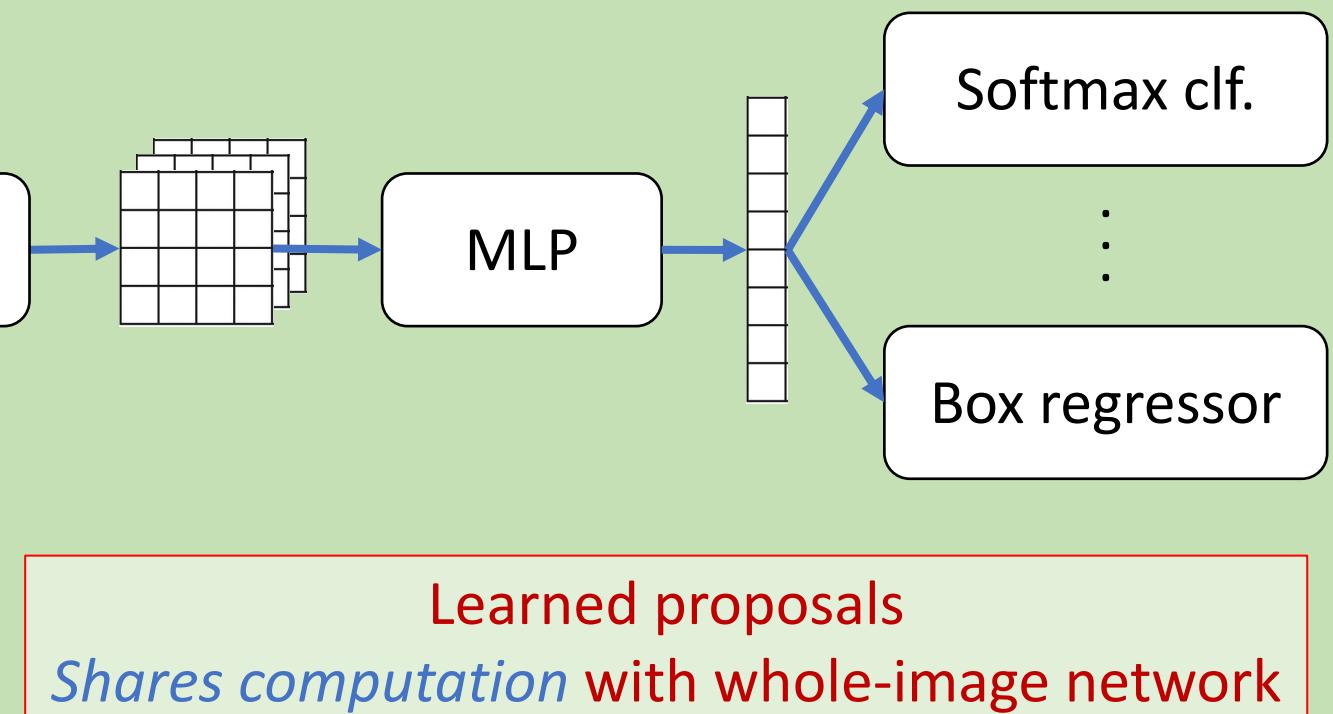
Anchor boxes: K anchors per location with different scales and aspect ratios

Faster R-CNN

Per-image computation



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

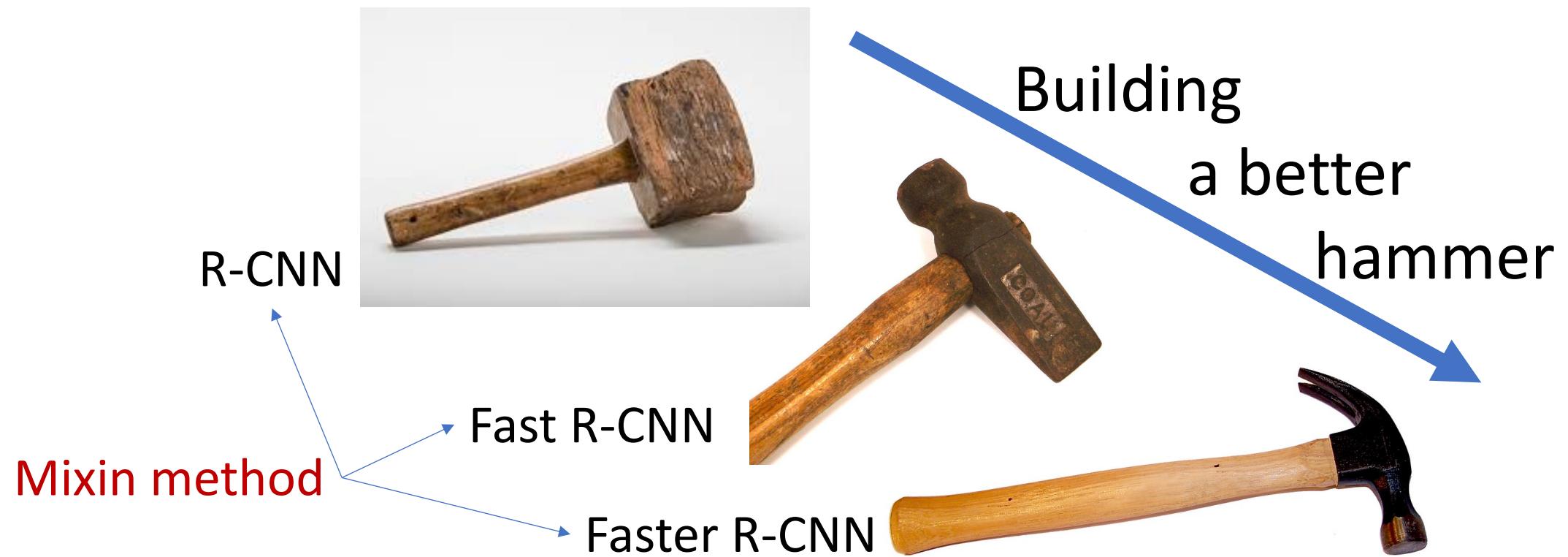


Faster R-CNN

	R-CNN	Fast R-CNN	Faster R-CNN
Test time per image (with proposals)	50 seconds	2 seconds	0.2 seconds
(Speedup)	1x	25x	250x
mAP (VOC 2007)	66.0	66.9	66.9

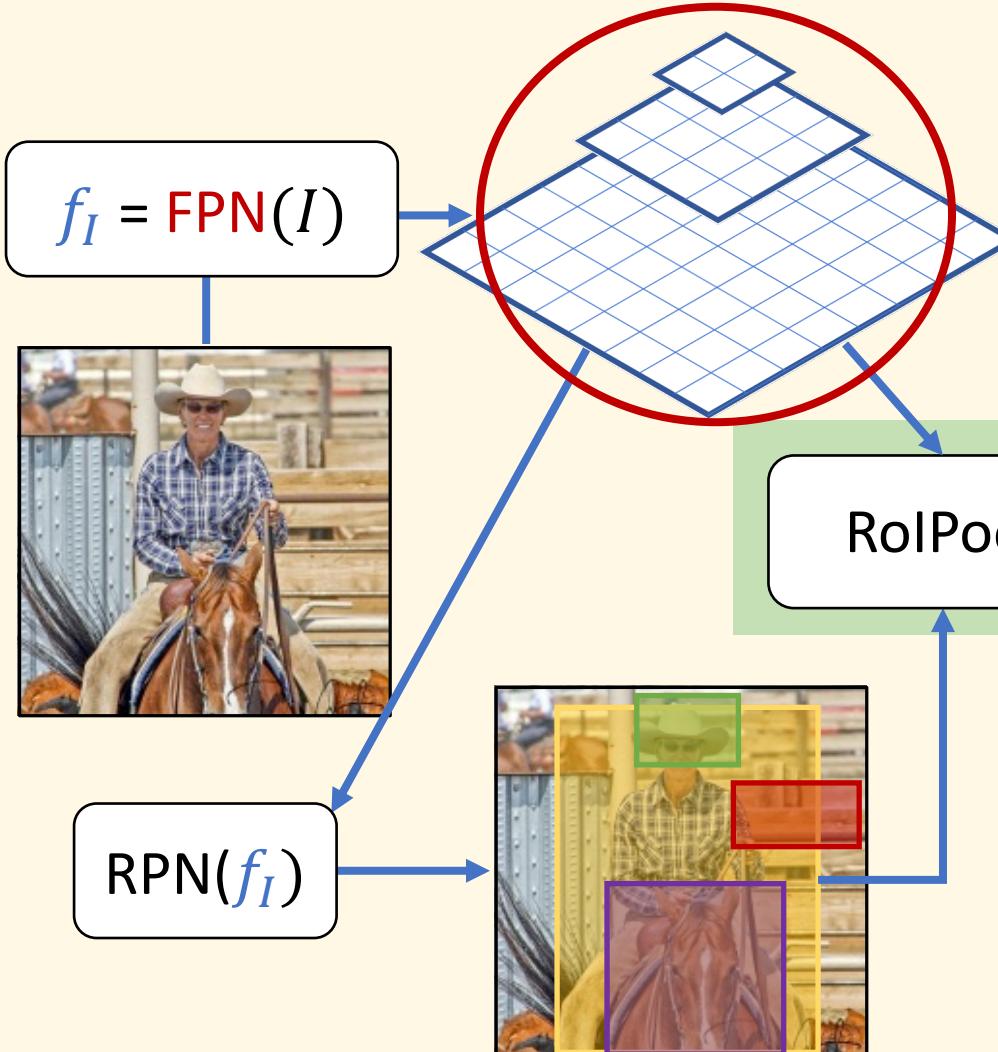
Improvements from “Mixin Methods”

- Largely orthogonal to the base detector
- Can be added to many different detectors as a “mixin”
(A better backbone network is one example)

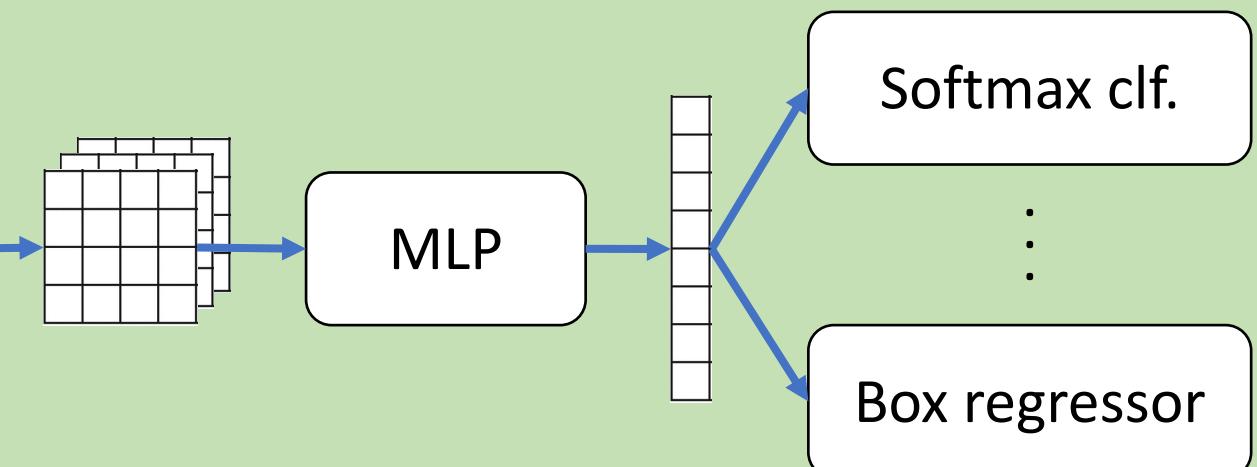


Faster R-CNN with a Feature Pyramid Network

Per-image computation

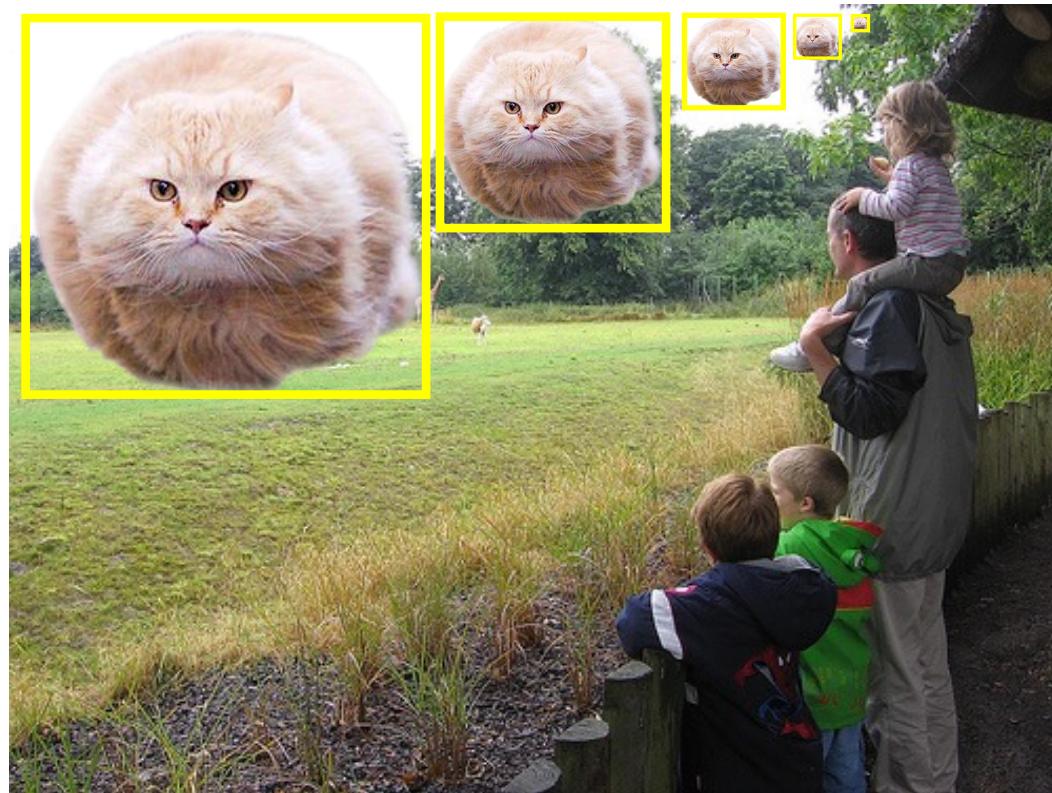


Per-region computation for each $r_i \in r(I)$



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

FPN: Improving Scale Invariance and Equivariance



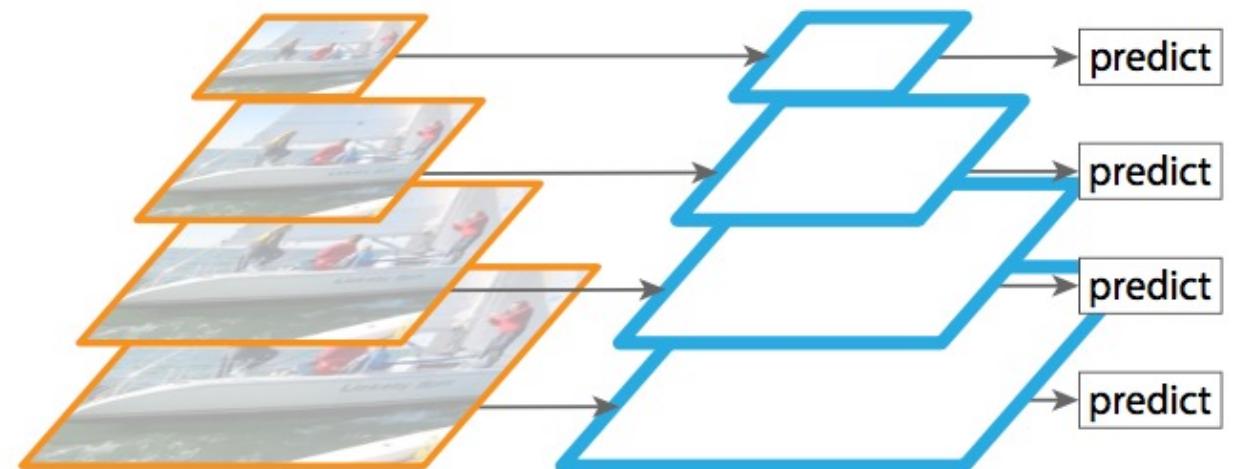
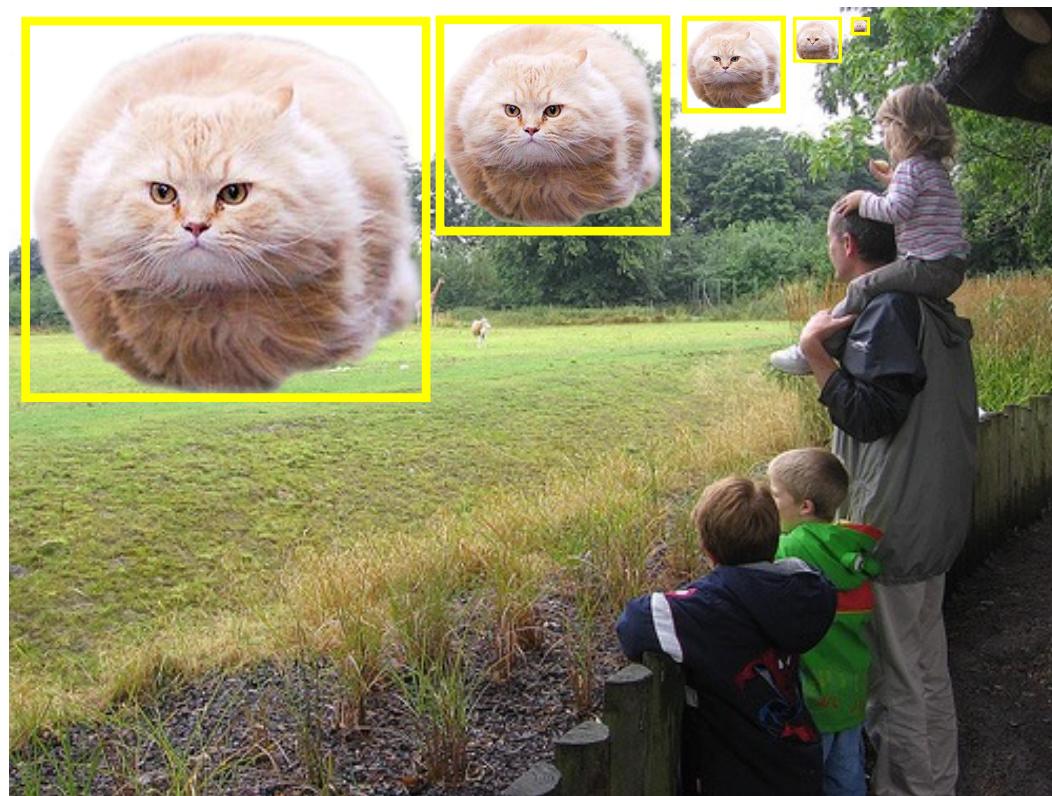
Detectors need to

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

objects over a **wide range of scales**

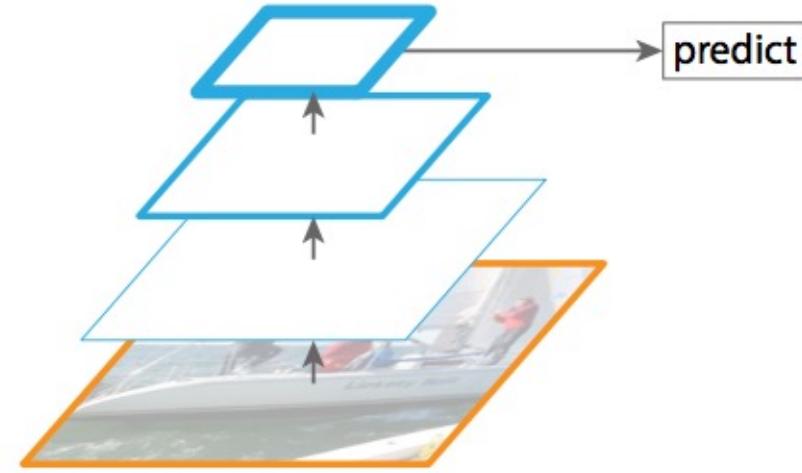
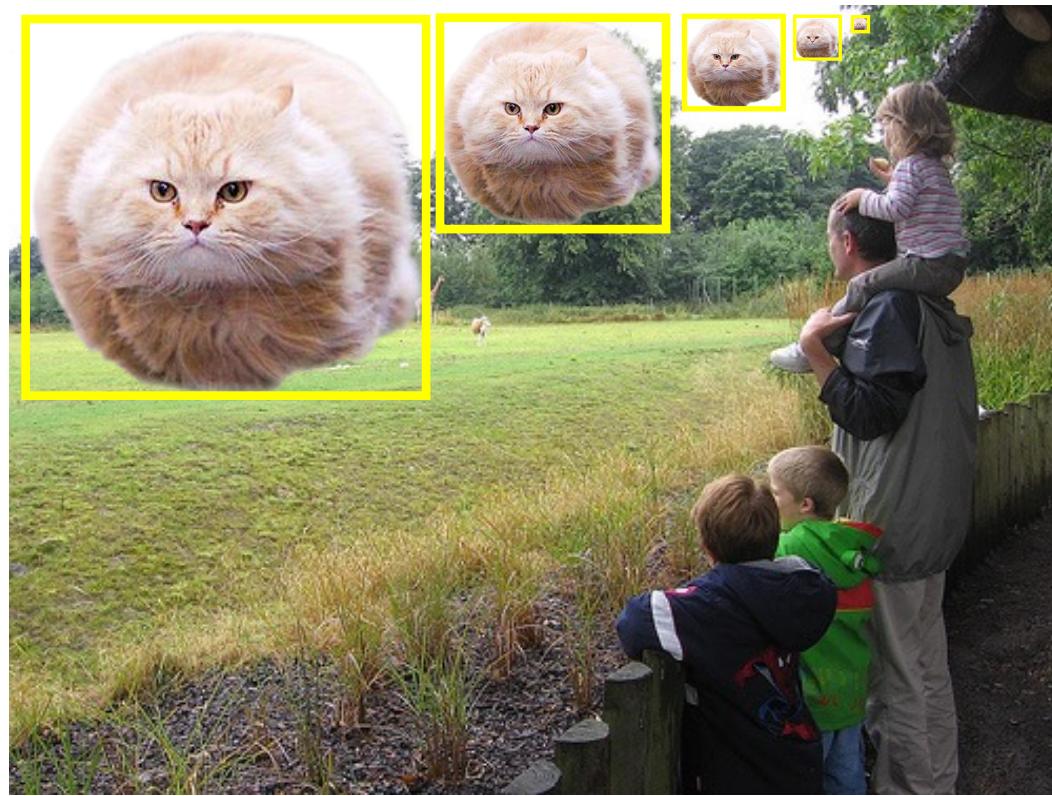
FPN improves this ability

Strategy 1: Image Pyramid



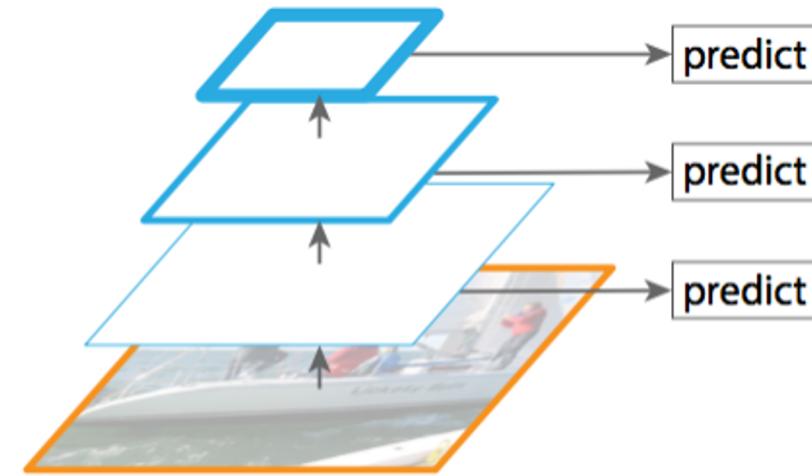
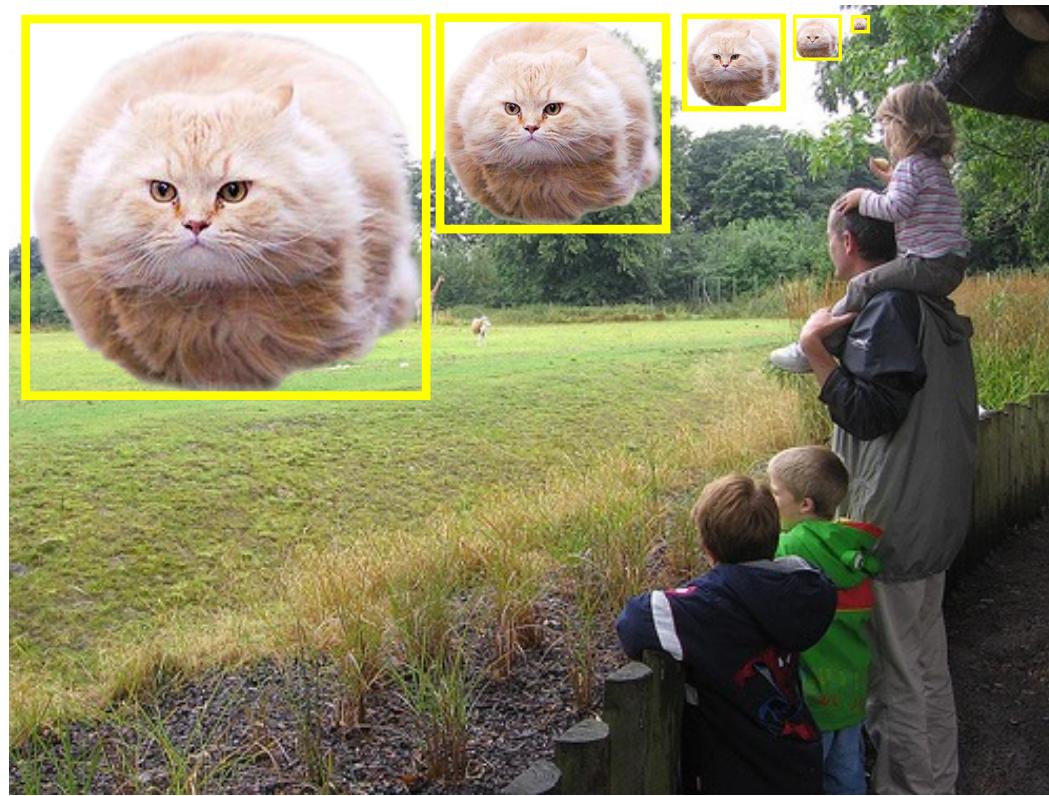
(a) Featurized image pyramid
Standard solution – *slow!*
(E.g., Viola & Jones, HOG, DPM, SPP-net,
multi-scale Fast R-CNN, ...)

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



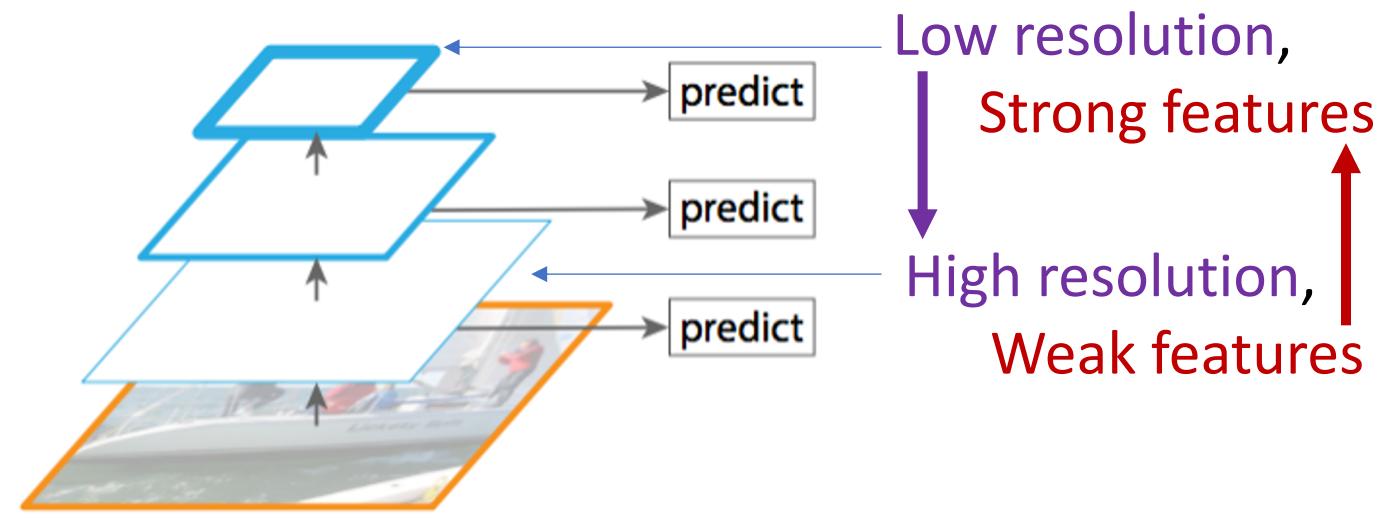
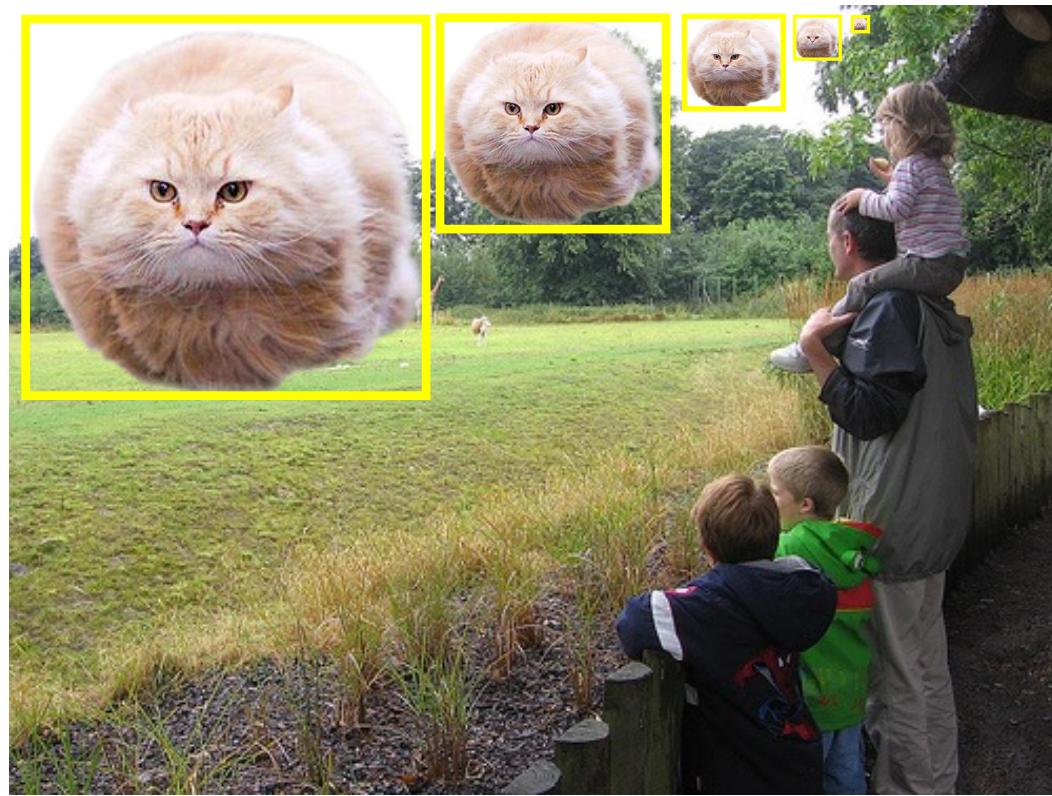
(b) Single feature map
Leave it all to the features – fast, suboptimal
(E.g., Fast/er R-CNN, YOLO, ...)

Strategy 3: Naïve In-network Pyramid



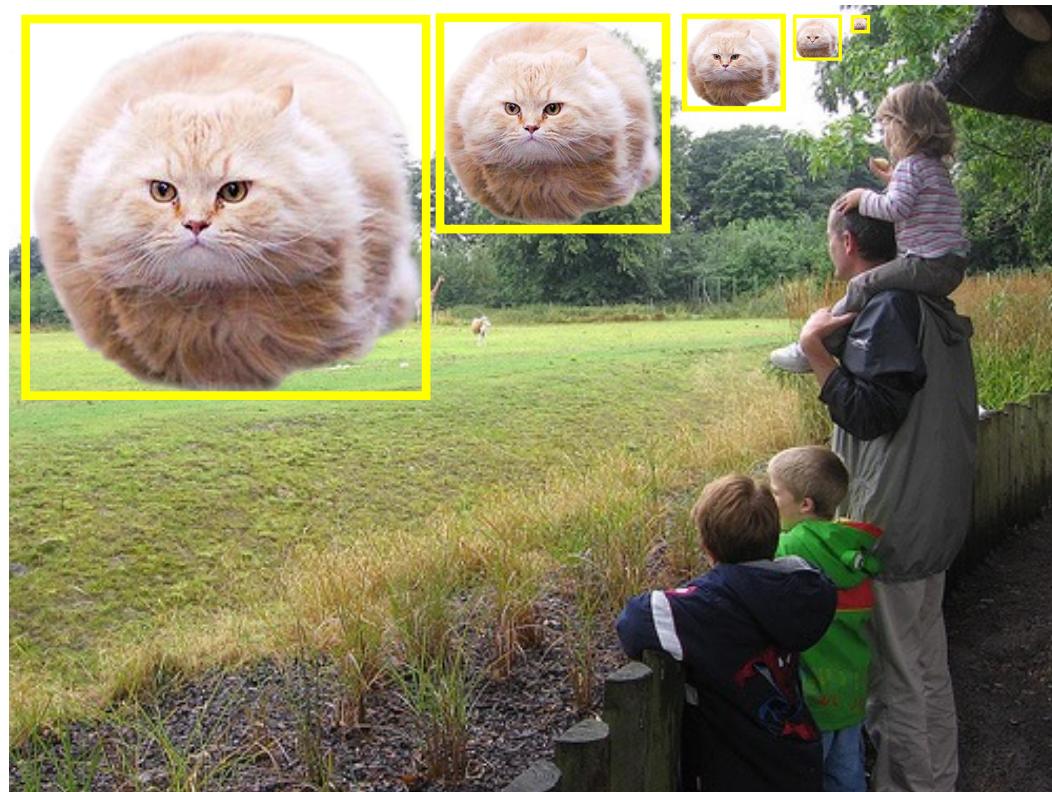
(c) Pyramidal feature hierarchy
Use the internal pyramid – *fast, suboptimal*
(E.g., \approx SSD, ...)

Strategy 3: Naïve In-network Pyramid



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

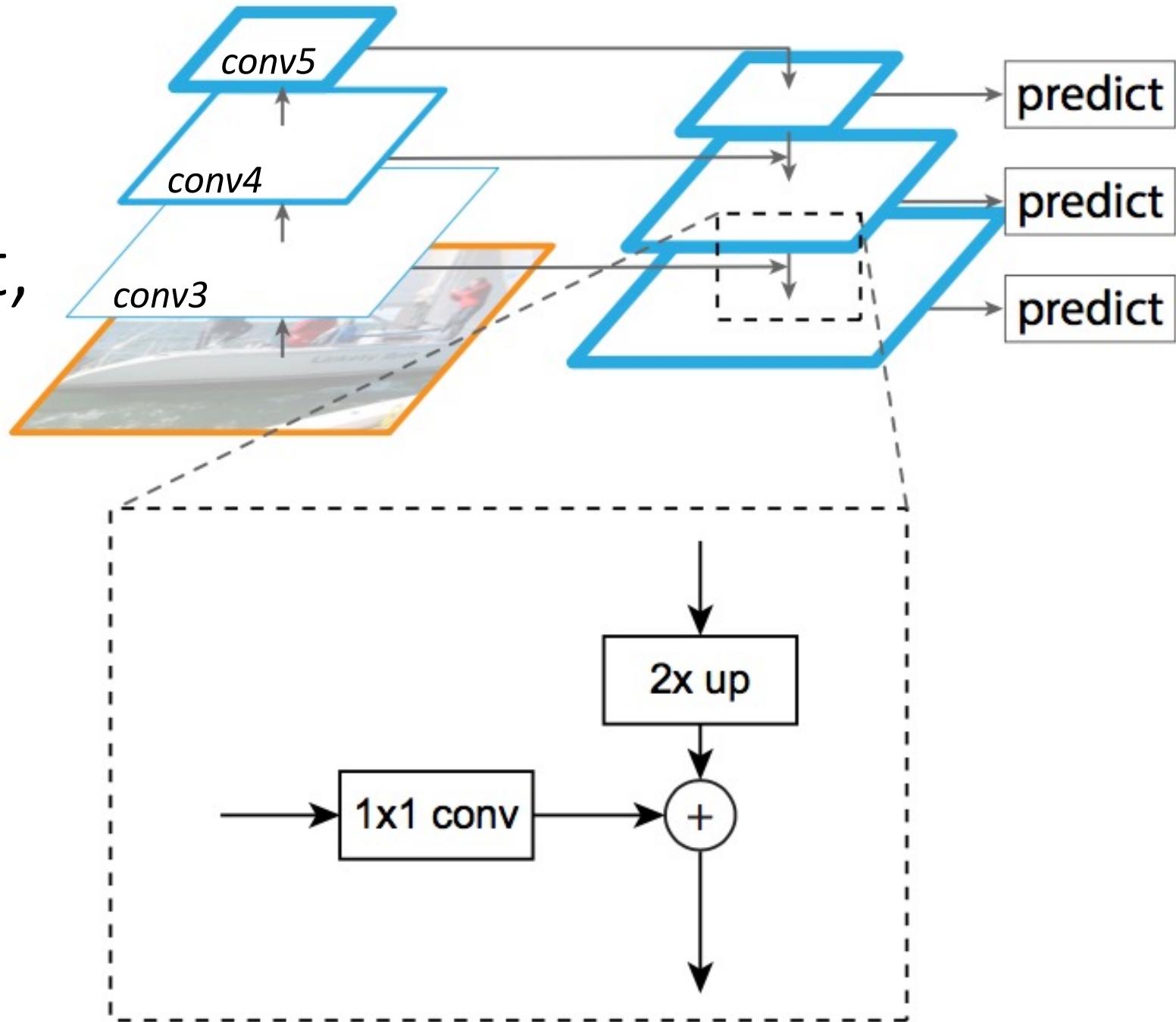
Strategy 4: Feature Pyramid Network



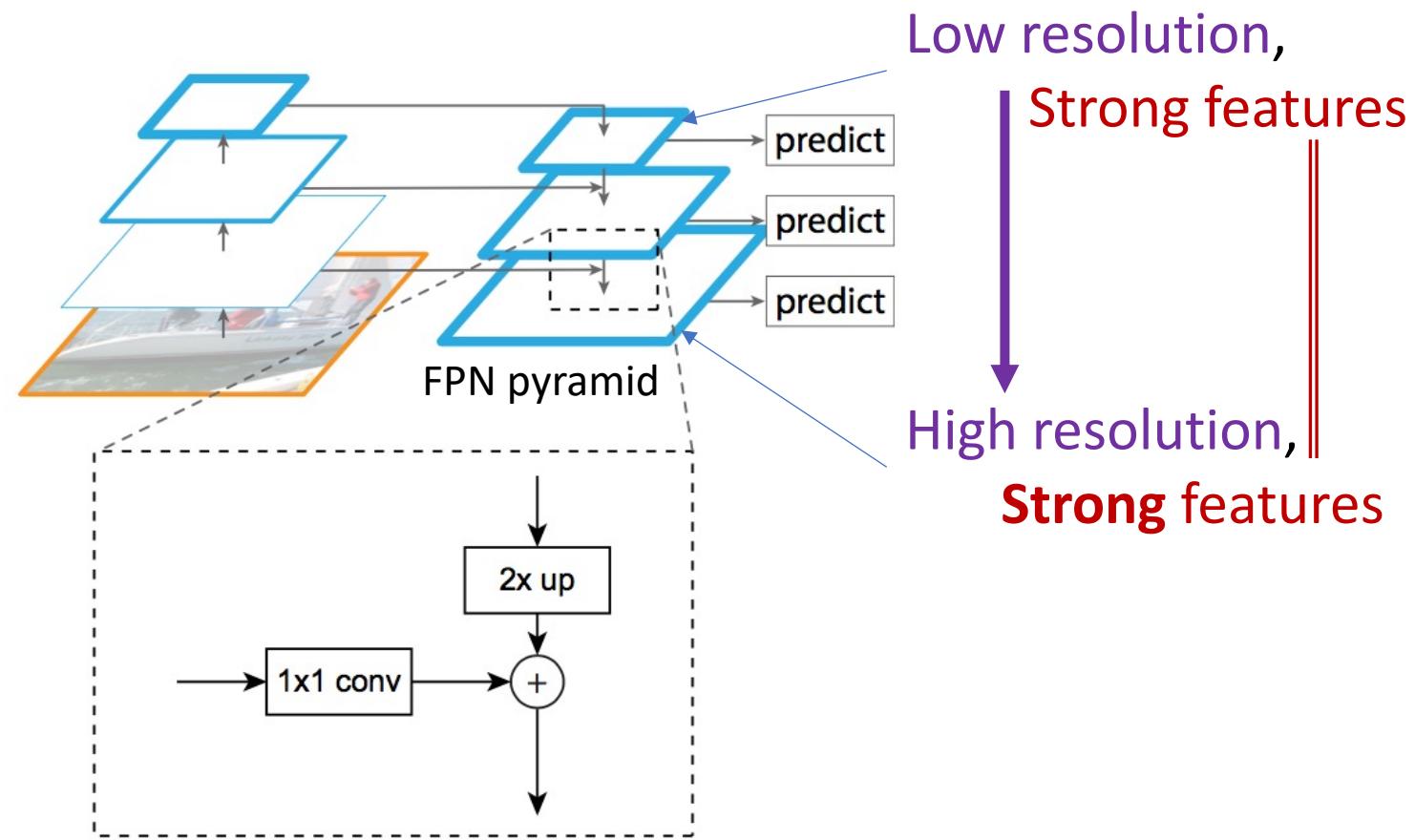
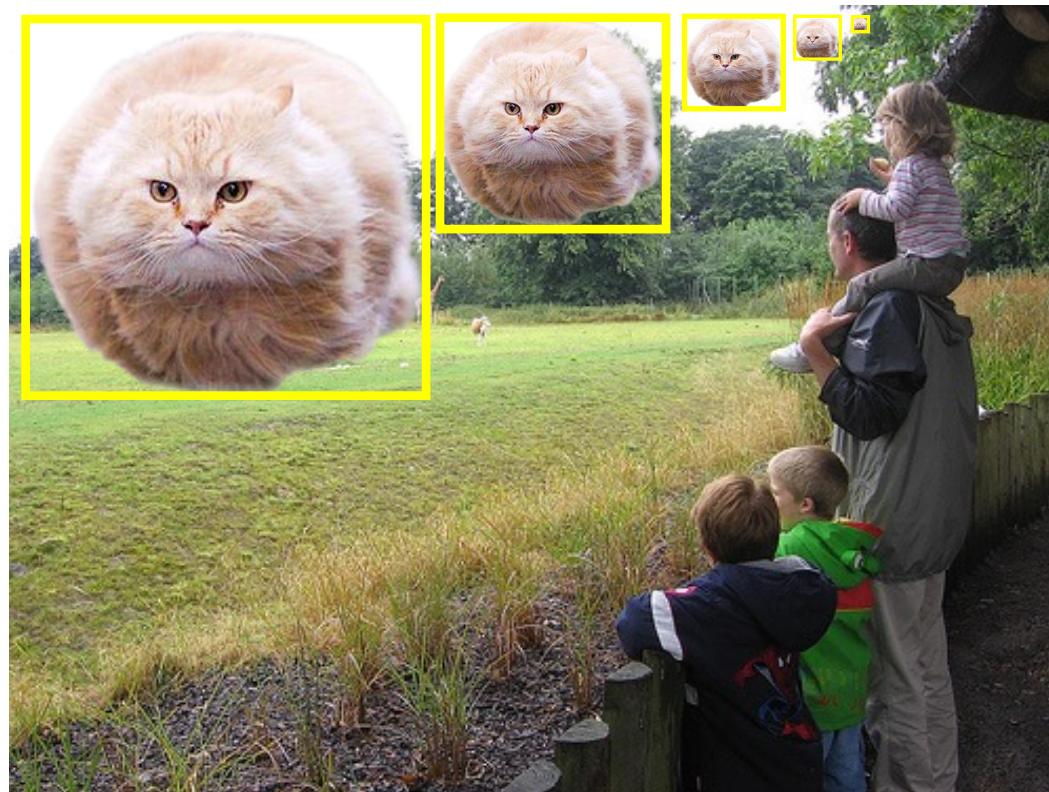
(d) Feature Pyramid Network

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

FPN:
Light-weight,
Top-down
Refinement
Module



No Compromise on Feature Quality, Still Fast



Mask R-CNN: The Final Hammer

R-CNN



Fast R-CNN



Faster R-CNN



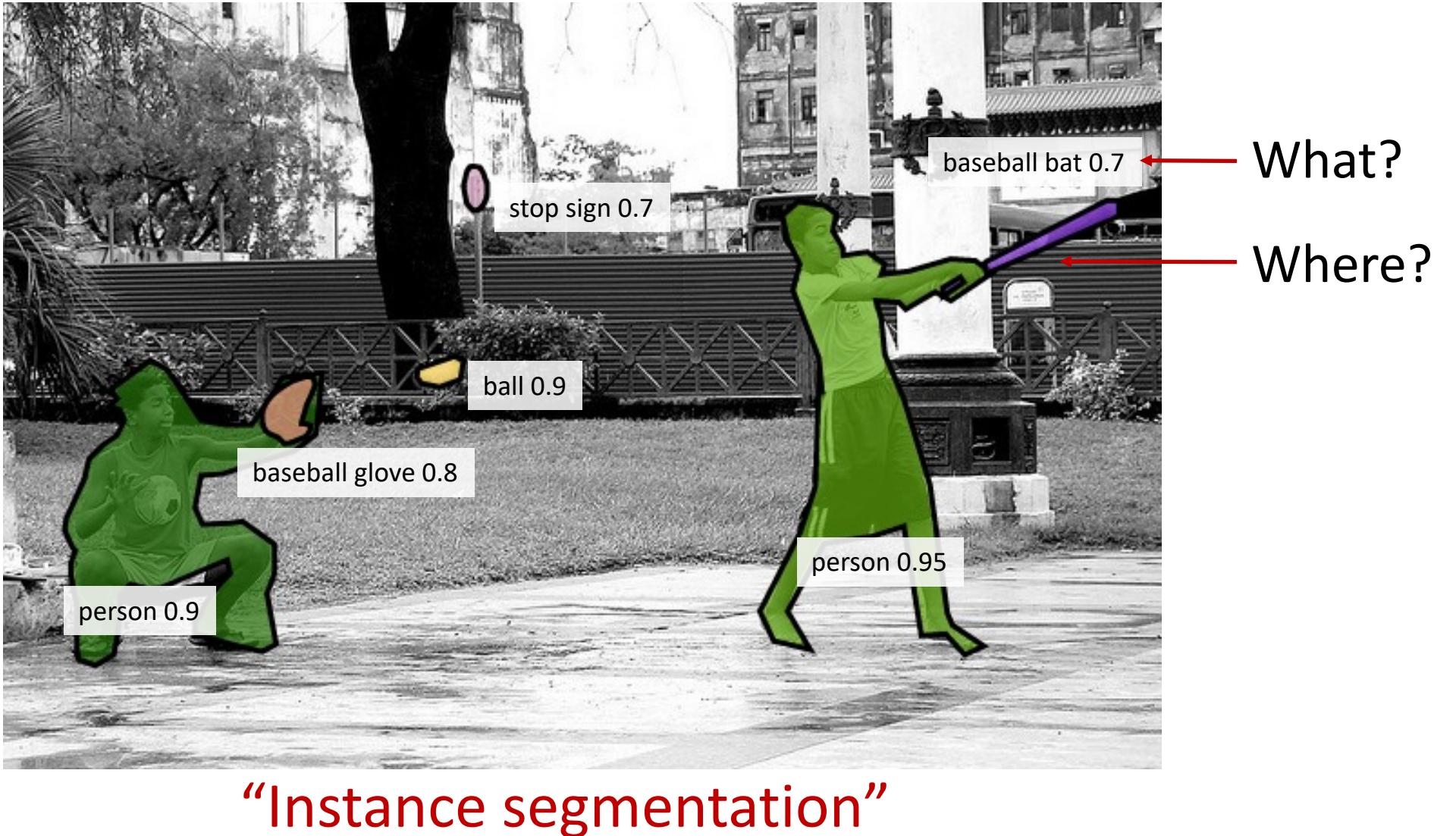
Mask R-CNN



Building
a better
hammer

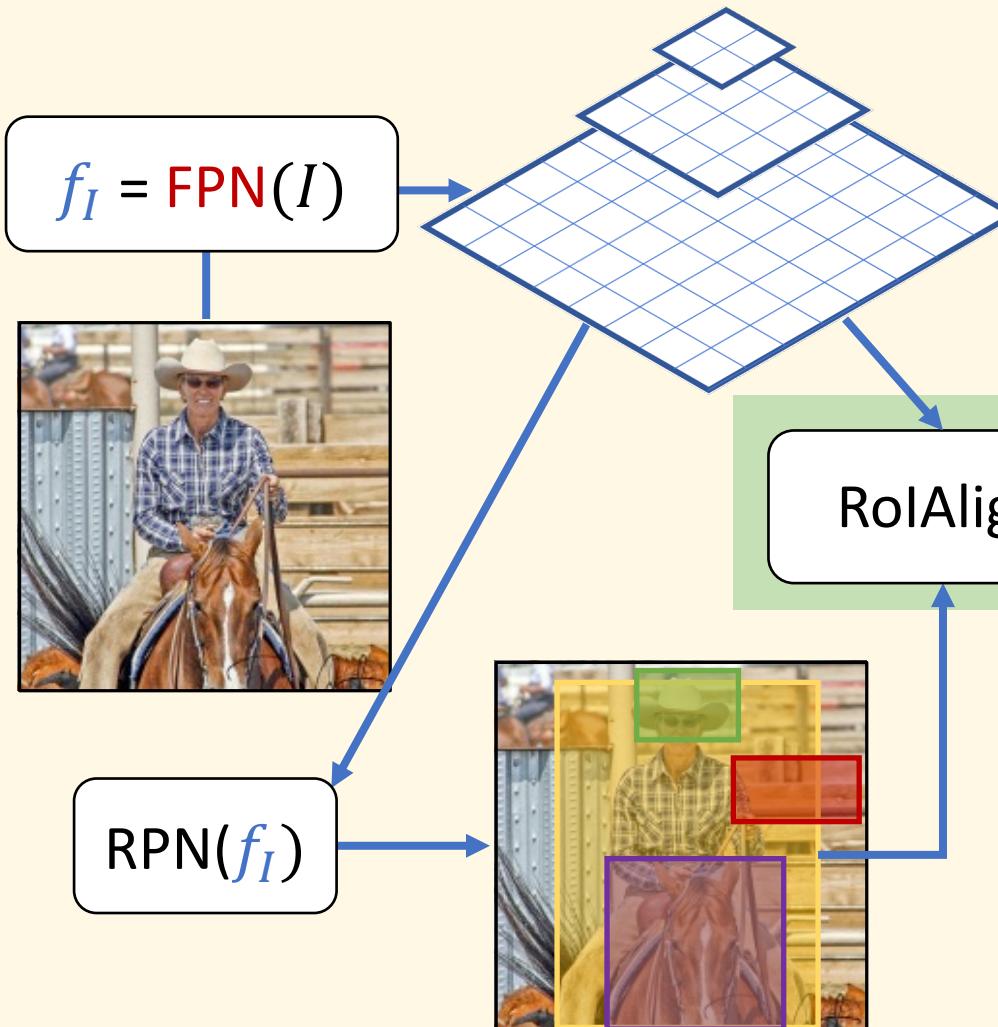
A blue arrow points diagonally from the "Fast R-CNN" hammer towards the "Mask R-CNN" hammer, indicating a progression or evolution.

Object Detection with Segmentation Masks

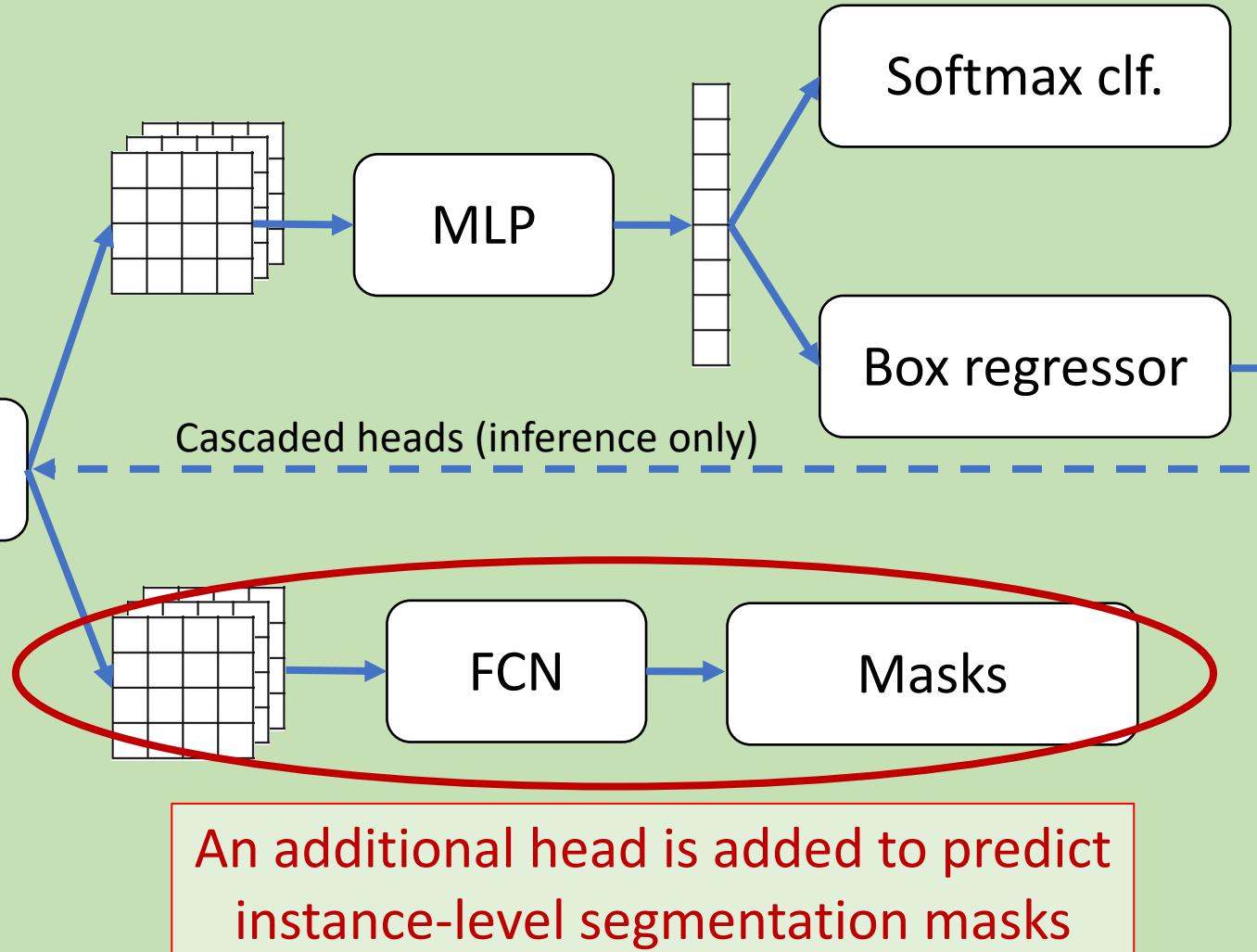


Mask R-CNN

Per-image computation

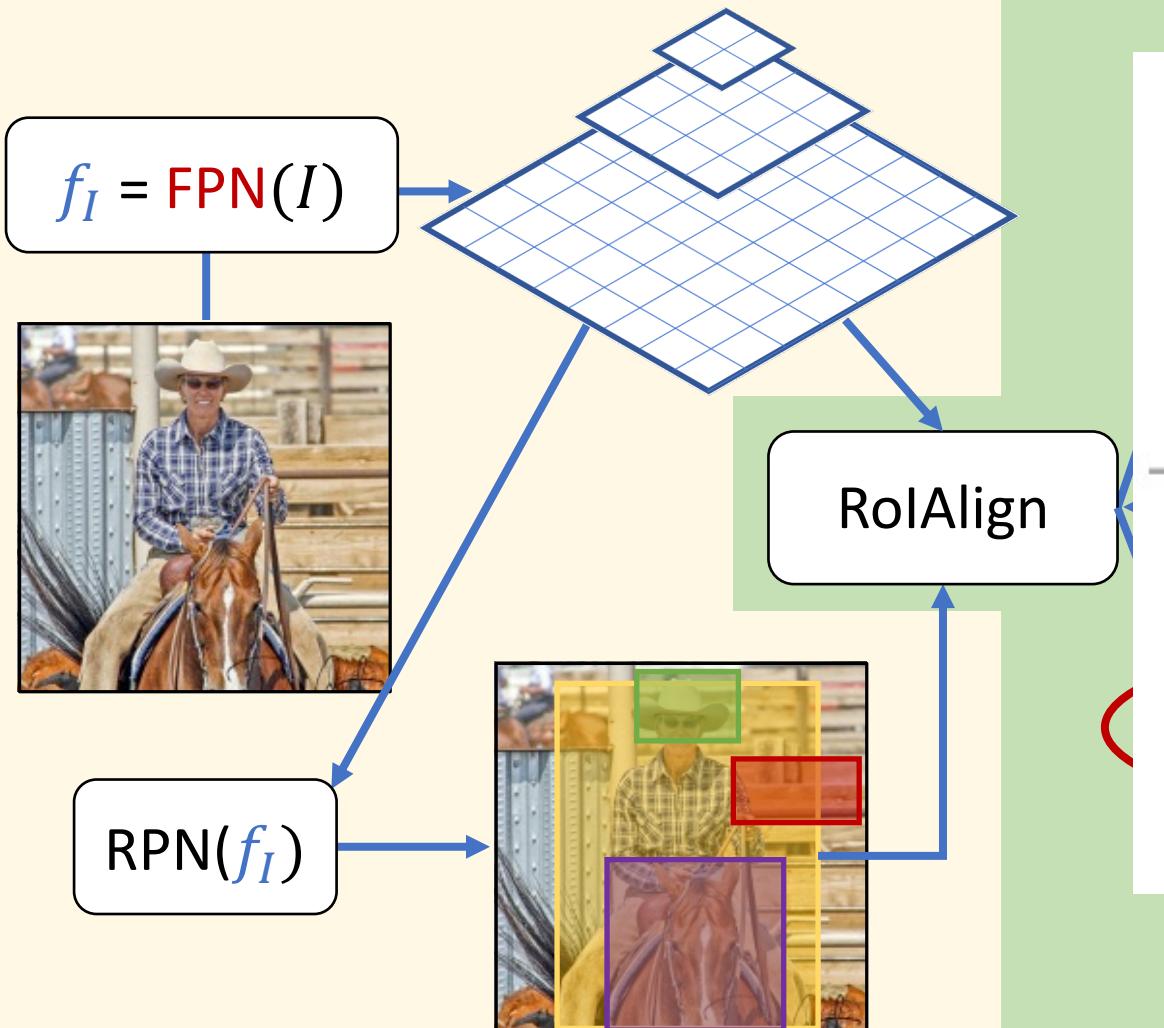


Per-region computation for each $r_i \in r(I)$

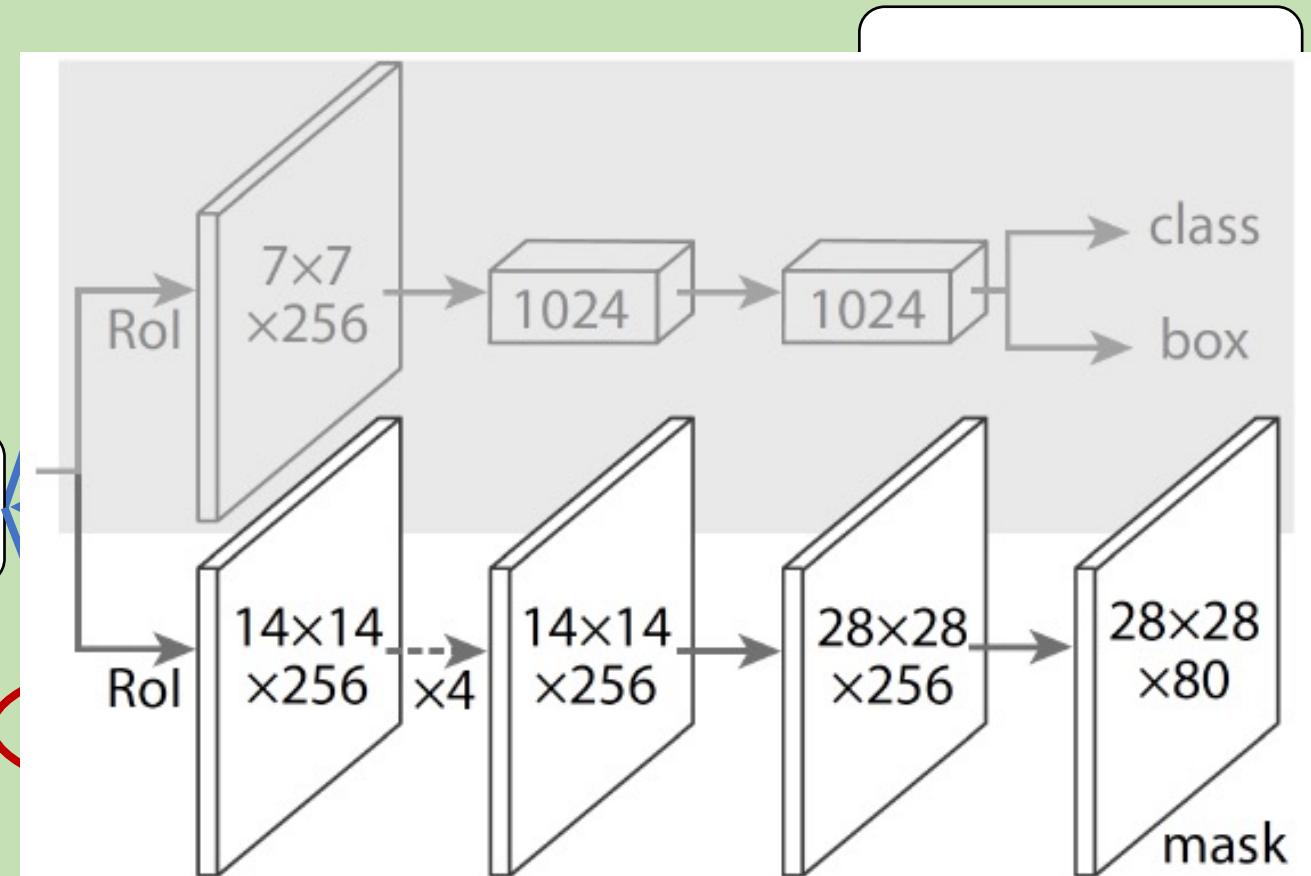


Mask R-CNN

Per-image computation



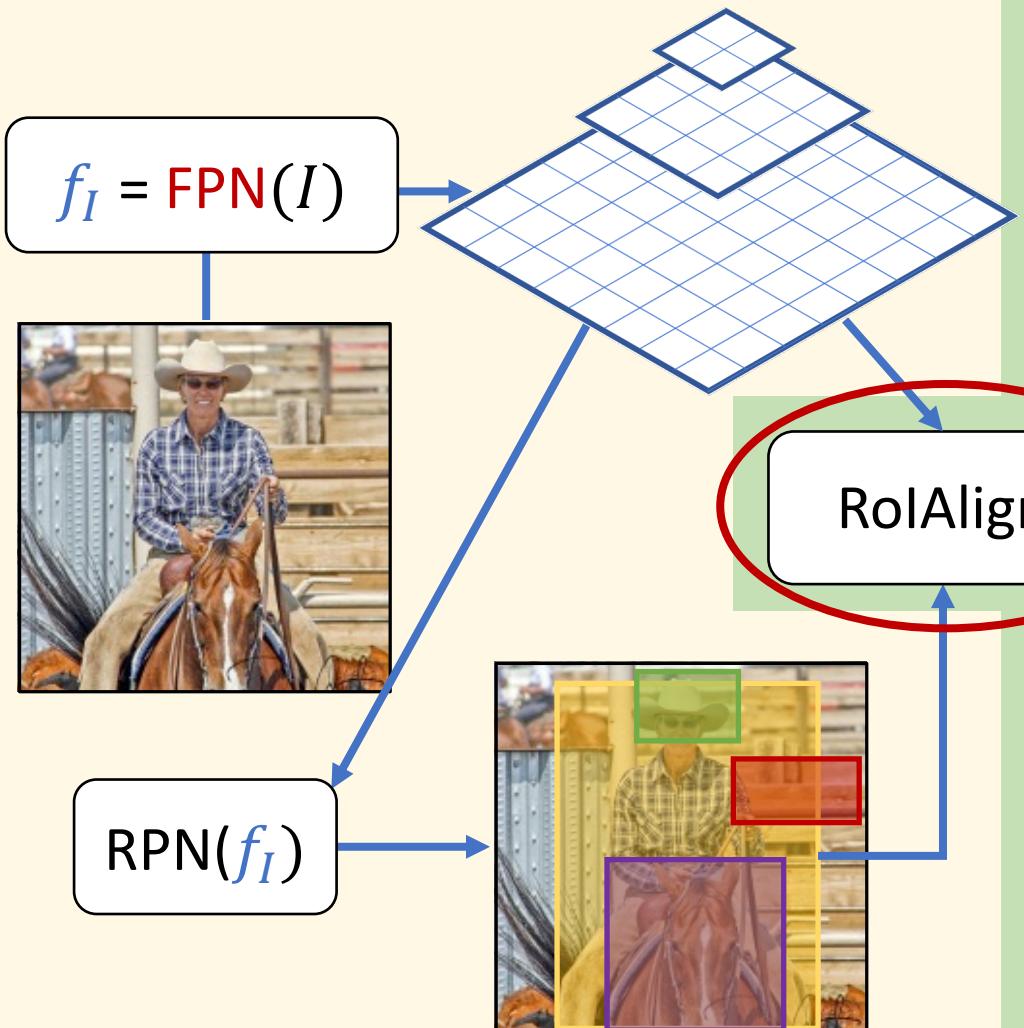
Per-region computation for each $r_i \in r(I)$



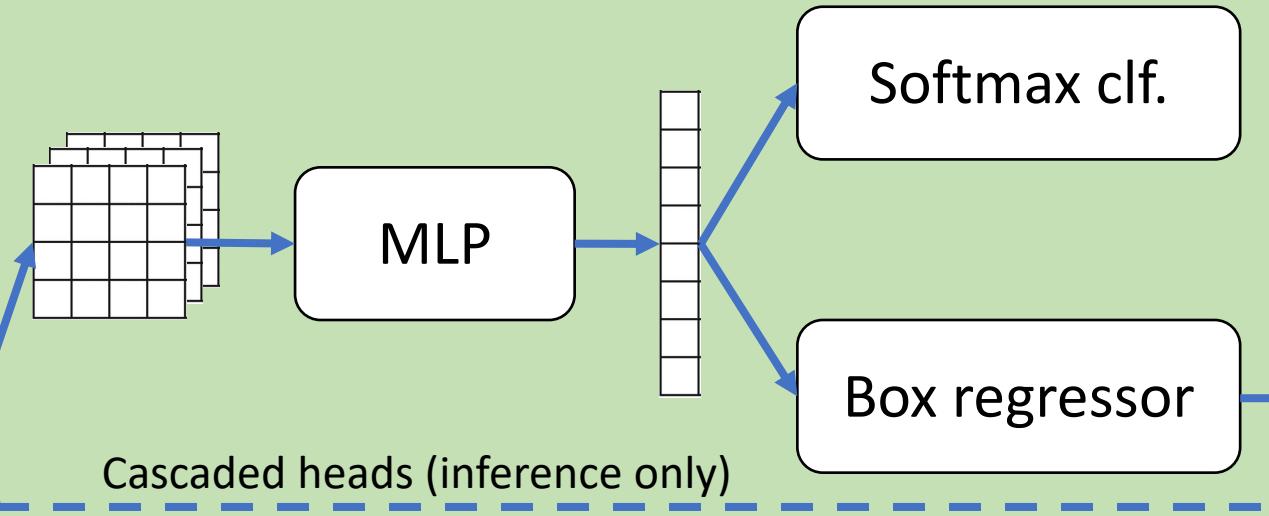
An additional head is added to predict instance-level segmentation masks

Mask R-CNN

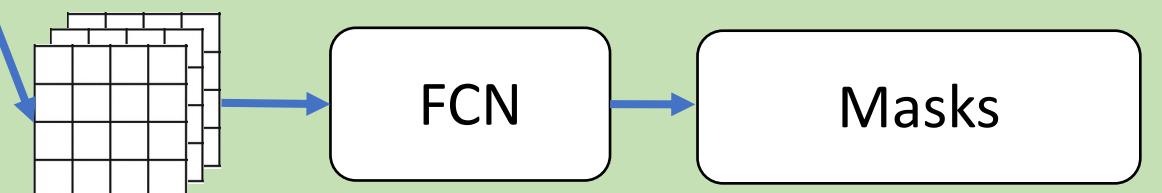
Per-image computation



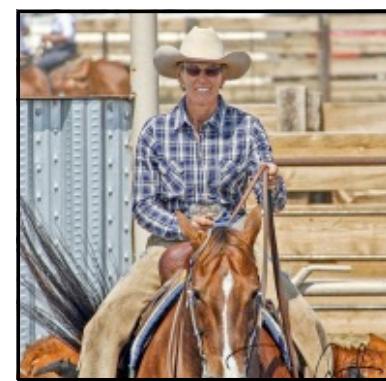
Per-region computation for each $r_i \in r(I)$



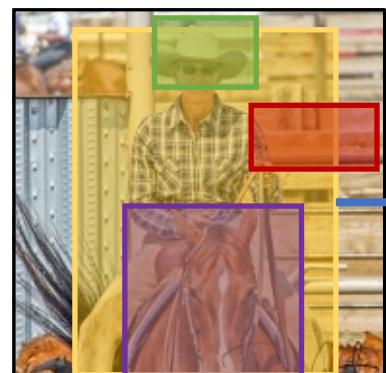
RoIPool is replaced with RoIAlign



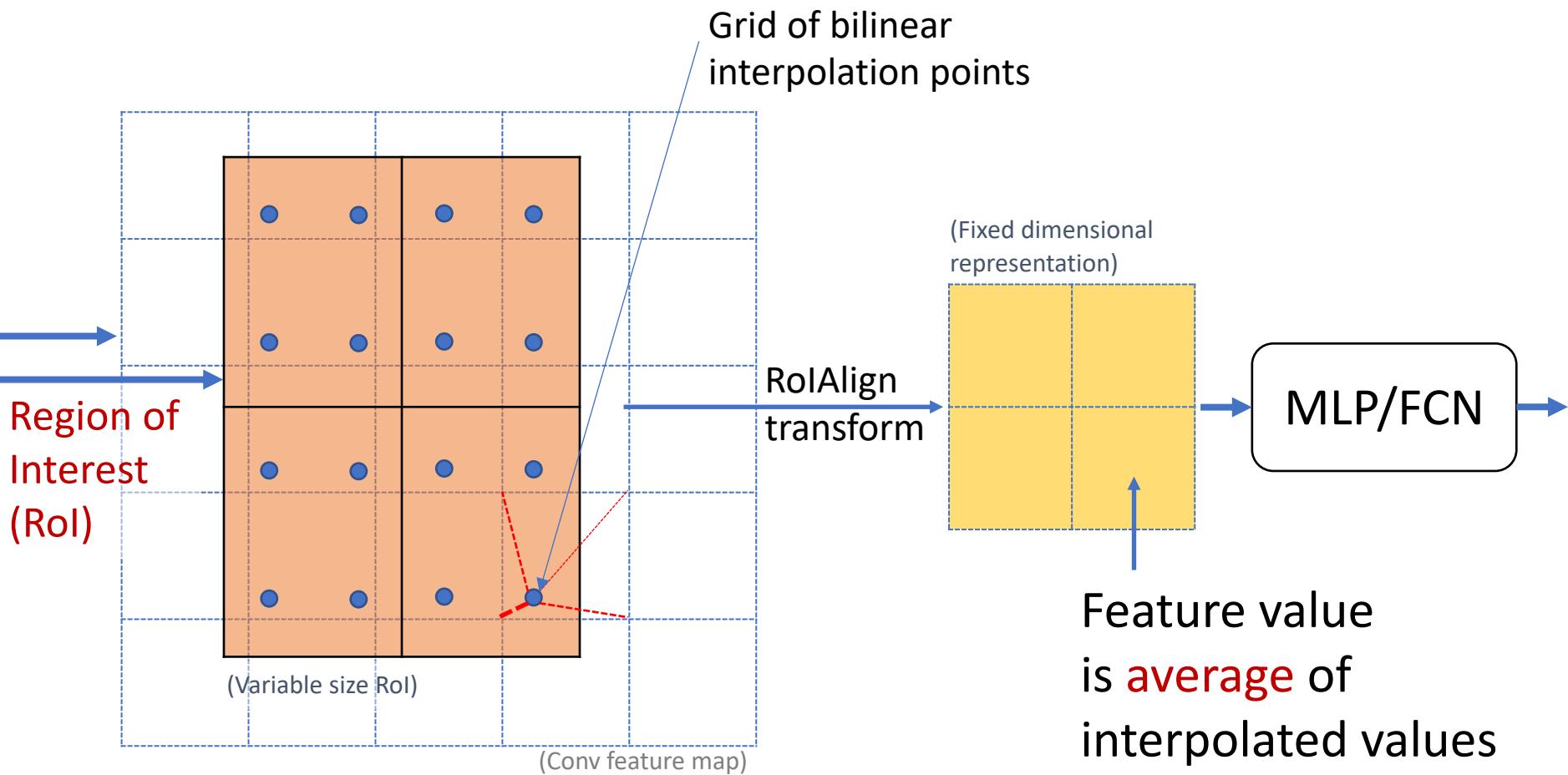
RoIAlign Operation (on each Proposal)



$$f_I = \text{FCN}(I)$$

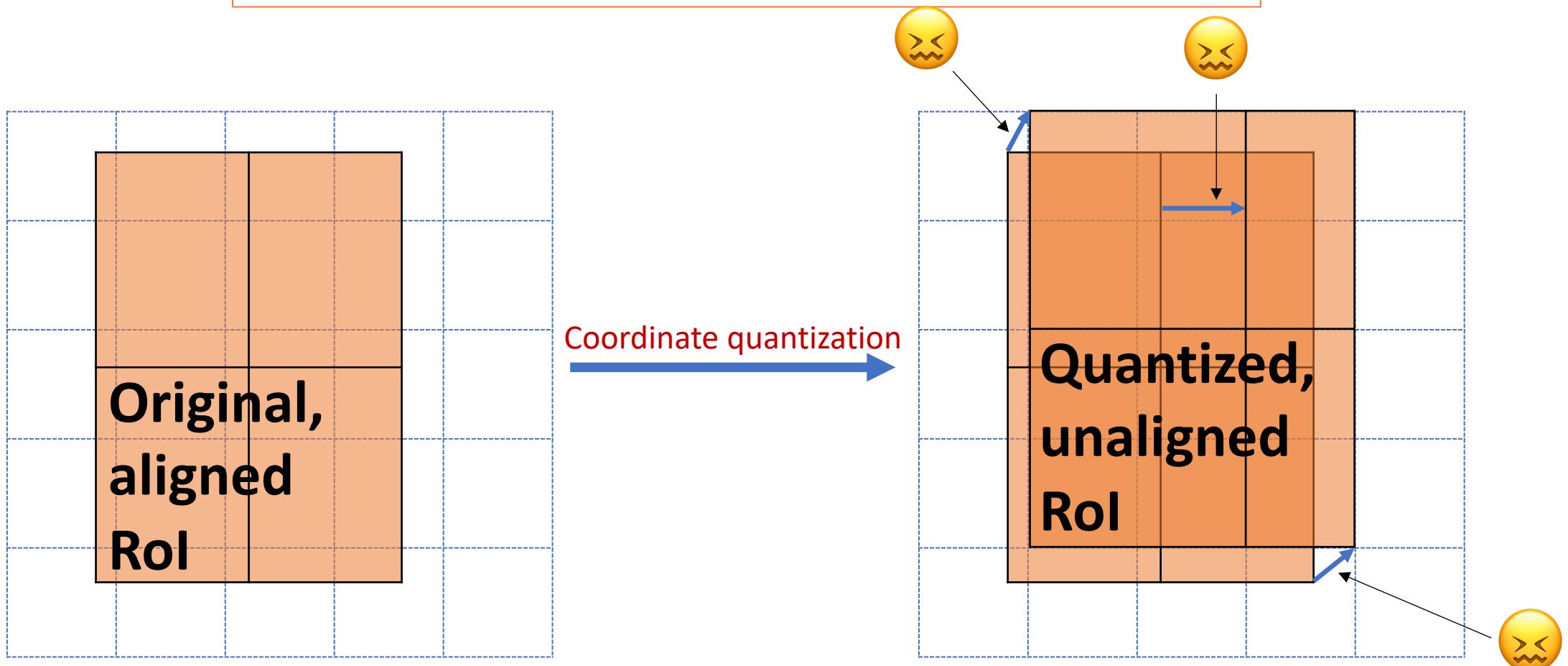


Transform **arbitrary size proposal** into a **fixed-dimensional representation** (e.g., 2x2)



Compare to ROI Pool and ROI Warp

Quantization breaks pixel-to-pixel alignment
between input and output

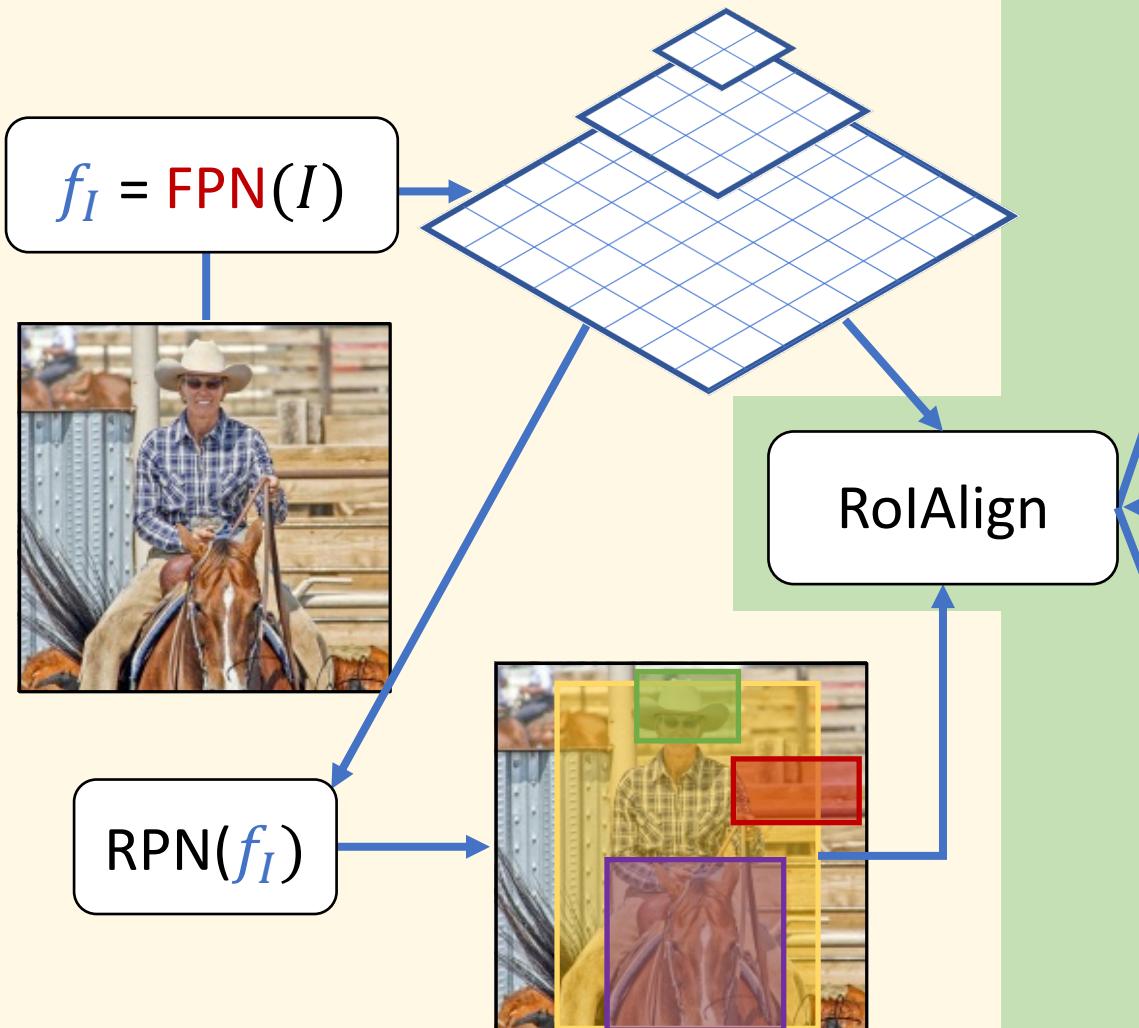


Effects of RoIAlign Operation

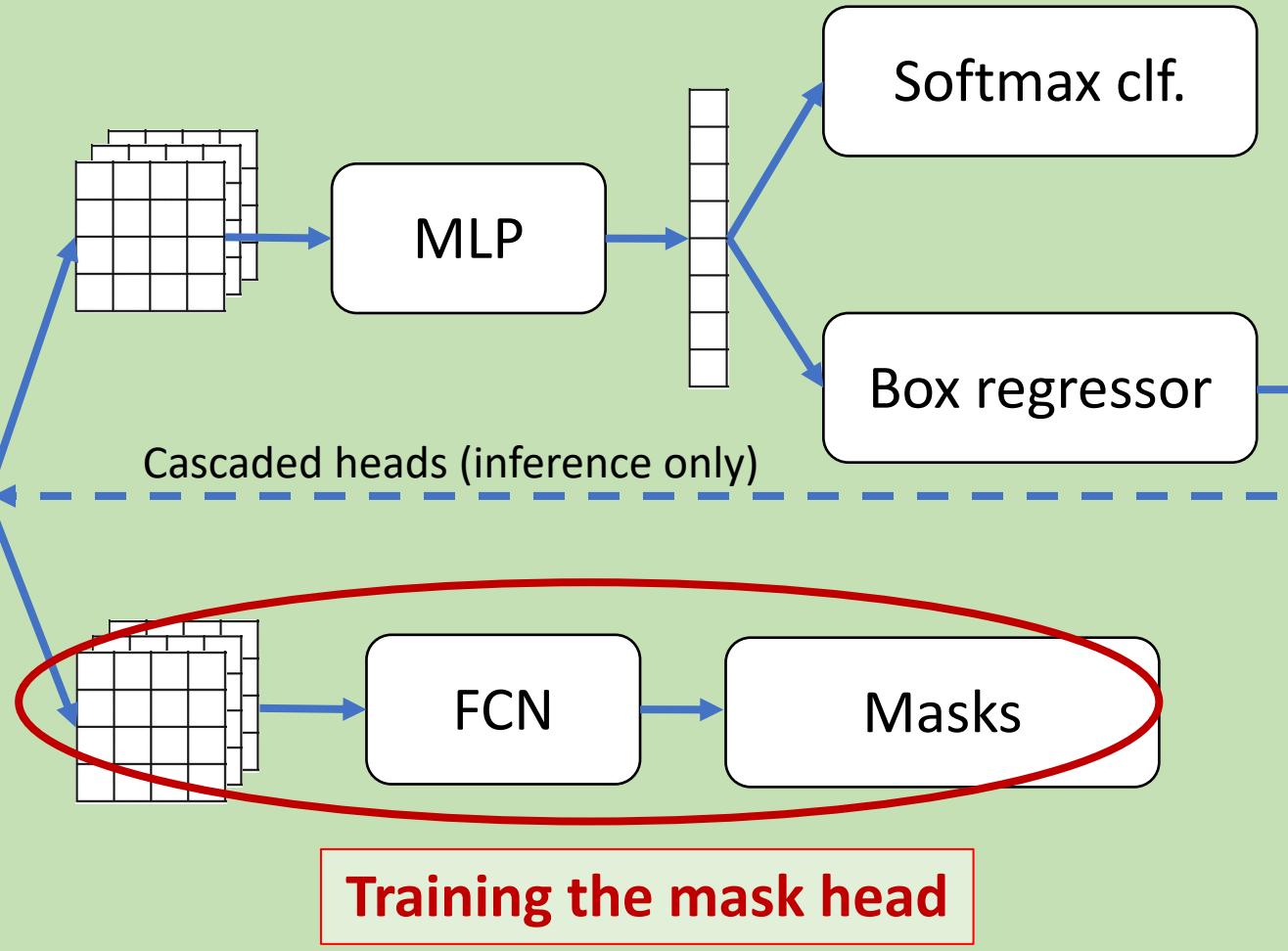
	AP	AP ₅₀	AP ₇₅	AP ^{bb}	AP ^{bb} ₅₀	AP ^{bb} ₇₅
<i>RoIPool</i>	23.6	46.5	21.6	28.2	52.7	26.9
<i>RoIAlign</i>	30.9	51.8	32.1	34.0	55.3	36.4
	+7.3	+ 5.3	+10.5	+5.8	+2.6	+9.5

Mask R-CNN

Per-image computation



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



Mask R-CNN: Inference

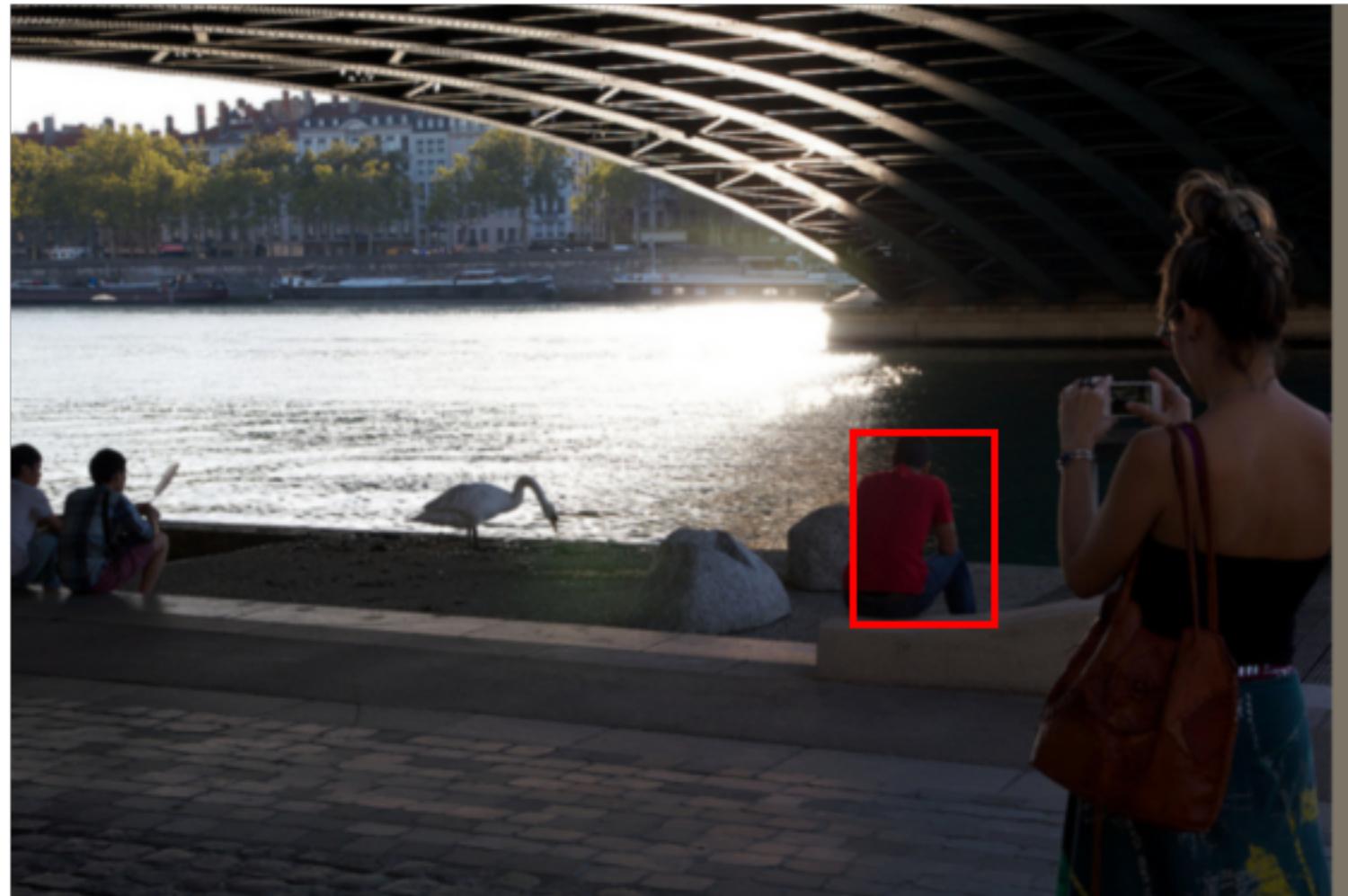
1. Perform Faster R-CNN inference

- Generate proposals (RPN)
- Score the proposals
- Regress from proposals to refined detection boxes
- Apply NMS and take the top K ($= 100$, e.g.)

2. Run RoIAlign and mask head on top- K refined, post-NMS boxes

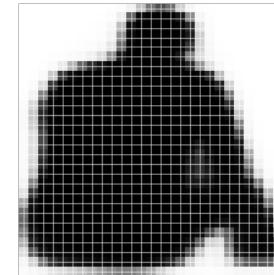
- Fast (only compute masks for top- K detections)
- Improves accuracy (uses *refined* detection boxes, not proposals)

Mask Prediction



Validation image with box detection shown in red

28x28 soft prediction from Mask R-CNN
(enlarged)



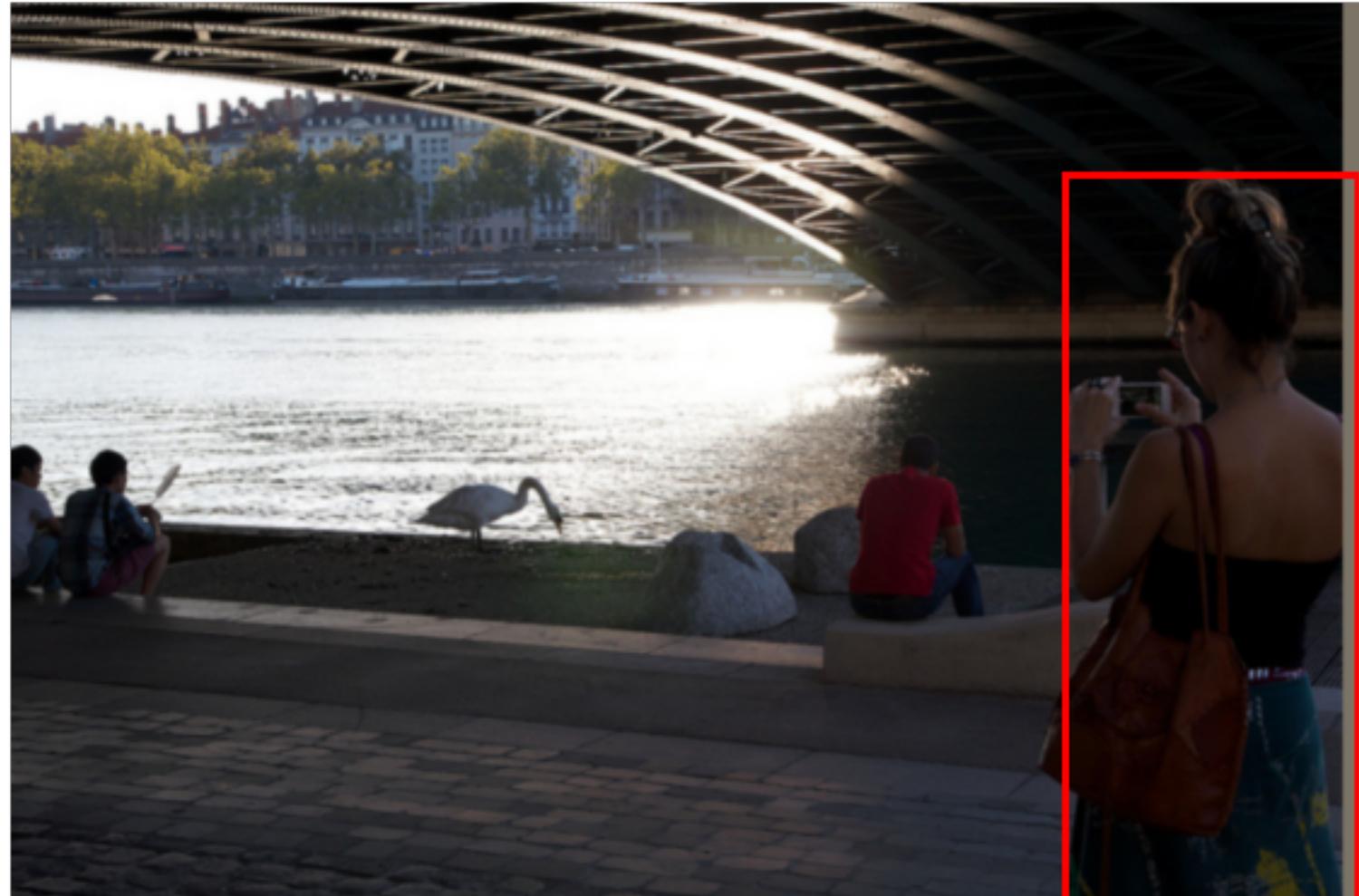
Soft prediction **resampled to image coordinates**
(bilinear and bicubic interpolation work equally well)



Final prediction (threshold at 0.5)

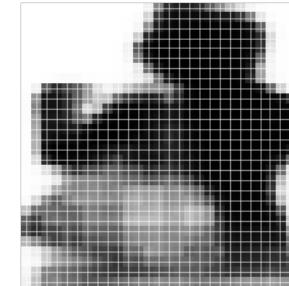


Mask Prediction



Validation image with box detection shown in red

28x28 soft prediction



Resized soft prediction



Final mask

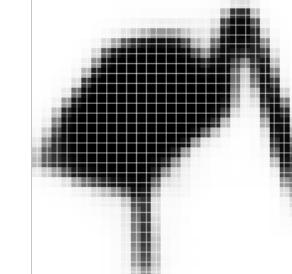


Mask Prediction



Validation image with box detection shown in red

28x28 soft prediction



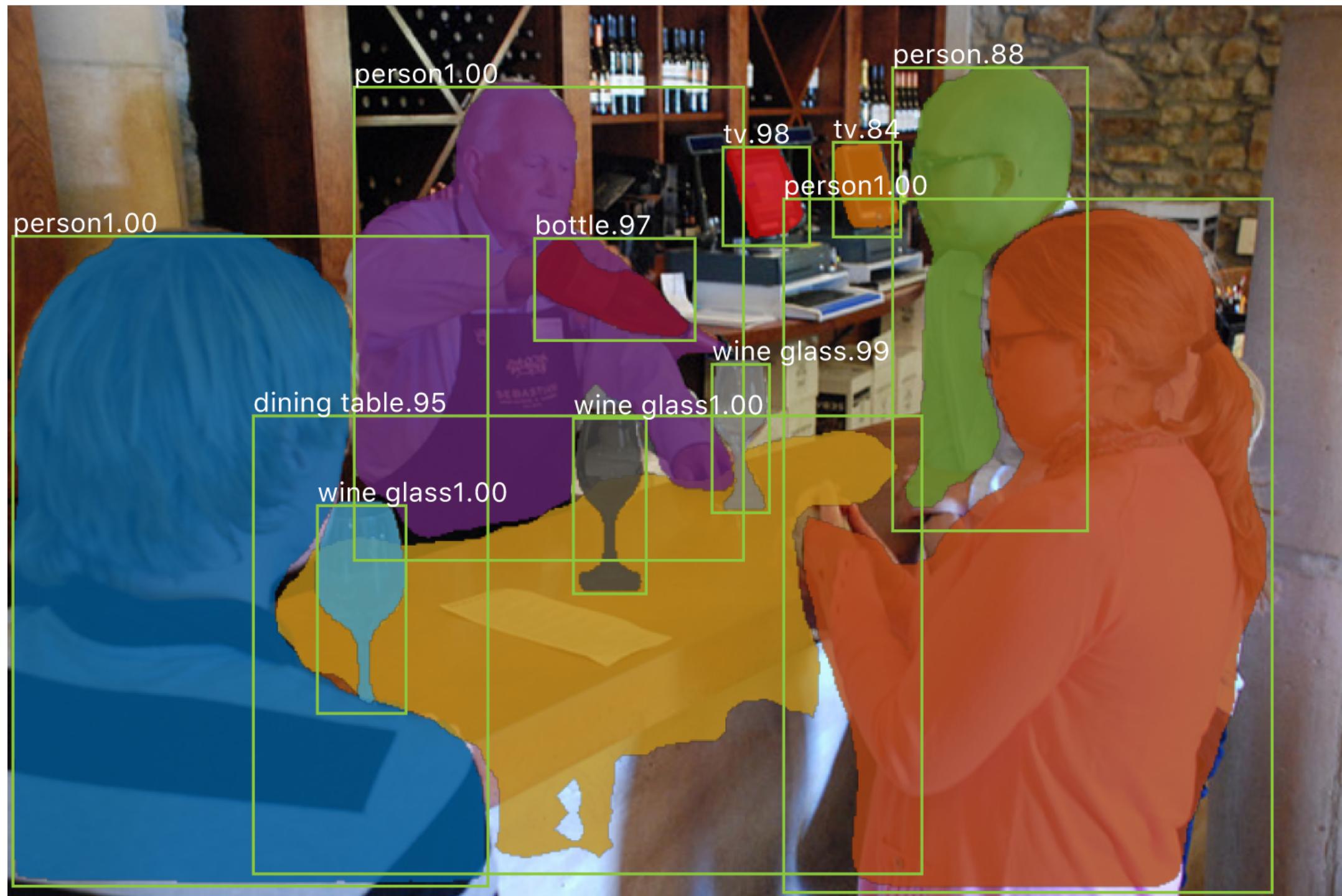
Resized Soft prediction



Final mask





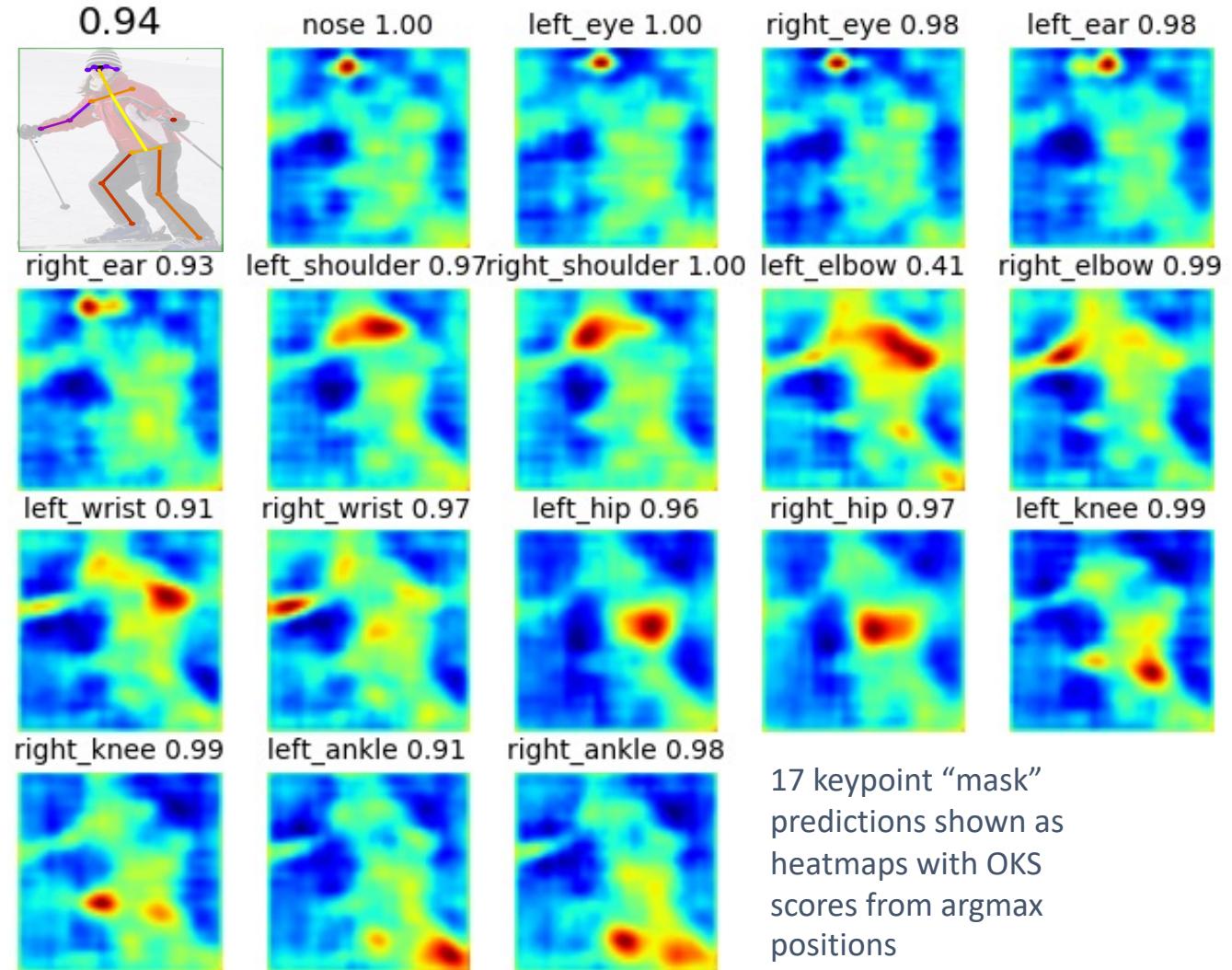
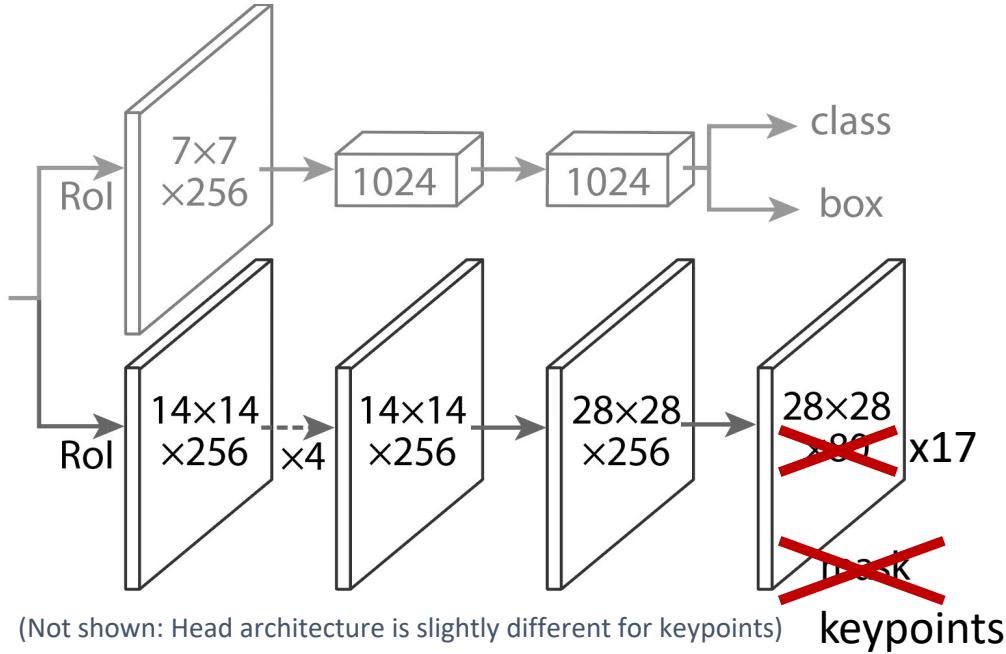


Instance Segmentation Results on COCO

	backbone	AP	AP ₅₀
MNC [7]	ResNet-101-C4	24.6	44.3
FCIS [20] +OHEM	ResNet-101-C5-dilated	29.2	49.5
FCIS+++ [20] +OHEM	ResNet-101-C5-dilated	33.6	54.5
Mask R-CNN	ResNet-101-C4	33.1	54.9
Mask R-CNN	ResNet-101-FPN	35.7	58.0

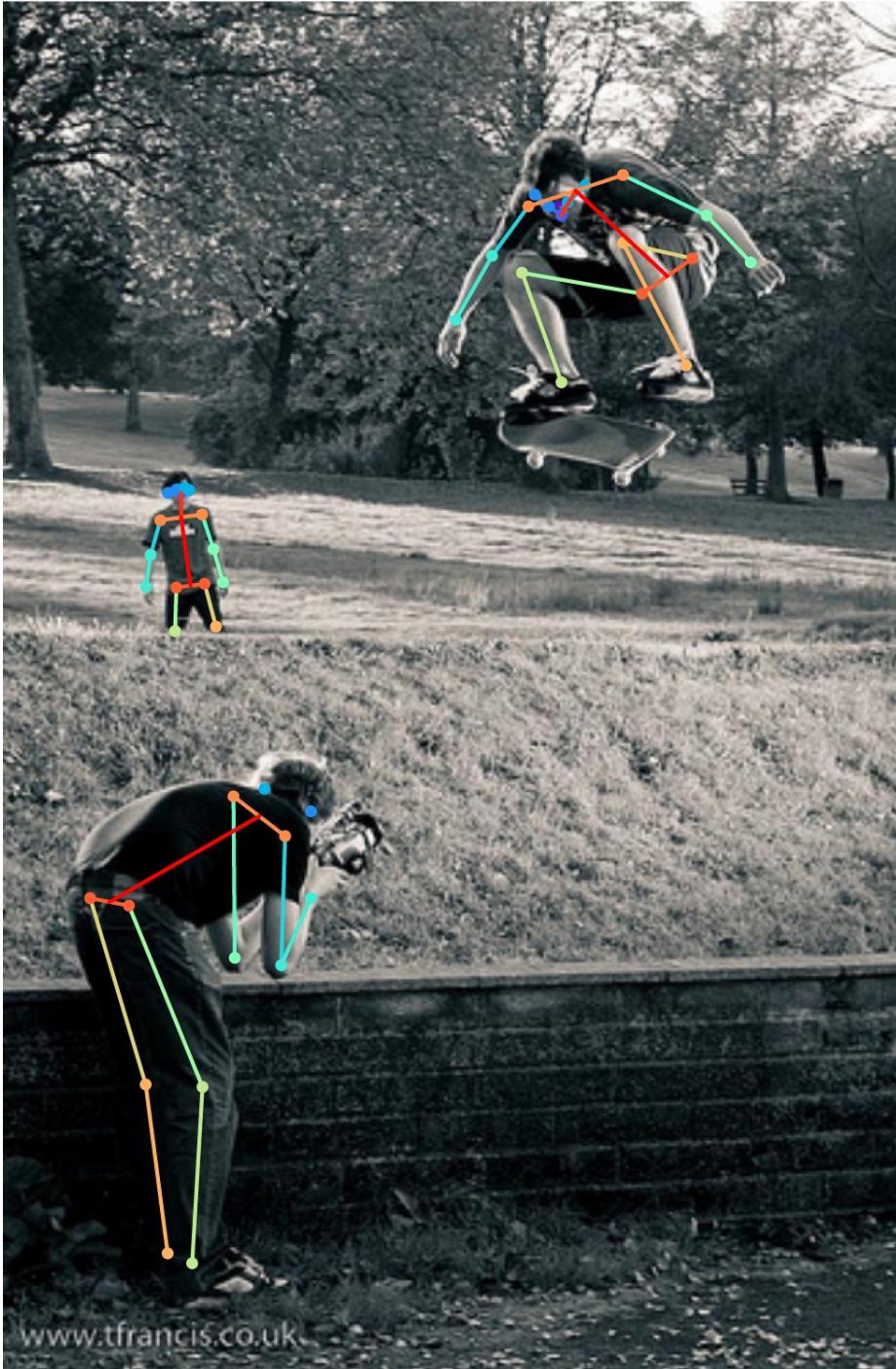
- without bells and whistles, **2 AP better** than 2016 winner
- **200ms / img**

Human Pose



17 keypoint “mask”
predictions shown as
heatmaps with OKS
scores from argmax
positions

- Add keypoint head ($28 \times 28 \times 17$)
- Predict one “mask” for each keypoint
- Softmax over spatial locations (encodes one keypoint per mask “prior”)







Overview

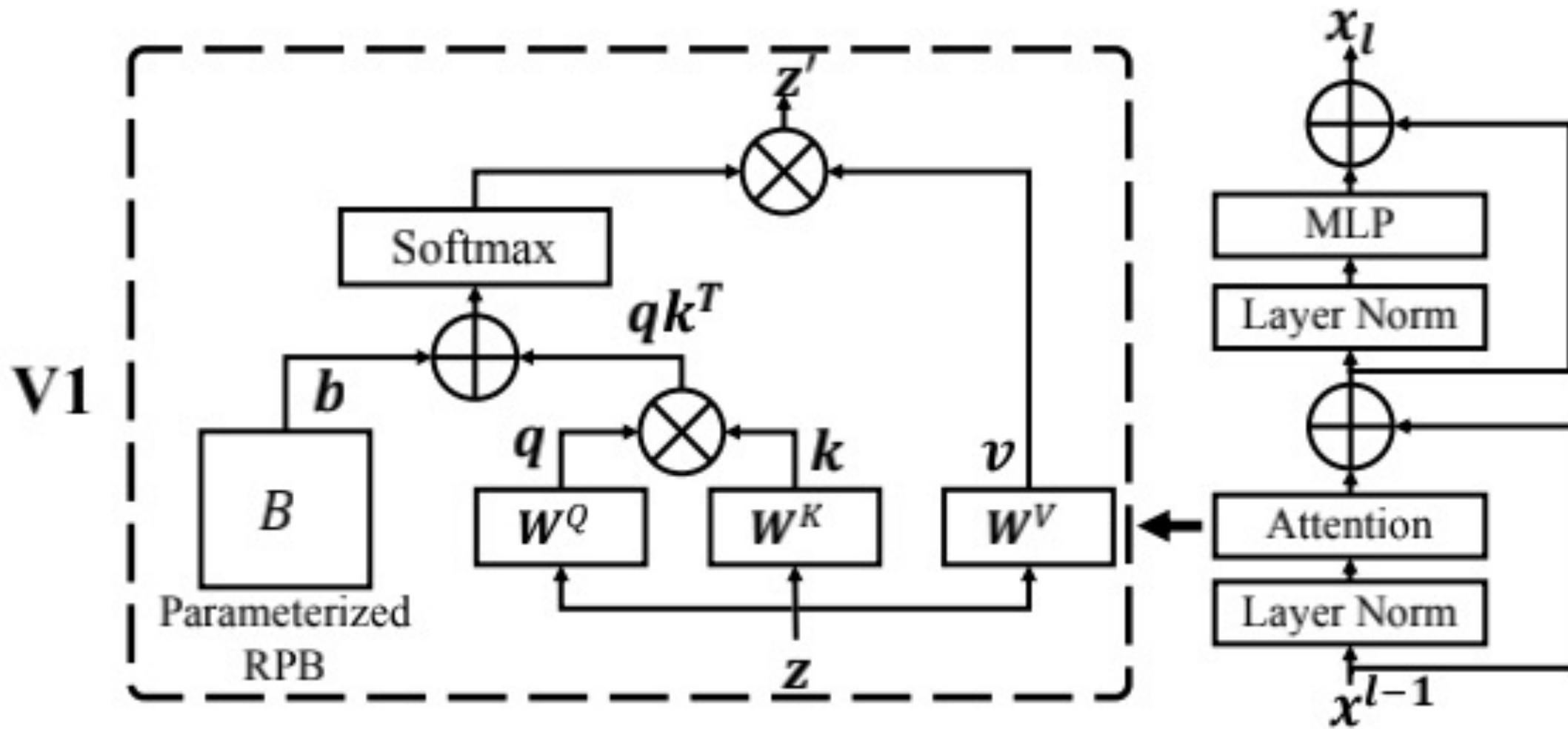
Topics to cover

- Object detection intro
- The Generalized R-CNN framework
 - **R-CNN**
 - Fast R-CNN
 - Faster R-CNN
 - **Mask R-CNN**
- Future works + Discussion
 - **SwinV2**

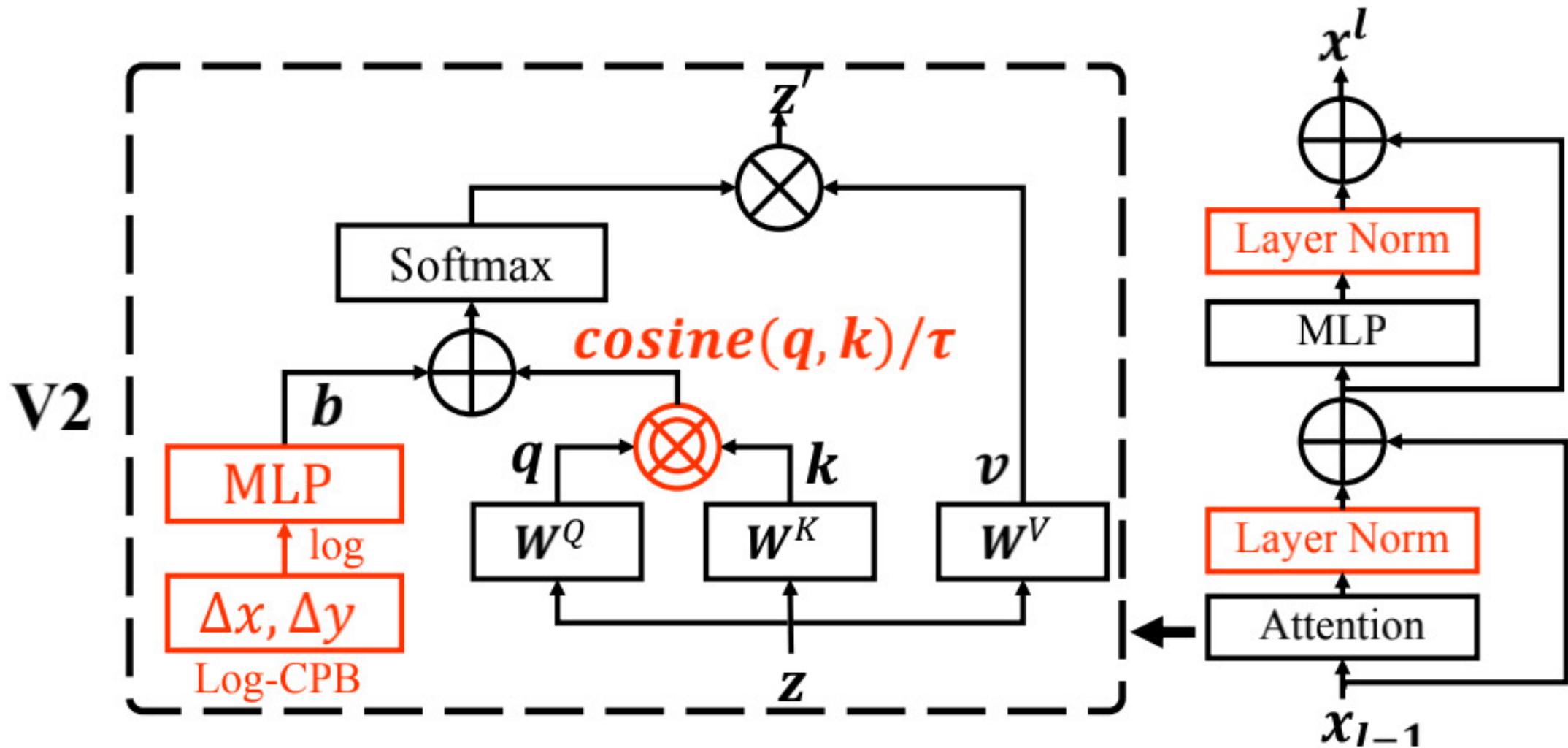
Vision Transformer



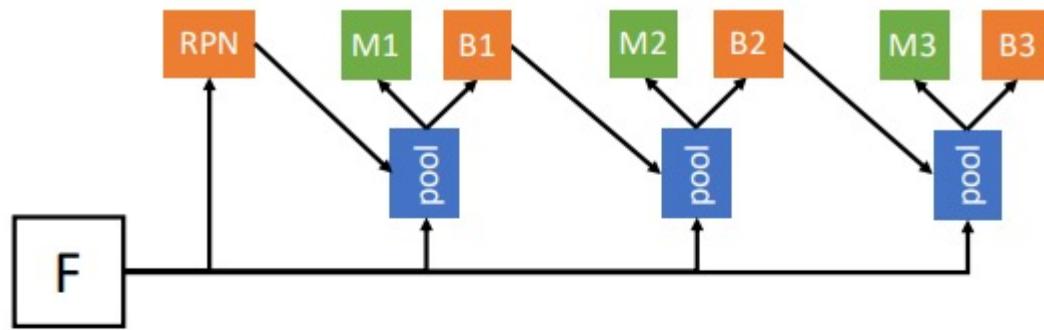
Swin



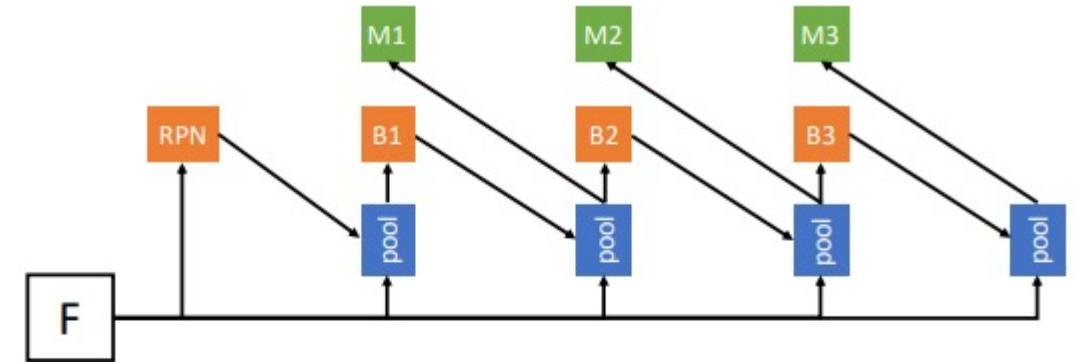
SwinV2



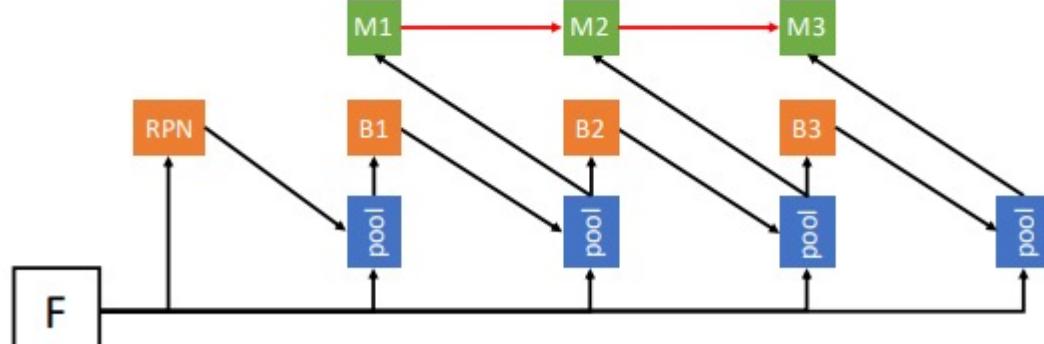
Hybrid Task Cascade



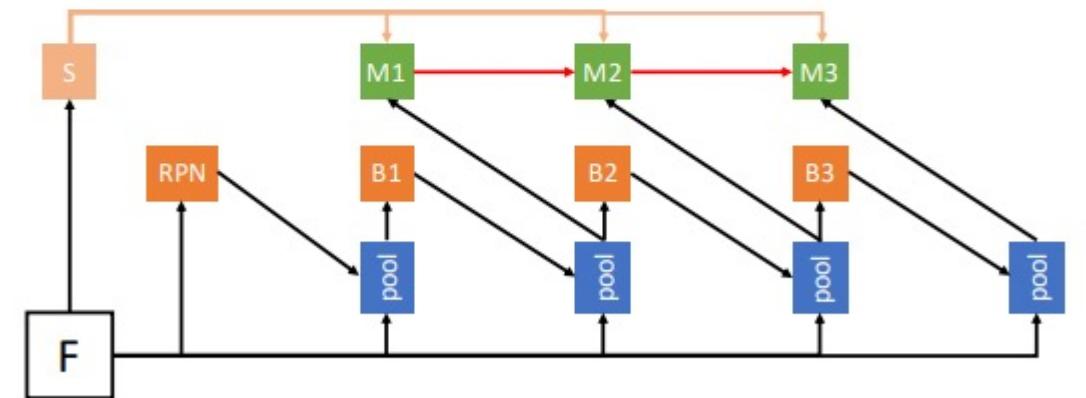
(a) Cascade Mask R-CNN



(b) Interleaved execution

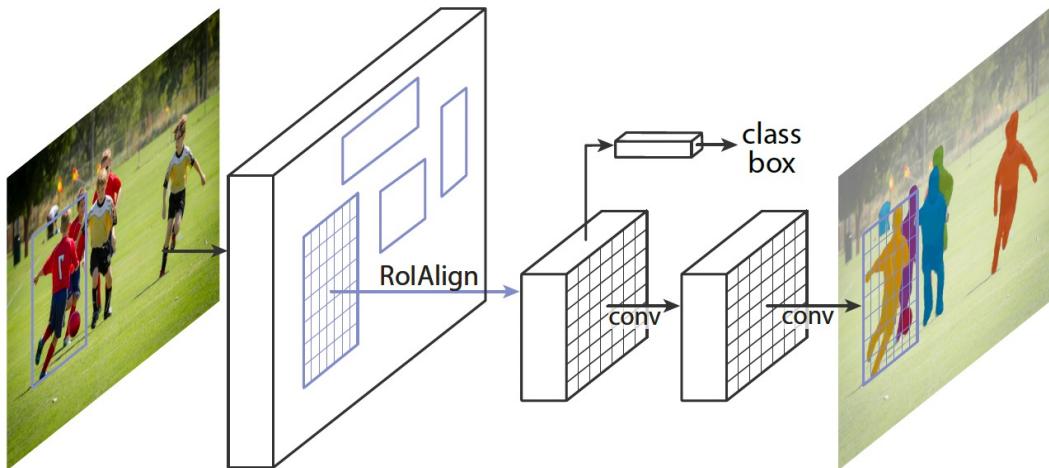


(c) Mask information flow



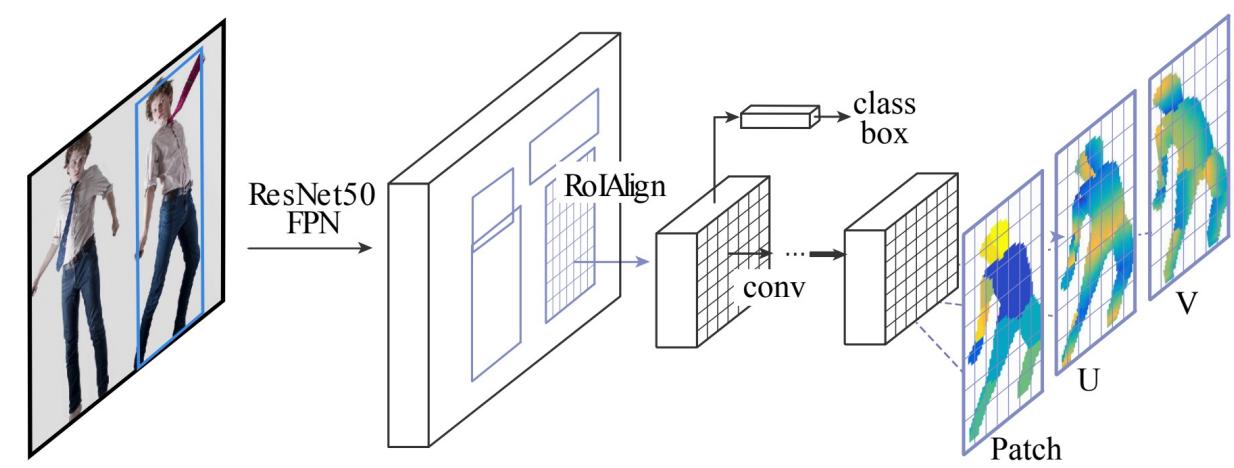
(d) Hybrid Task Cascade (semantic feature fusion with box branches is not shown on the figure for neat presentation.)

Generalized R-CNN: Adding More Heads



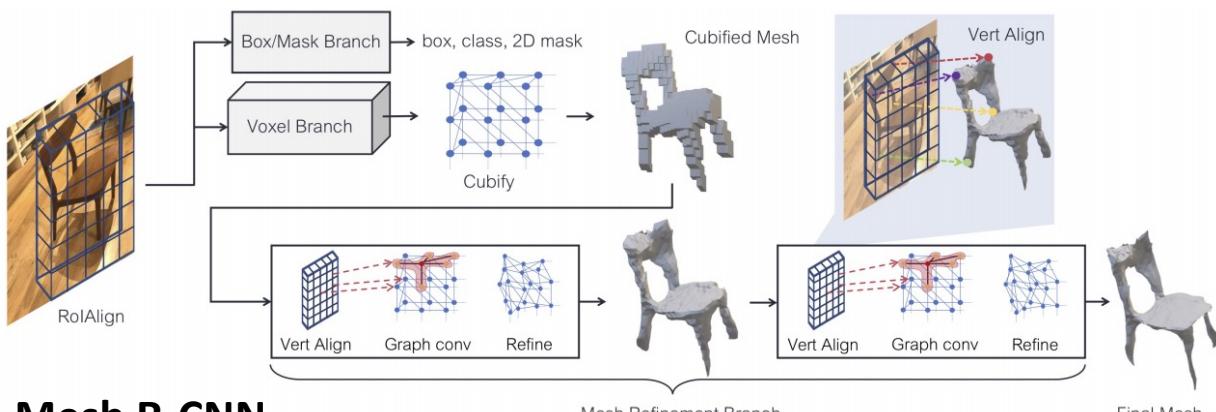
Mask R-CNN

[He, Gkioxari, Dollár, Girshick]



DensePose

[Güler, Neverova, Kokkinos]



Mesh R-CNN

[Gkioxari, Malik, Johnson. ICCV 2019]

Size

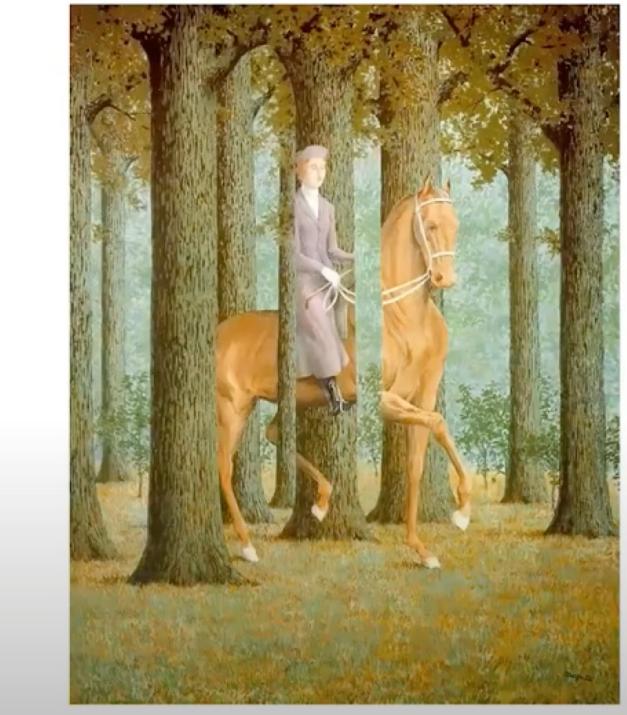


Rene Magritte, *The Listening Room*

Open Question Time

~

Interposition



Position, Probability, Size



Elephant in the Fridge

Thank You!