

# The Generalized R-CNN Framework for Object Detection

CVPR 2018 Tutorial

Ross Girshick  
Facebook AI Research (FAIR)

# Overview of this Tutorial

## Topics to cover

- Object detection intro (very brief)
- The Generalized R-CNN framework

# Overview of this Tutorial

## Topics to cover

- Object detection intro (very brief)
- The Generalized R-CNN framework

# Bounding-Box Object Detection



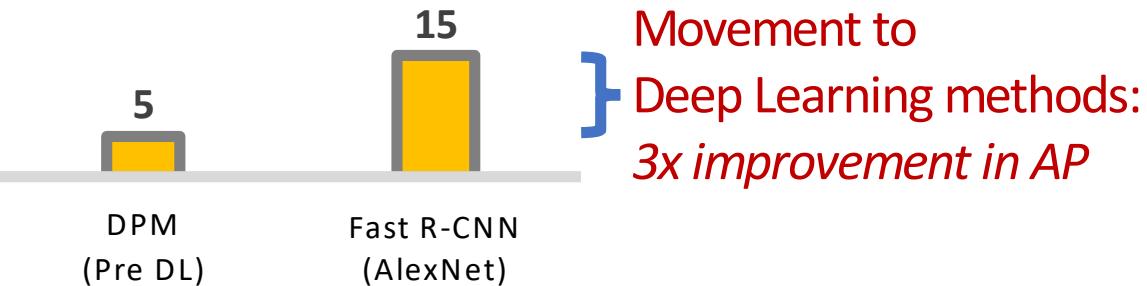
Object detection

# How Well are Detectors Working?

# COCO Object Detection Average Precision (%)

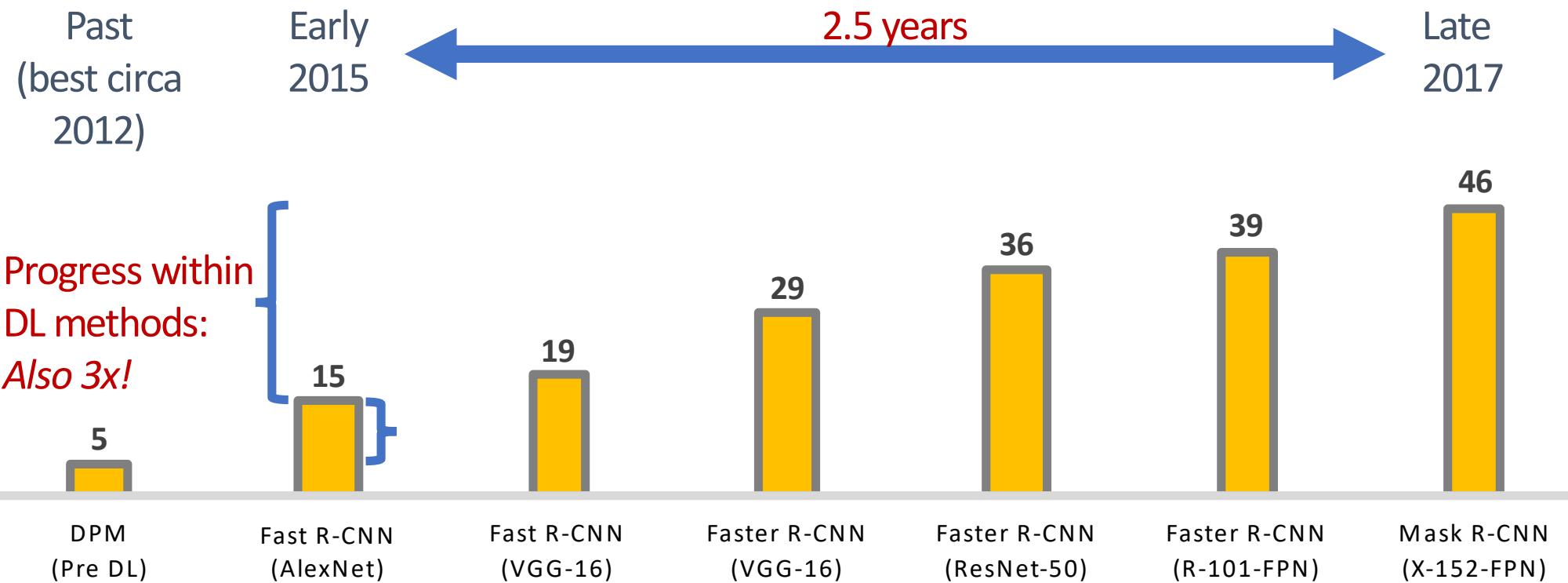
Past  
(best circa  
2012)

Early  
2015



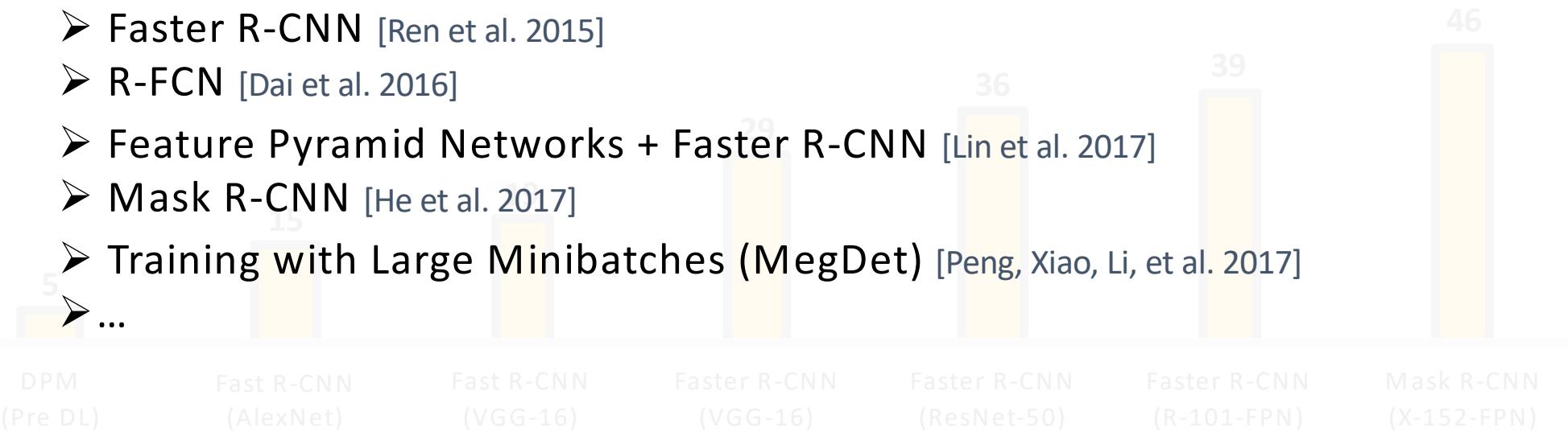
Movement to  
Deep Learning methods:  
*3x improvement in AP*

# COCO Object Detection Average Precision (%)

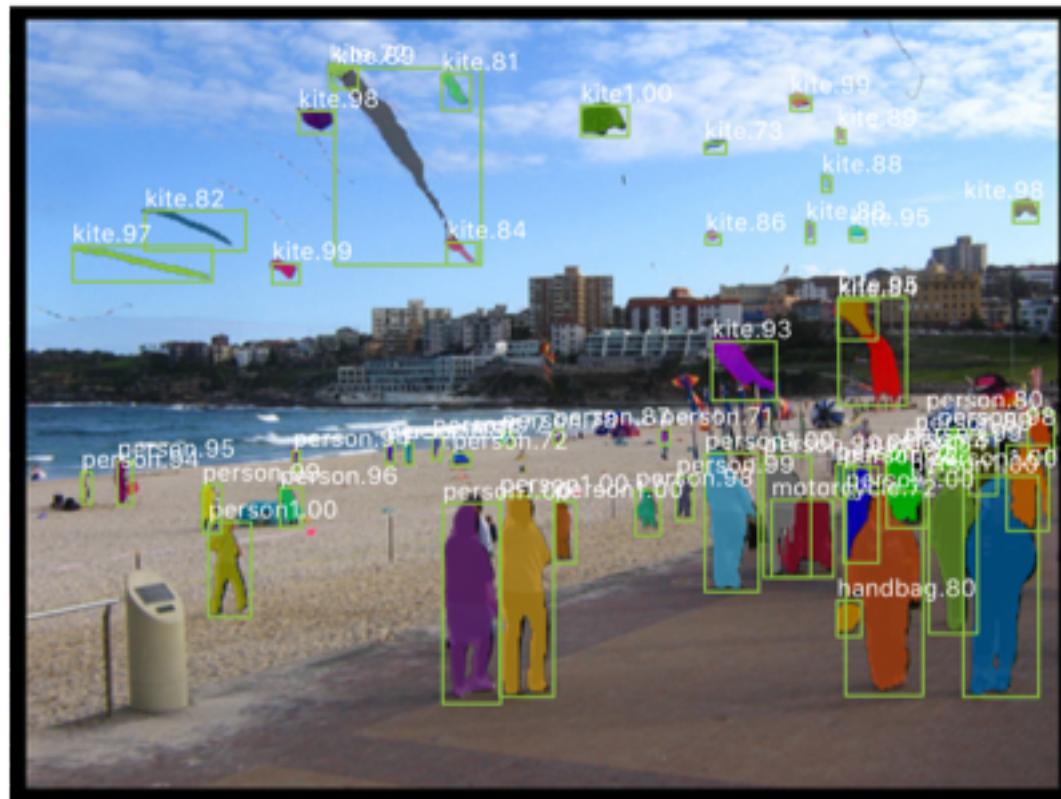


# Steady Progress

- 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]
- ...



# Detection Beyond Bounding Boxes



**Mask R-CNN**

[He, Gkioxari, Dollár, Girshick]

# Detection Beyond Bounding Boxes



**DensePose: Dense Human Pose Estimation In The Wild**  
[Güler, Neverova, Kokkinos]

# Overview of this Tutorial

## Topics to cover

- Object detection intro (very brief)
- The Generalized R-CNN framework

# R-CNN

## References

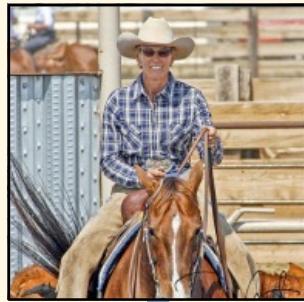
- B. Alexe, T. Deselaers, and V. Ferrari. Measuring the objectness of image windows. TPAMI, 2012.
- I. Endres and D. Hoiem. Category independent object proposals. In ECCV, 2010.
- J. Uijlings, K. van de Sande, T. Gevers, and A. Smeulders. Selective search for object recognition. IJCV, 2013.
- A. Krizhevsky, I. Sutskever, and G. Hinton. ImageNet classification with deep convolutional neural networks. In NIPS, 2012.
- X. Wang, M. Yang, S. Zhu, and Y. Lin. Regionlets for generic object detection. In ICCV, 2013.
- R. Girshick, J. Donahue, T. Darrell, and J. Malik. Rich feature hierarchies for accurate object detection and semantic segmentation. In CVPR, 2014.

Girshick, Donahue, Darrell, Malik. Rich Feature Hierarchies for Accurate Object Detection and Semantic Segmentation. CVPR 2014.

# R-CNN (Region-based Convolutional Neural Net)

Per-image computation

$I:$



Selective search,  
Edge Boxes,  
MCG, ...

1



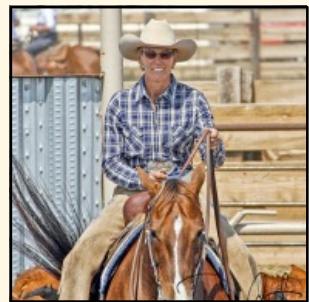
Use an off-the-shelf region/object/detection proposal algorithm (~2k proposals per image)

Girshick, Donahue, Darrell, Malik. Rich Feature Hierarchies for Accurate Object Detection and Semantic Segmentation. CVPR 2014.

# R-CNN

Per-image computation

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



$I:$

Selective search,  
Edge Boxes,  
MCG, ...

1



Crop &  
warp

2



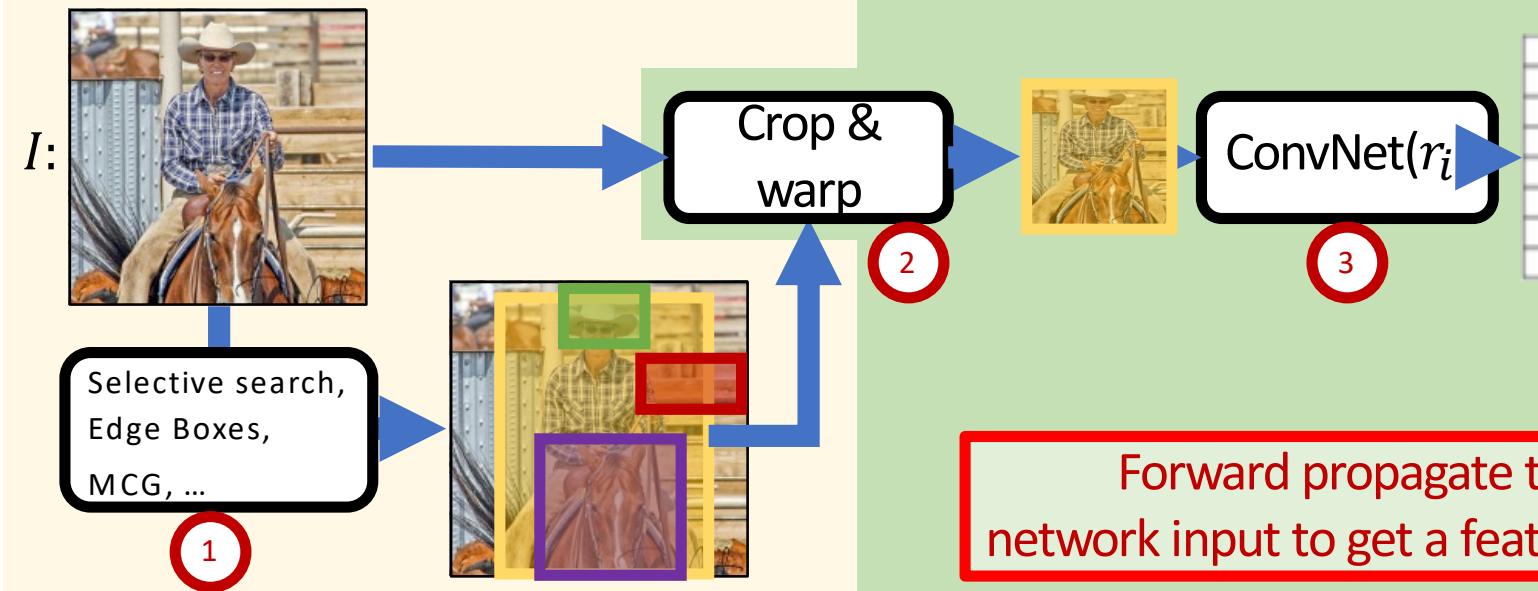
Crop and warp each proposal image window  
to obtain a fixed-size network input

Girshick, Donahue, Darrell, Malik. Rich Feature Hierarchies for Accurate Object Detection and Semantic Segmentation. CVPR 2014.

# R-CNN

Per-image computation

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

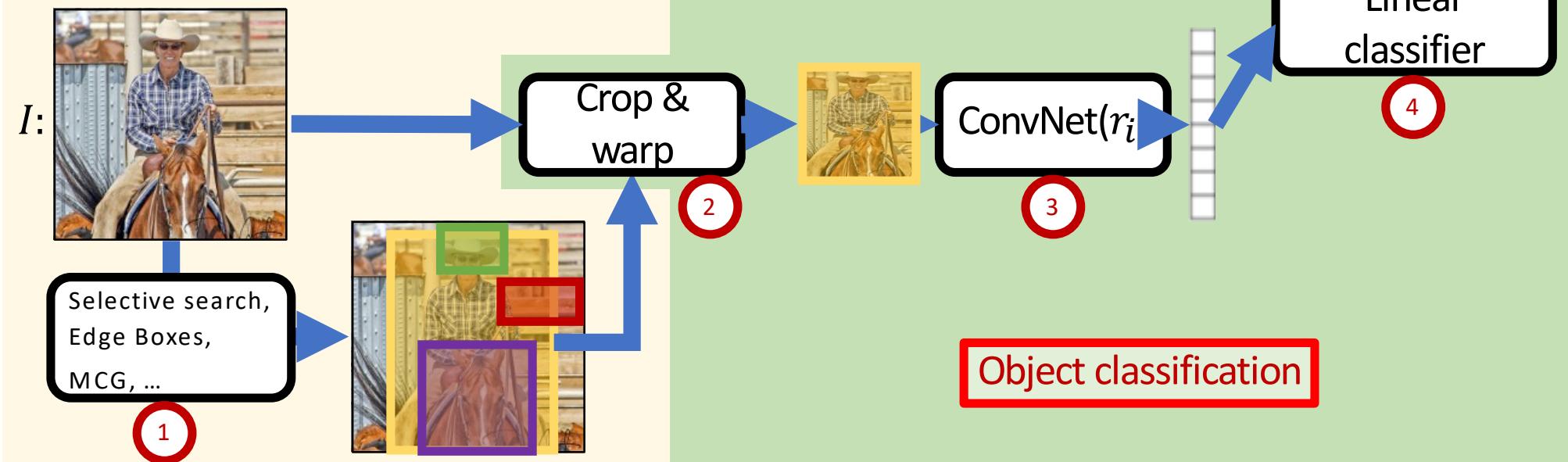


Girshick, Donahue, Darrell, Malik. Rich Feature Hierarchies for Accurate Object Detection and Semantic Segmentation. CVPR 2014.

# R-CNN

Per-image computation

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

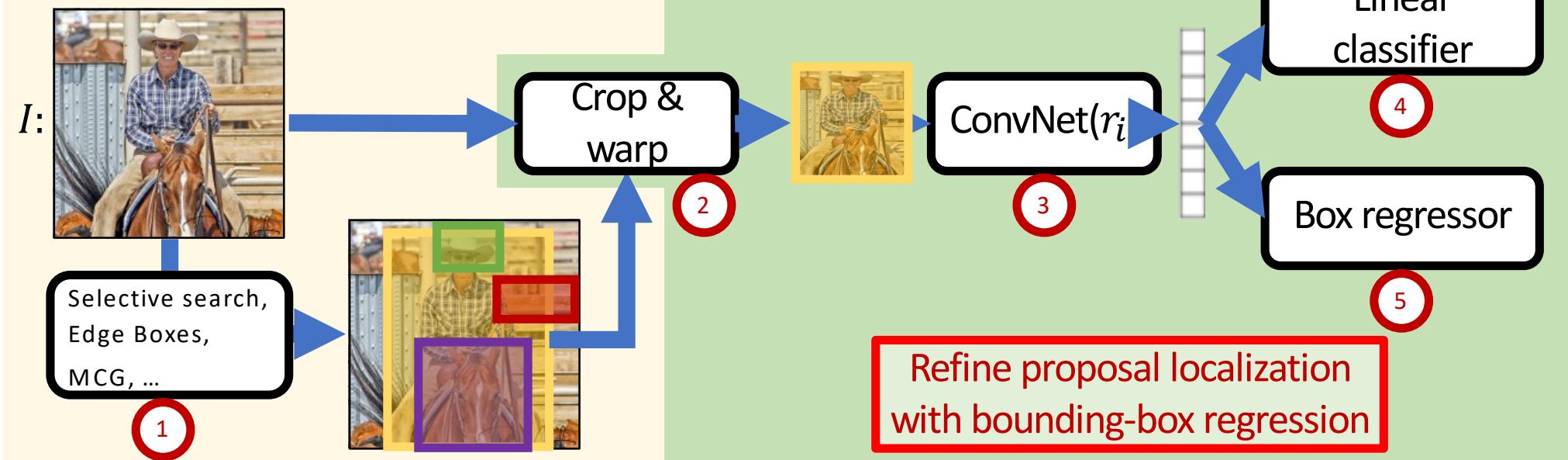


Girshick, Donahue, Darrell, Malik. Rich Feature Hierarchies for Accurate Object Detection and Semantic Segmentation. CVPR 2014.

# R-CNN

Per-image computation

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



# Generalized R-CNN Framework

A common framework for understanding

- R-CNN
- Fast R-CNN
- Faster R-CNN
- Feature Pyramid Networks (FPN) + Faster R-CNN
- Mask R-CNN
- ... and more (e.g., DensePose)

<https://github.com/facebookresearch/Detectron>

# Generalized R-CNN Framework

Per-image computation



$I:$

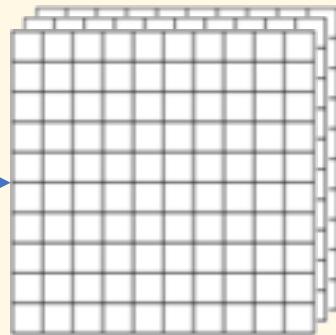
Per-region computation for each  $r_i \in r(I)$

Input image  
per-image operations | per-region operations

# Generalized R-CNN Framework

Per-image computation

$$f_I = f(I)$$

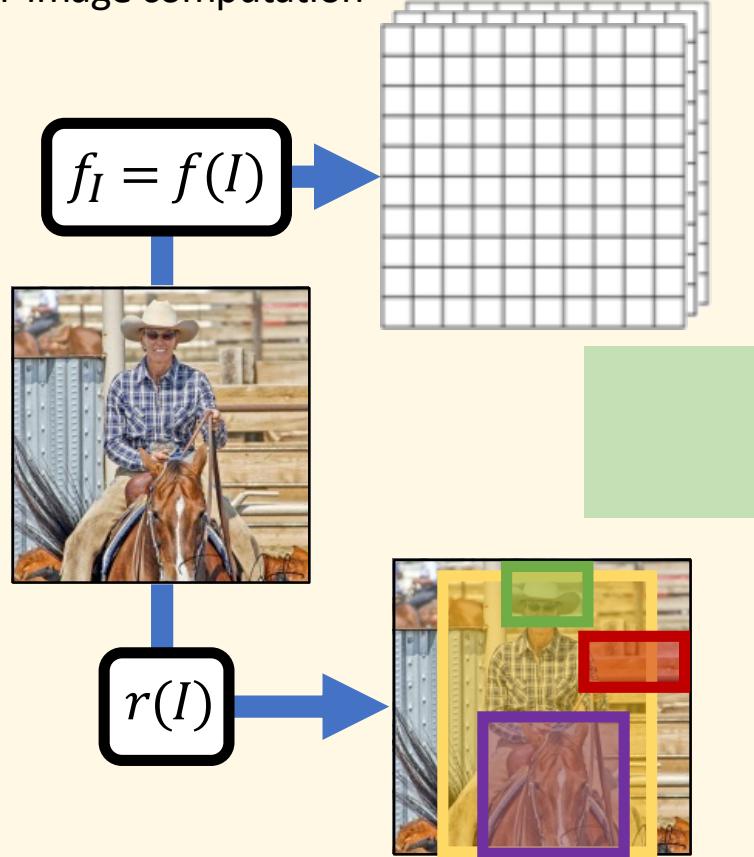


Per-region computation for each  $r_i \in r(I)$

Transformation of the input image  
into a featurized representation

# Generalized R-CNN Framework

Per-image computation

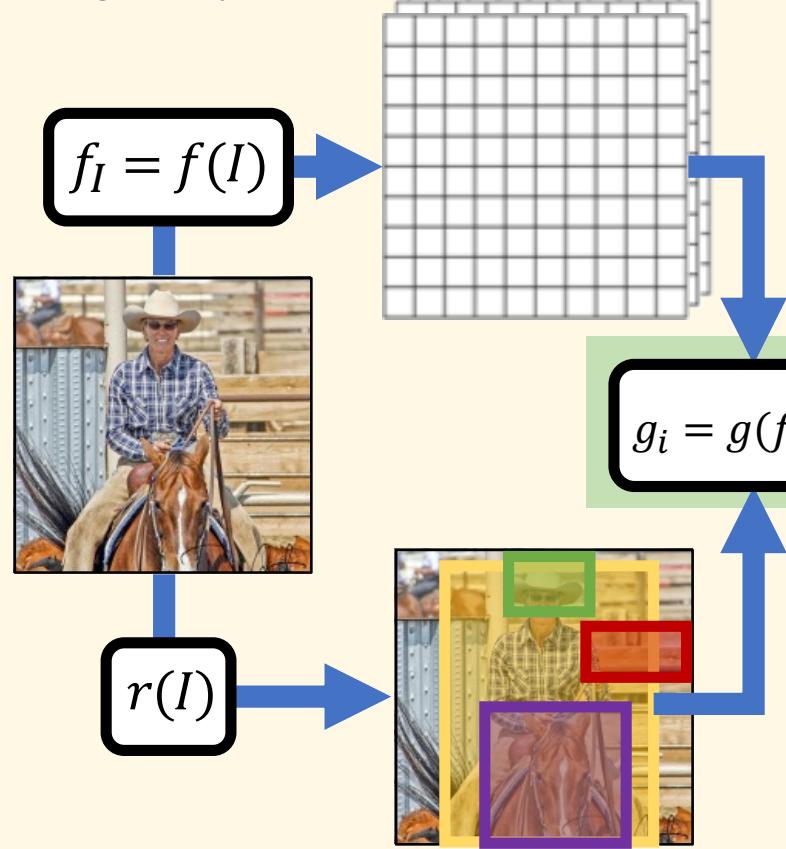


Per-region computation for each  $r_i \in r(I)$

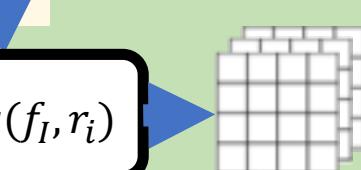
Object/region/detection proposals  
computed for the image

# Generalized R-CNN Framework

Per-image computation

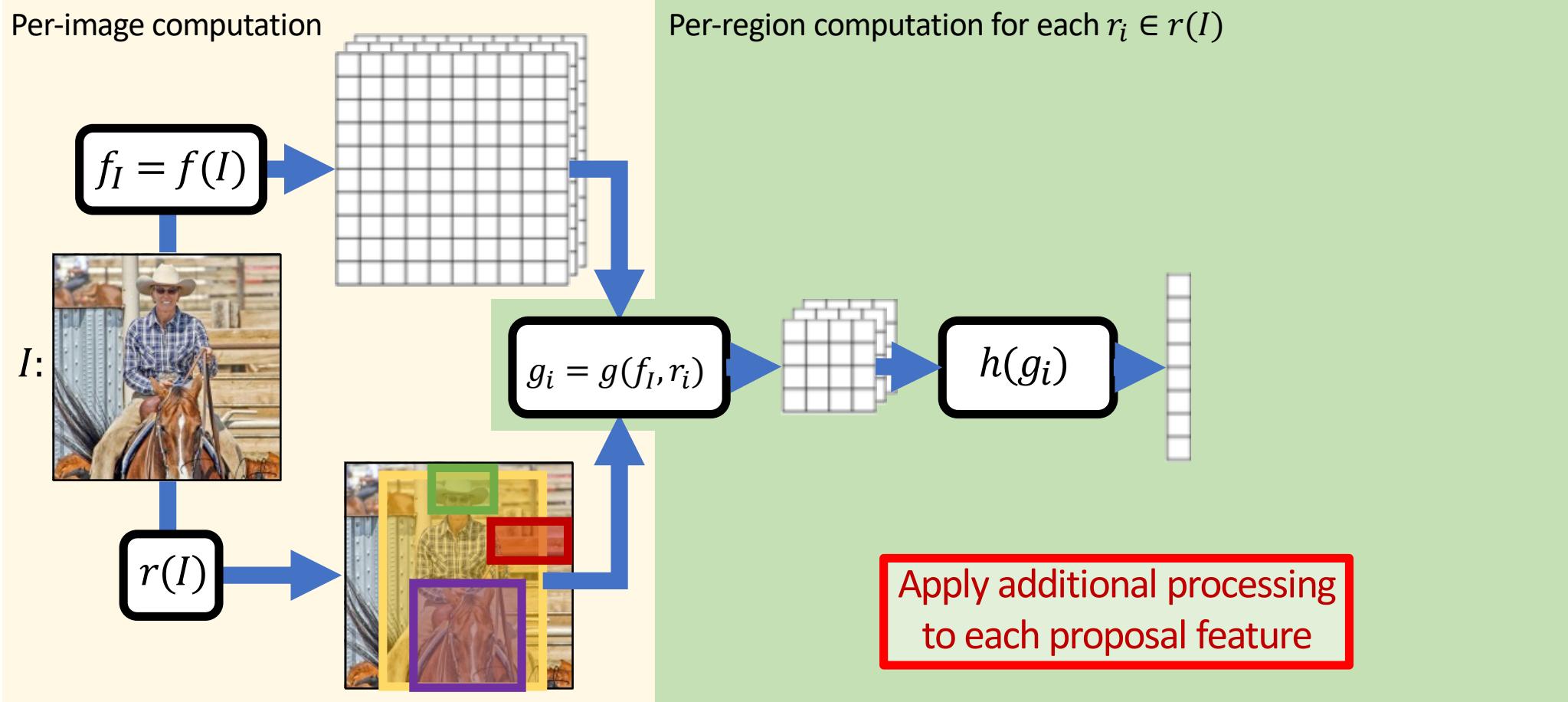


Per-region computation for each  $r_i \in r(I)$

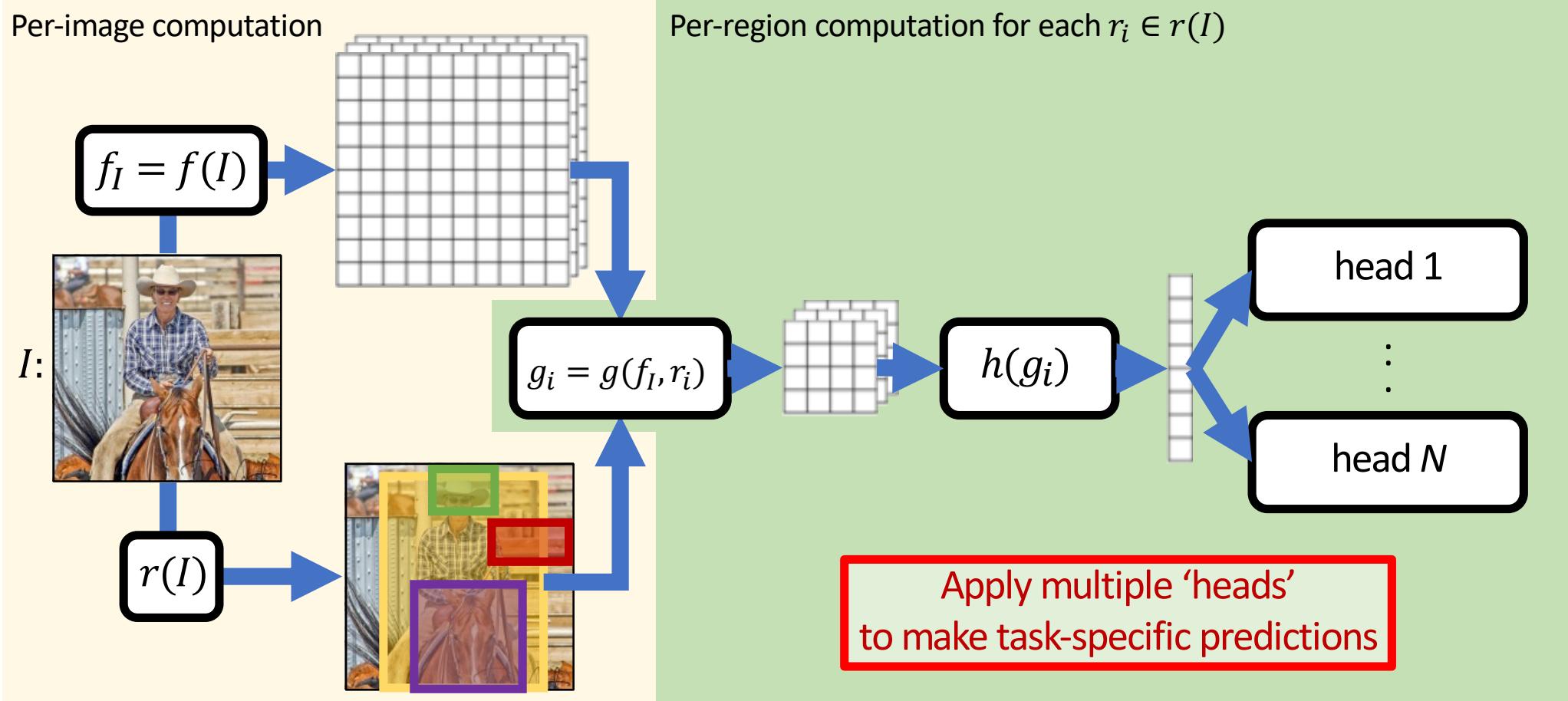


Compute a featurized representation of each proposal

# Generalized R-CNN Framework



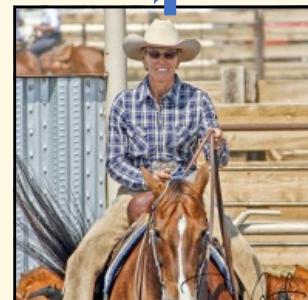
# Generalized R-CNN Framework



# R-CNN in the Generalized Framework

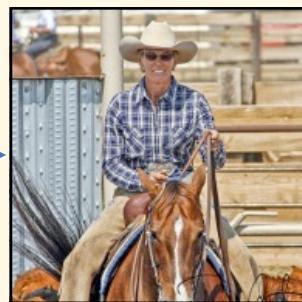
Per-image computation

$$f_I = I$$



$I:$

Selective search,  
Edge Boxes,  
MCG, ...



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

Crop & warp



ConvNet( $r_i$ )



Linear classifier

Box regressor

R-CNN in the generalized framework

# The Problem with R-CNN

# “Slow” R-CNN

Per-image computation

$$f_I = I$$



$I:$

Selective search,  
Edge Boxes,  
MCG, ...



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

Crop & warp



ConvNet( $r_i$ )



Linear classifier

Box regressor

Very heavy *per-region* computation  
E.g., 2000 full network evaluations

# Fast R-CNN

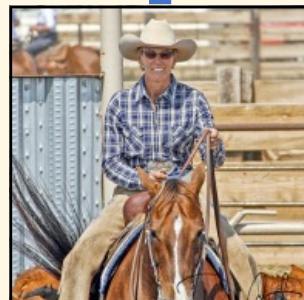
## References

- 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.

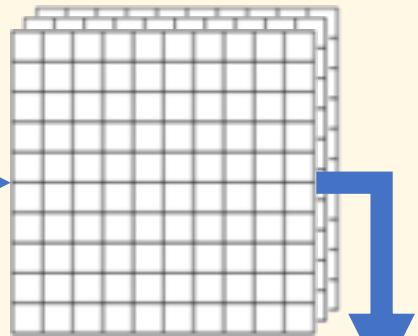
# Generalized R-CNN → Fast R-CNN

Per-image computation

$$f_I = f(I)$$

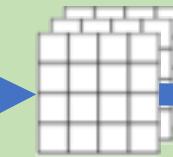


Selective search,  
Edge Boxes,  
MCG, ...



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

$$g_i = g(f_I, r_i)$$



$$h(g_i)$$

head 1

⋮

head  $N$

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

# Generalized R-CNN → Fast R-CNN

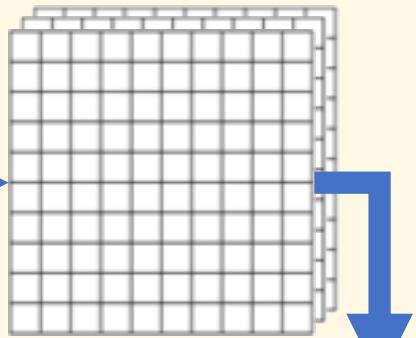
Per-image computation



$I$ :

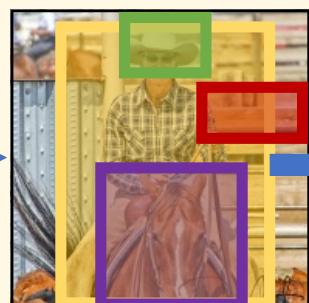


Selective search,  
Edge Boxes,  
MCG, ...

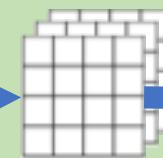


Per-region computation for each  $r_i \in r(I)$

$$g_i = g(f_I, r_i)$$



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



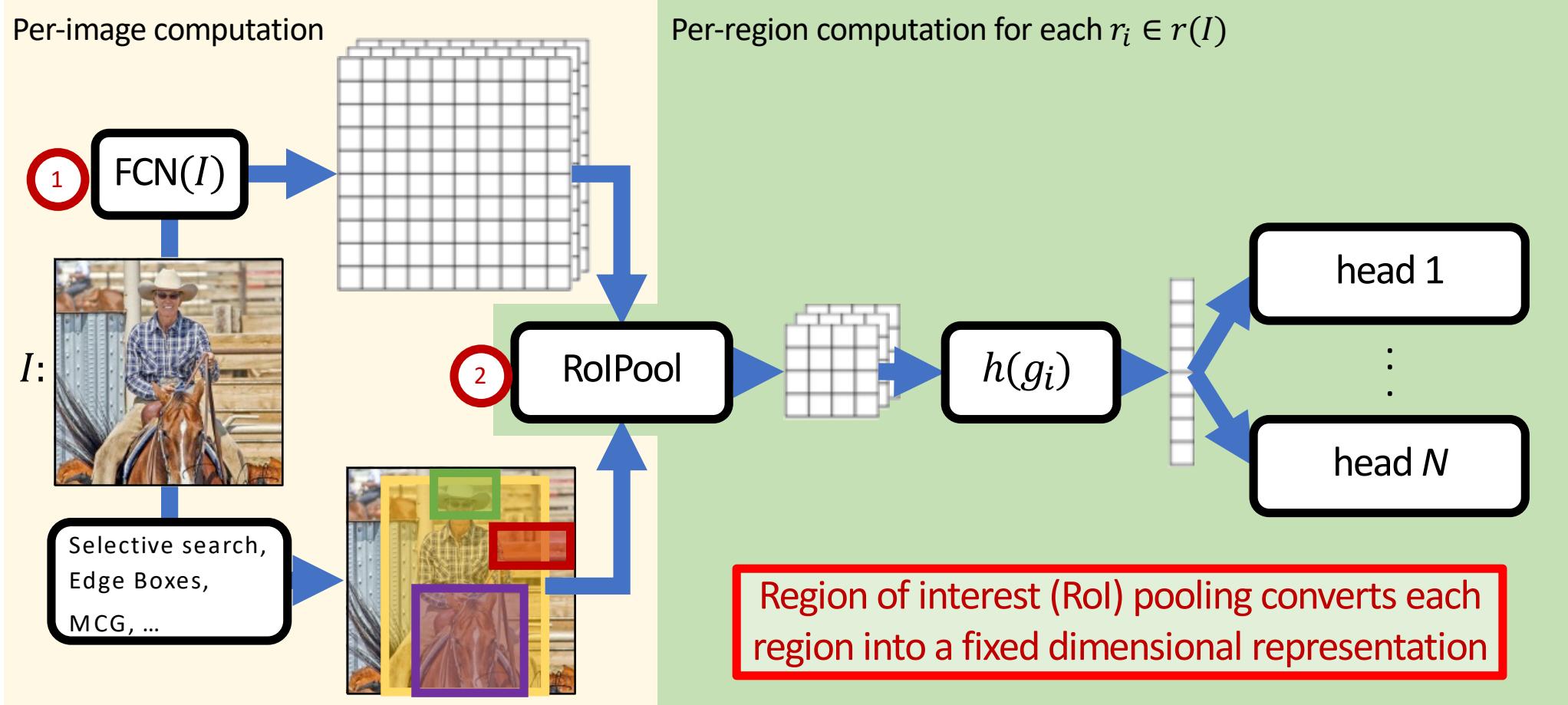
$$h(g_i)$$

head 1

⋮

head  $N$

# Generalized R-CNN → Fast R-CNN

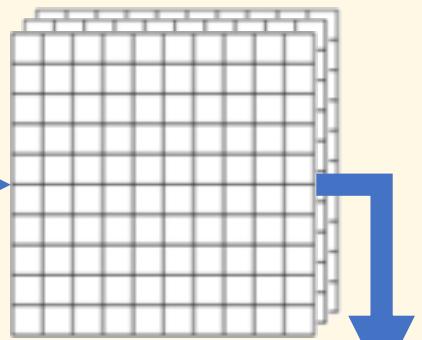


# Generalized R-CNN → Fast R-CNN

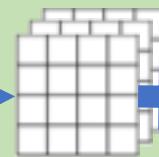
Per-image computation



Selective search,  
Edge Boxes,  
MCG, ...



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



MLP

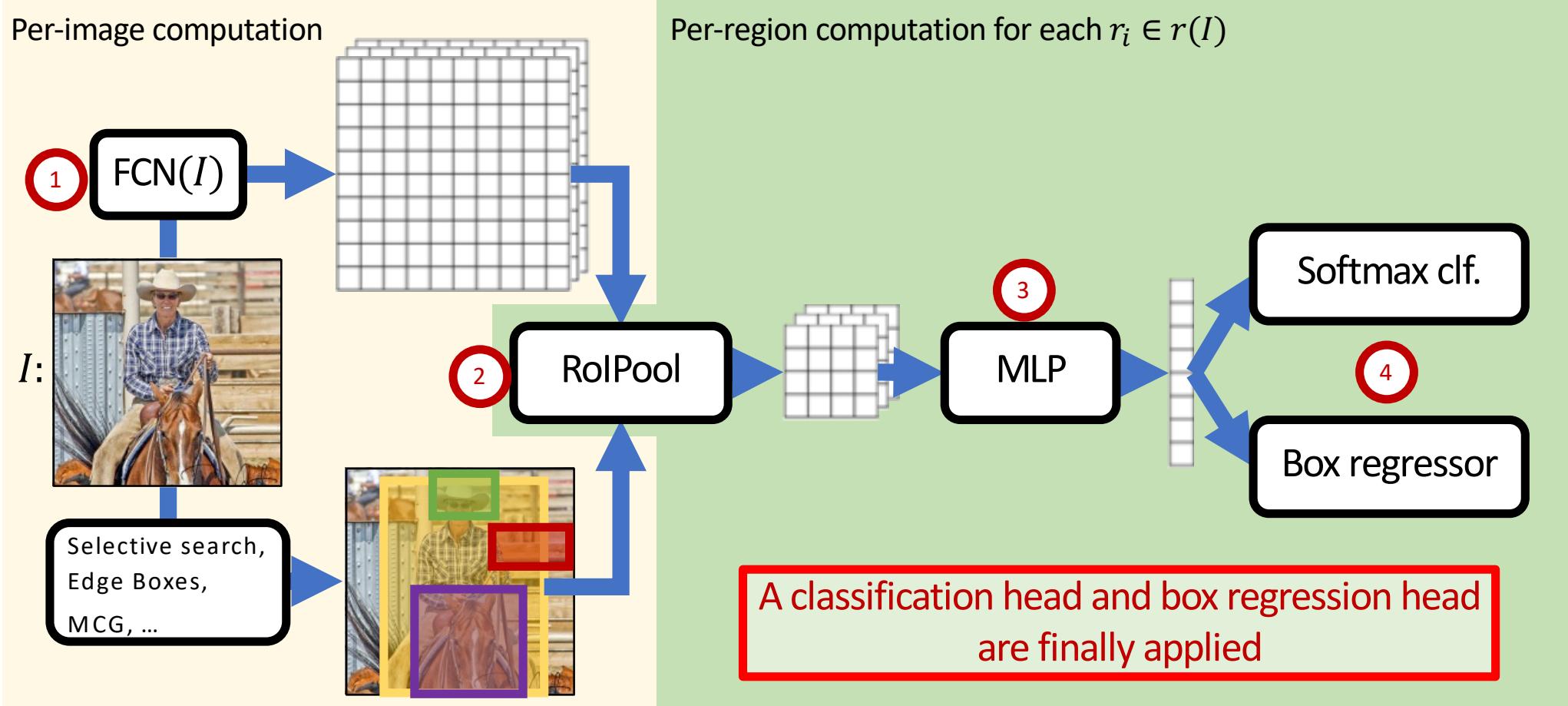
head 1

:

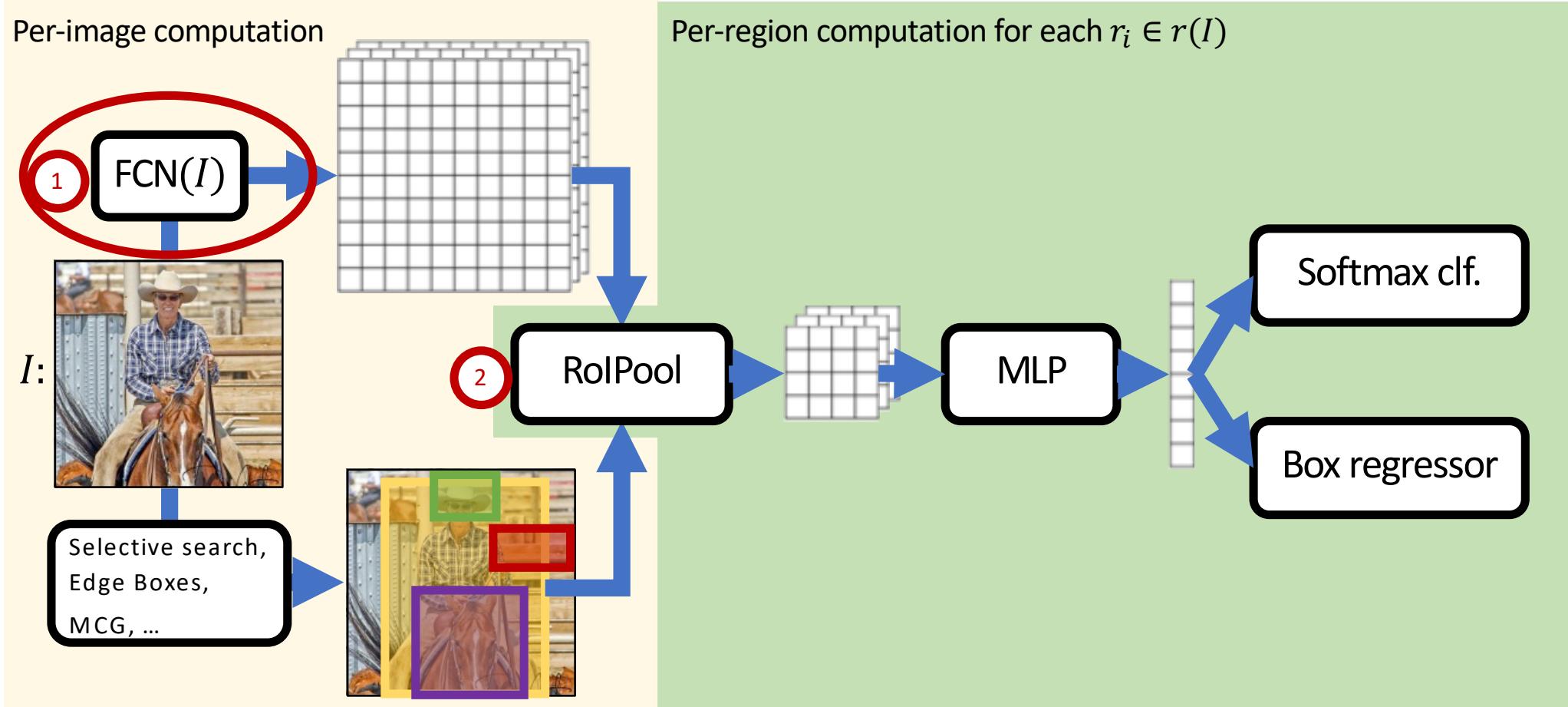
head  $N$

A lightweight MLP processes each region feature map

# Generalized R-CNN → Fast R-CNN



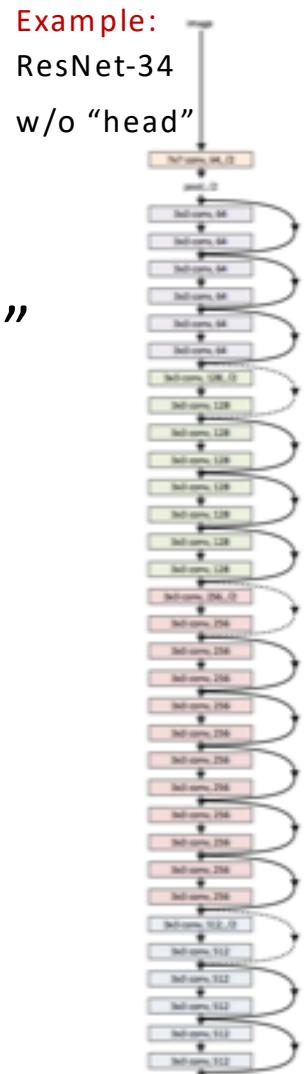
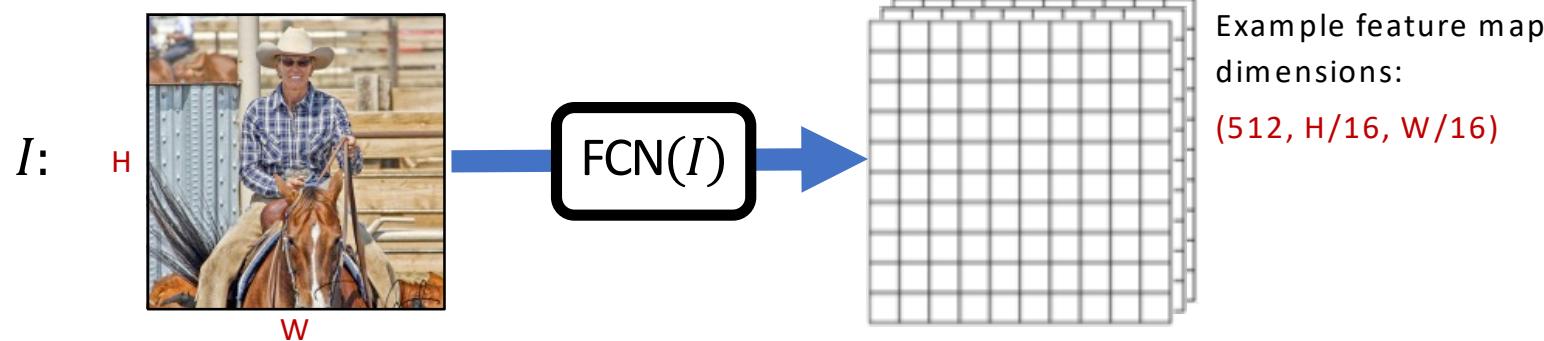
# Fast R-CNN



# Whole-image, Fully Conv. Network (FCN)

Use **any standard ConvNet** as the “backbone architecture”

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

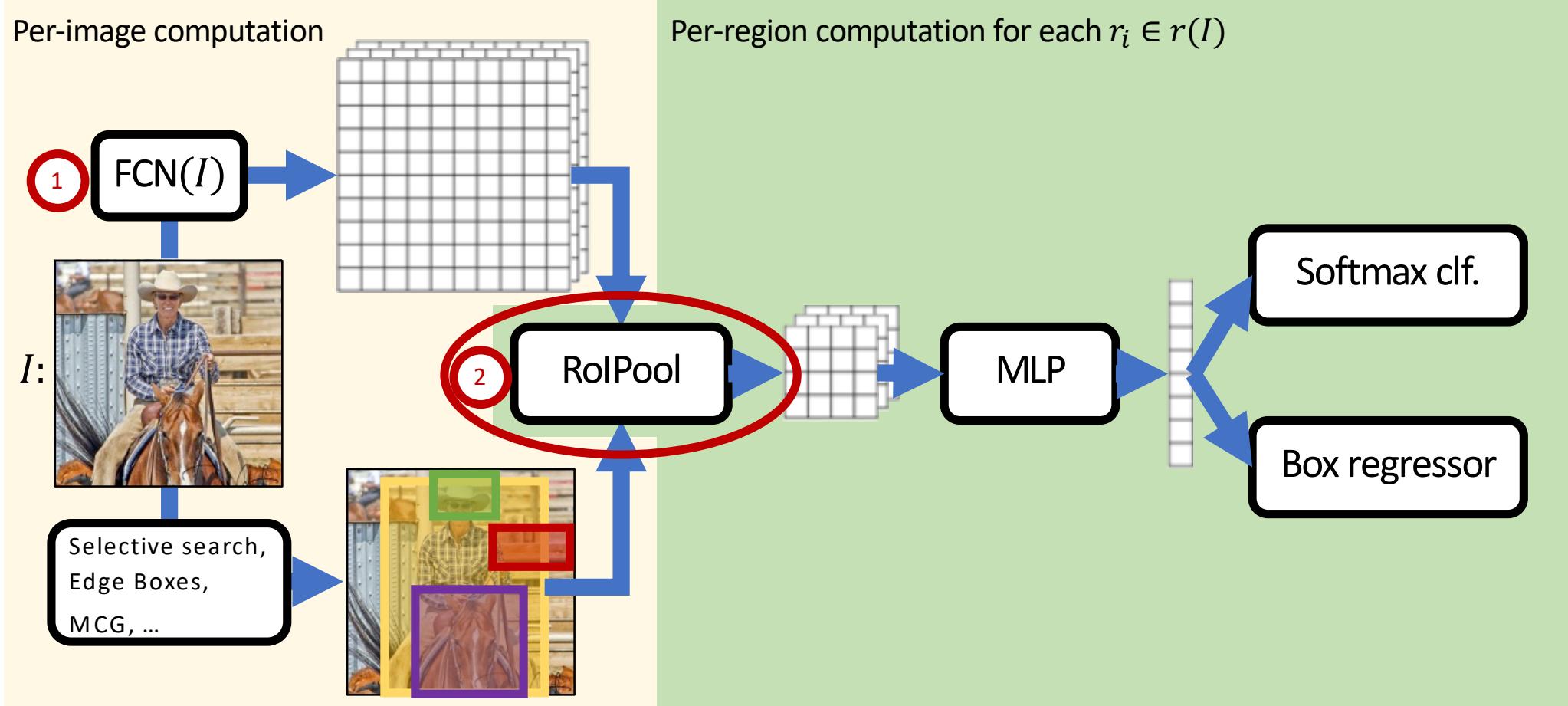


# The Engine of Recognition

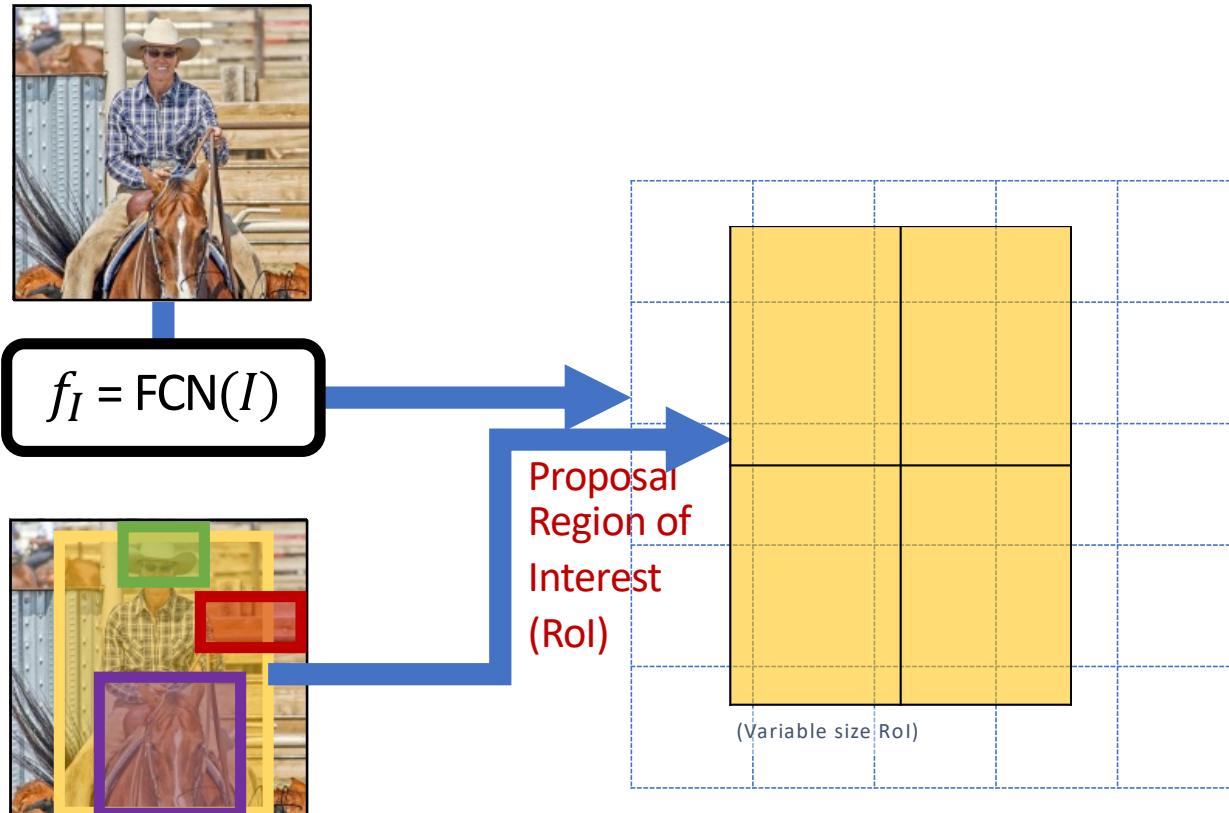
A good network lift all boats (“features matter”)

- AlexNet
- VGG
- GoogleNet / Inception
- ResNet
- ResNeXt
- ...

# Fast R-CNN

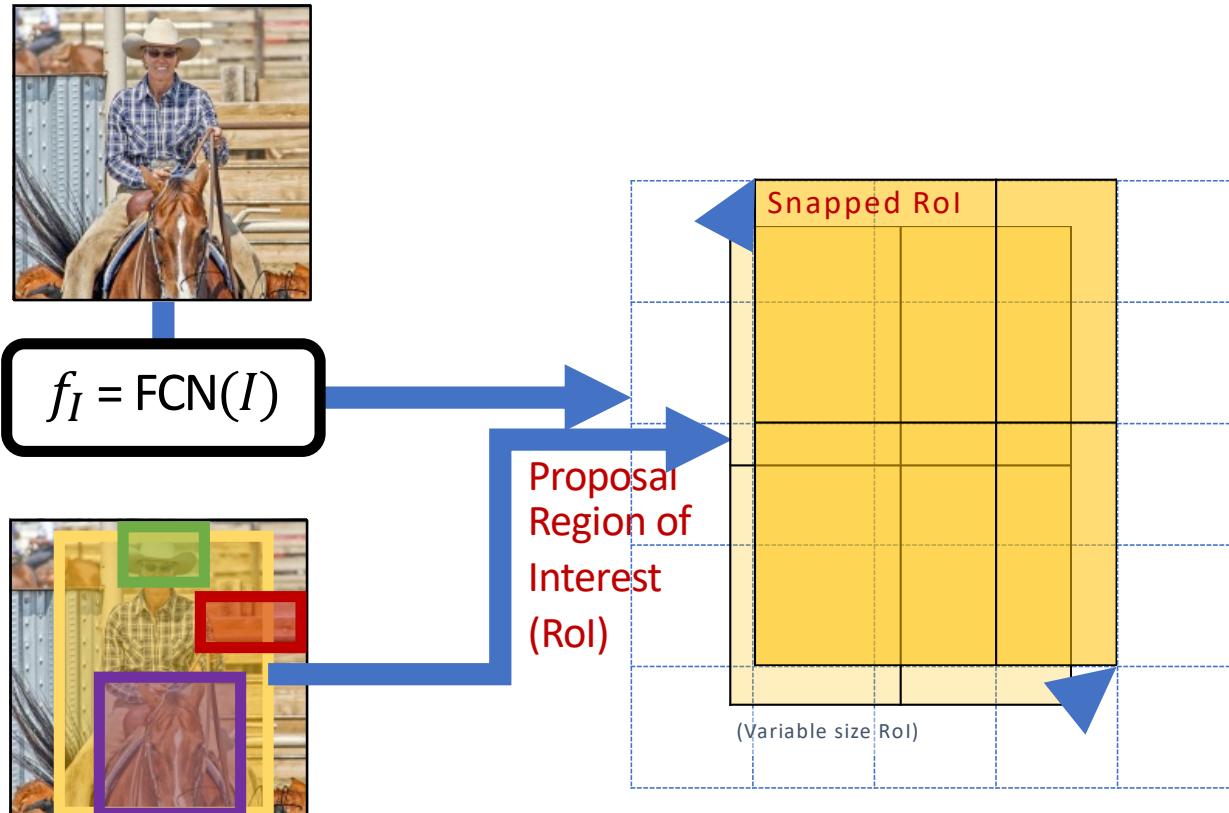


# RoIPool (on each Proposal)



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

# RoIPool (on each Proposal)



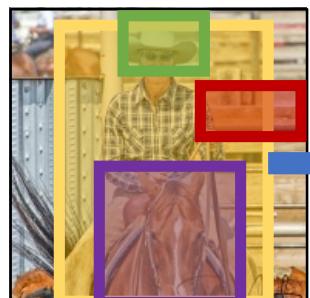
Key innovation in SPP-net  
[He et al. 2014]

# RoIPool (on each Proposal)

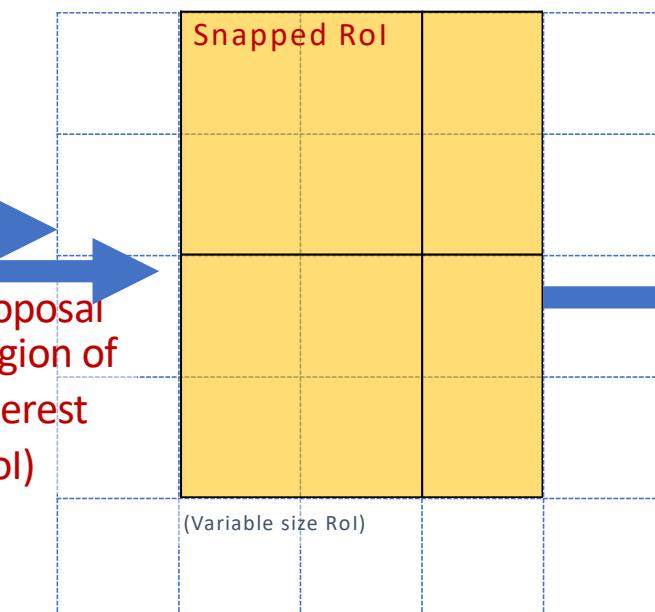
Transform **arbitrary size proposal** into a **fixed-dimensional representation** (e.g., 2x2)



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



Proposal  
Region of  
Interest  
(RoI)



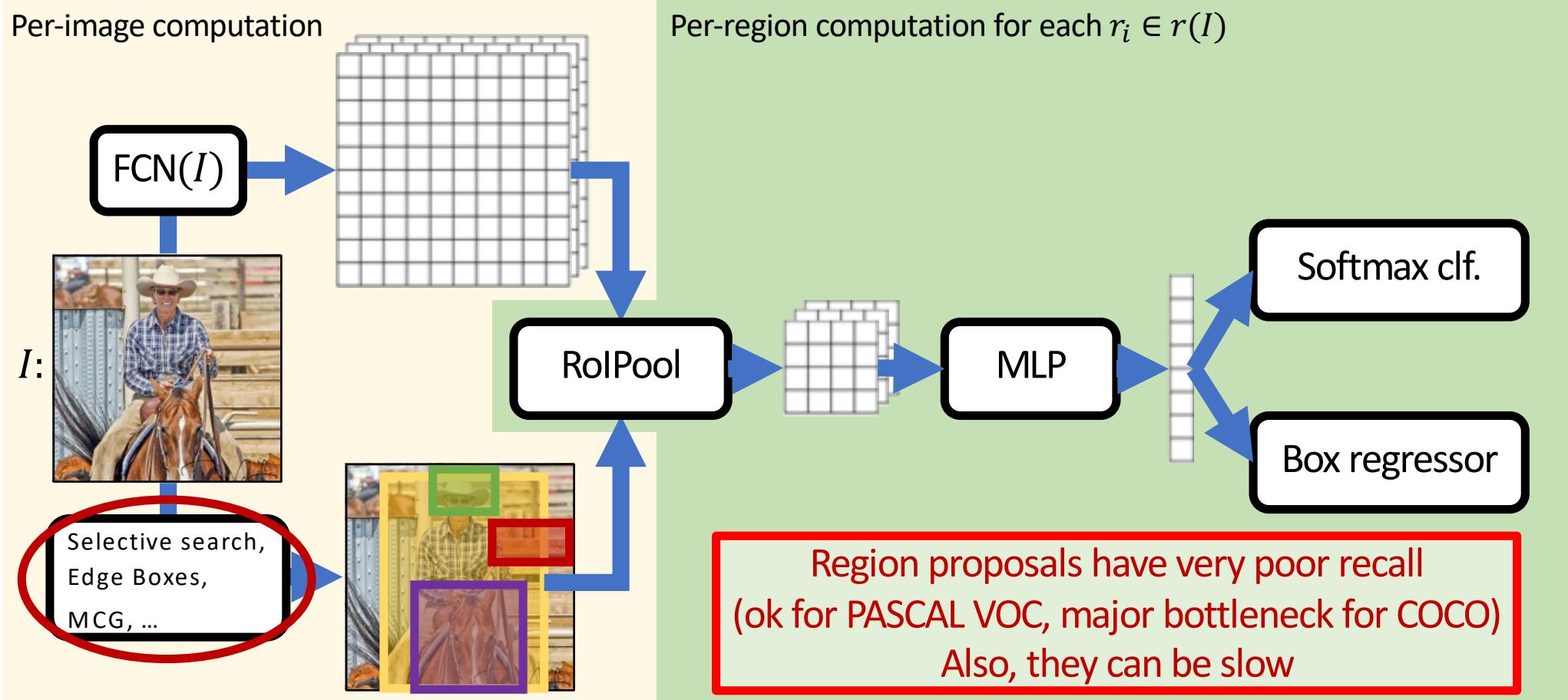
(Fixed dimensional  
representation)

MLP

Feature value  
is **max** over input  
cells

# The Problem with Fast R-CNN

# The Problem with Fast R-CNN



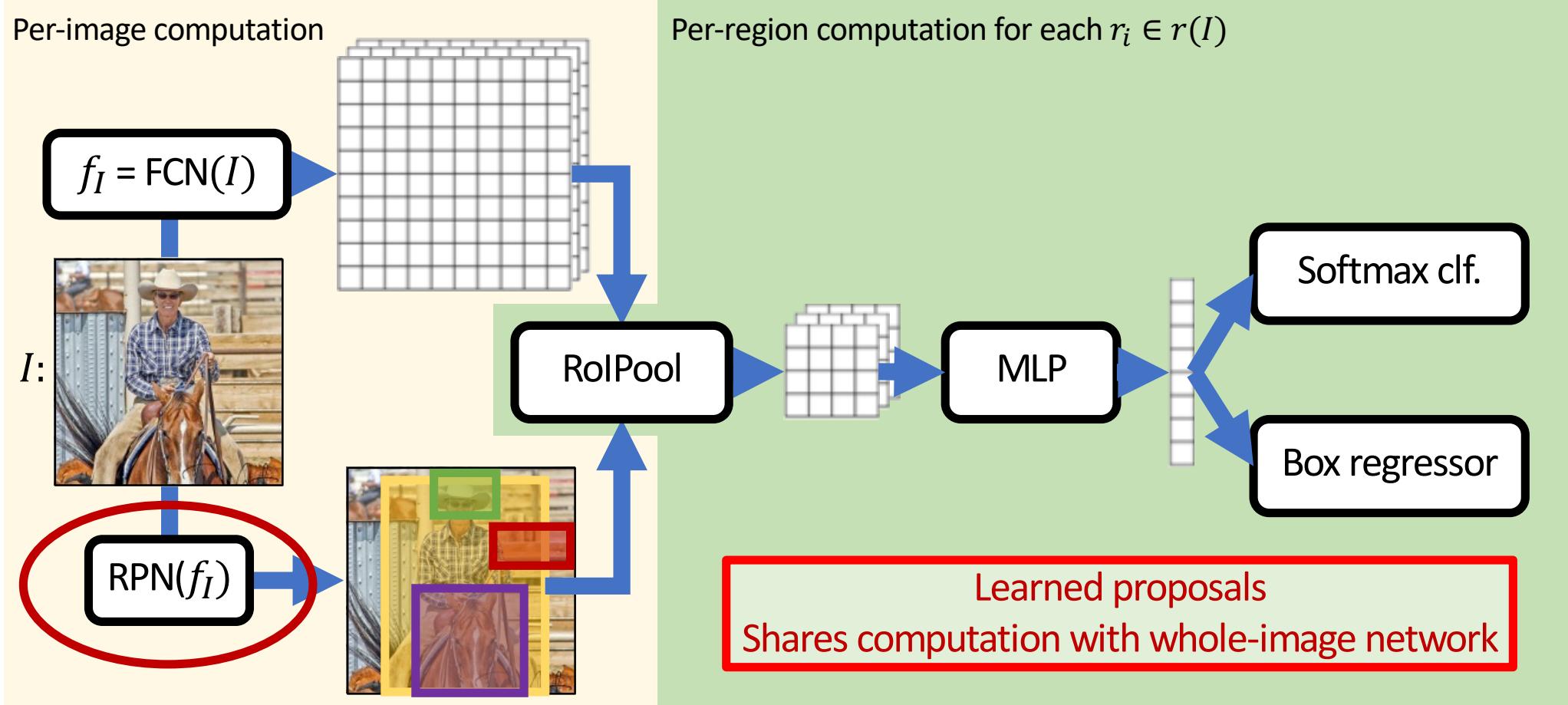
# Faster R-CNN

## 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.

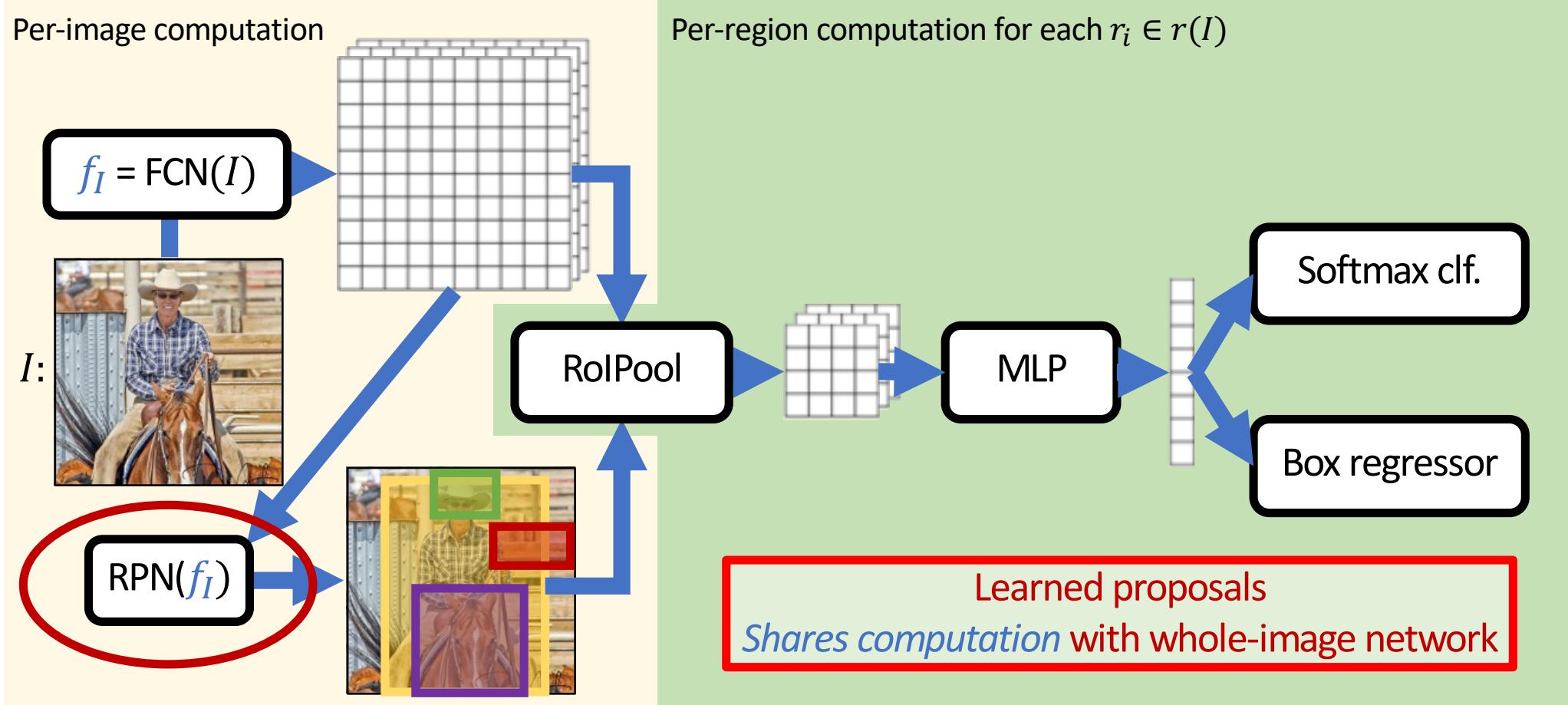
Ren, He, Girshick, Sun. Faster R-CNN:  
Towards Real-Time Object Detection. NIPS 2015.

# Faster R-CNN

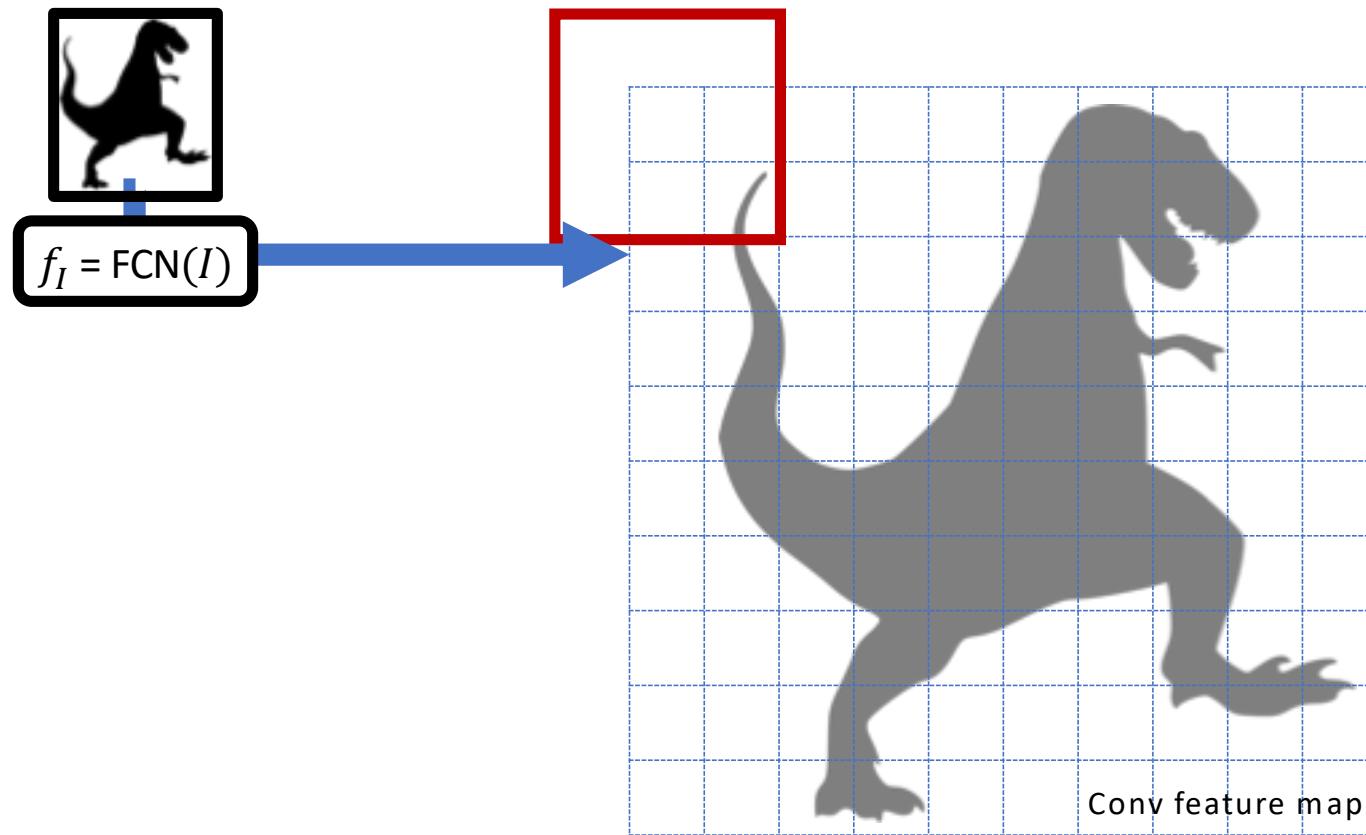


Ren, He, Girshick, Sun. Faster R-CNN:  
Towards Real-Time Object Detection. NIPS 2015.

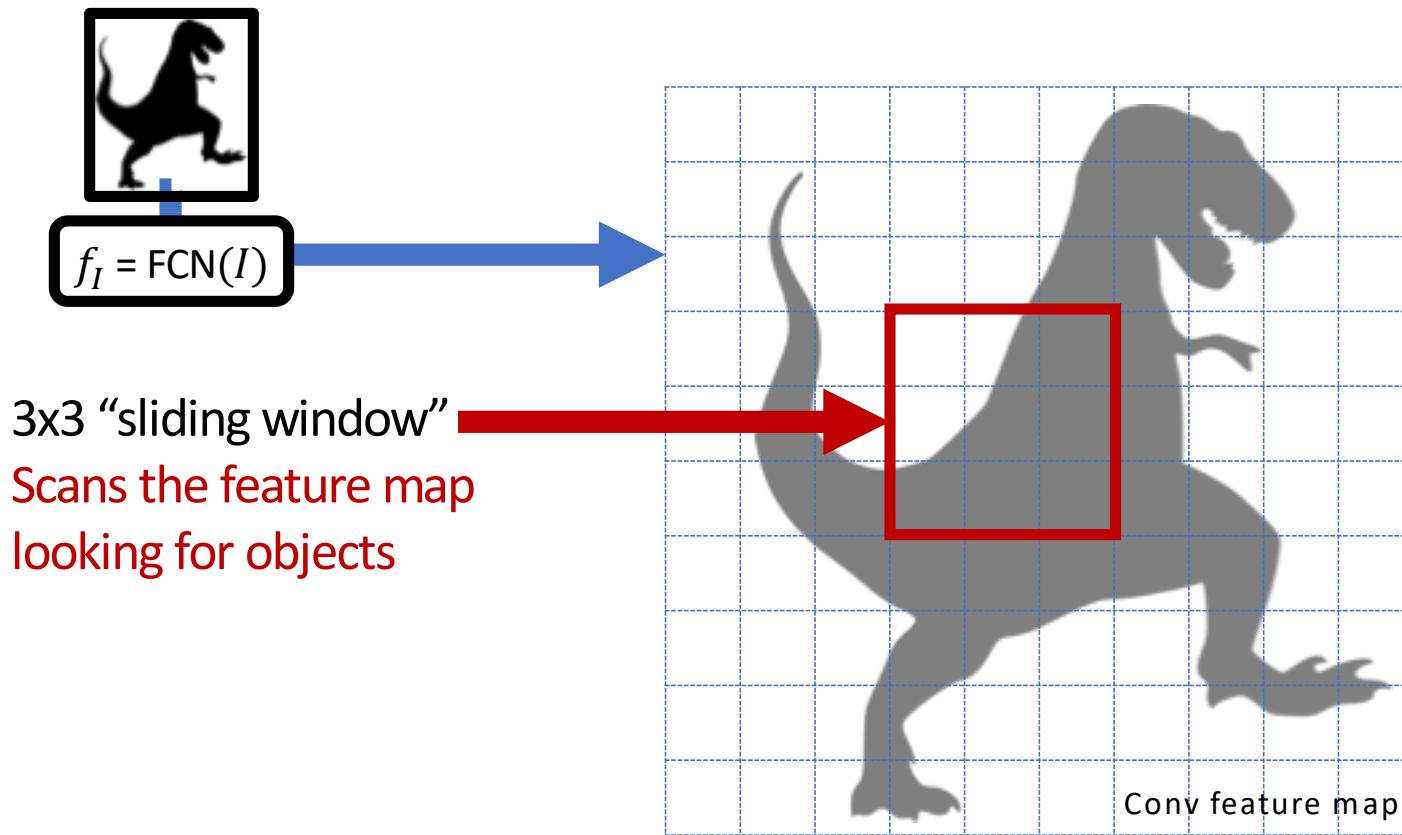
# Faster R-CNN



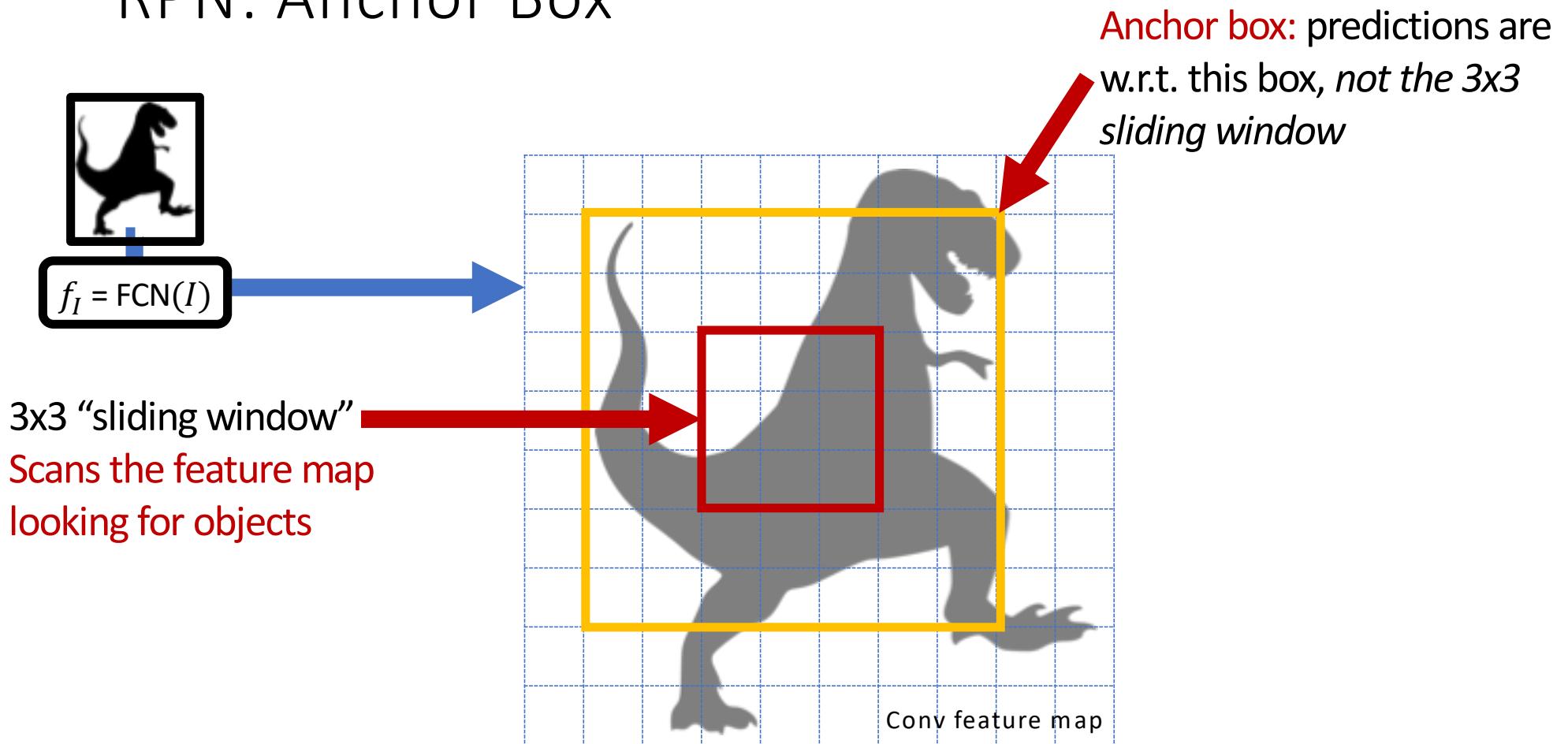
# RPN: Region Proposal Network



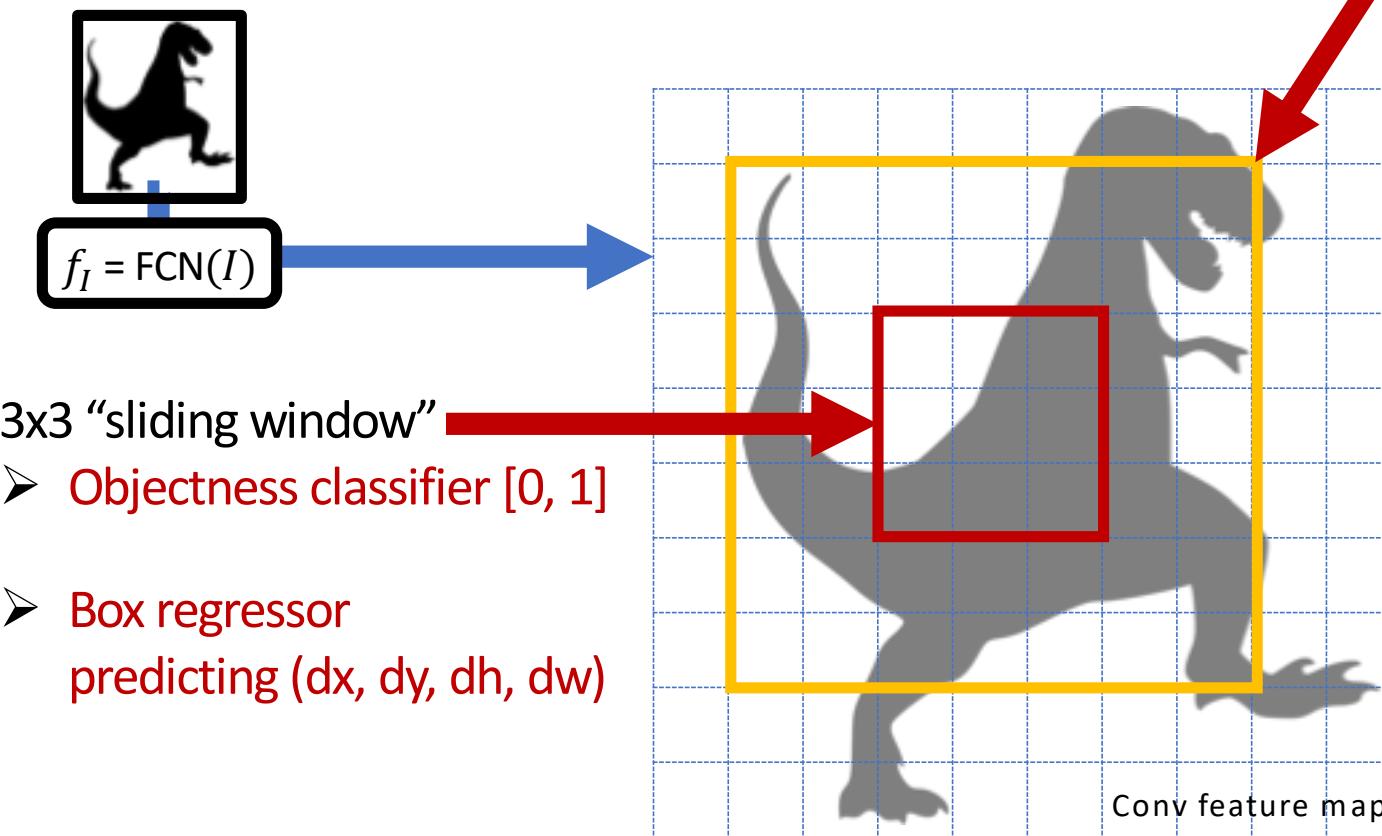
# RPN: Region Proposal Network



# RPN: Anchor Box



# RPN: Anchor Box



## RPN: Prediction (on object)

Objectness score

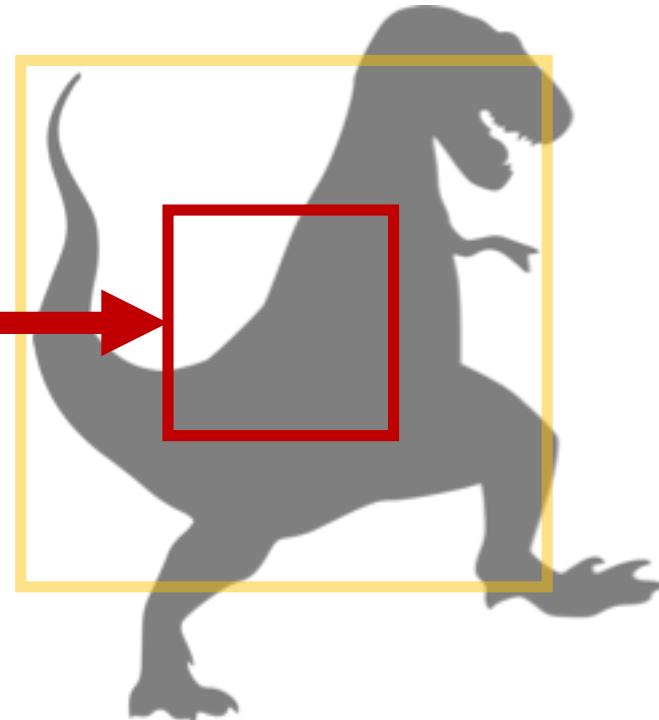


$$P(\text{object}) = 0.94$$

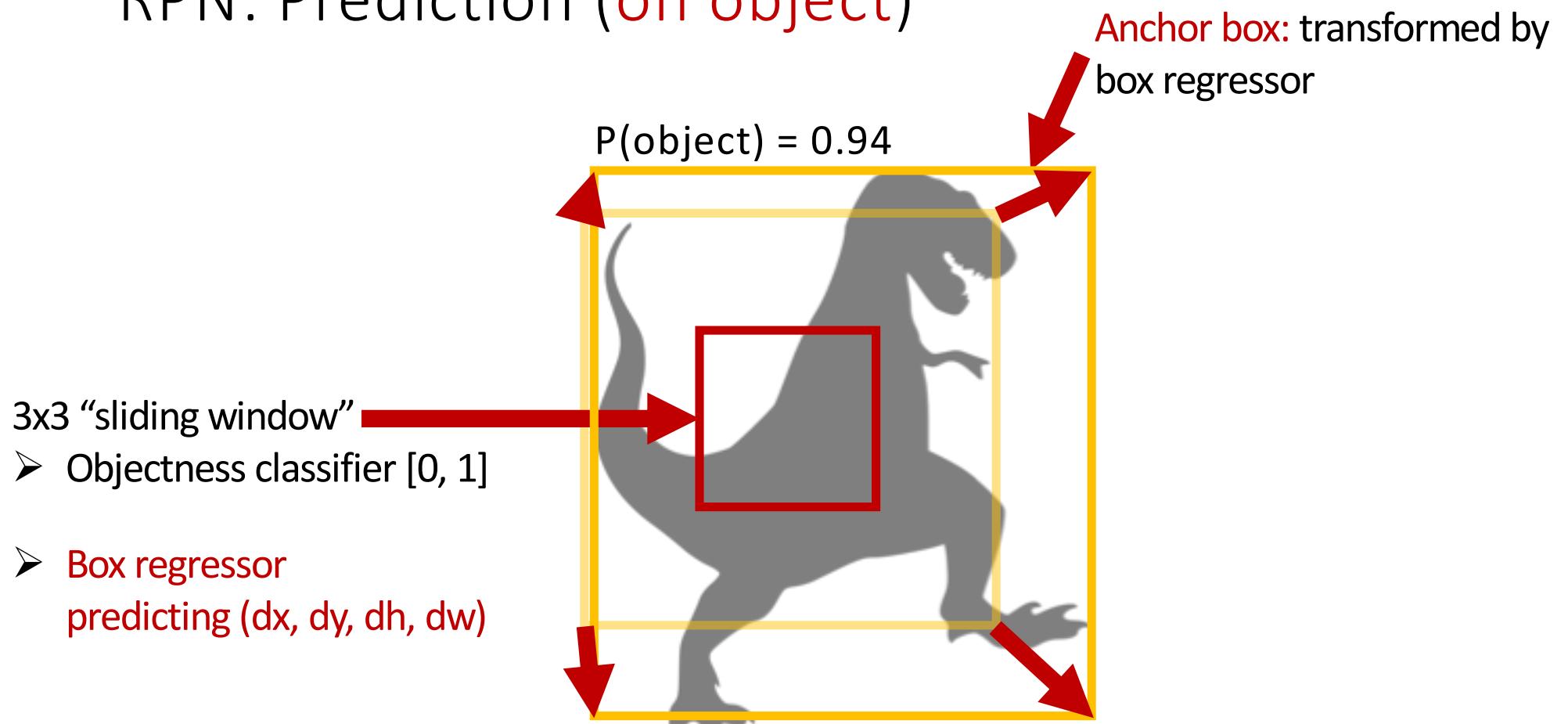
3x3 “sliding window”

➤ Objectness classifier [0, 1]

➤ Box regressor  
predicting (dx, dy, dh, dw)



## RPN: Prediction (on object)



# RPN: Prediction (off object)

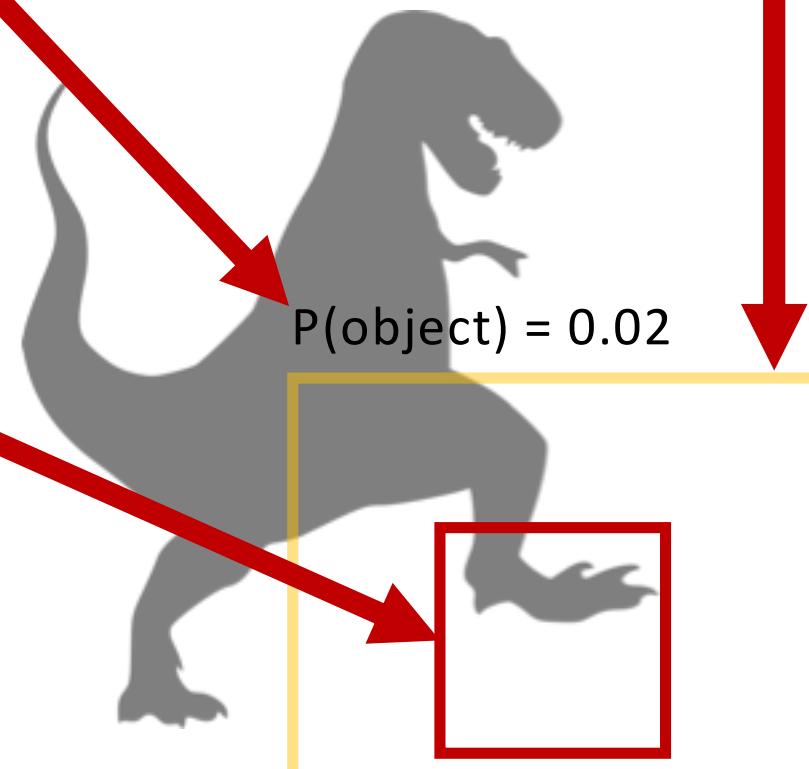
Objectness score

3x3 “sliding window”

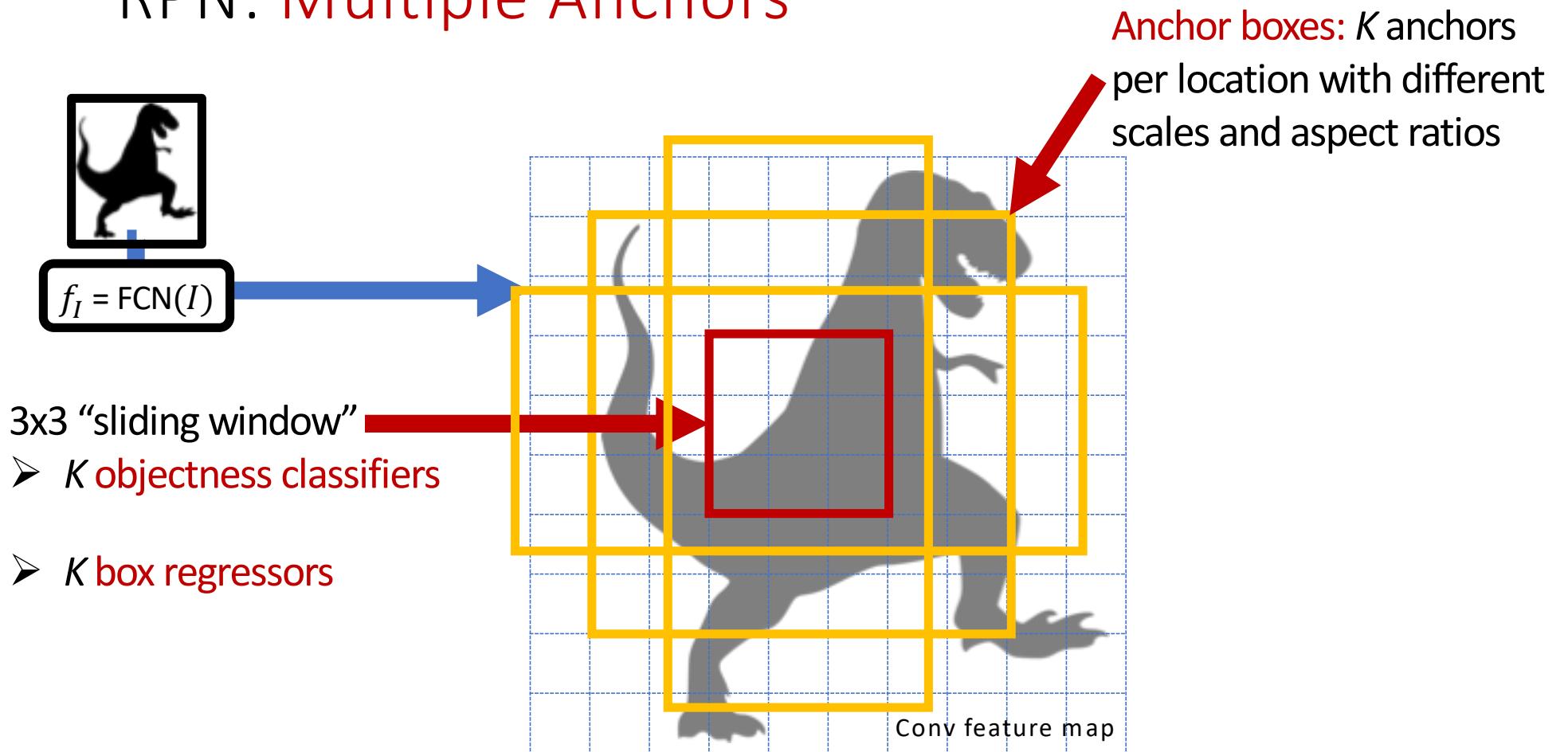
➤ Objectness classifier

➤ Box regressor  
predicting  $(dx, dy, dh, dw)$

Anchor box: transformed by  
box regressor

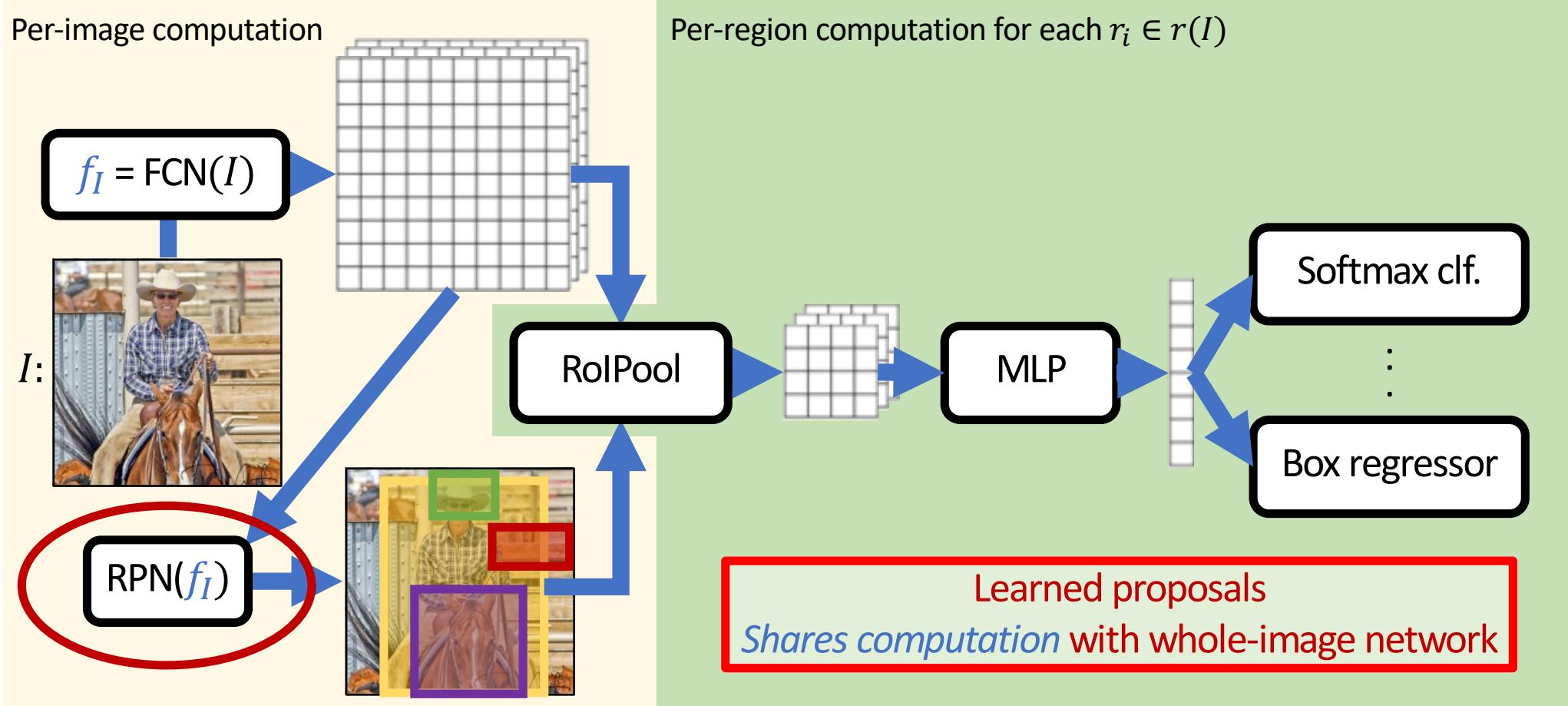


## RPN: Multiple Anchors



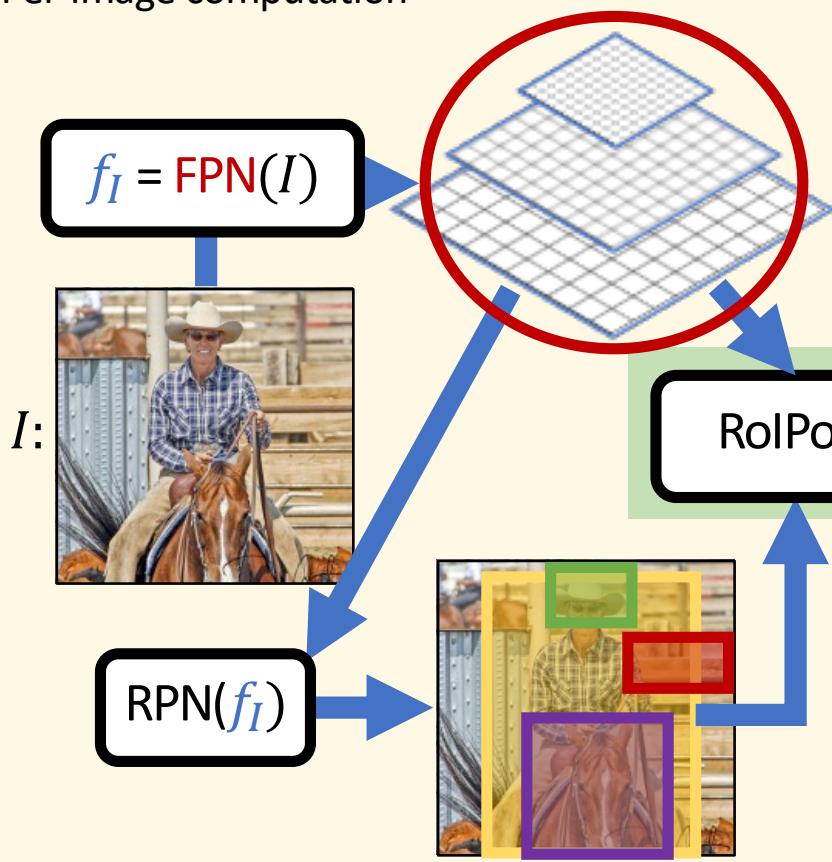
Ren, He, Girshick, Sun. Faster R-CNN:  
Towards Real-Time Object Detection. NIPS 2015.

# Faster R-CNN

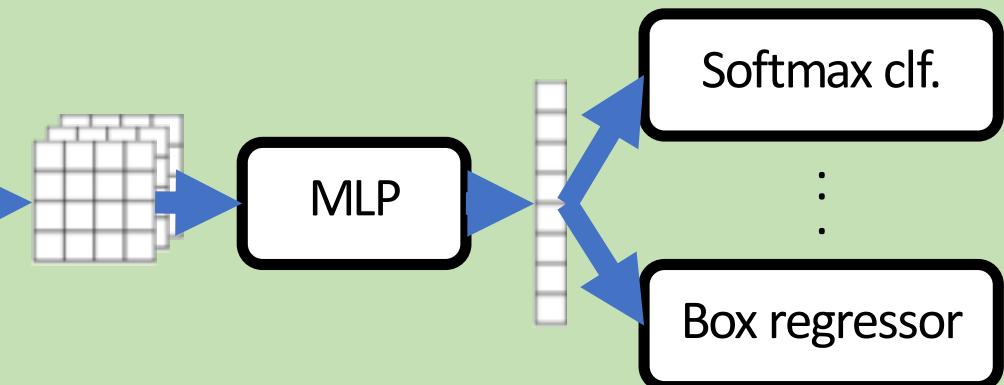


# Faster R-CNN with a Feature Pyramid Network

Per-image computation



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



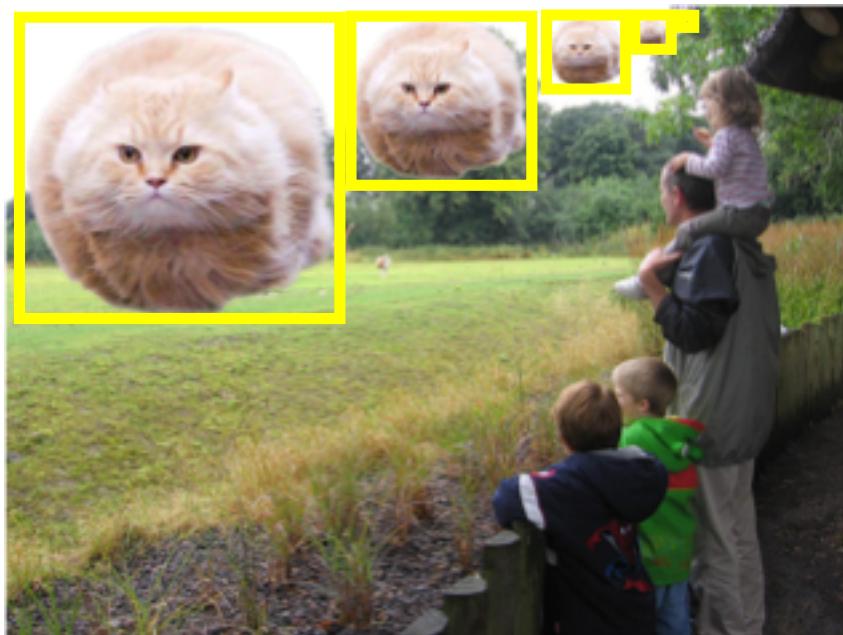
The whole-image feature representation can be improved by making it *multi-scale*

# Feature Pyramid Network (FPN)

## References

- O. Ronneberger, P. Fischer, and T. Brox. U-Net: Convolutional networks for biomedical image segmentation. In MIC- CAI, 2015.
- P. O. Pinheiro, T.-Y. Lin, R. Collobert, and P. Dollár. Learning to refine object segments. In ECCV, 2016.
- W. Liu, D. Anguelov, D. Erhan, C. Szegedy, and S. Reed. SSD: Single shot multibox detector. In ECCV, 2016.
- A. Shrivastava, R. Sukthankar, J. Malik, and A. Gupta. Beyond skip connections: Top-down modulation for object detection. arXiv:1612.06851, 2016.
- T.-Y. Lin, P. Dollár, R. Girshick, K. He, B. Hariharan, S. Belongie. Feature Pyramid Networks for Object Detection. In CVPR, 2017.

# FPN: Improving Scale Invariance and Equivariance



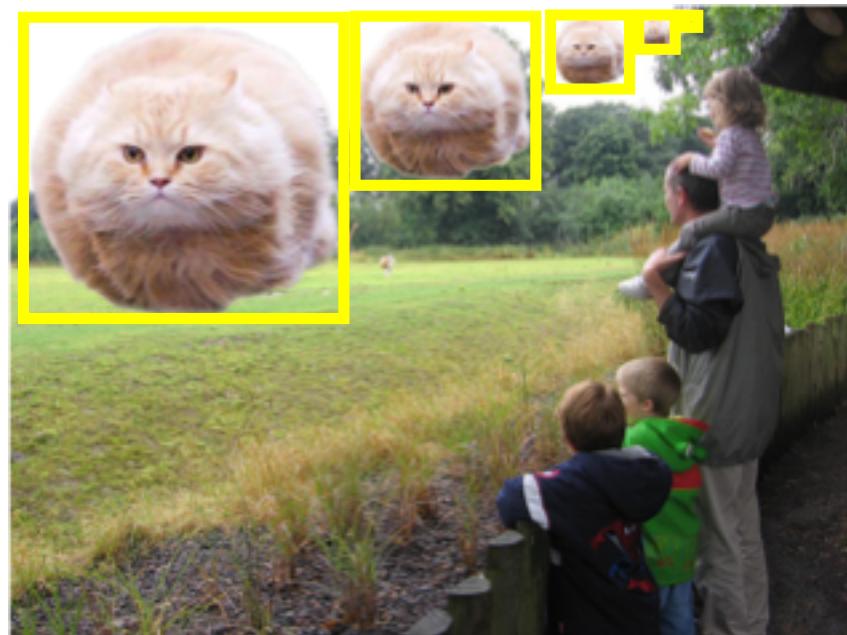
Detectors need to

1. classify and
2. localize

objects over a **wide range of scales**

FPN improves this ability

# Strategy 1: Image Pyramid

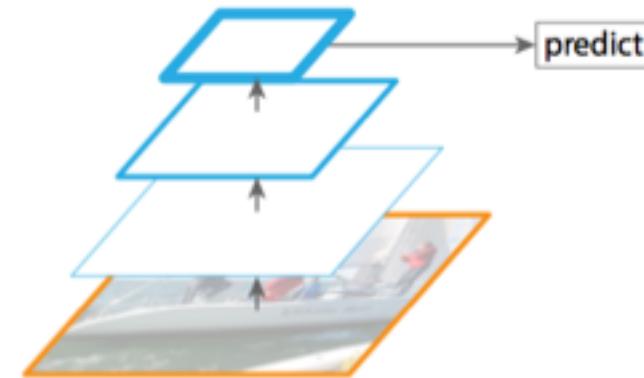
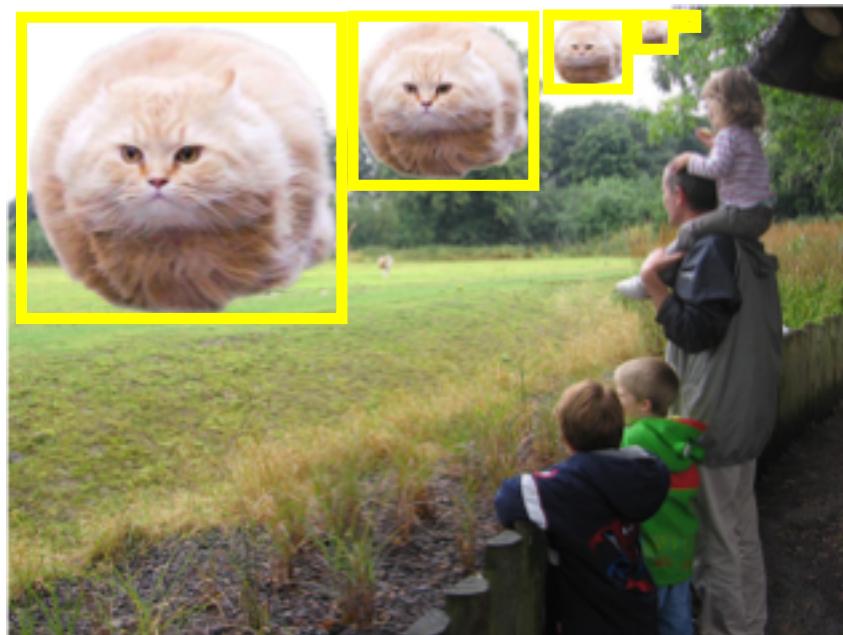


(a) Featurized image pyramid

Standard solution – *slow!*

(E.g., Viola & Jones, HOG, DPM, SPP-net, multi-scale Fast R-CNN, ...)

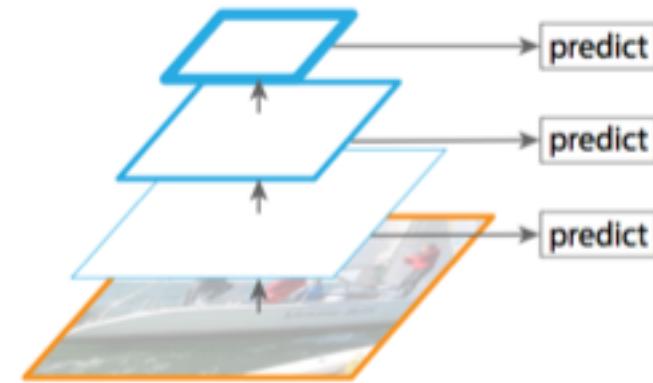
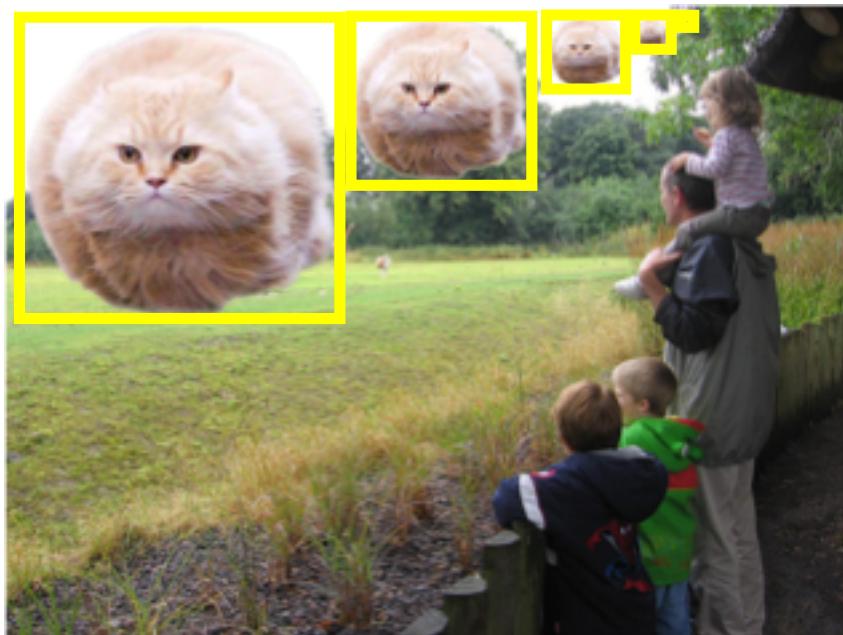
## Strategy 2: Multi-scale Features (Single-scale Map)



(b) Single feature map

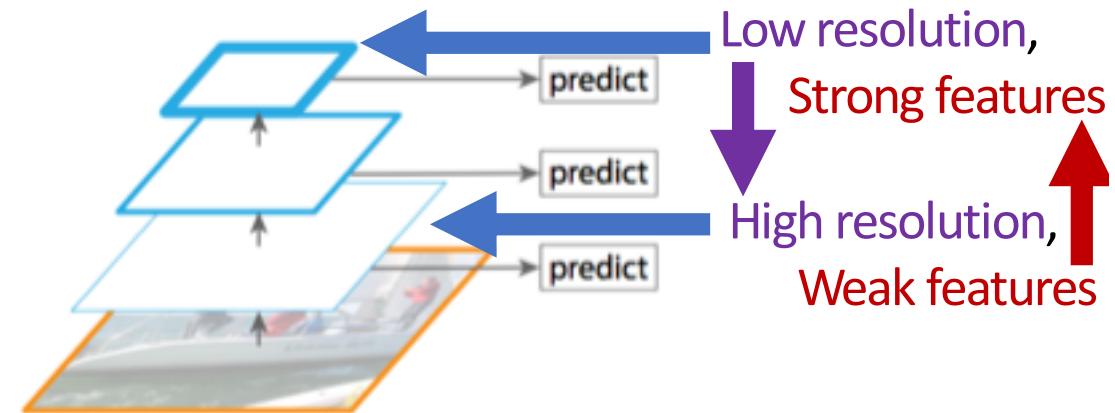
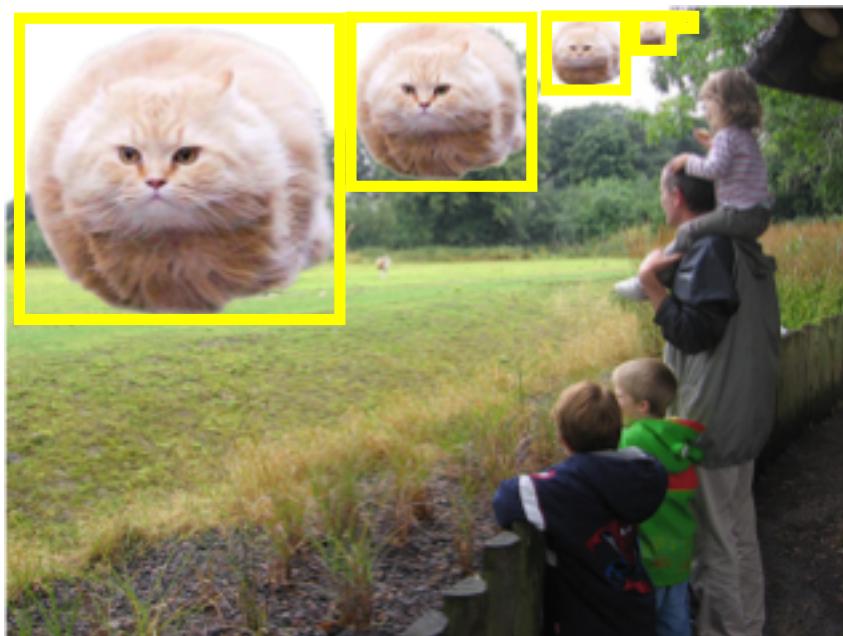
Leave it all to the features – *fast, suboptimal*  
(E.g., Fast/er R-CNN, YOLO, ...)

## Strategy 3: Naïve In-network Pyramid



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

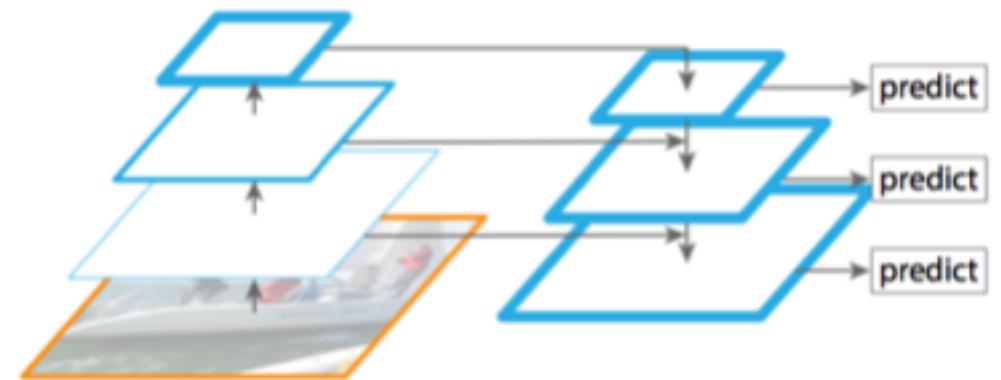
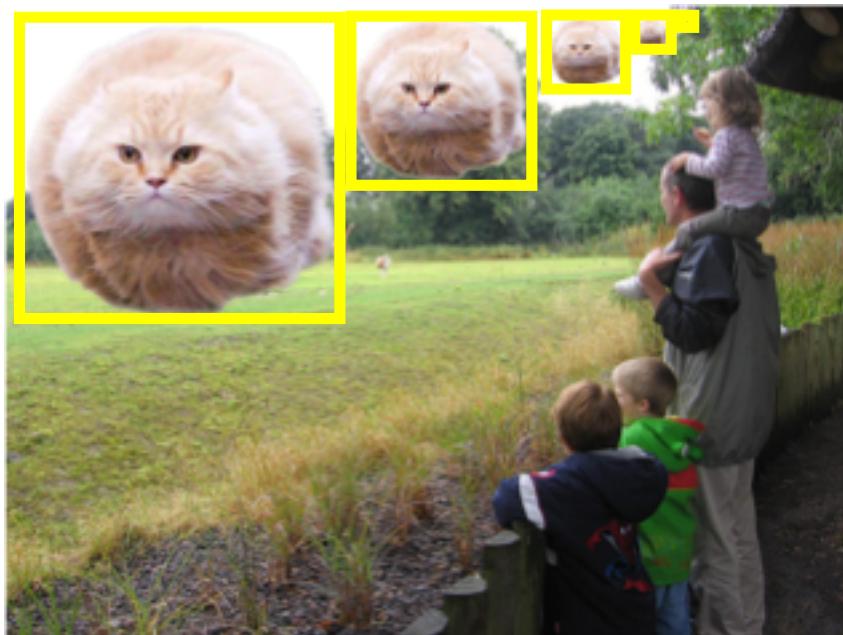
## Strategy 3: Naïve In-network Pyramid



(c) Pyramidal feature hierarchy

Use the internal pyramid – *fast, suboptimal*  
(E.g.,  $\approx$  SSD, ...)

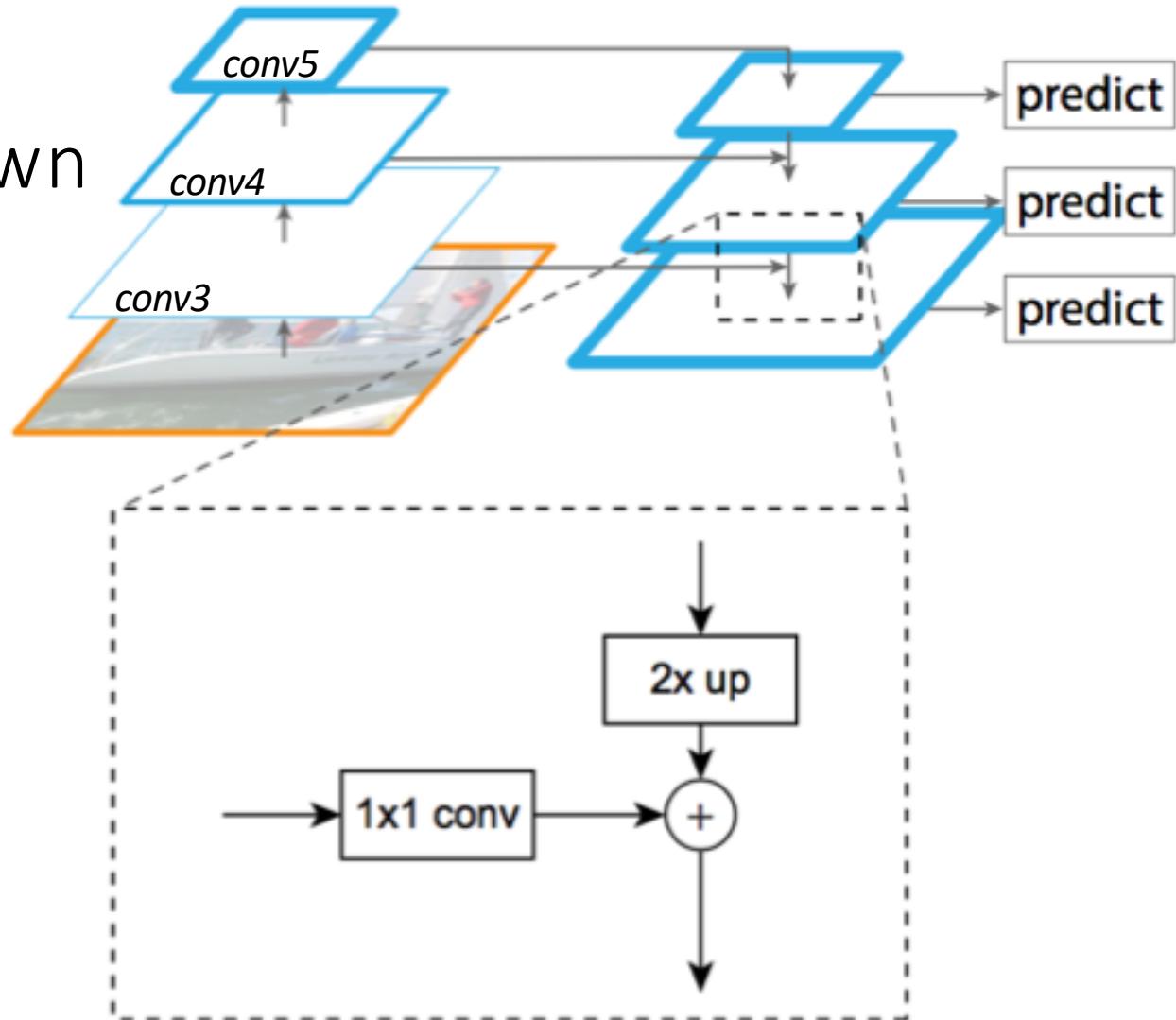
## Strategy 4: Feature Pyramid Network



(d) Feature Pyramid Network

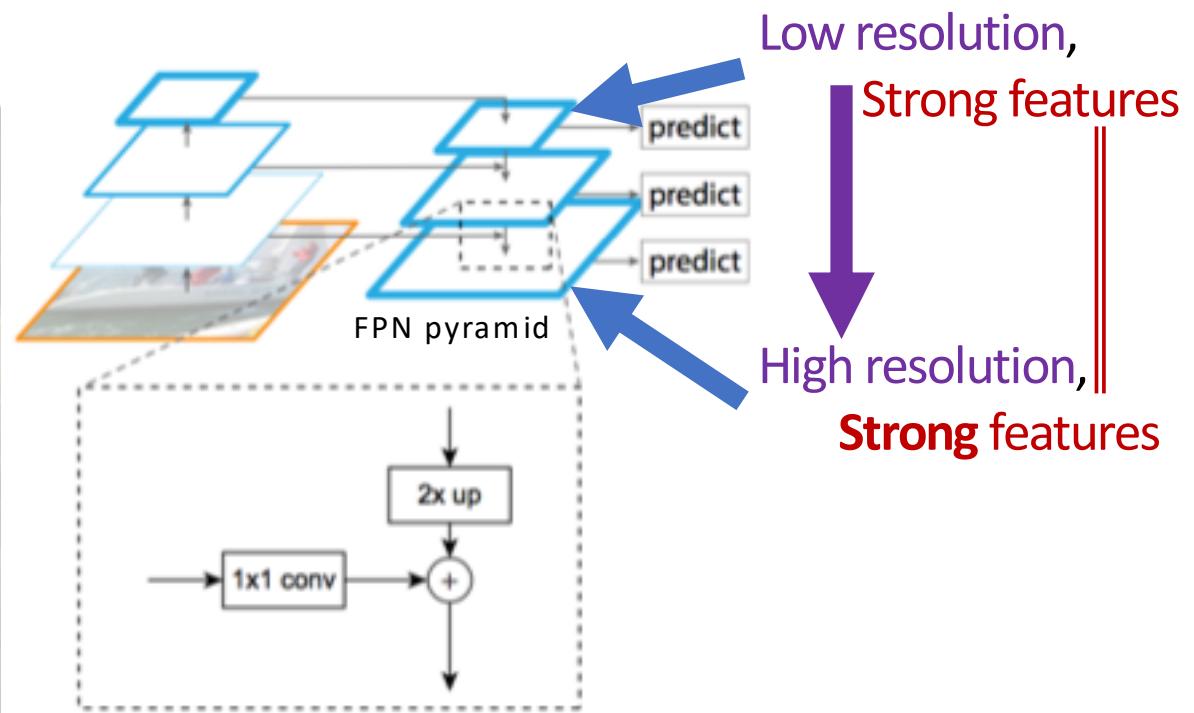
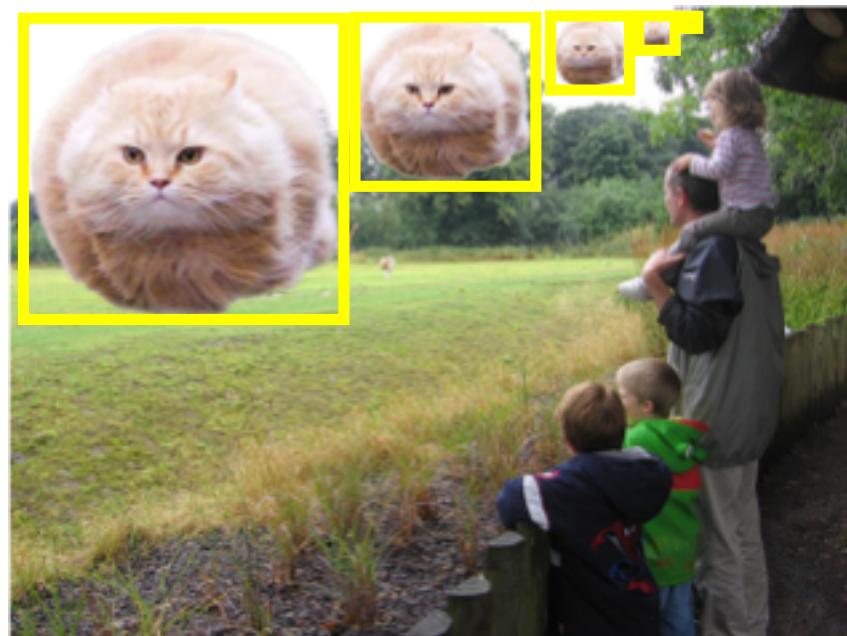
Top-down enrichment of high-res features –  
*fast, less suboptimal*

# FPN Top-down Refinement Module



Lin et al. Feature Pyramid Networks for Object Detection. CVPR 2017.

# No Compromise on Feature Quality, still Fast



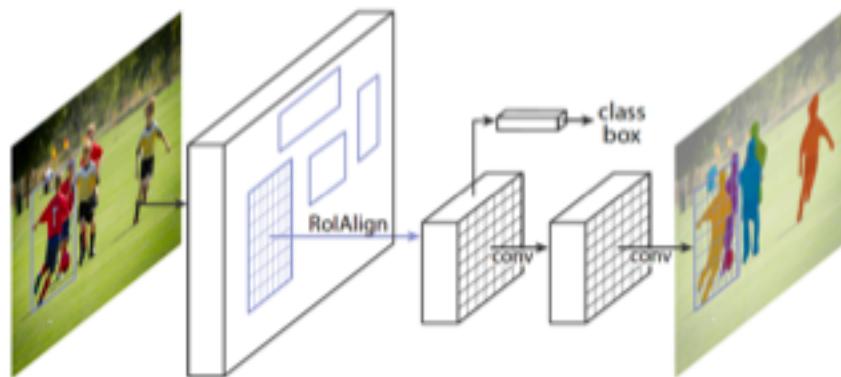
Lin et al. Feature Pyramid Networks for Object Detection. CVPR 2017. See also: Shrivastava's TDM.

# FPN – A Generic Backbone Modification

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

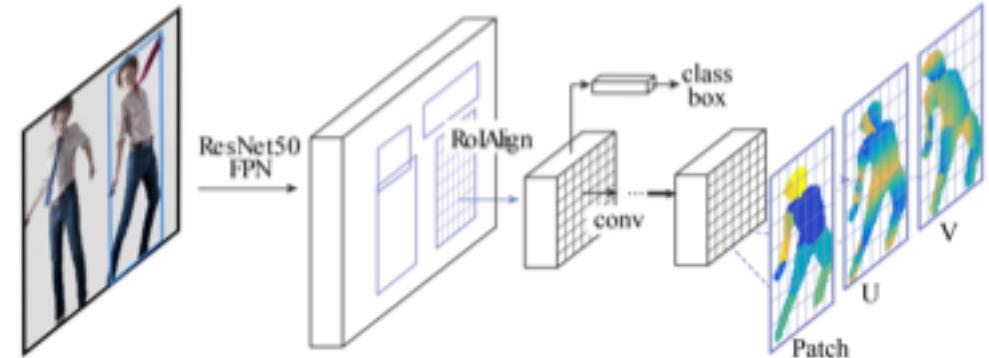
- RPN
- Fast/er R-CNN
- Mask R-CNN
- RetinaNet
- FPN-based semantic segmentation FCN

# Generalized R-CNN: Adding More Heads



**Mask R-CNN**

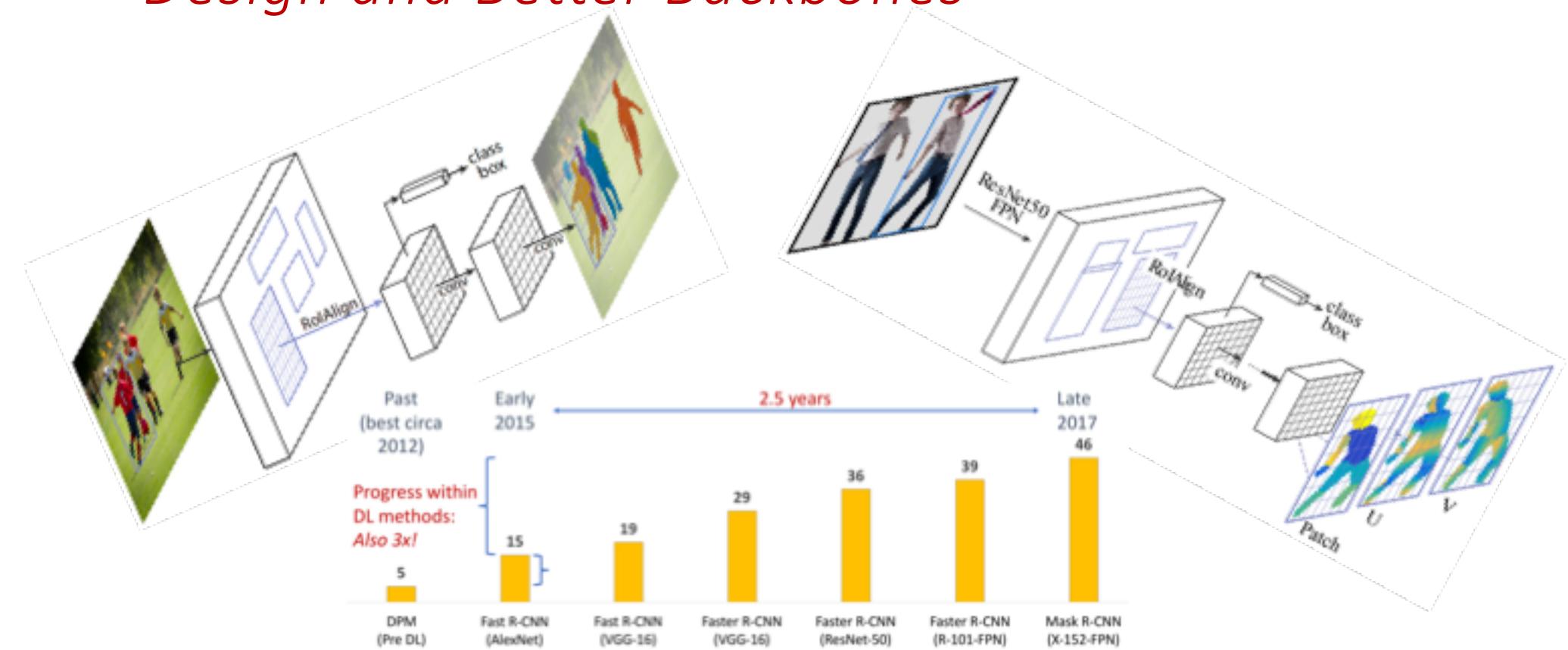
[He, Gkioxari, Dollár, Girshick]



**DensePose**

[Güler, Neverova, Kokkinos]

# *Summary: Progress through Better Detector Design and Better Backbones*



## Takeaway Points

1. Detection is not solved – *datasets are just a model of the world*
2. Features matter – *the backbone net is the engine of recognition*
3. Detector design matters – *progress from domain knowledge*
4. Flexible / extensible framework with OSS code

**<https://github.com/facebookresearch/Detectron>**

# Additional Material Not Covered in 2018

## Topics to cover

- Object detection intro (very brief)
- The Generalized R-CNN framework
- One-stage vs. multi-stage detectors & speed/accuracy tradeoffs

## Stages: What and Why?

Detection output space:  $N = H \times W$  pixel image has  $O(N^2)$  boxes

Output space is **HUGE**, even for small images it's *billions of boxes*

Number of target objects is small, say **10**

**Massive** foreground / background class imbalance

# How do we Deal with Class Imbalance?

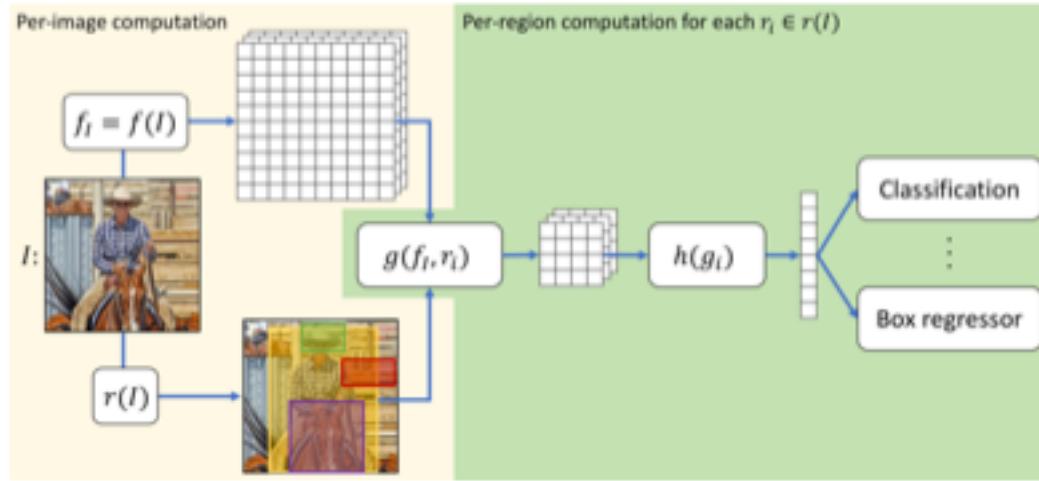
## 1. Subsample the output space while maintaining *high recall*

- Sliding window
- Object proposals

## 2. Classify boxes in a *cascade of stages*

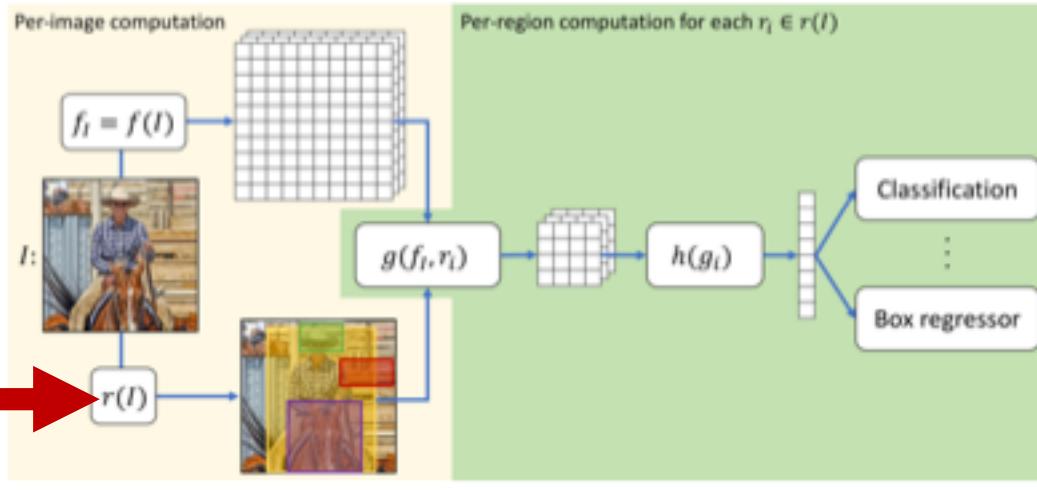
- Most famous example: Viola-Jones detector
- Typically combined with subsampling
- *Note that subsampling is already an implicit stage*

## Example of two stages using proposals: Generalized R-CNN Framework



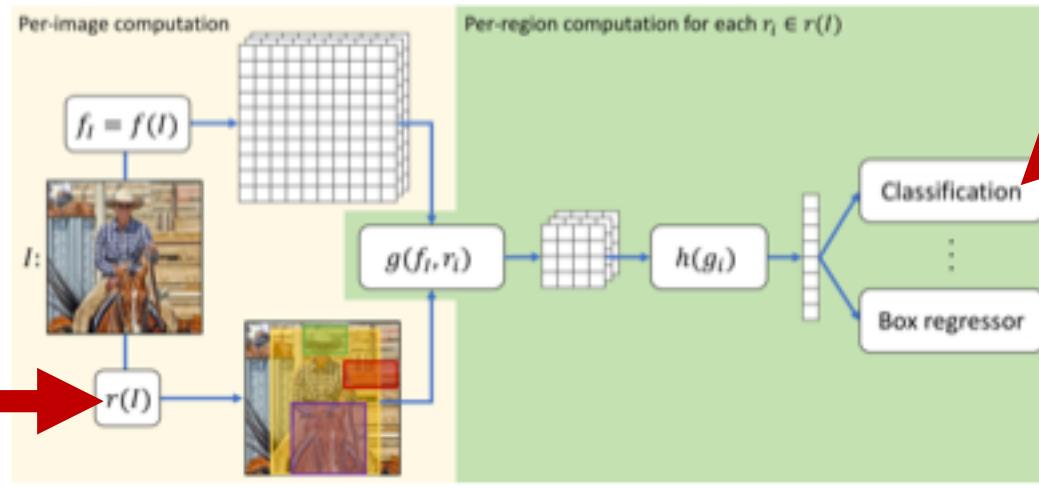
## Example of two stages using proposals: Generalized R-CNN Framework

1. Proposals  
eliminate most  
of the output space



## Example of two stages using proposals: Generalized R-CNN Framework

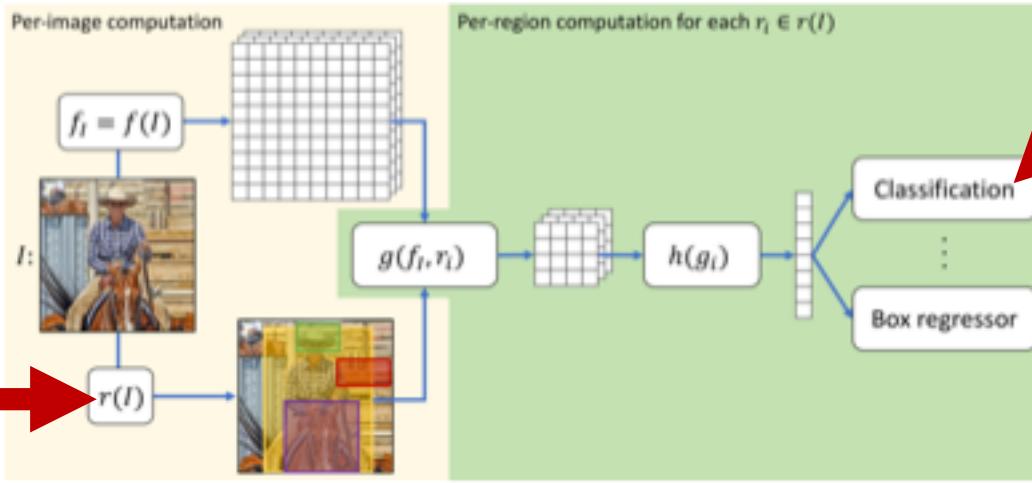
1. Proposals  
eliminate most  
of the output space



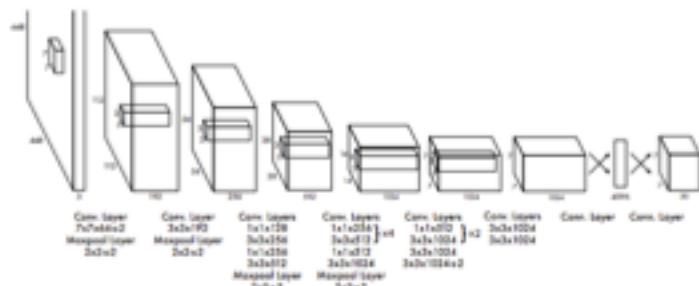
2. Classification  
on small subset  
of output space

Example of two stages using proposals: Generalized R-CNN Framework

1. Proposals  
eliminate most  
of the output space



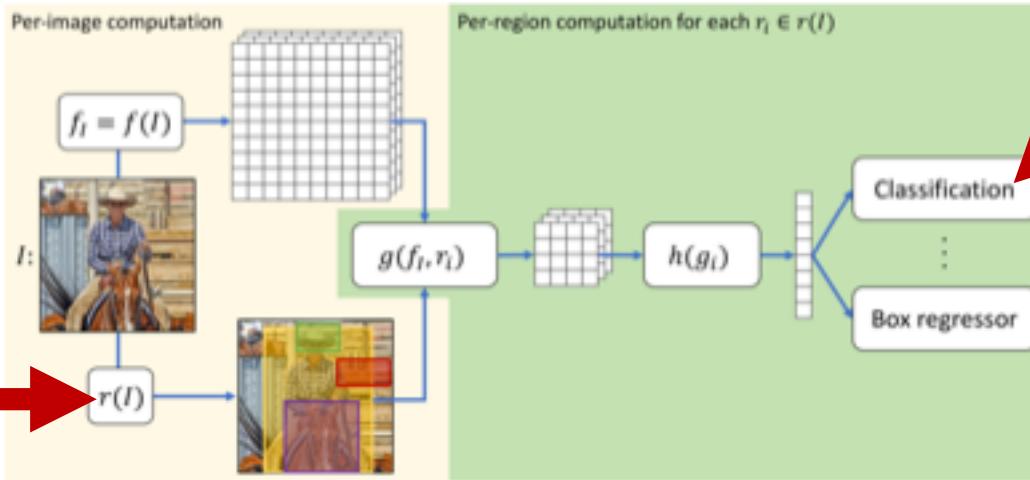
Example one stage using *dramatic* subsampling



Redmond et al. You Only Look Once:  
Unified Real-time Object Detection. CVPR 2016.

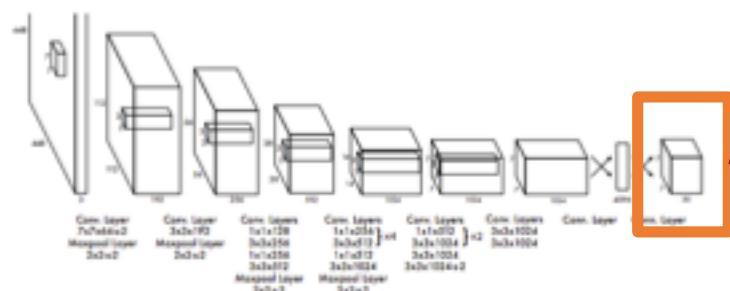
Example of two stages using proposals: Generalized R-CNN Framework

1. Proposals  
eliminate most  
of the output space



2. Classification  
on small subset  
of output space

Example one stage using *dramatic* subsampling



1. Consider a *tiny* subset of the output space  
by design; directly classify this small set of boxes

“You only look once” = “Single shot”  
= “One stage”

Redmond et al. You Only Look Once:  
Unified Real-time Object Detection. CVPR 2016.

# Subsampled Output Space [+ Hard Example Mining]

Central challenge for one-stage detectors: class imbalance

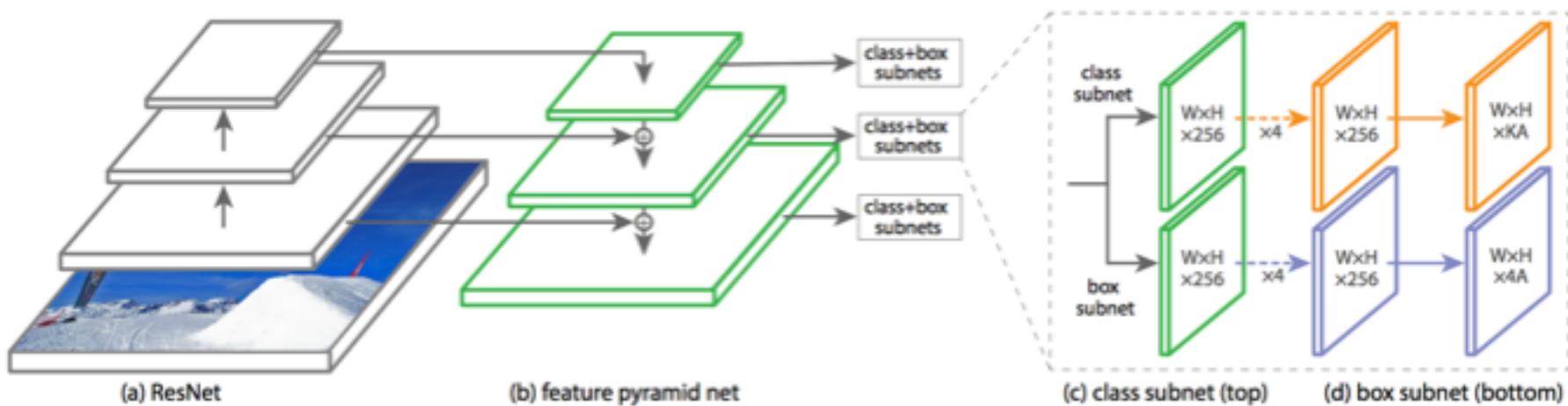
- YOLOv1 – 98 boxes
- YOLOv2 – ~ 1k
- OverFeat – ~ 1-2k
- SSD – ~ 8–26k (hard-example mining)
  
- RetinaNet – ~ 100-200k (“soft” hard-example mining)

# RetinaNet

## References

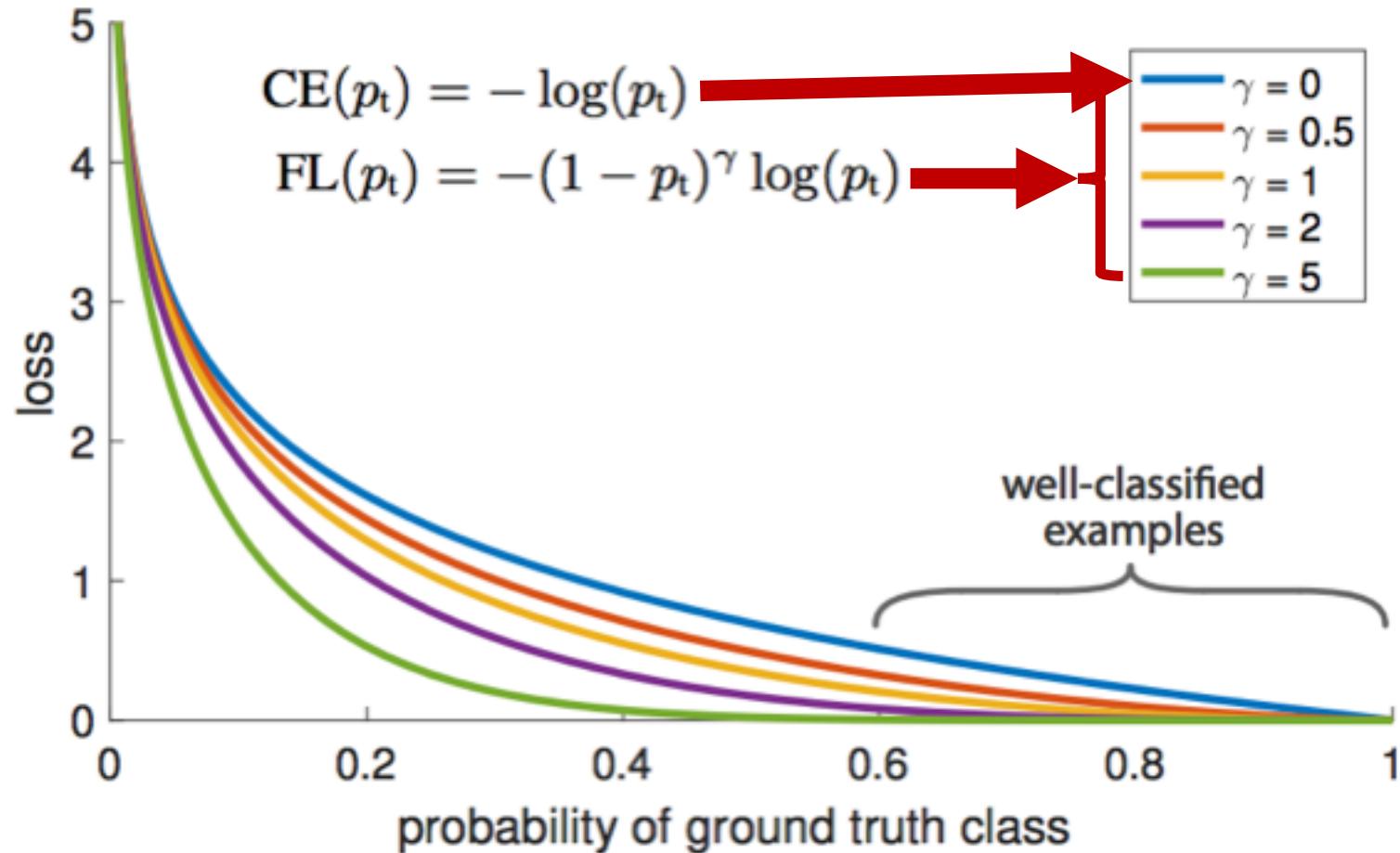
- P. Sermanet, D. Eigen, X. Zhang, M. Mathieu, R. Fergus, and Y. LeCun. Overfeat: Integrated recognition, localization and detection using convolutional networks. In ICLR, 2014.
- A. Shrivastava, A. Gupta, and R. Girshick. Training region-based object detectors with online hard example mining. In CVPR, 2016.
- P. O. Pinheiro, R. Collobert, and P. Dollar. Learning to segment object candidates. In NIPS, 2015.

# RetinaNet Model Description



- Backbone with FPN + class-specific RPN (final detections)
- 6 anchors per location (2 scales  $\times$  3 aspect ratios)
- 100 – 200k anchor boxes to classify per image → “dense” detection

## Focal Loss: “Soft” Hard-Example Mining

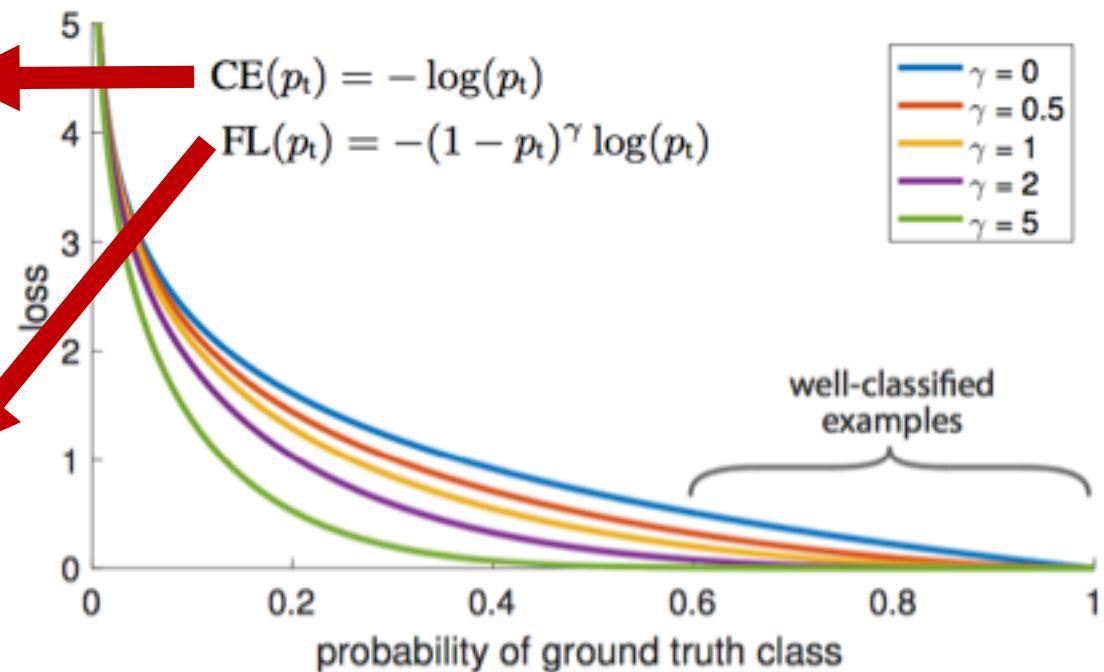


# Impact of Focal Loss (FL)

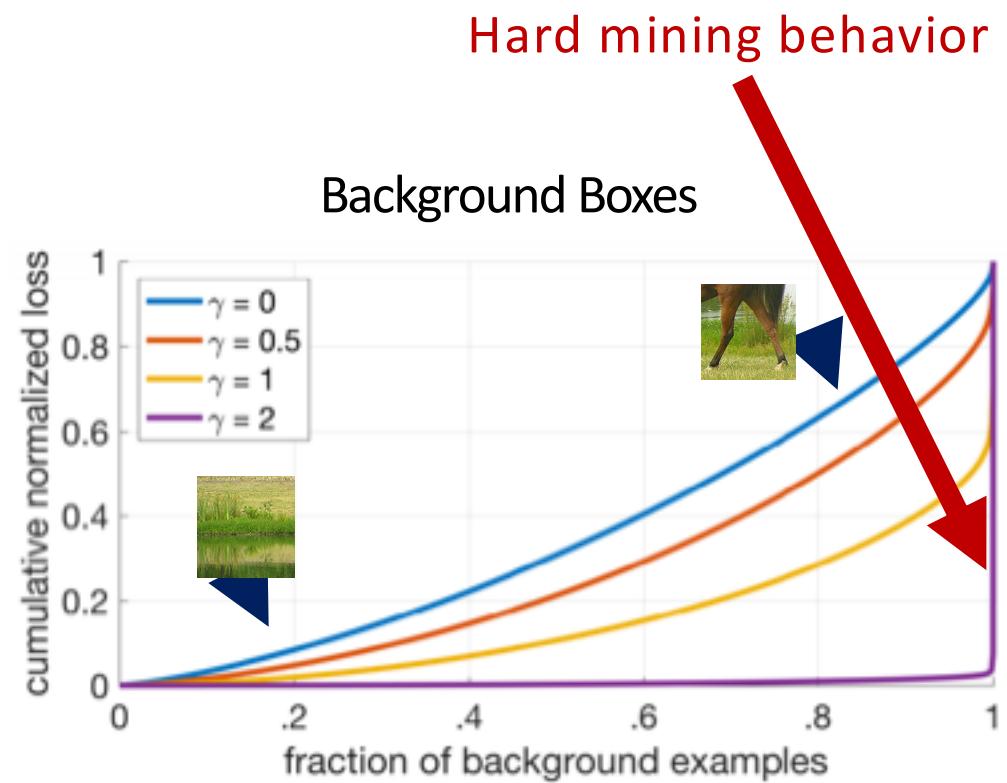
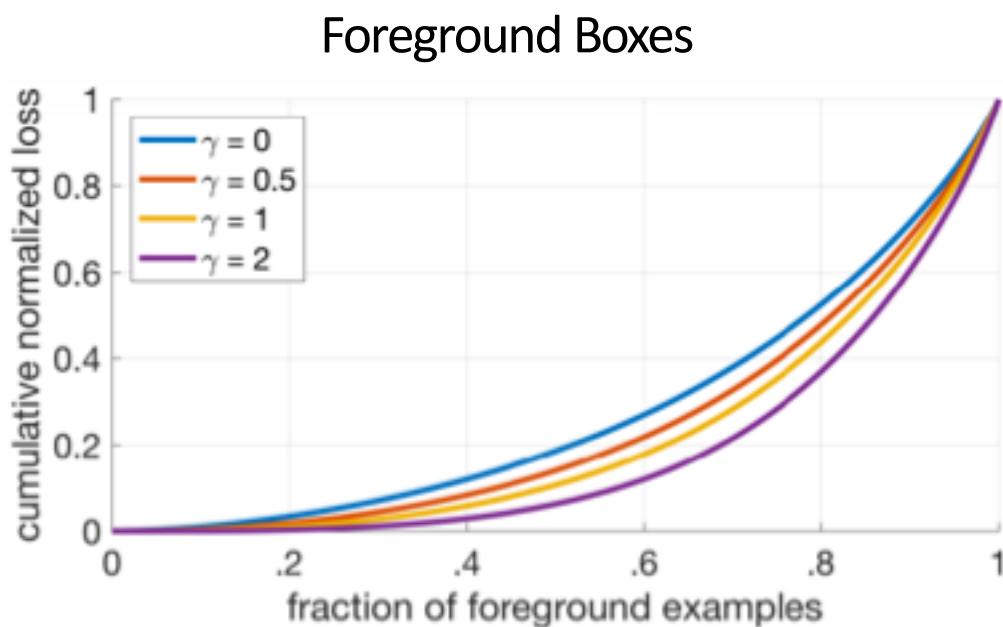
$\gamma$	$\alpha$	AP	AP <sub>50</sub>	AP <sub>75</sub>
0	.75	31.1	49.4	33.0
0.1	.75	31.4	49.9	33.1
0.2	.75	31.9	50.7	33.4
0.5	.50	32.9	51.7	35.2
1.0	.25	33.7	52.0	36.2
2.0	.25	<b>34.0</b>	<b>52.5</b>	<b>36.5</b>
5.0	.25	32.2	49.6	34.8

(b) Varying  $\gamma$  for FL (w. optimal  $\alpha$ )

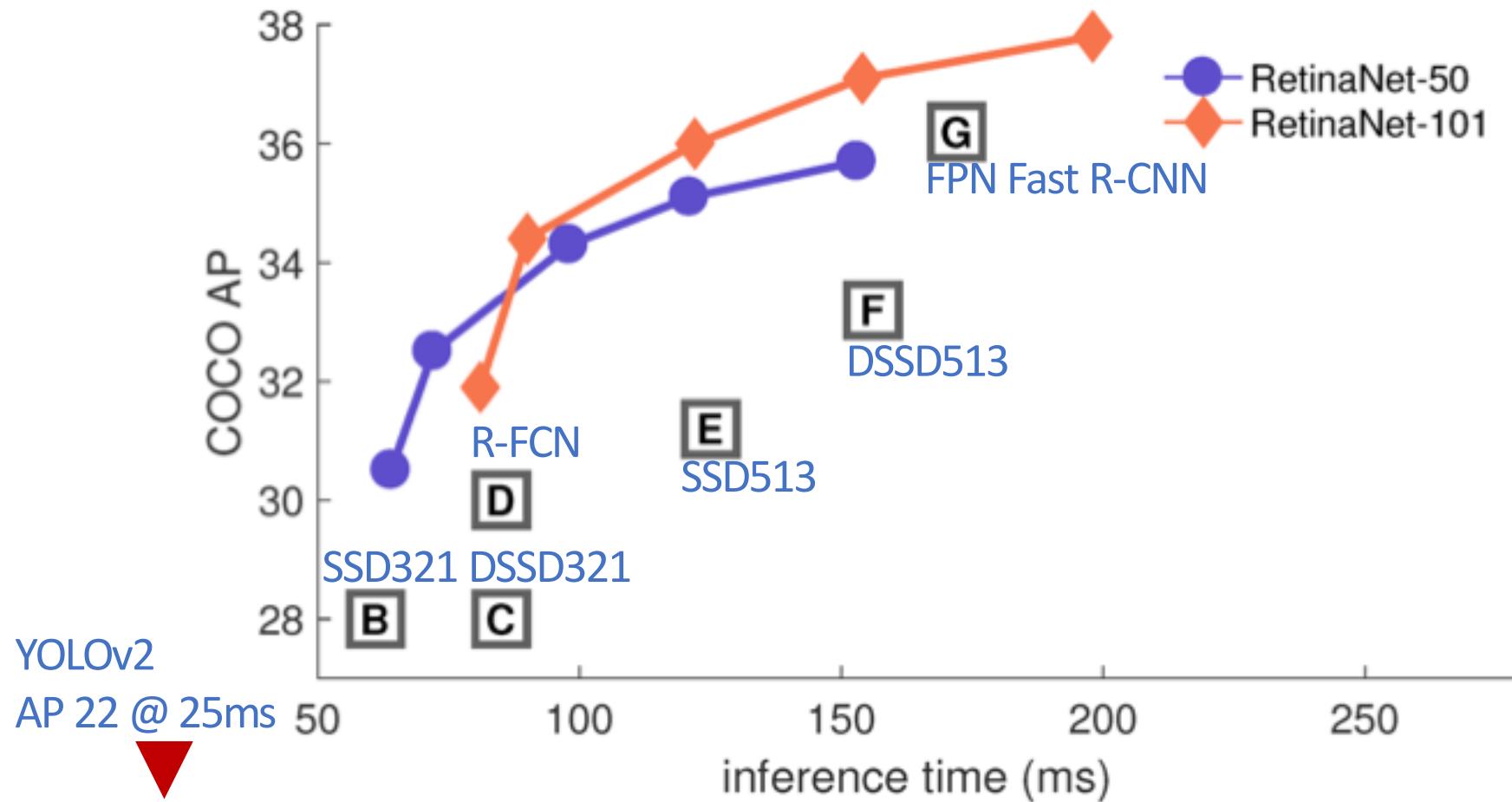
(ResNet-50-FPN 600px input image)



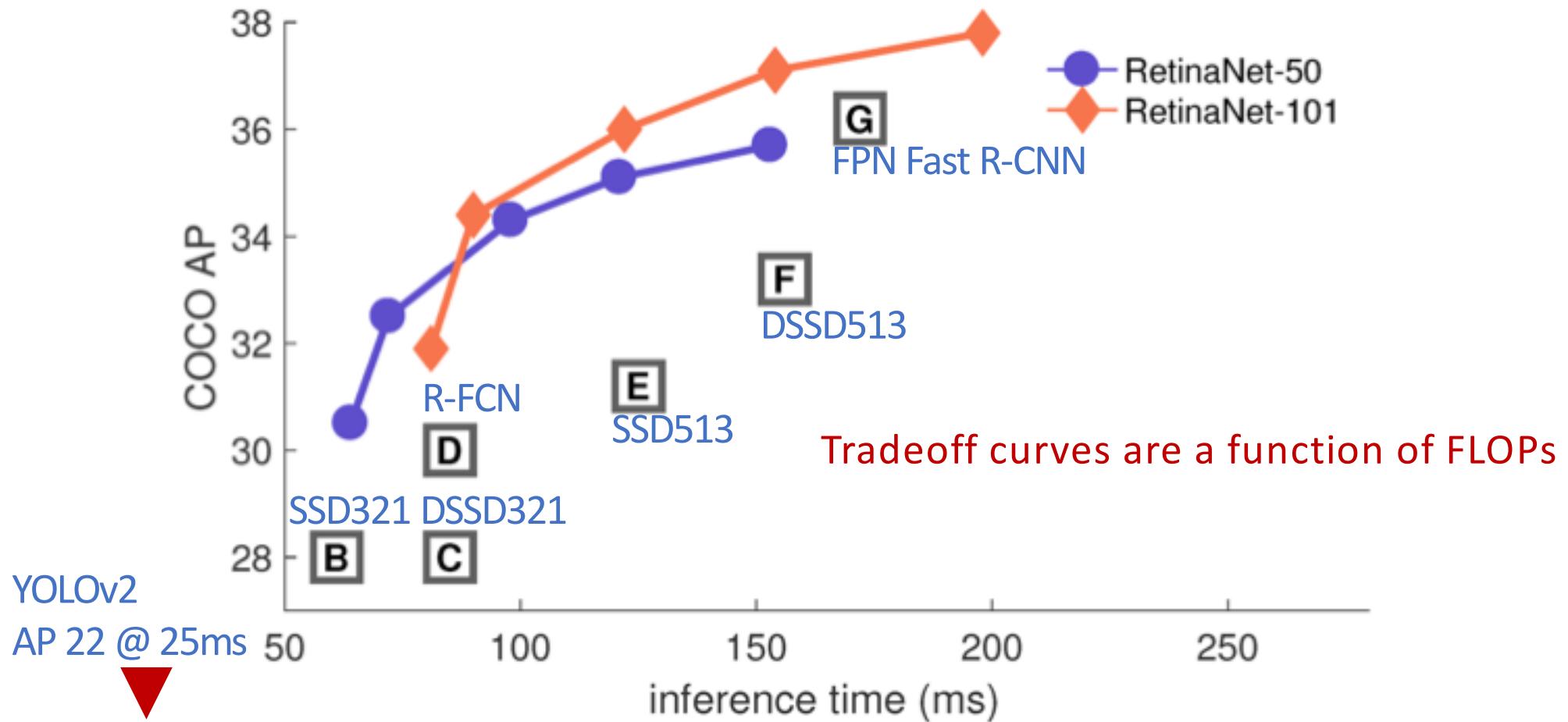
# Loss Distribution under Focal Loss



# Speed/Accuracy Tradeoff



# Speed/Accuracy Tradeoff

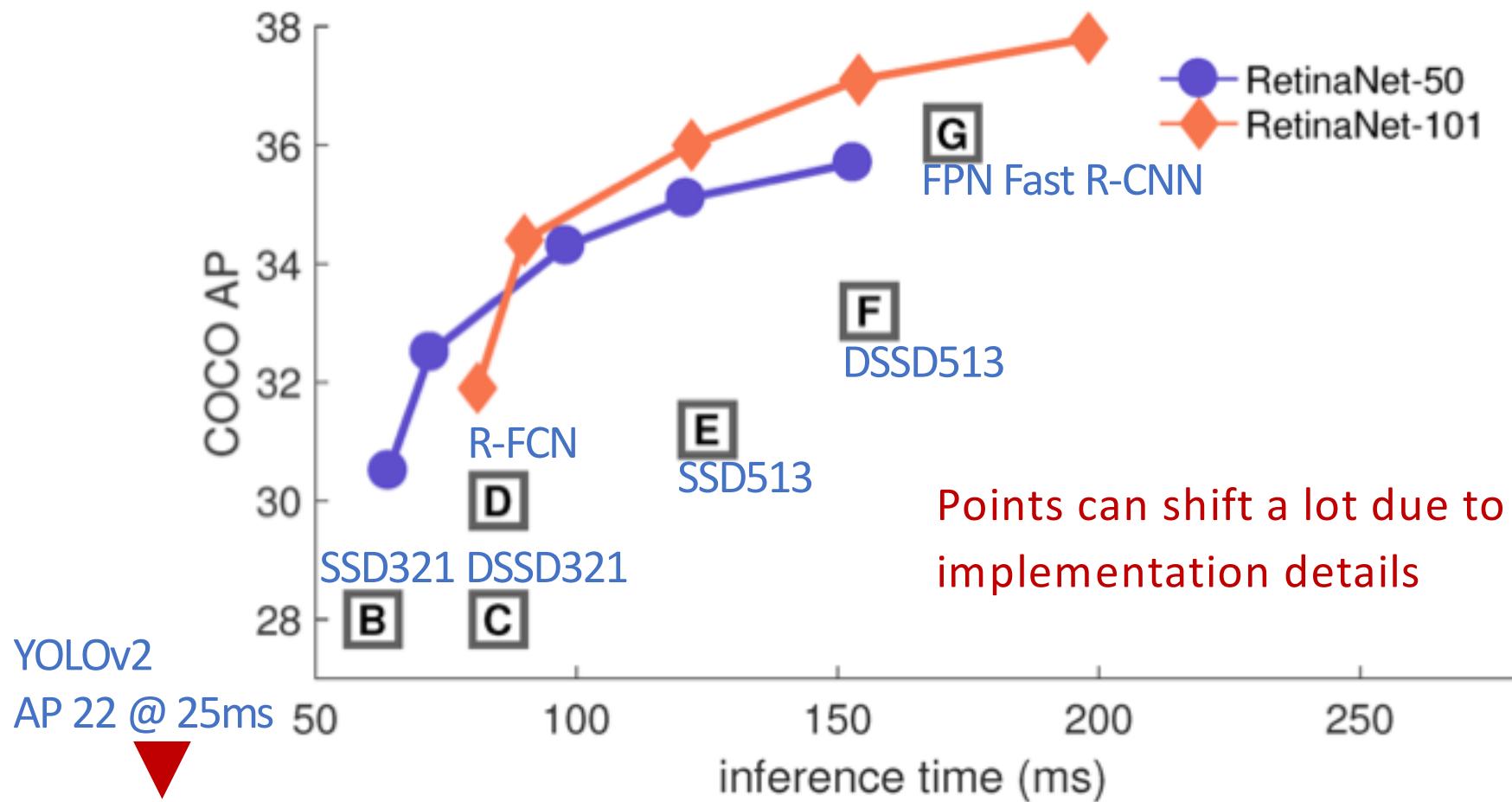


# Warning about Speed/Accuracy Tradeoffs

Comparisons across publications are *completely uncontrolled*

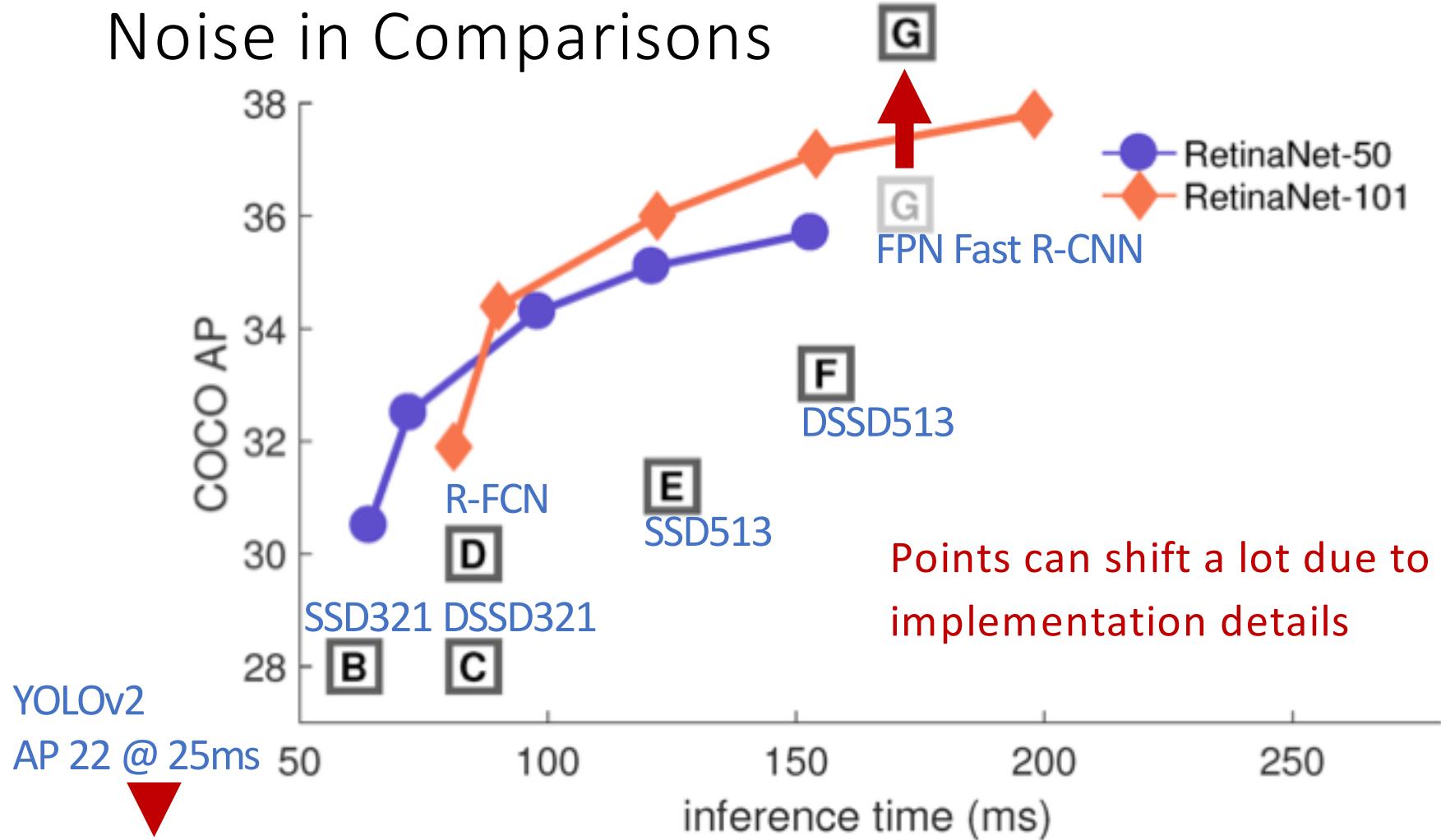
- Results should be taken with a grain of salt
- Accuracy varies with hyper-parameters ('recipe')
- Speed varies with low-level optimization (perf tuning)
- Speed varies with 'tricks' (e.g., batching during inference)

## Noise in Comparisons



## Noise in Comparisons

Improved recipe  
(longer sched., RoIAlign, tuned NMS, e2e training)



## *Fast Detectors*

**Key design factors (*it's all about the FLOPs*)**

- Low resolution input
- Lightweight backbone network
- Top- $k$  proposals with small  $k = 10-50$  (if applicable)

**“Anti-factors”**

- One-stage instead of two-stage

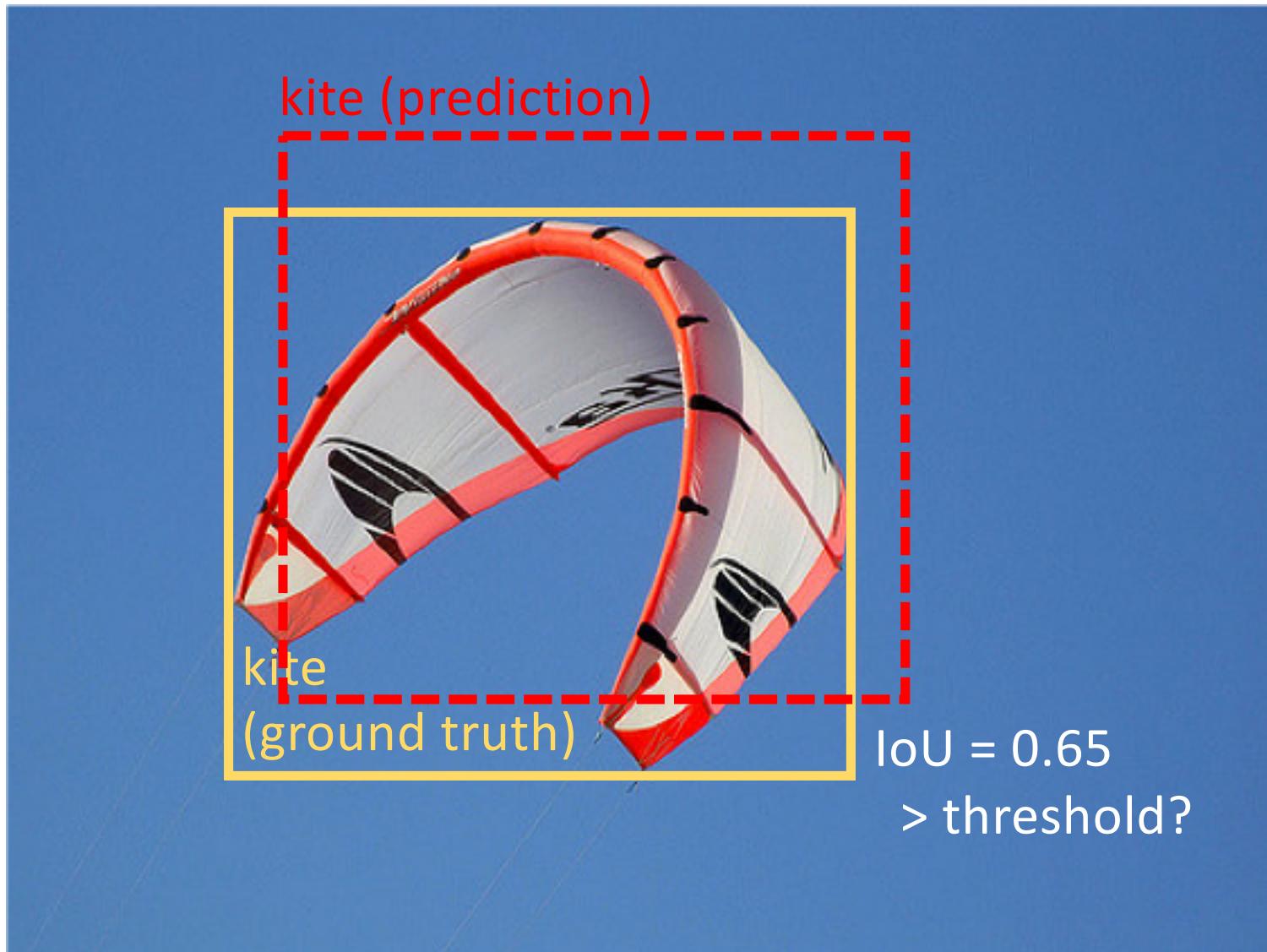
## In Conclusion, Four Takeaway Points

1. Detection is not solved – *datasets are just a model of the world*
2. Features matter – *the backbone net is the engine of recognition*
3. Detector design matters – *progress from domain knowledge*
4. Speed/accuracy tradeoff is *not* a point – *it's a curve/surface*

End of Slides



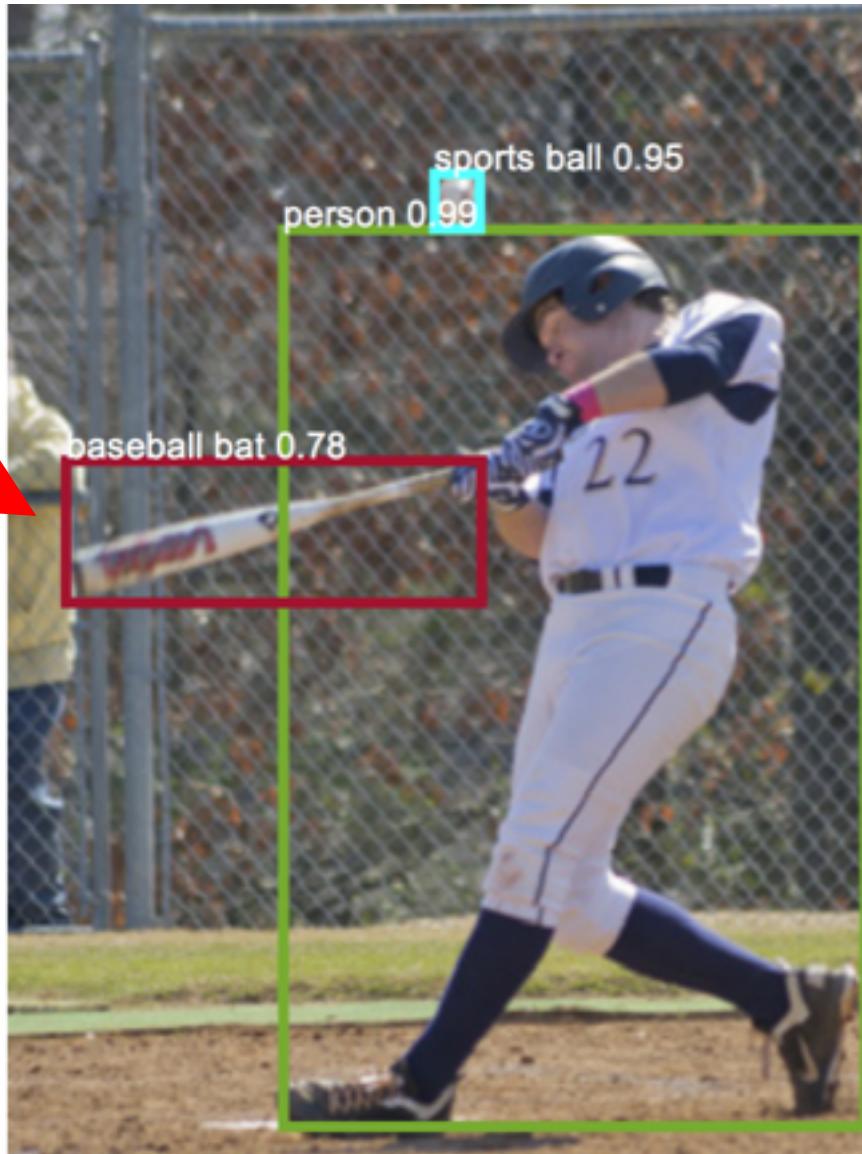
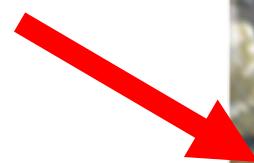


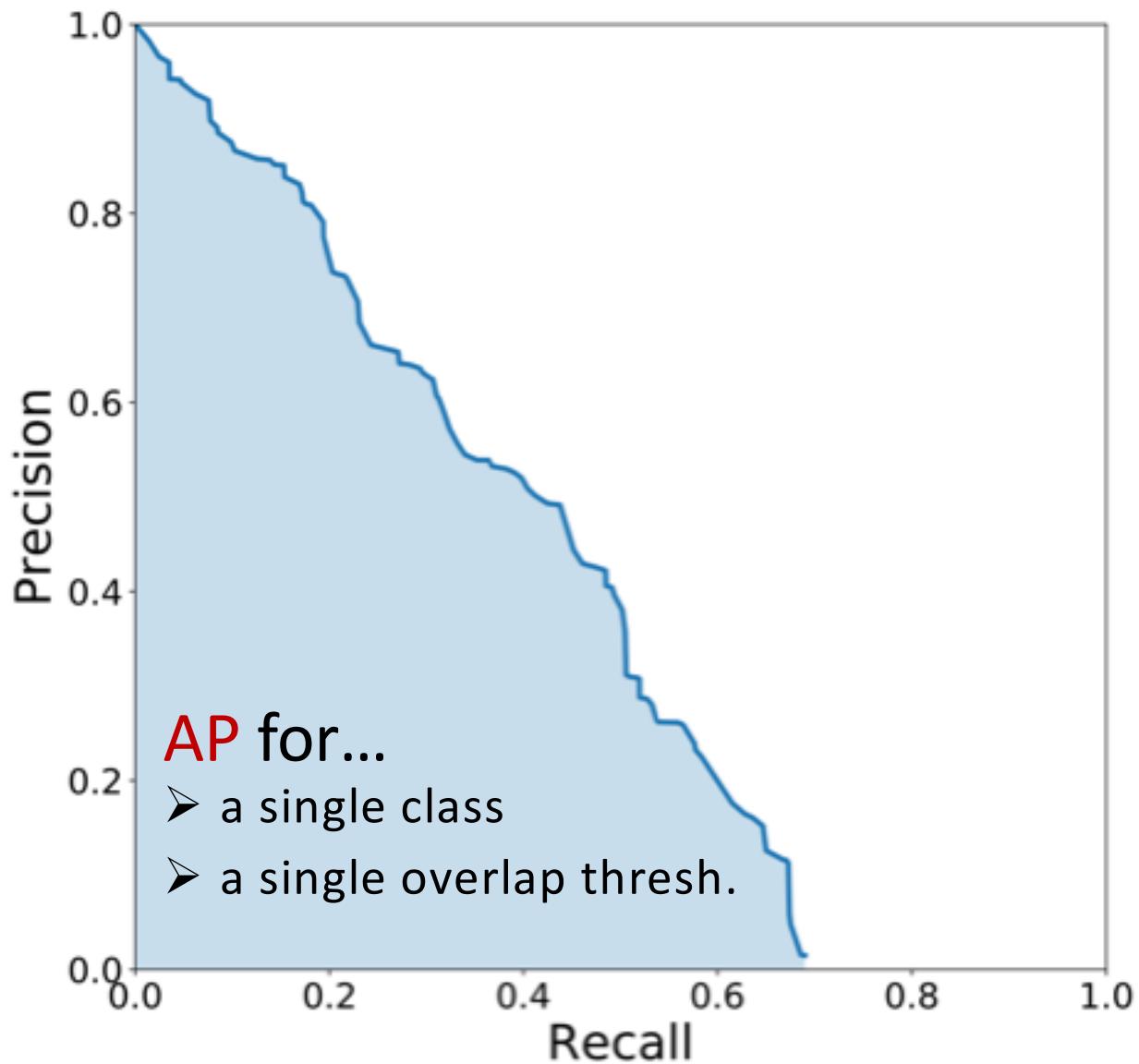




False Positive  
IoU too low  
(or duplicate)

**False Negative  
(missing detection)**





# COCO Average Precision (%)

AP averaged over...

- Object categories
- True positive overlap (IoU) requirement (from 0.5 to 0.95)

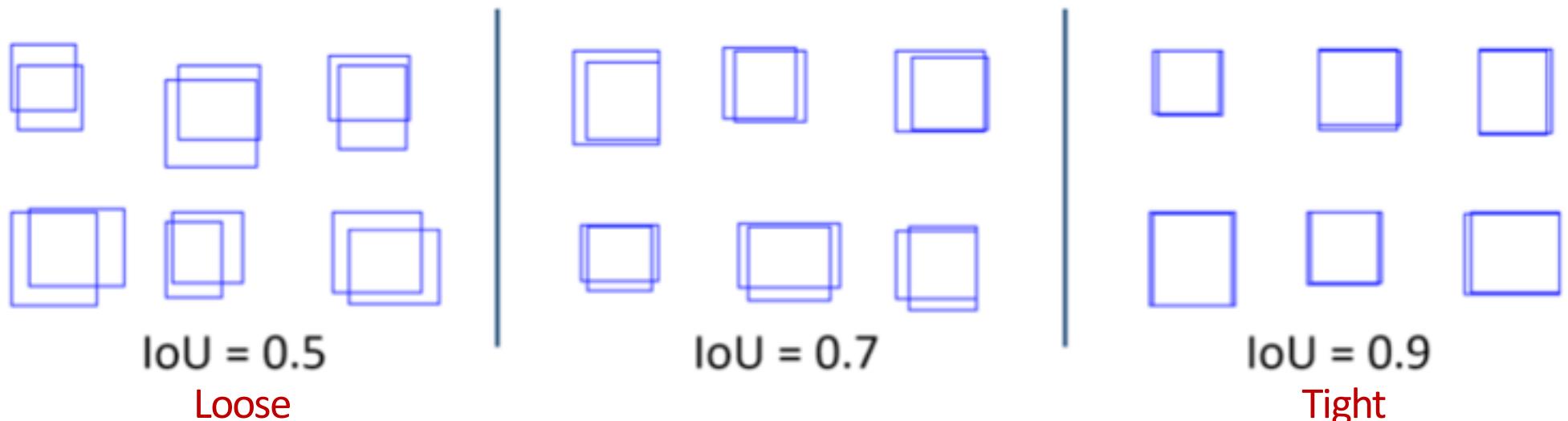


Figure credits: Dollár and Zitnick

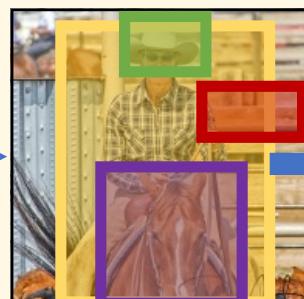
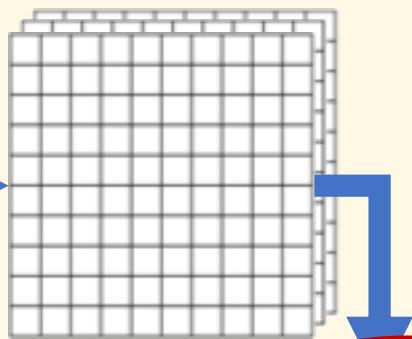
## R-FCN

Per-image computation

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



$$\text{RPN}(f_I)$$



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

PSRoIPool cls  
PSRoIPool box

AvgPool

Softmax clf.

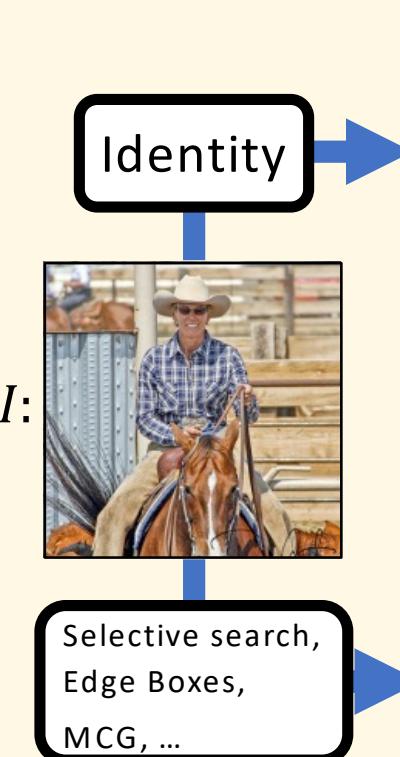
AvgPool

Box regressor

Extremely light-weight per-region computation

# Slow R-CNN in the Generalized Framework

Per-image computation



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



Crop & warp



$\text{ConvNet}(r_i)$



1-vs-rest SVMs

Box regressor