

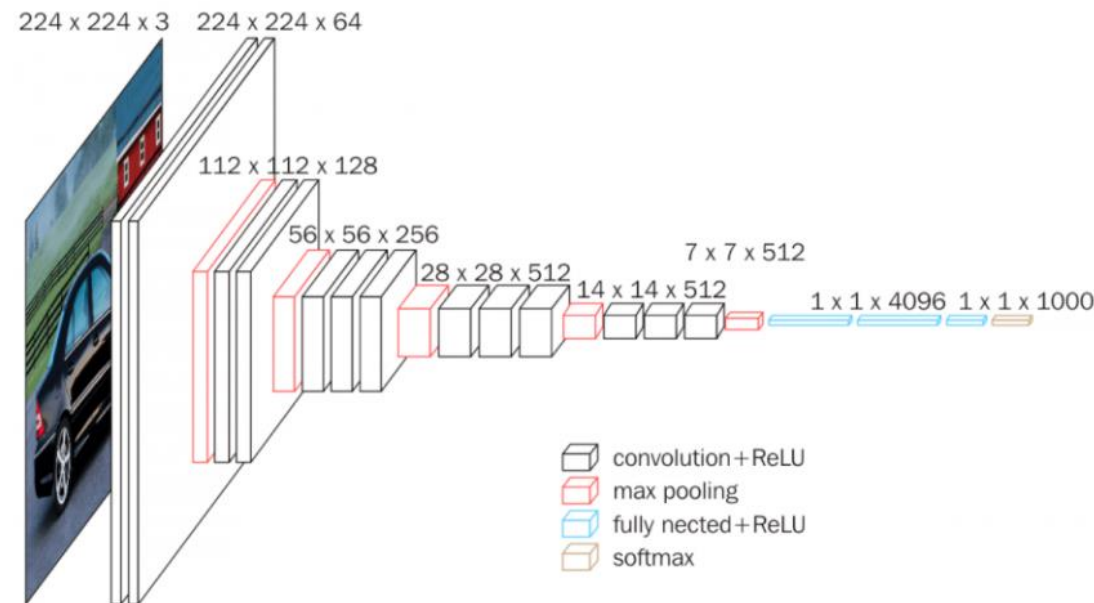
AI for Video Data Processing 3: R-CNN for Object Detection

Outline

1. R-CNN
2. Fast R-CNN
3. Faster R-CNN
4. Faster R-CNN with FPN
5. Mask R-CNN
6. RetinaNet

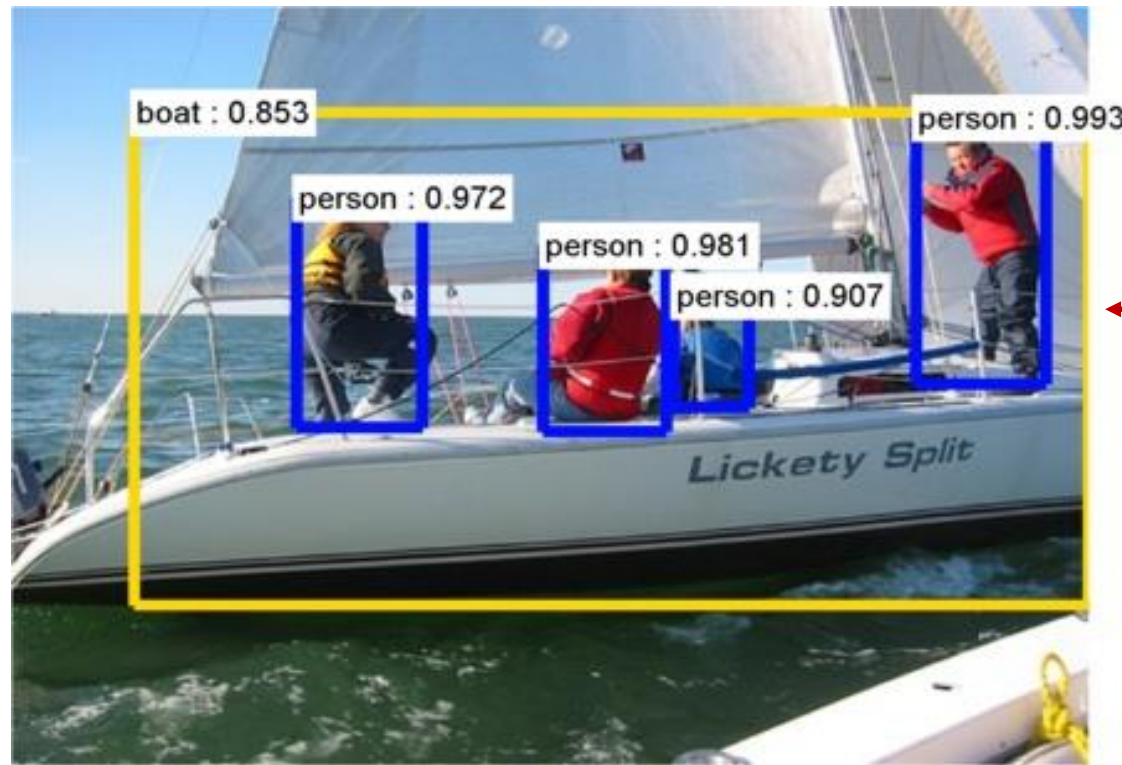
1. From CNN to R-CNN

Original CNNs are designed for Image Classification:
LeNet, AlexNet, VGG, ResNet, DenseNet, ...



1. From CNN to R-CNN

If we want to detect objects in an image, then ...



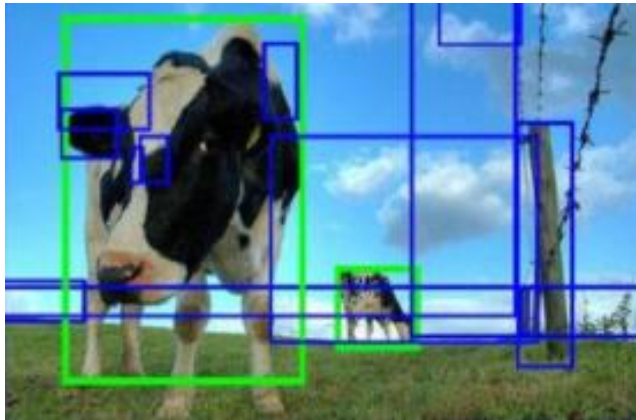
← What?

← Where?

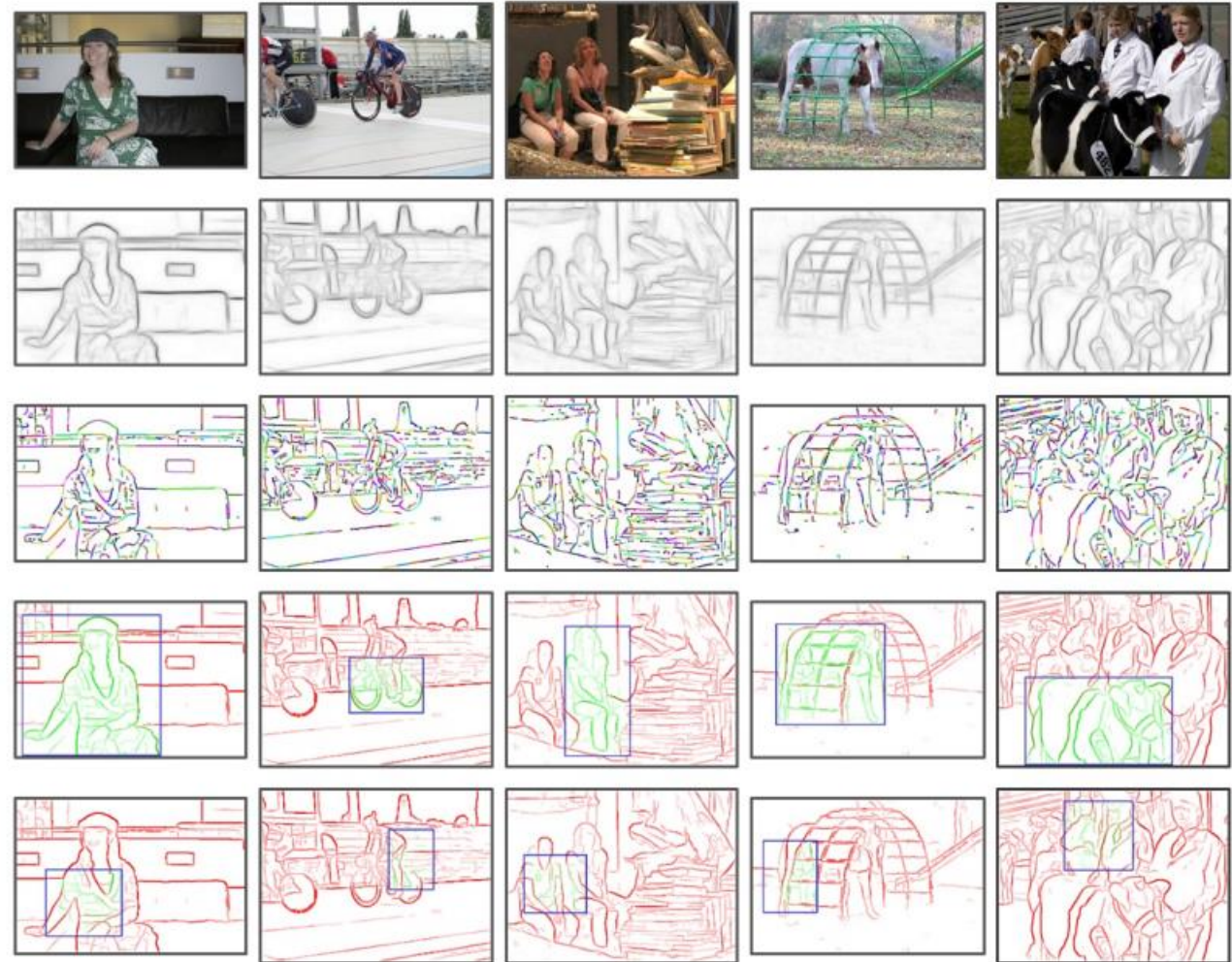
Object detection with bounding boxes

1. Object Proposal Algorithm in CV

Zitnick, C. Lawrence, and Piotr Dollár. "Edge boxes: Locating object proposals from edges." In *European Conference on Computer Vision*, pp. 391-405. Springer, Cham, 2014.



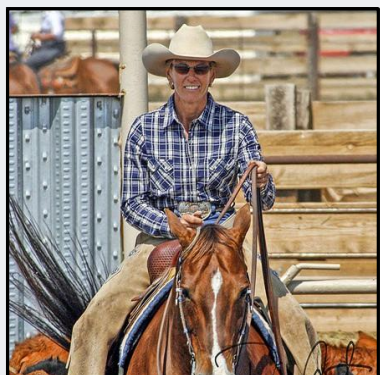
For example: 2k object proposals for one image



1.R-CNN (Region-based Convolutional Neural Net)

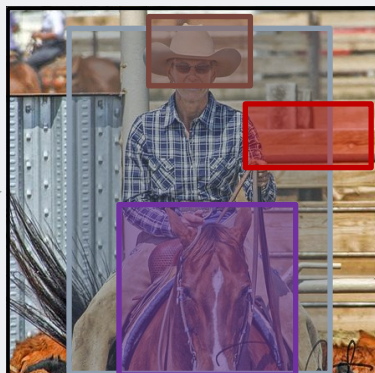
Per-image computation

Per-region computation



Selective search,
Edge Boxes,
MCG, ...

1

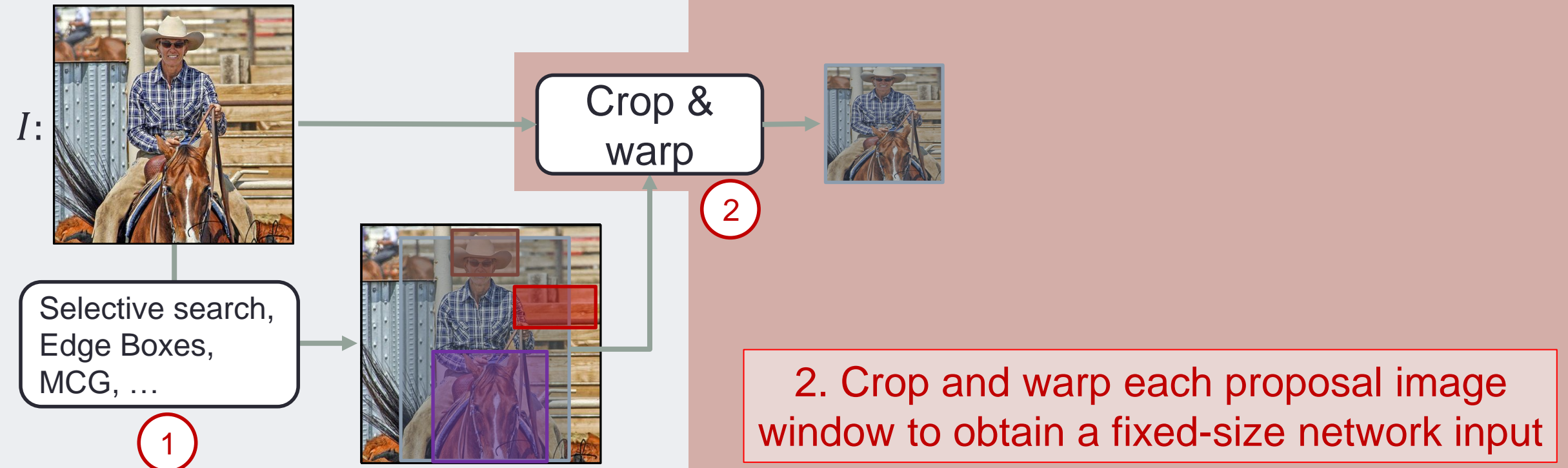


1. Use an off-the-shelf *Region of Interest* (RoI) proposal algorithm (~2k proposals per image)

1. R-CNN

Per-image computation

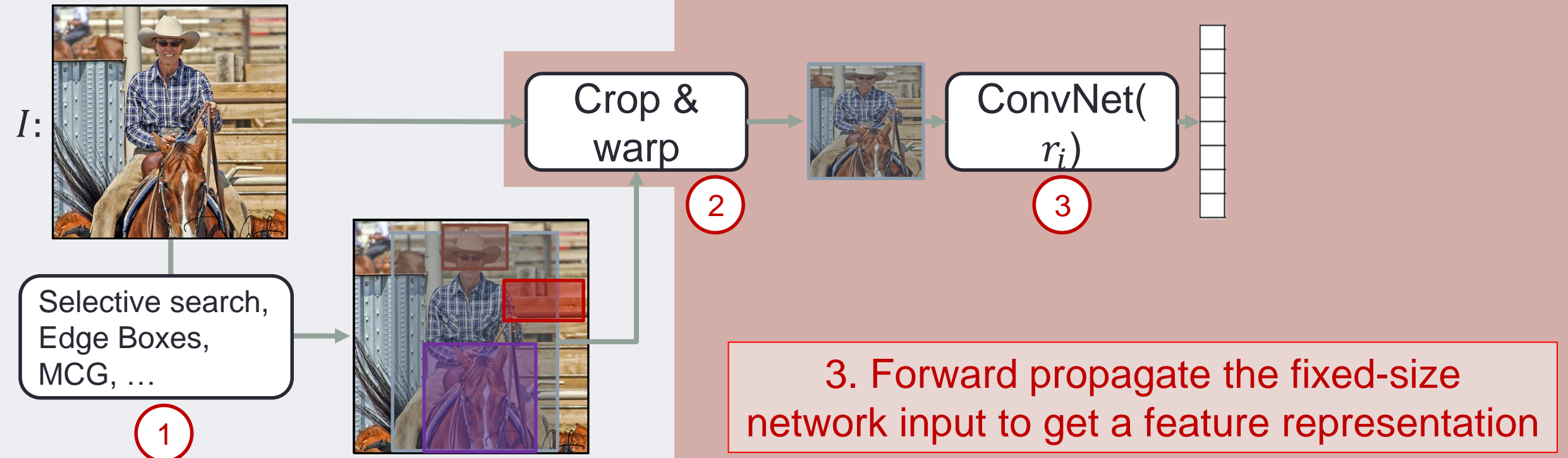
Per-region computation for each $r_i \in r(I)$



1. R-CNN

Per-image computation

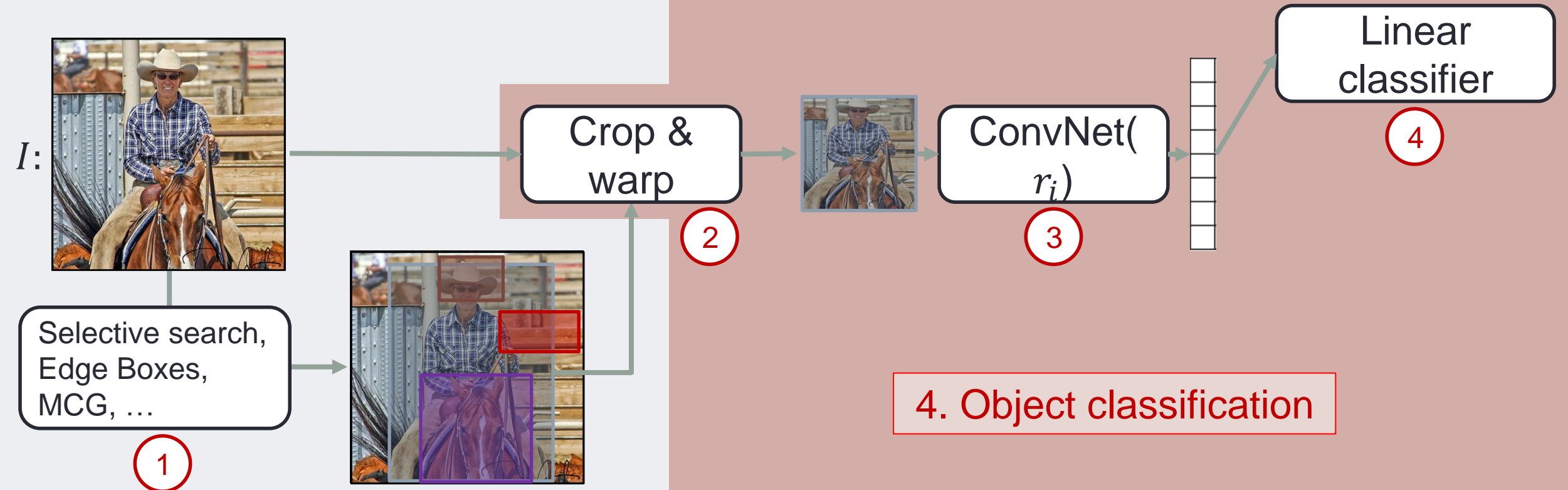
Per-region computation for each $r_i \in r(I)$



1. R-CNN

Per-image computation

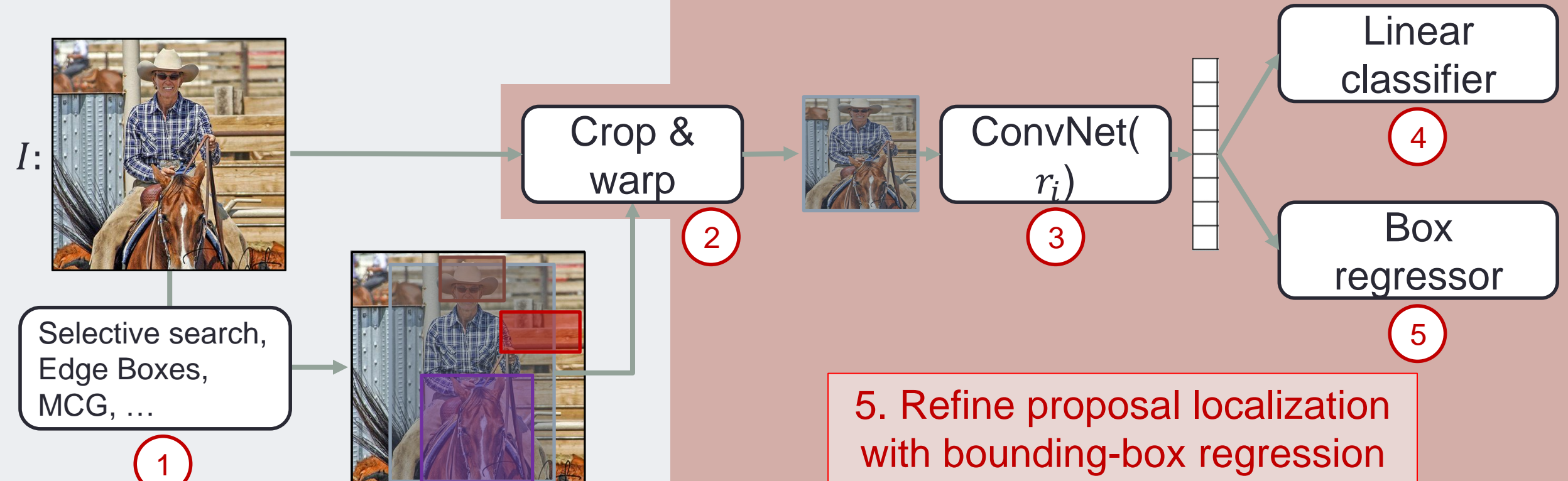
Per-region computation for each $r_i \in r(I)$



1. R-CNN

Per-image computation

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

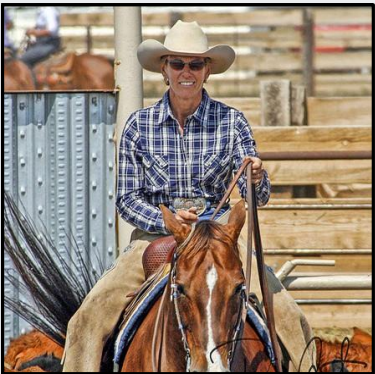


1. From R-CNN to Generalized R-CNN Framework

1. Generalized R-CNN Framework

Per-image computation

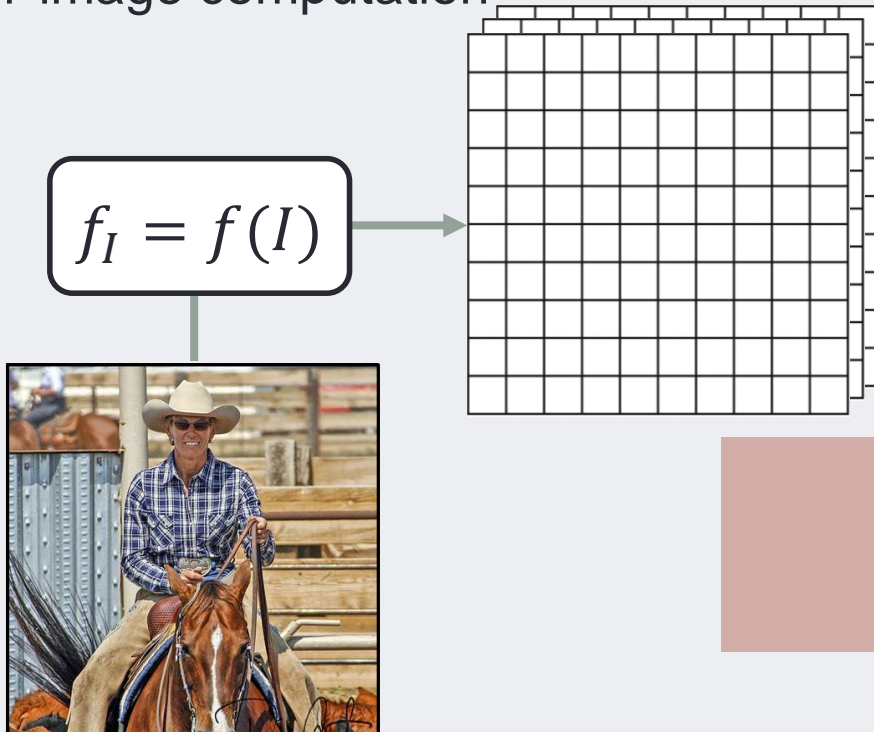
Per-region computation for each $r_i \in r(I)$



Input image
per-image operations | per-region operations

1. Generalized R-CNN Framework

Per-image computation



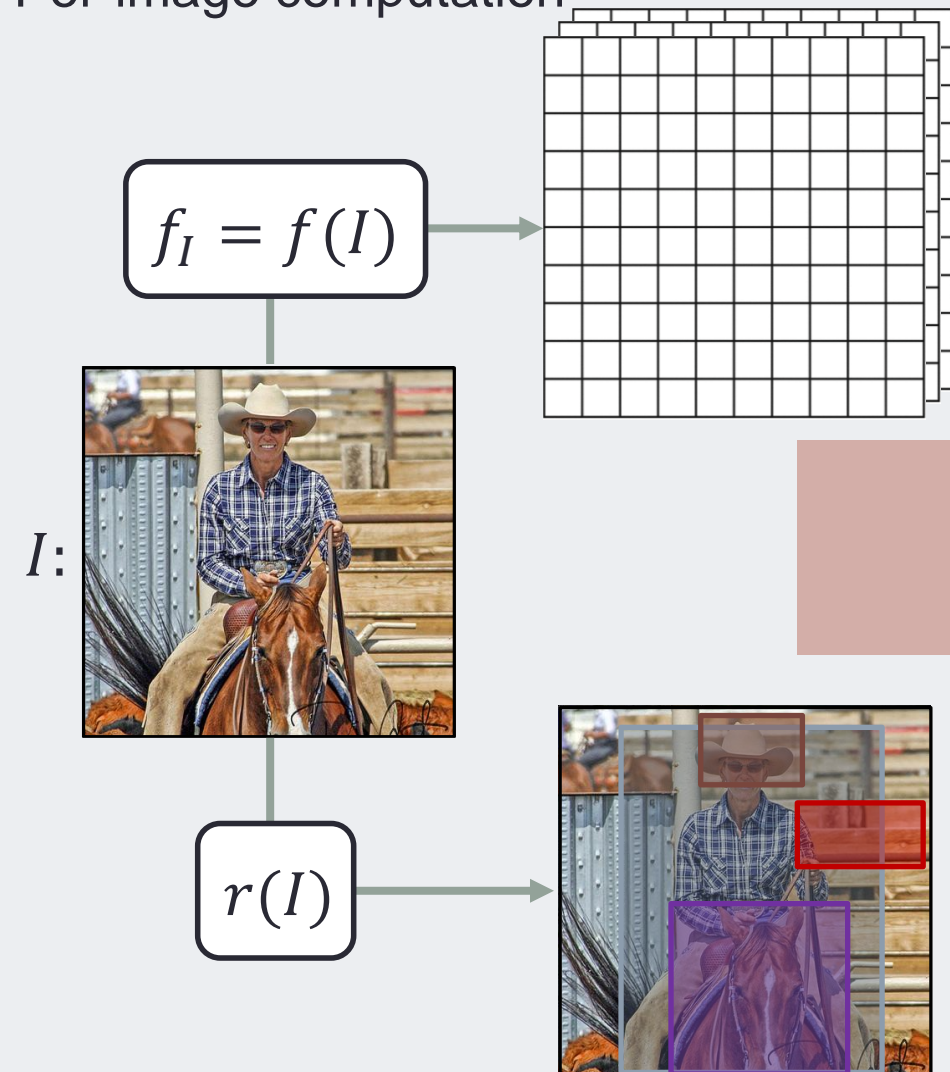
Per-region computation for each $r_i \in r(I)$

Transformation of the input image
into a featurized representation

1. Generalized R-CNN Framework

Per-image computation

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

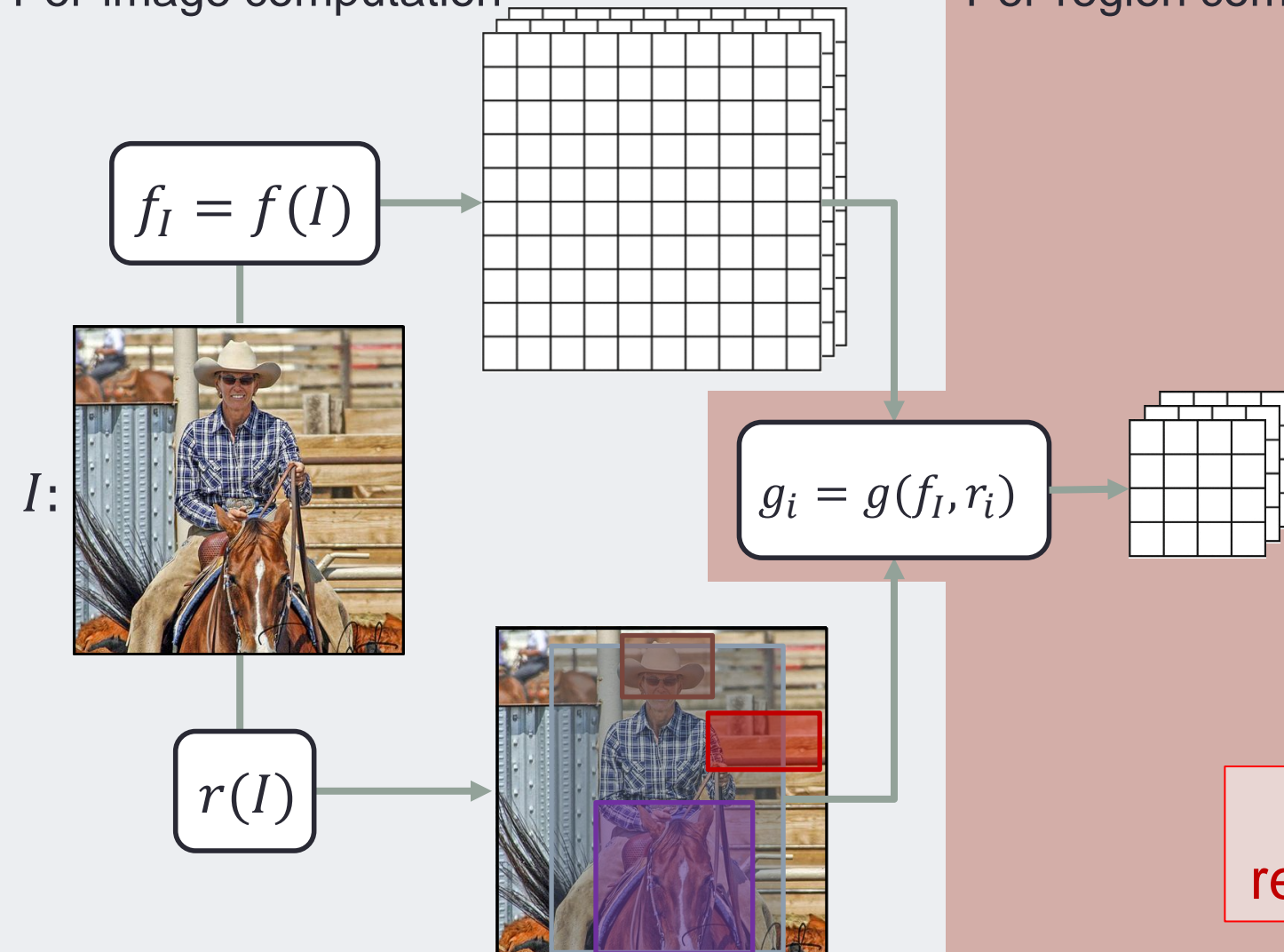


*Region of Interest proposals
computed for the image*

1. Generalized R-CNN Framework

Per-image computation

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

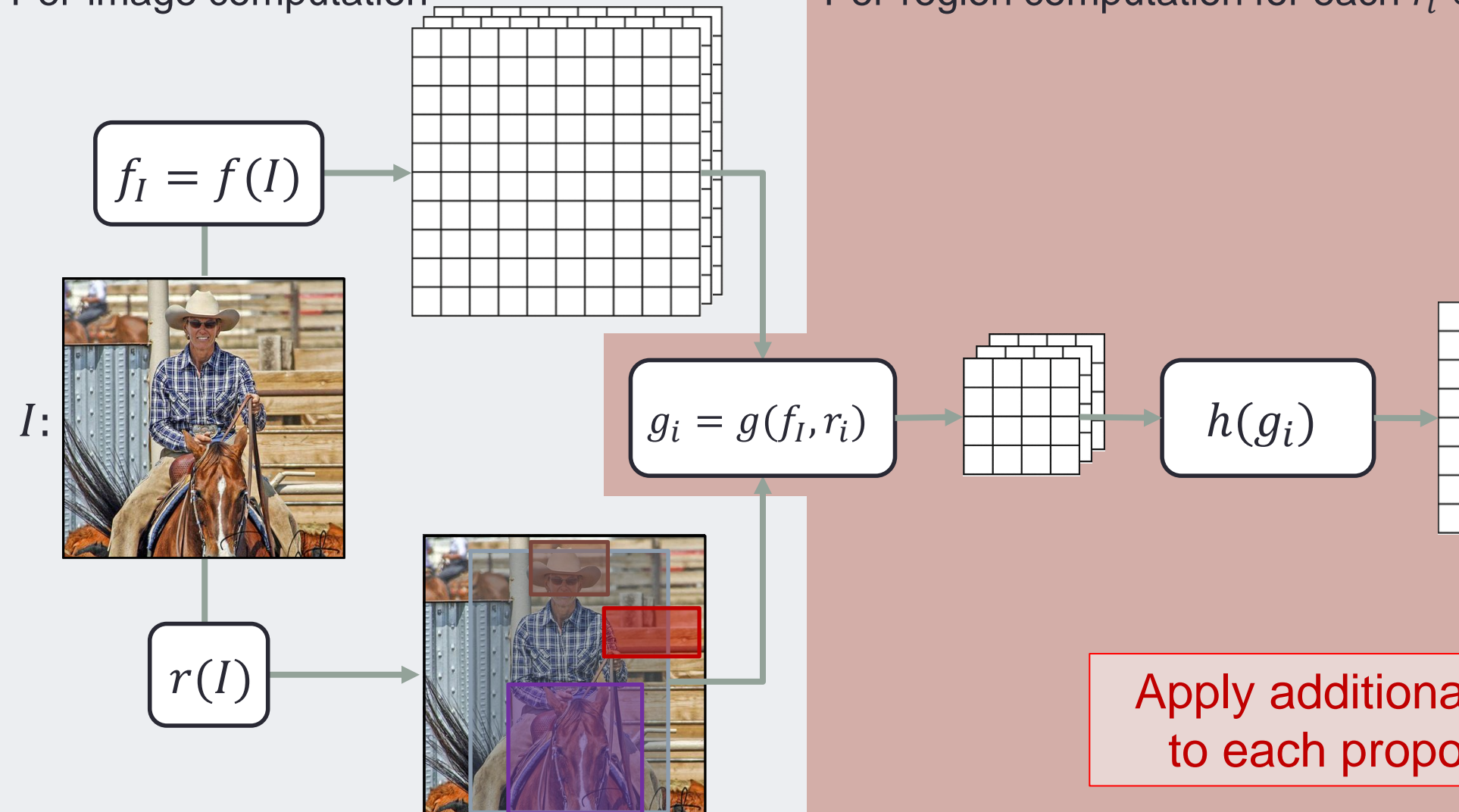


Compute a featurized representation of each proposal

1. Generalized R-CNN Framework

Per-image computation

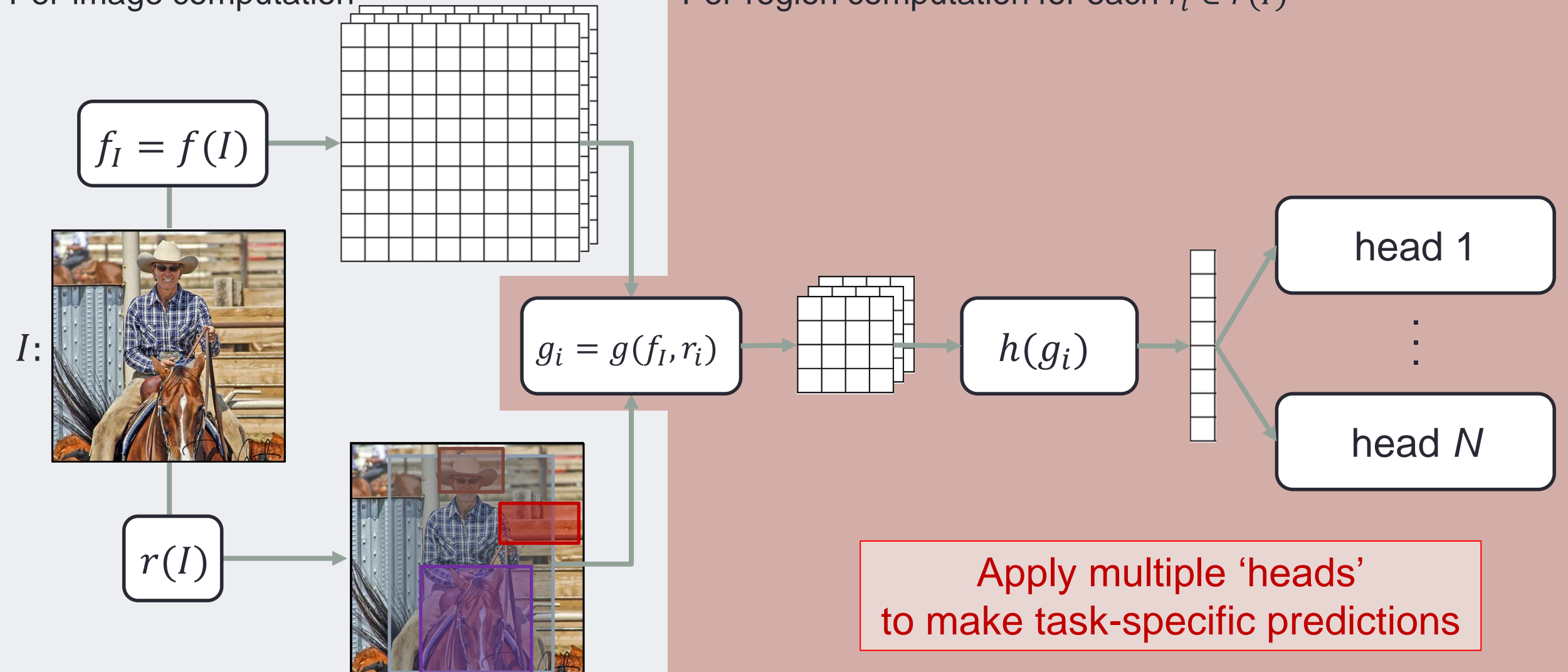
Per-region computation for each $r_i \in r(I)$



1. Generalized R-CNN Framework

Per-image computation

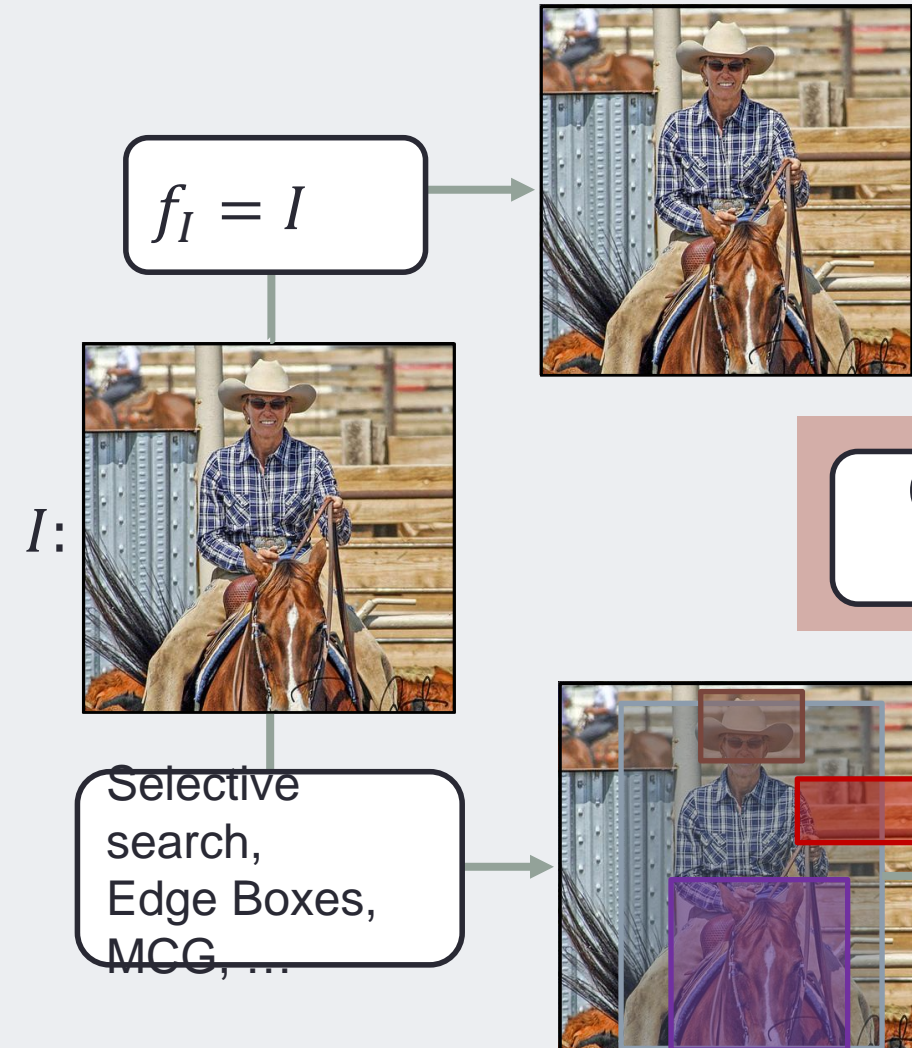
Per-region computation for each $r_i \in r(I)$



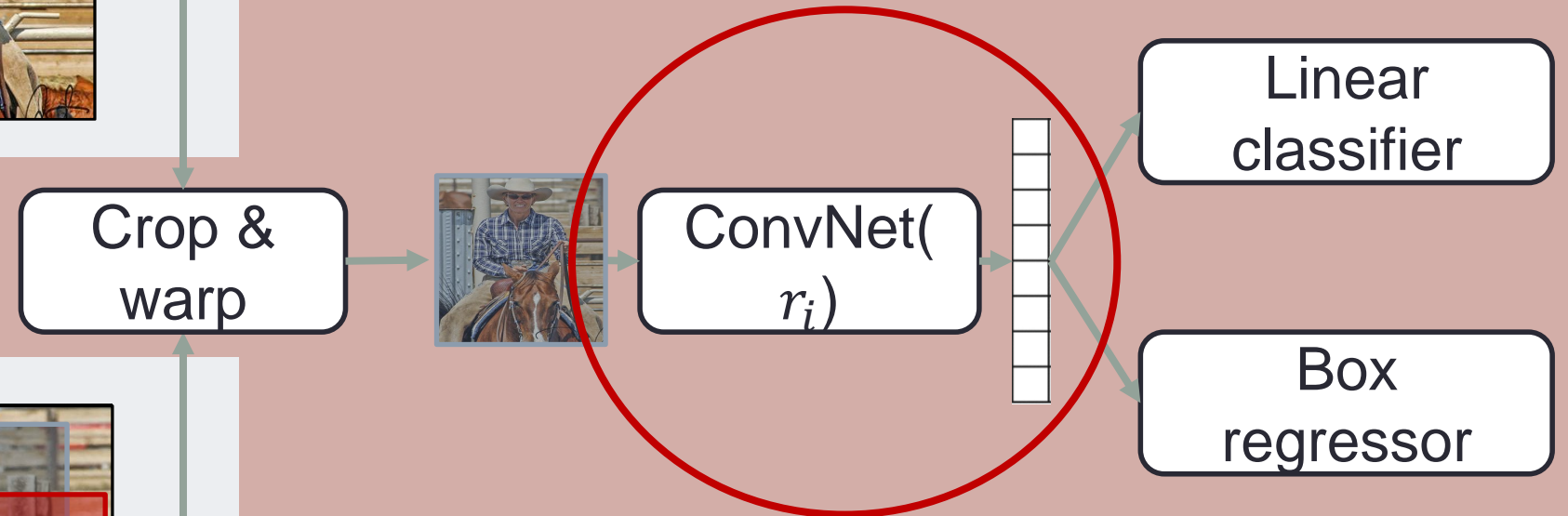
2. Fast R-CNN

“Slow” R-CNN

Per-image computation



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

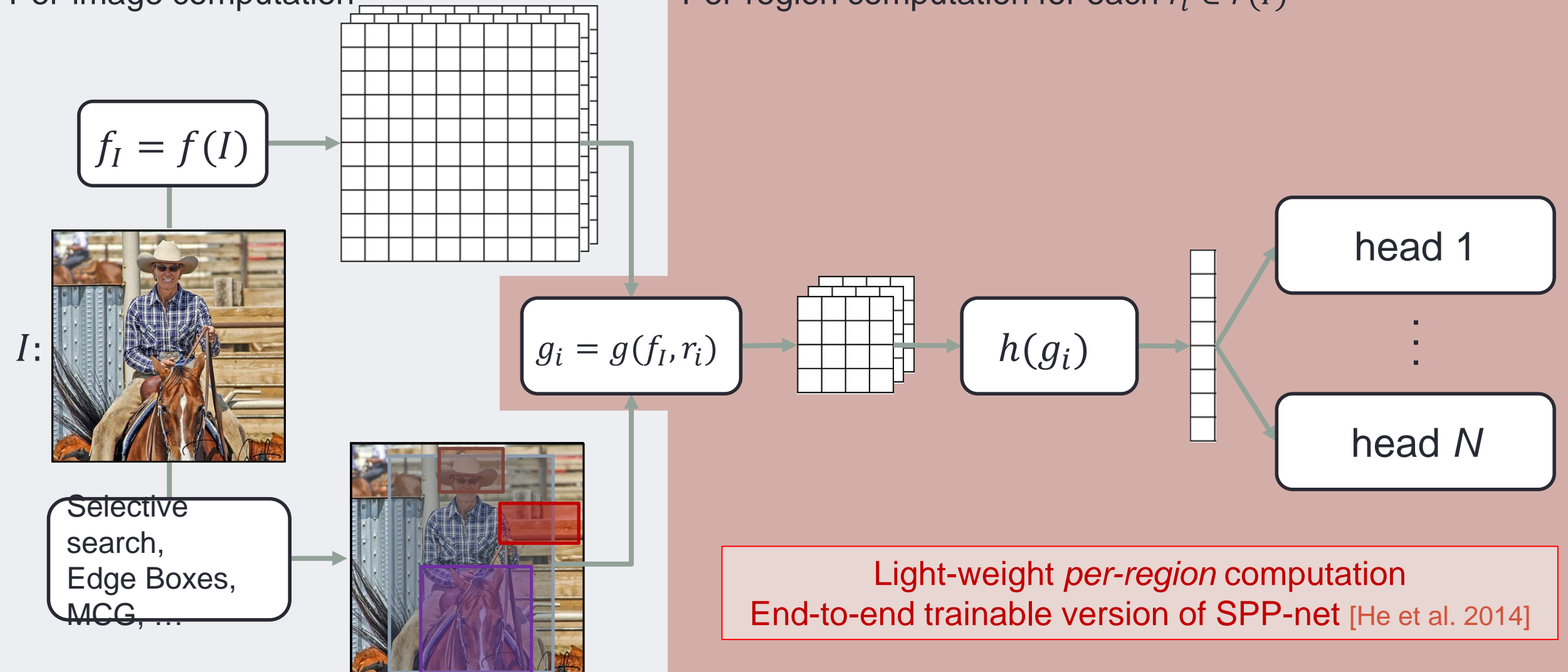


Very heavy *per-region* computation
E.g., 2000 full network evaluations

2. Generalized R-CNN → Fast R-CNN

Per-image computation

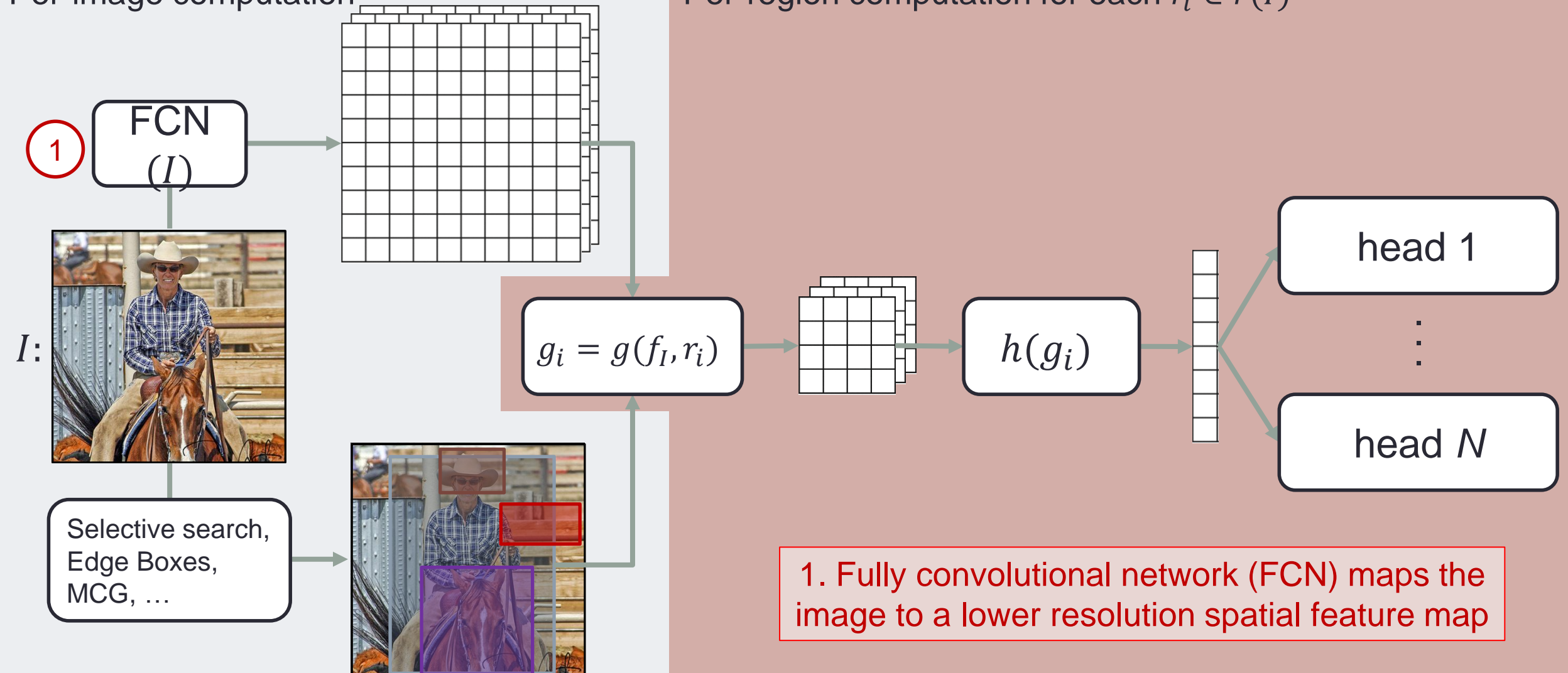
Per-region computation for each $r_i \in r(I)$



2. Generalized R-CNN → Fast R-CNN

Per-image computation

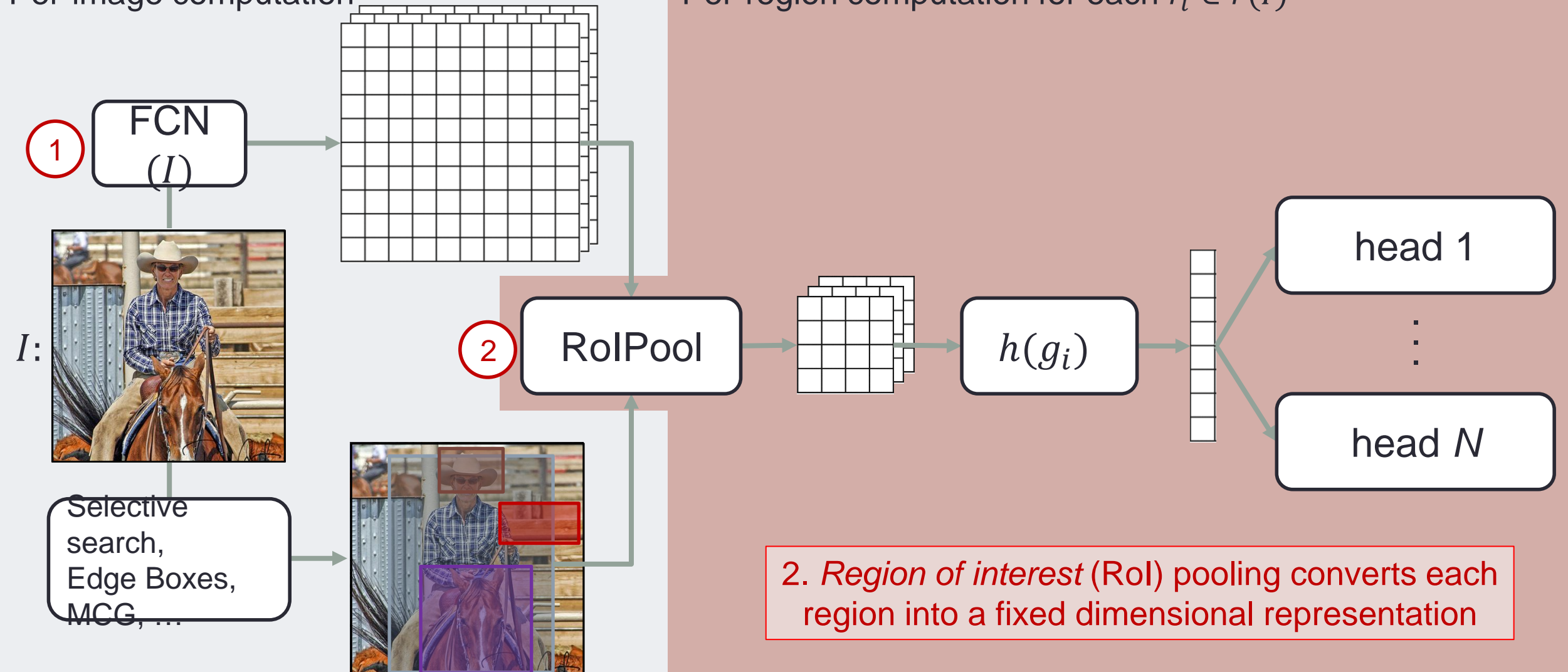
Per-region computation for each $r_i \in r(I)$



2. Generalized R-CNN → Fast R-CNN

Per-image computation

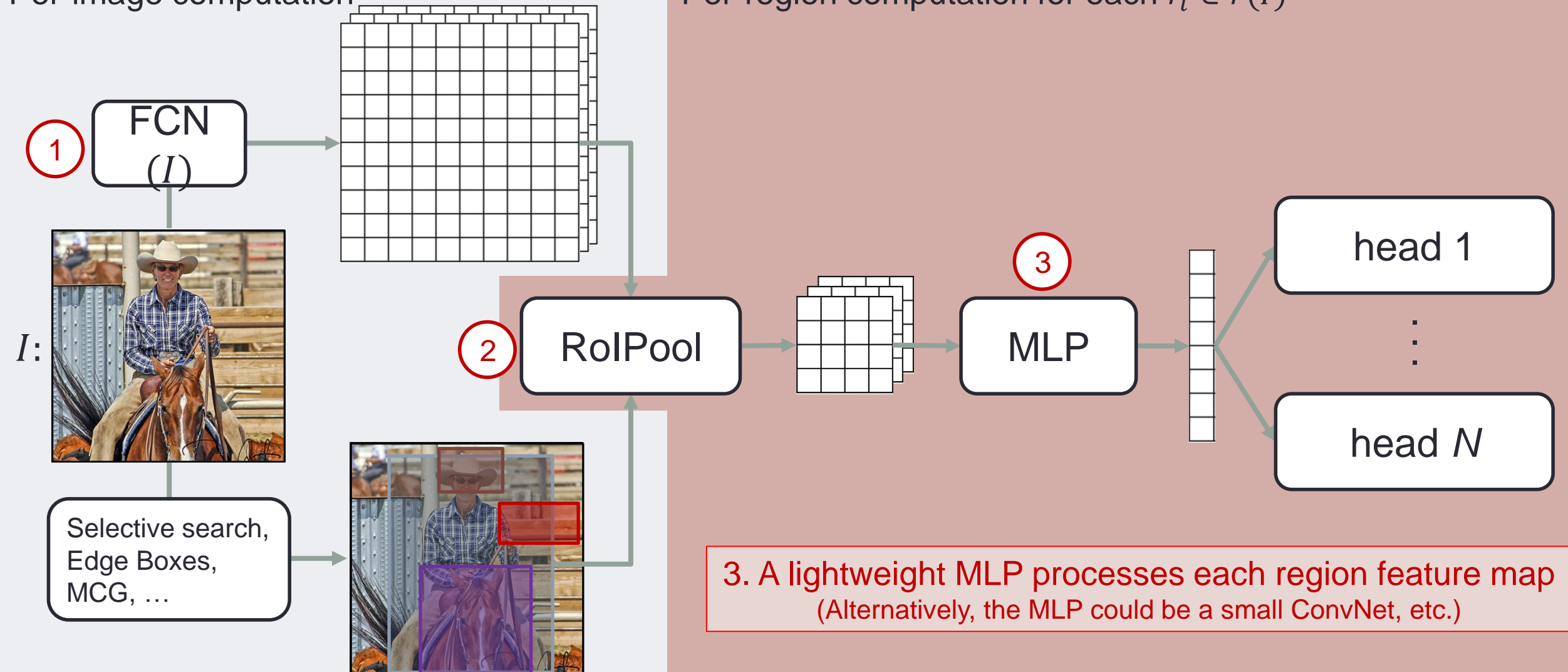
Per-region computation for each $r_i \in r(I)$



2. Generalized R-CNN → Fast R-CNN

Per-image computation

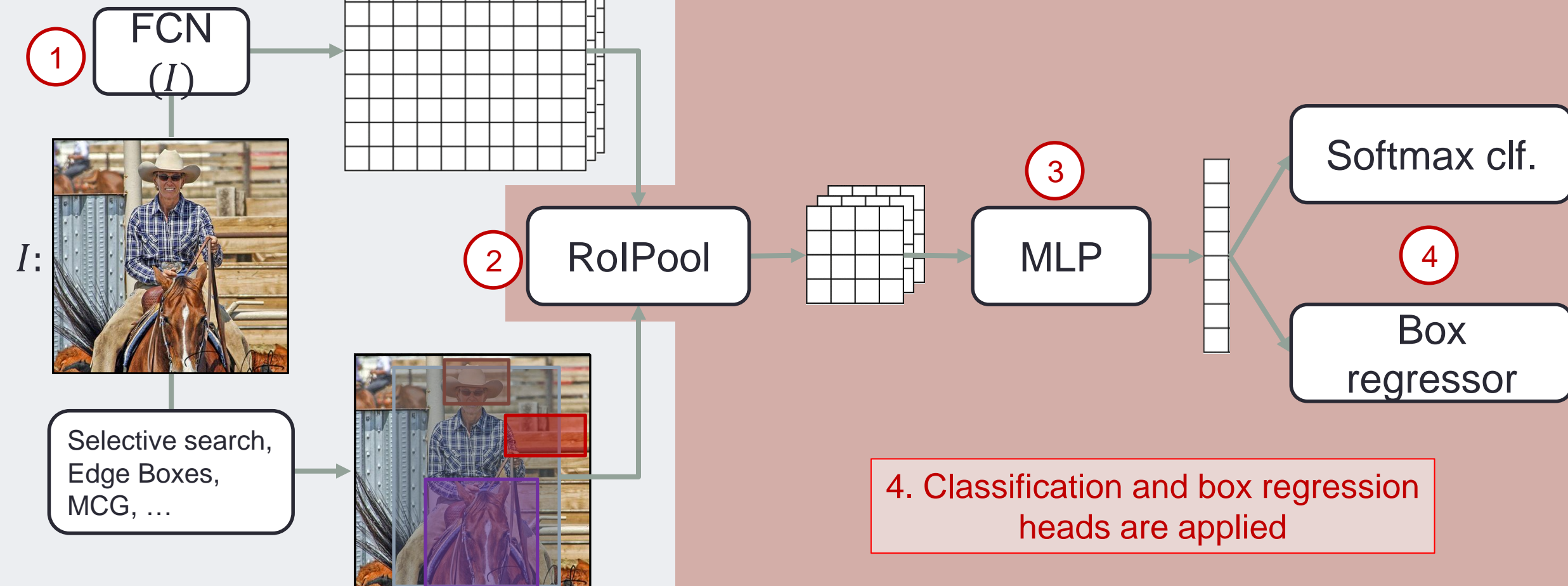
Per-region computation for each $r_i \in r(I)$



2. Generalized R-CNN → Fast R-CNN

Per-image computation

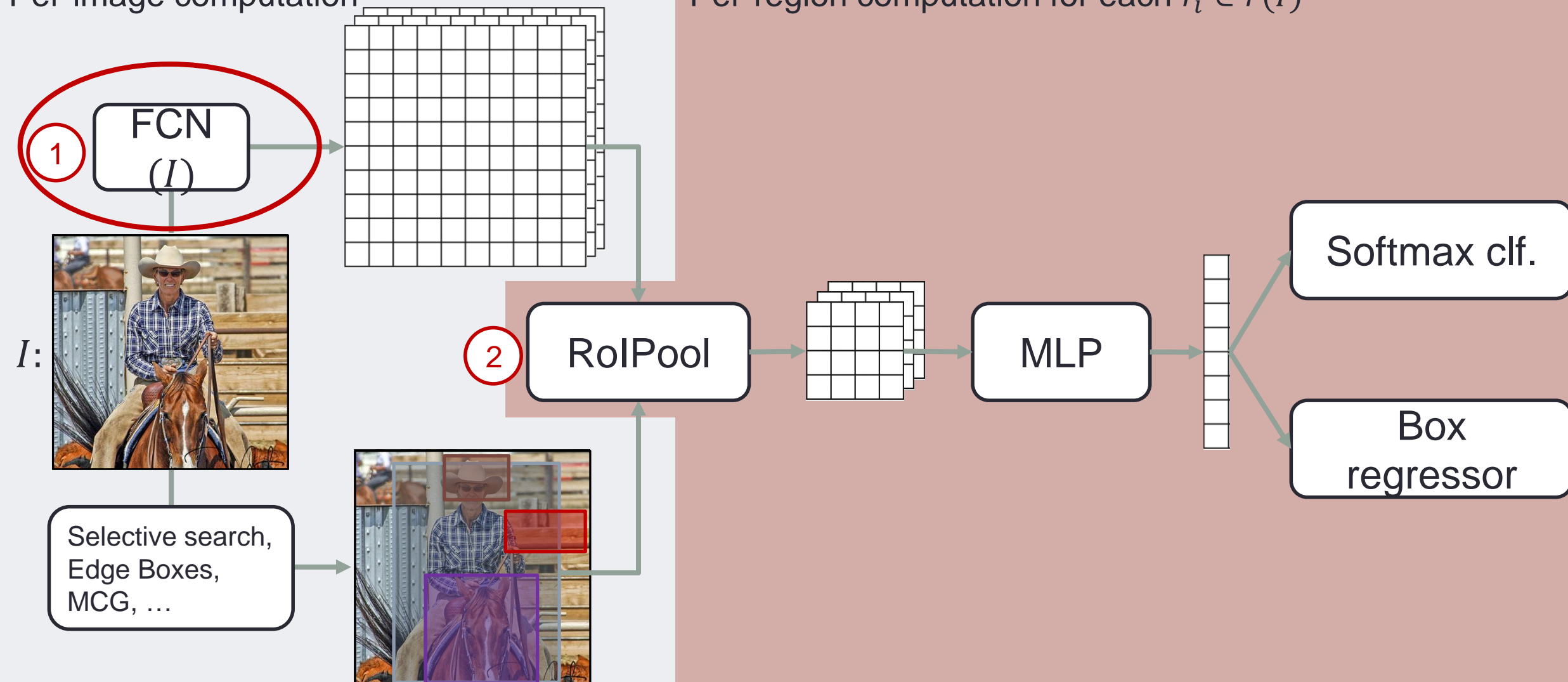
Per-region computation for each $r_i \in r(I)$



2. Fast R-CNN

Per-image computation

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



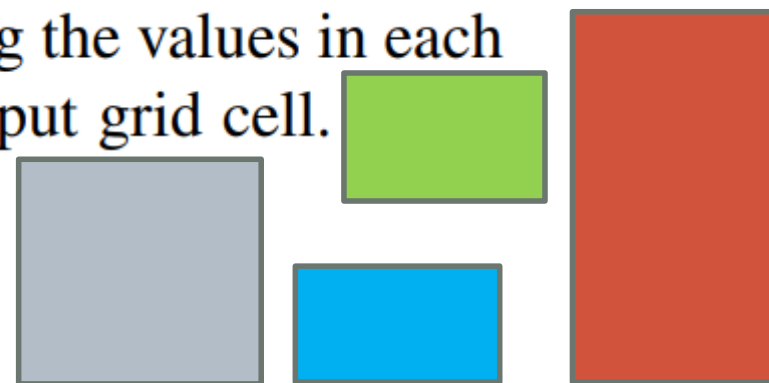
2. Fast R-CNN-RoI Pooling



2. Fast R-CNN-RoI Pooling

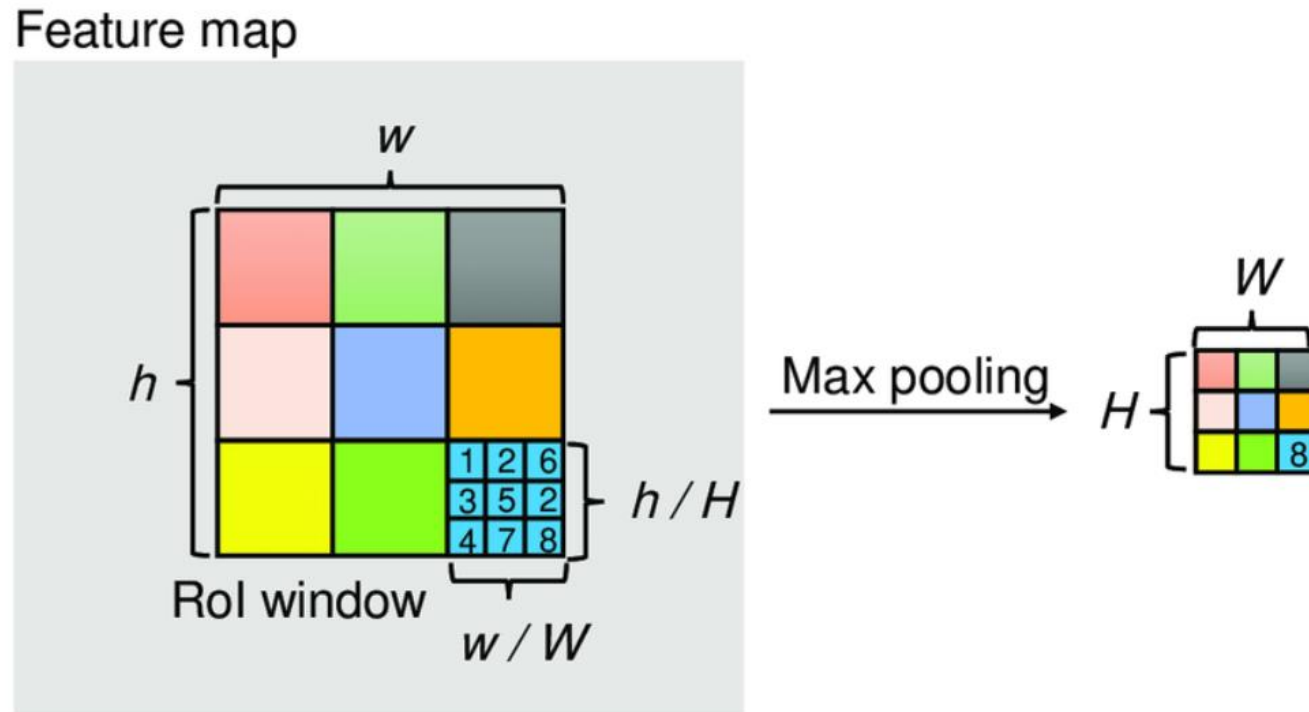
Transform **arbitrary size proposal** into a **fixed-dimensional** representation (e.g., $H=7 \times W=7$)

RoI max pooling works by dividing the $h \times w$ RoI window into an $H \times W$ grid of sub-windows of approximate size $h/H \times w/W$ and then max-pooling the values in each sub-window into the corresponding output grid cell.



2. Fast R-CNN-RoI Pooling

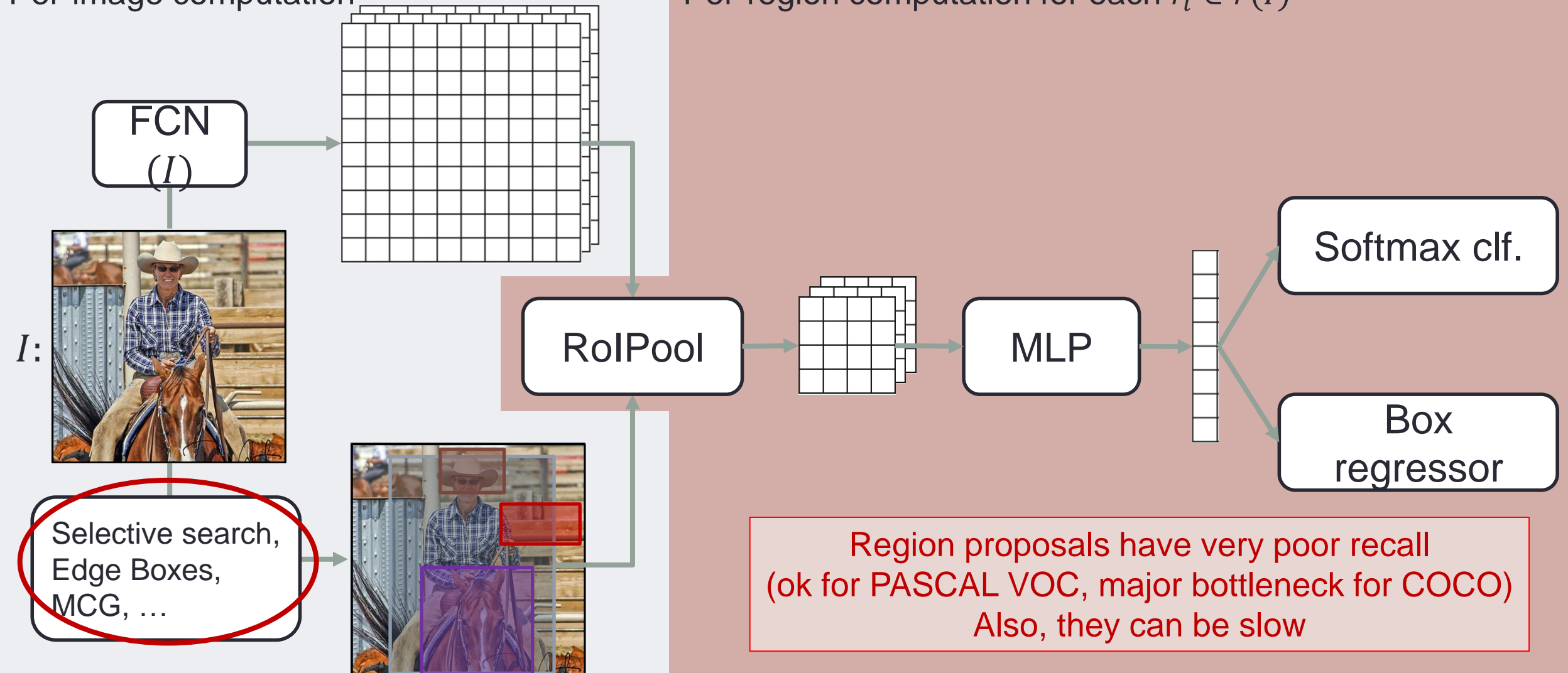
Transform **arbitrary size proposal** into a **fixed-dimensional** representation (e.g., $H=7 \times W=7$)



2. The Problem with Fast R-CNN

Per-image computation

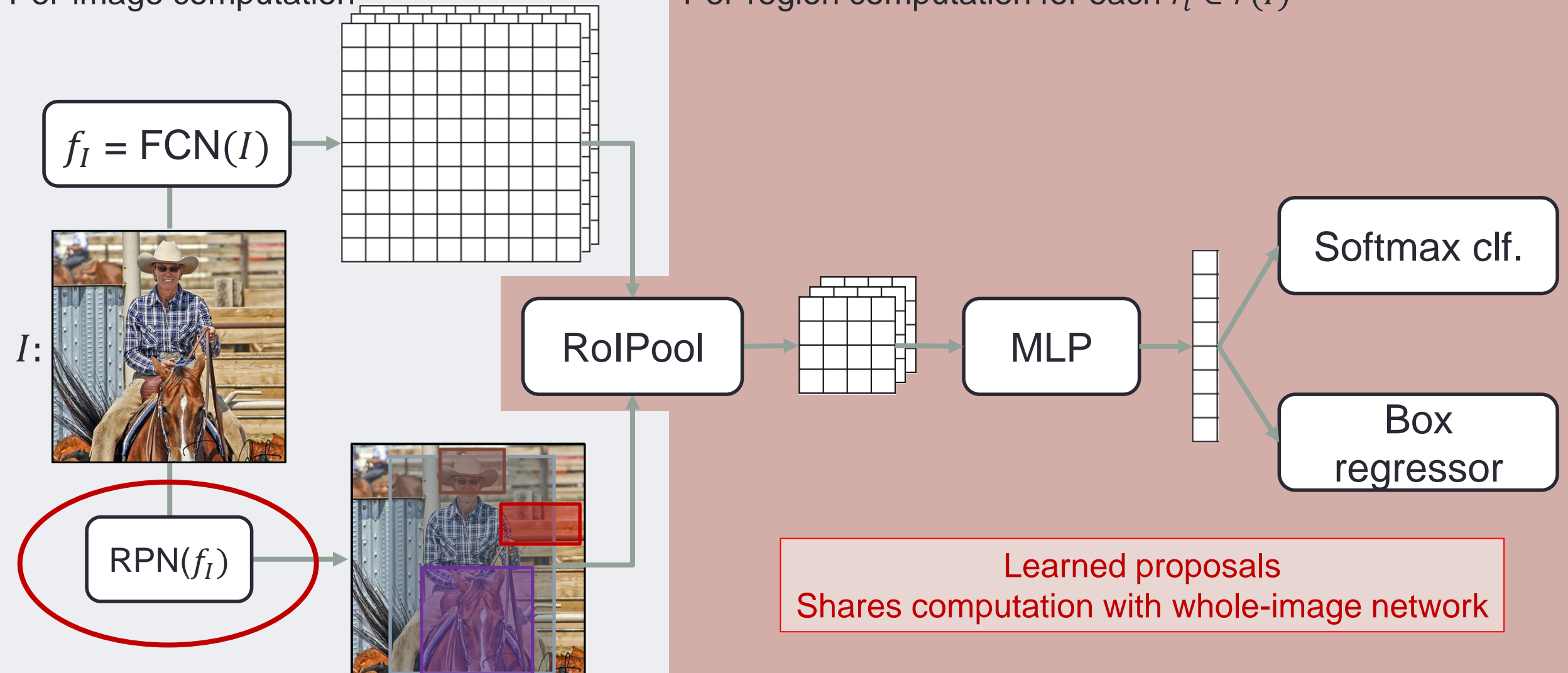
Per-region computation for each $r_i \in r(I)$



2. Faster R-CNN

Per-image computation

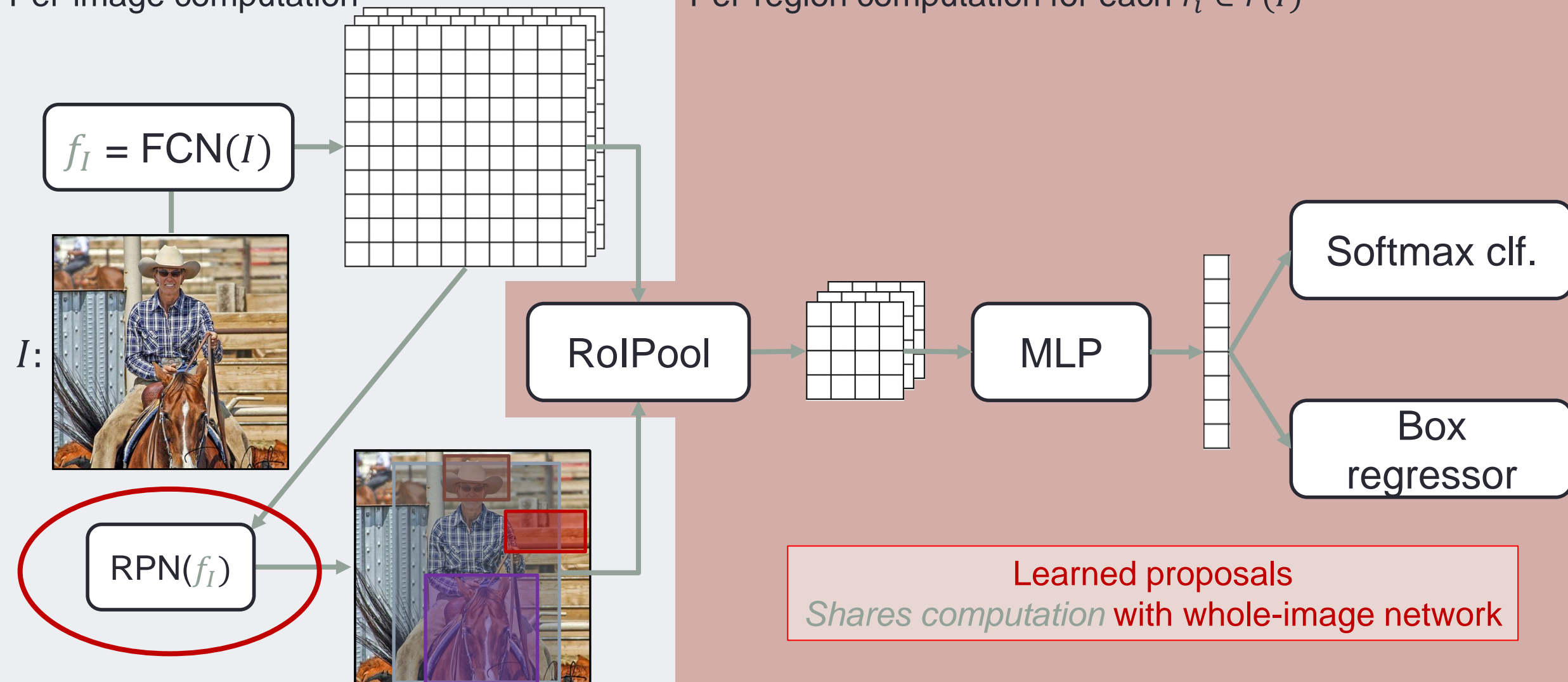
Per-region computation for each $r_i \in r(I)$



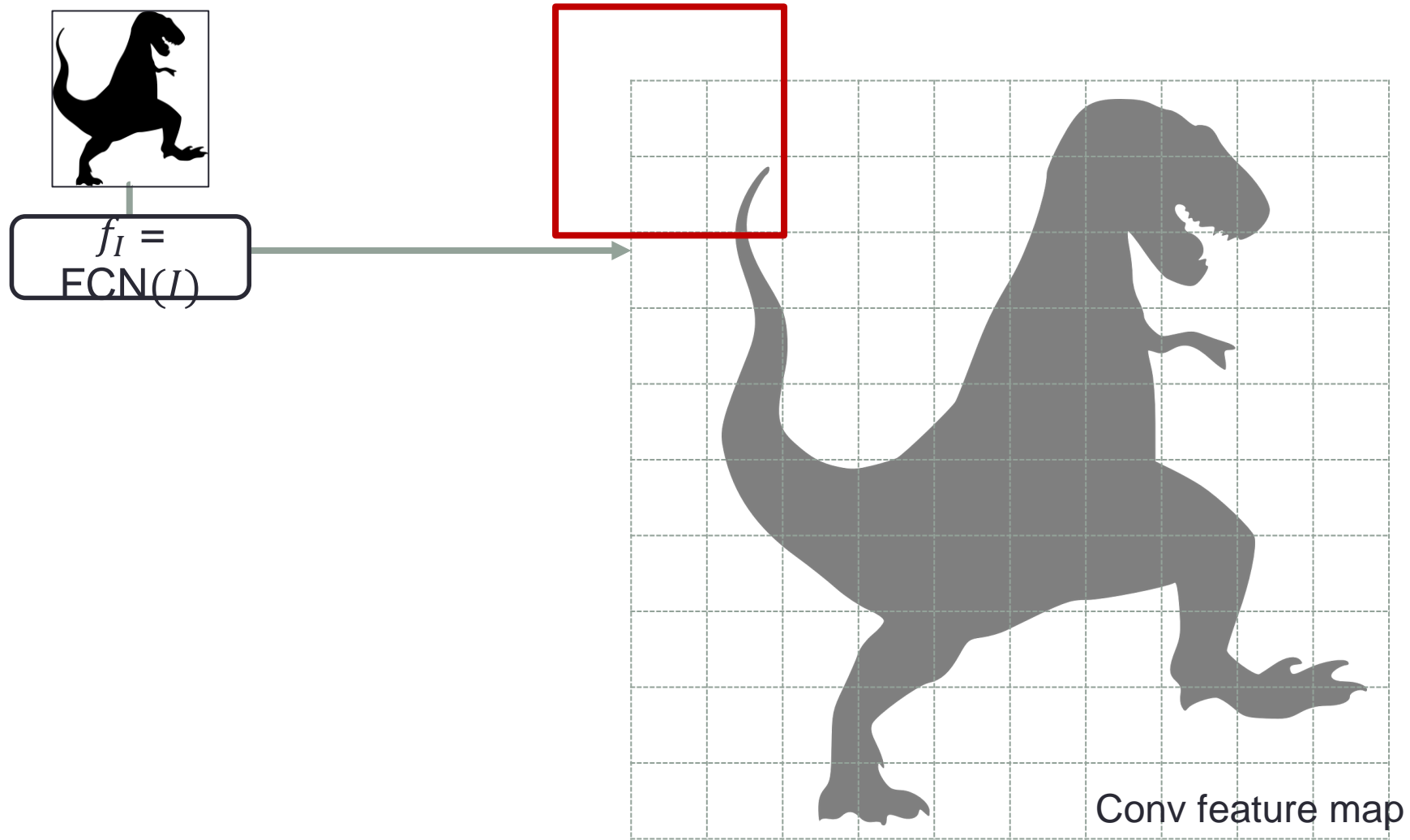
3. Faster R-CNN

Per-image computation

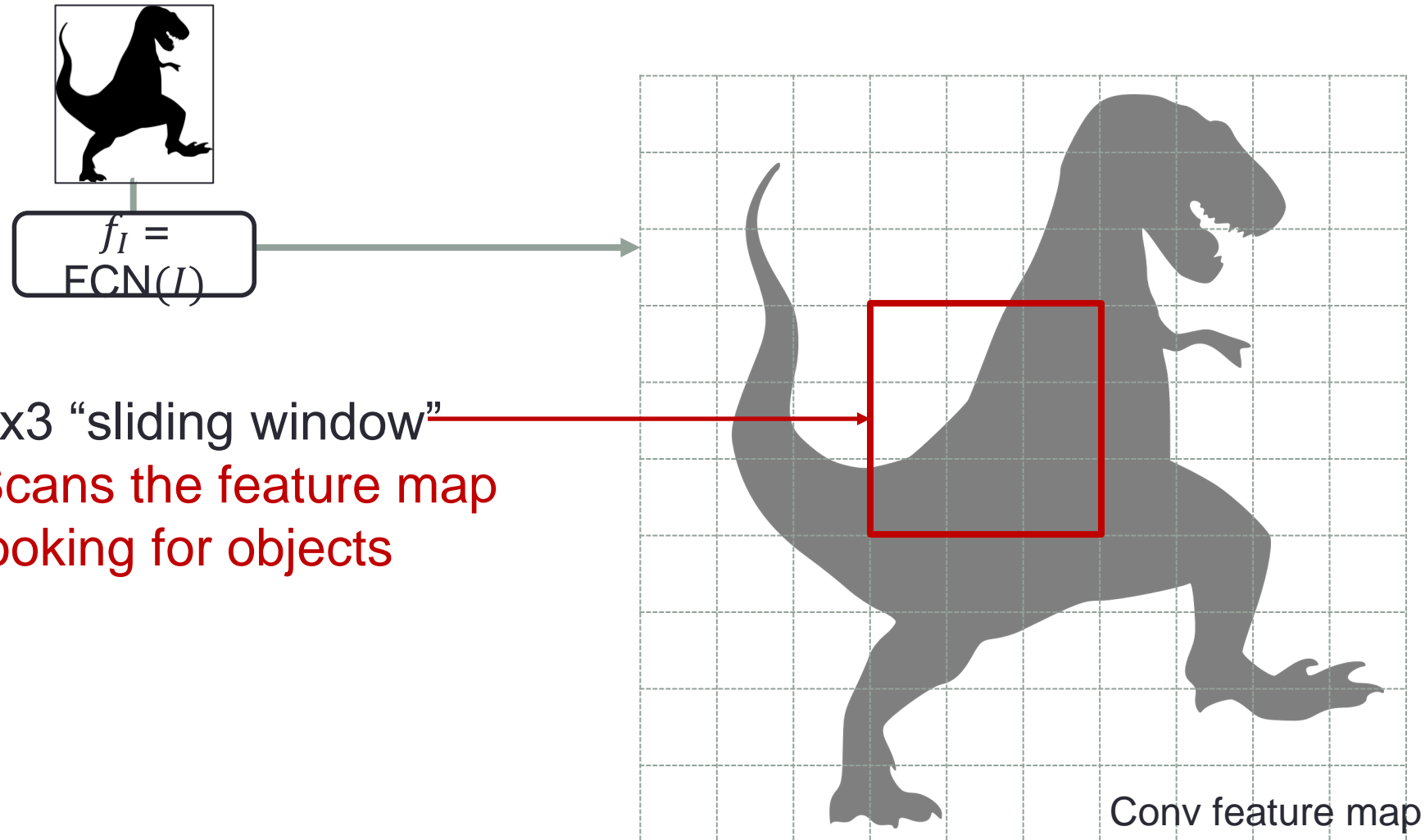
Per-region computation for each $r_i \in r(I)$



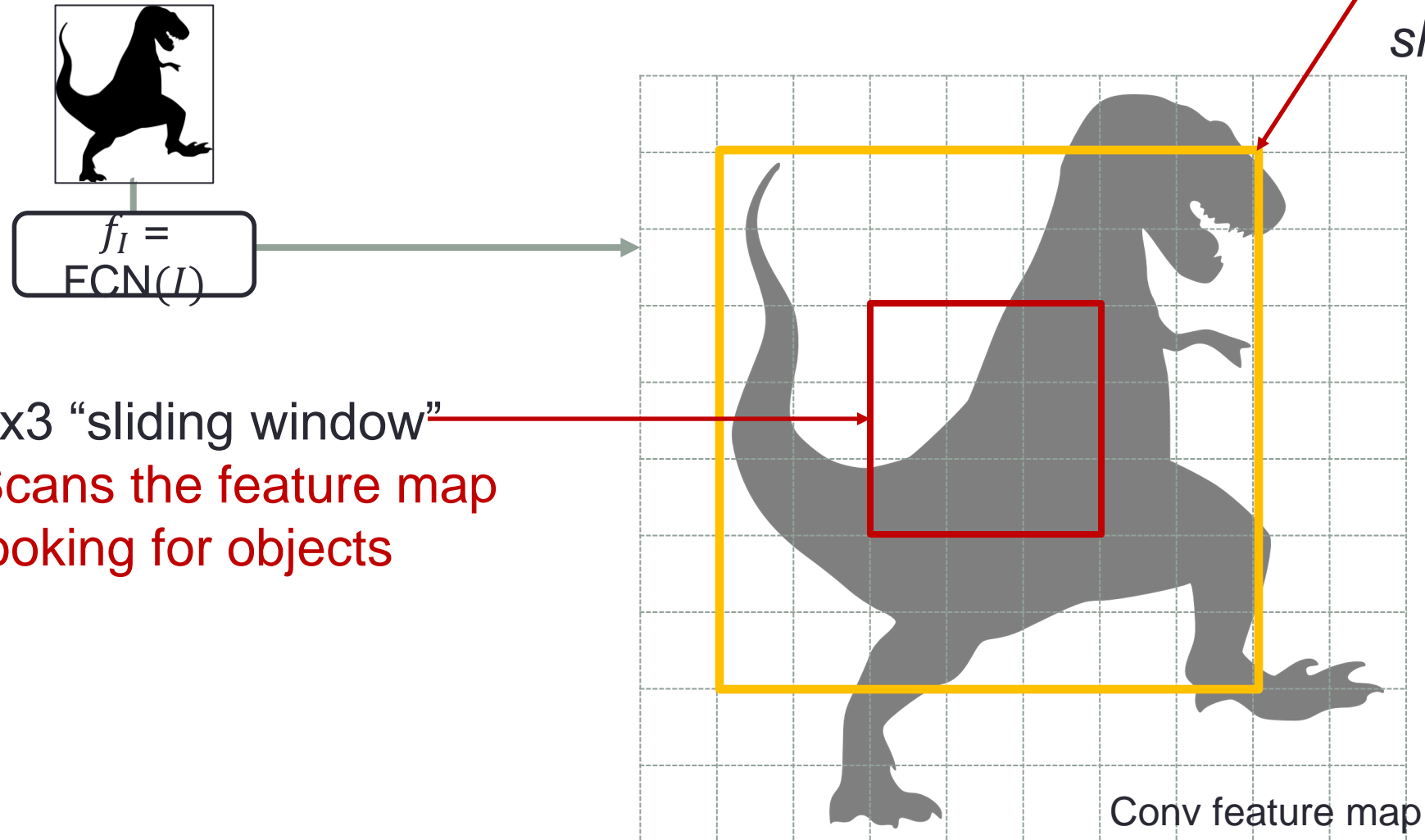
3. RPN: Region Proposal Network



3. RPN: Region Proposal Network



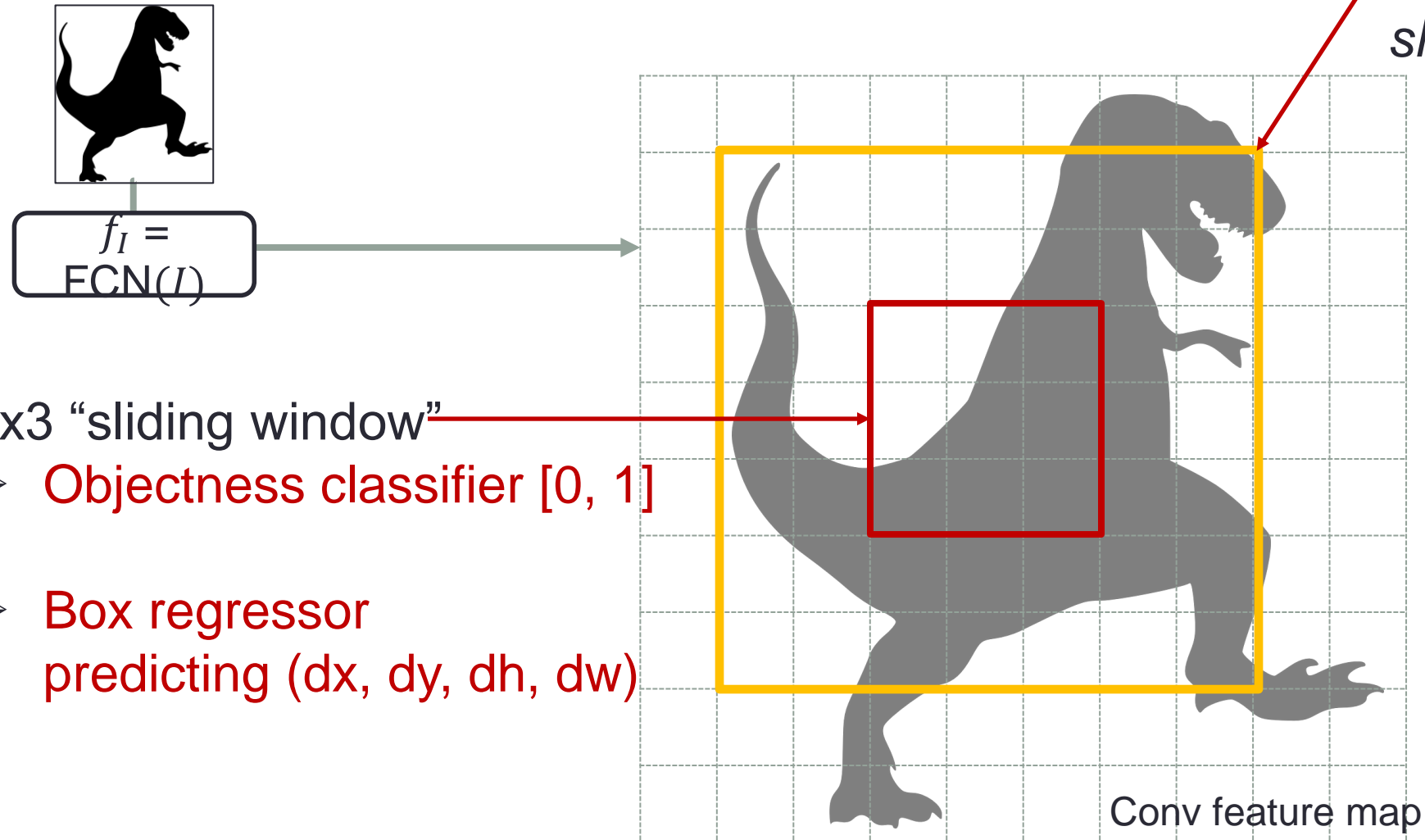
3. RPN: Anchor Box



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

3x3 "sliding window"
Scans the feature map
looking for objects

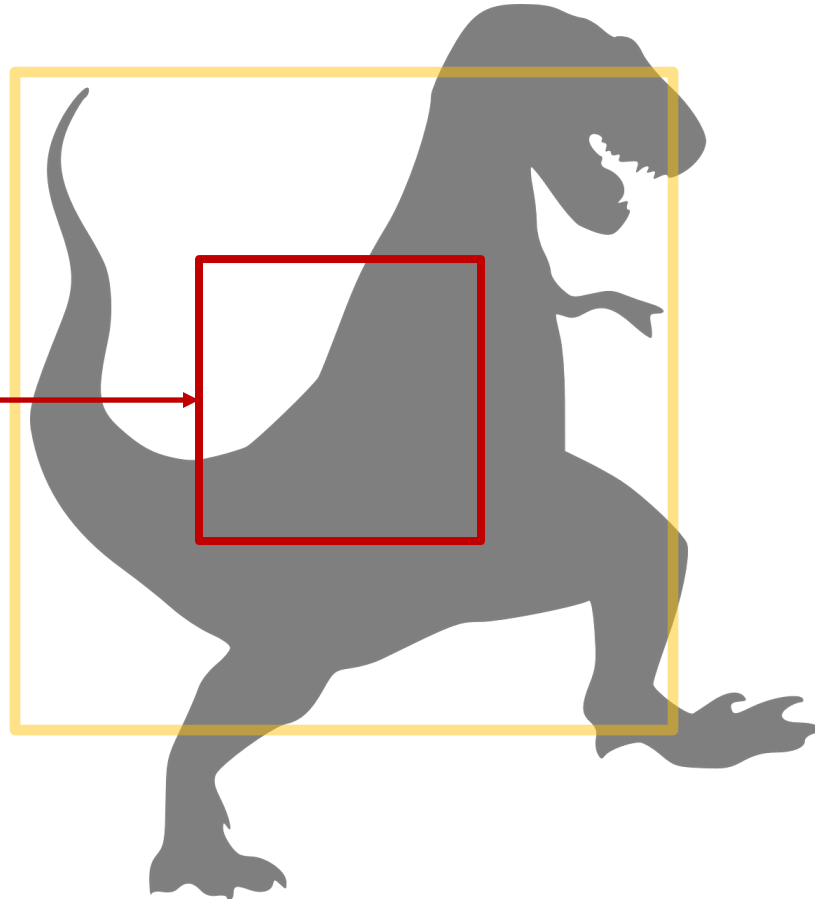
3. RPN: Anchor Box



3. RPN: Prediction (on object)

Objectness score

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



3x3 “sliding window”

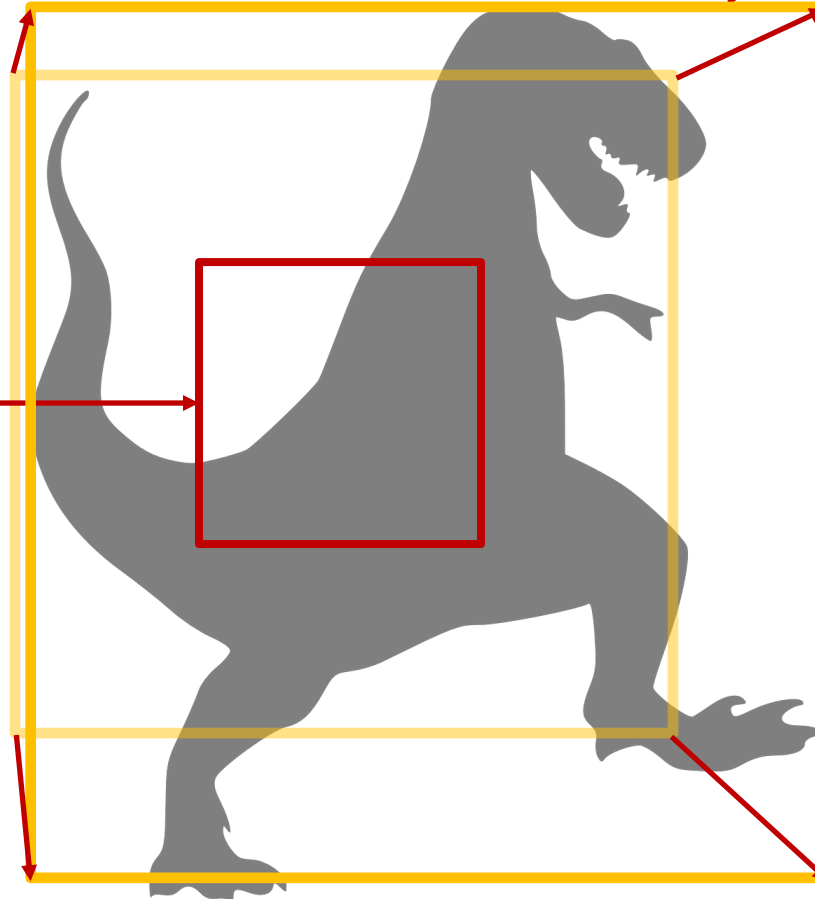
➤ Objectness classifier [0, 1]

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

3. RPN: Prediction (on object)

Anchor box: transformed by box regressor

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



3x3 “sliding window”

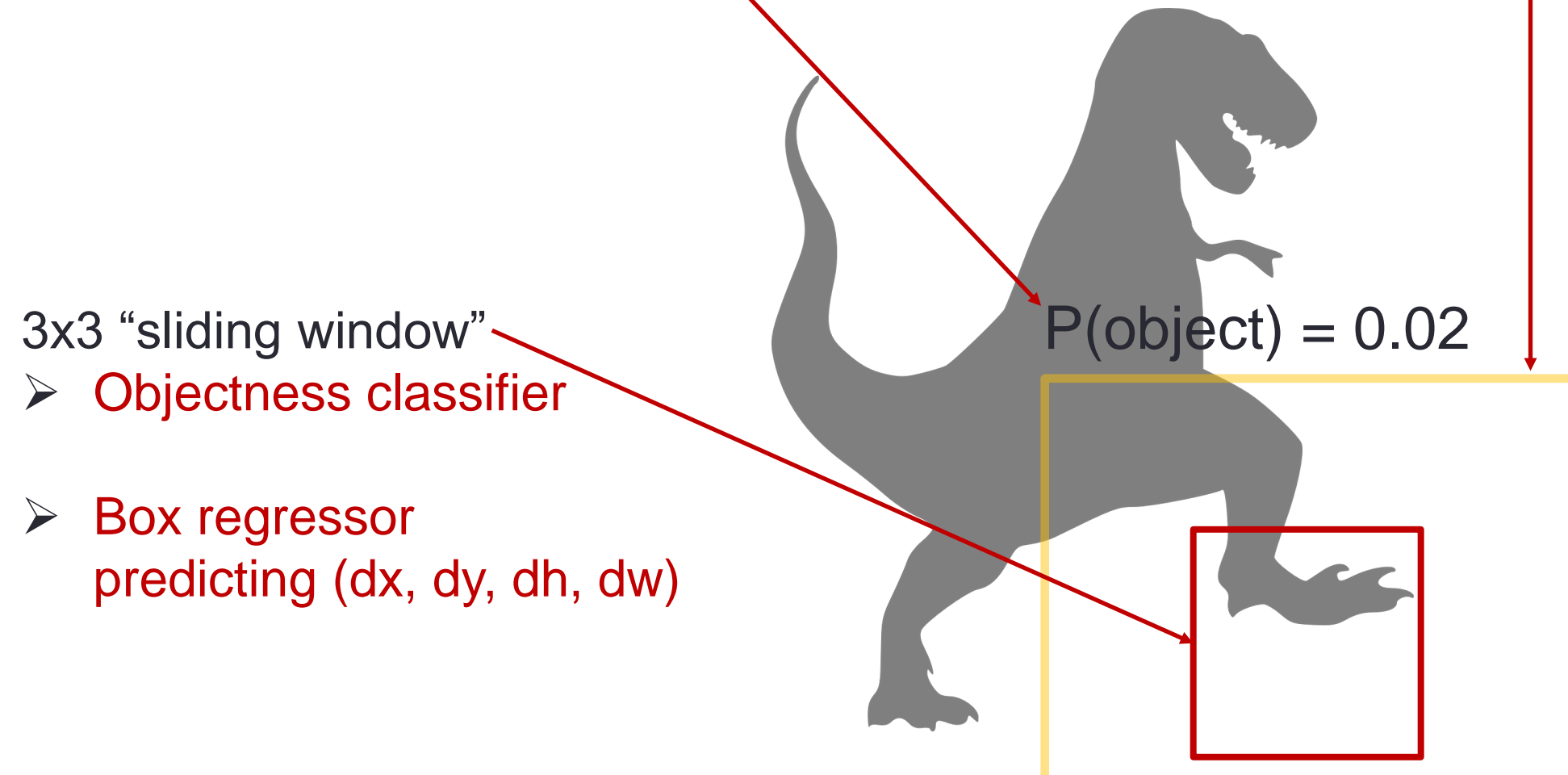
➤ Objectness classifier [0, 1]

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

3. RPN: Prediction (off object)

Objectness score

Anchor box: transformed by
box regressor



3. RPN: Multiple Anchors

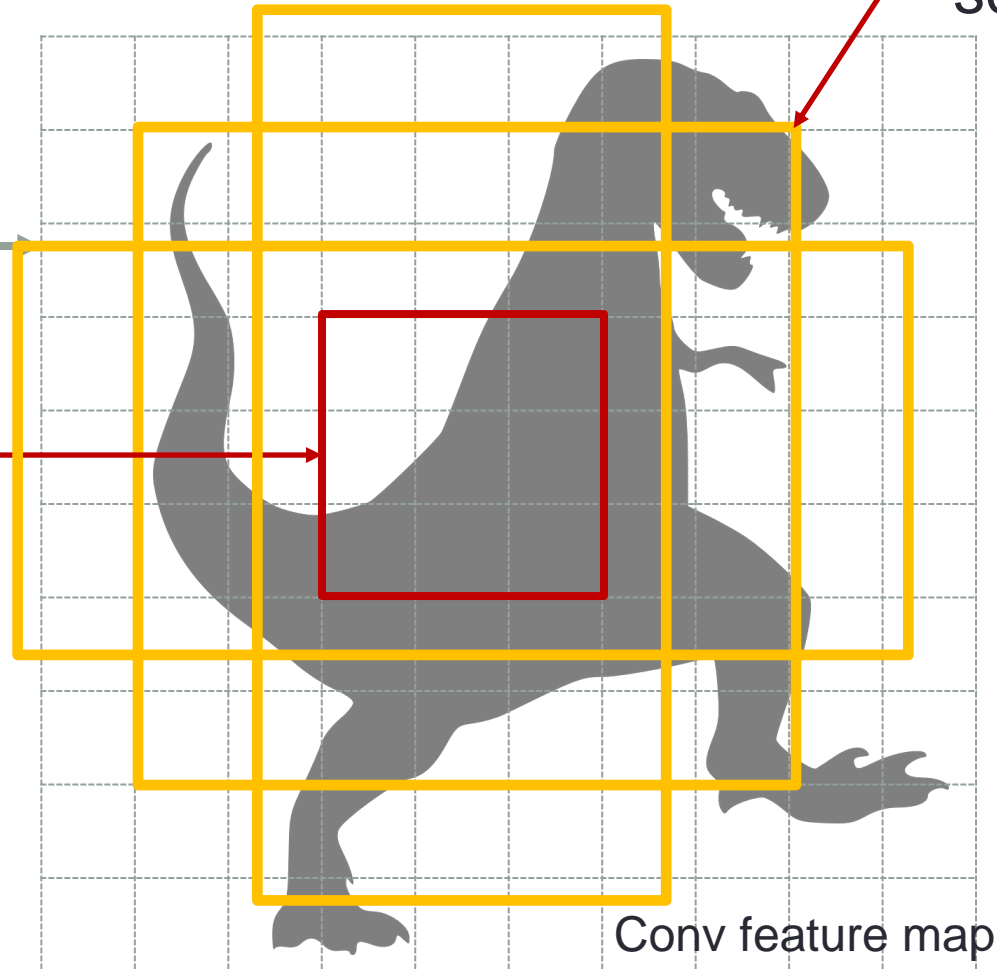


$$f_I = \text{FCN}(I)$$

3x3 “sliding window”

➤ K objectness classifiers

➤ K box regressors



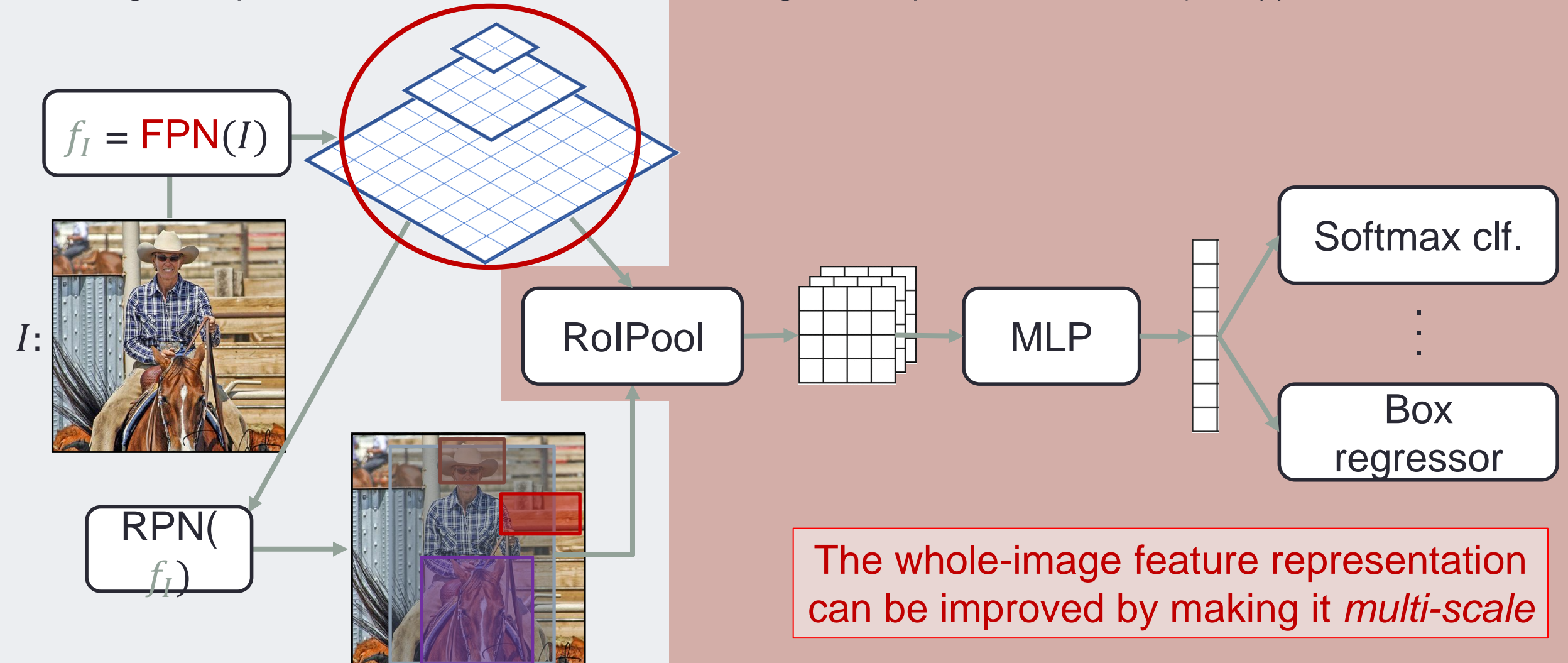
Anchor boxes: K anchors per location with different scales and aspect ratios

Conv feature map

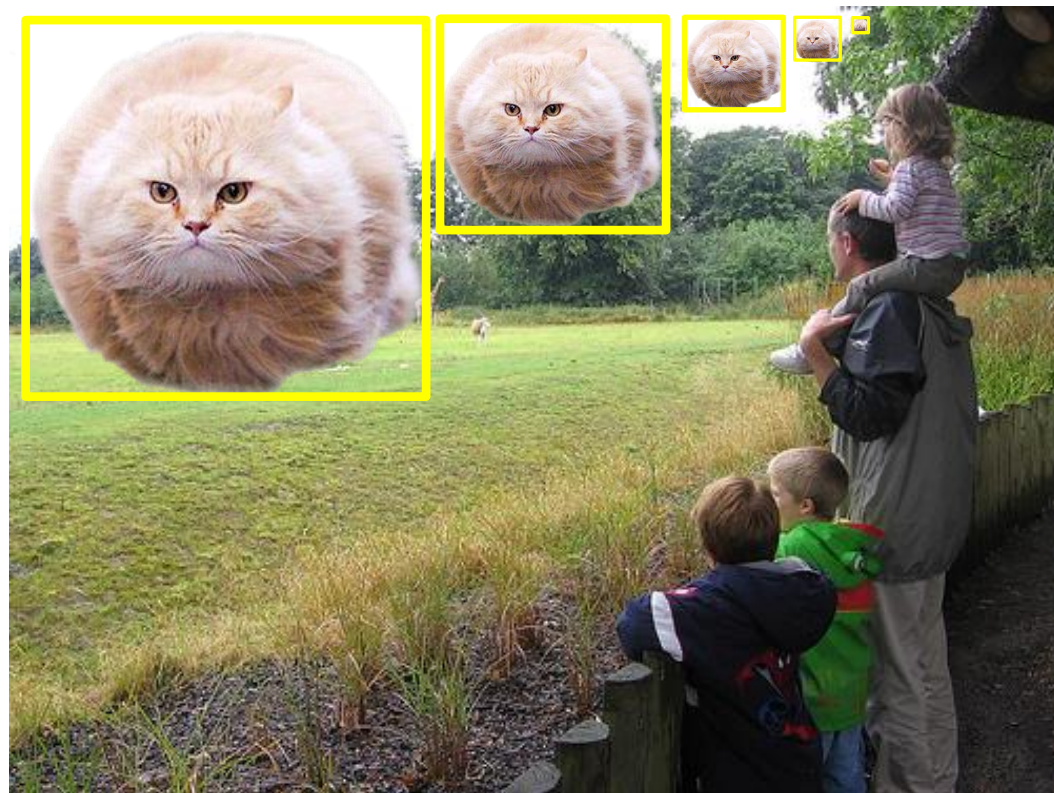
4. Faster R-CNN with a Feature Pyramid Network

Per-image computation

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



4. FPN: Improving Scale Invariance and Equivariance

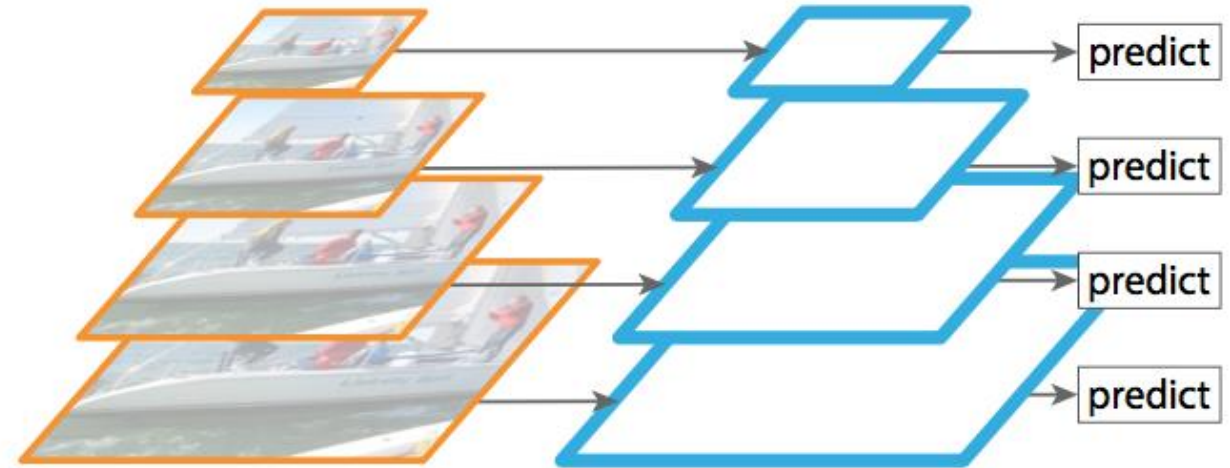
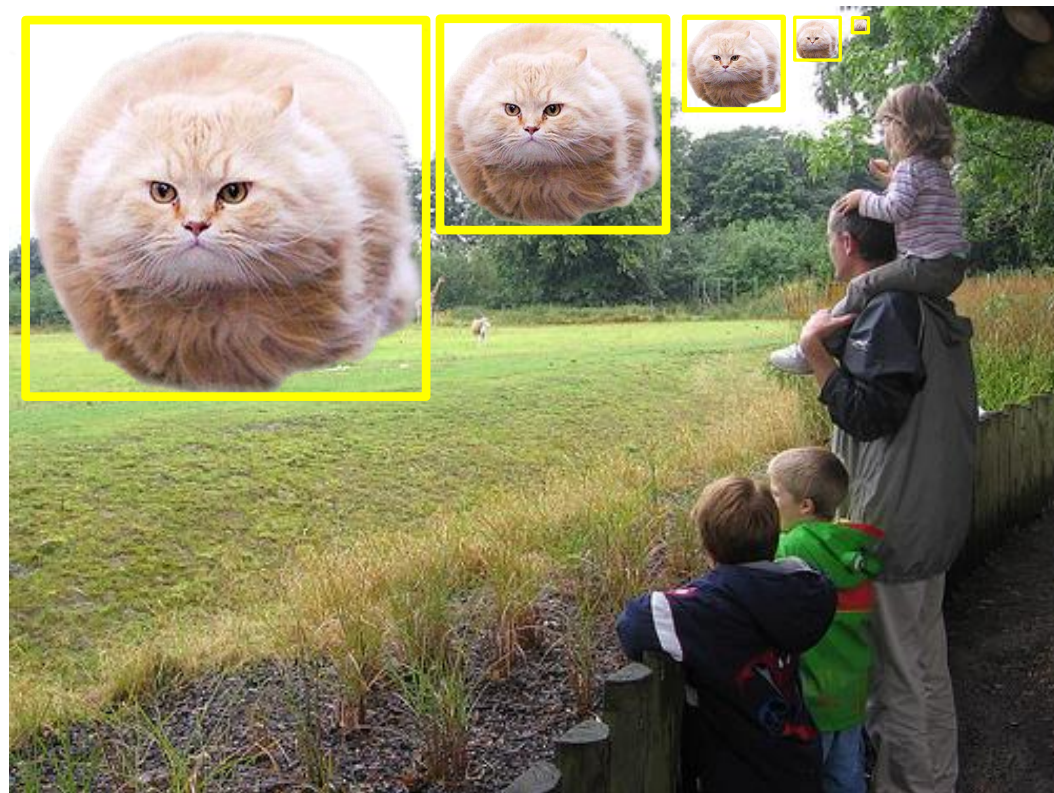


Detectors need to

1. classify and
 2. localize
- objects over a **wide range of scales**

FPN improves this ability

Strategy 1: Image Pyramid

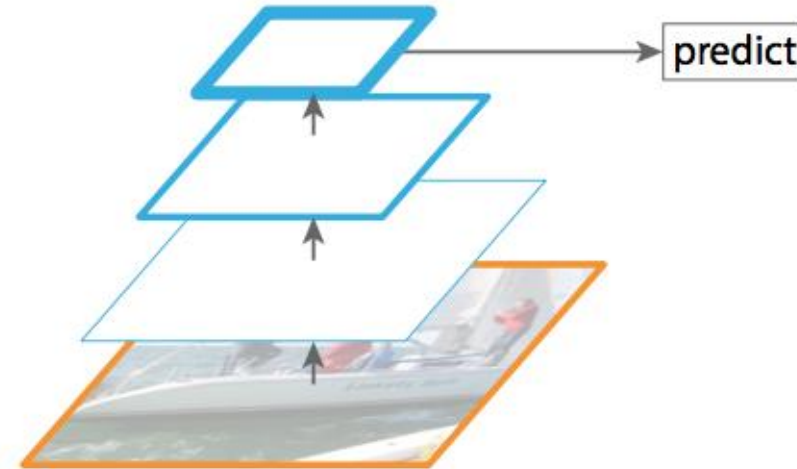
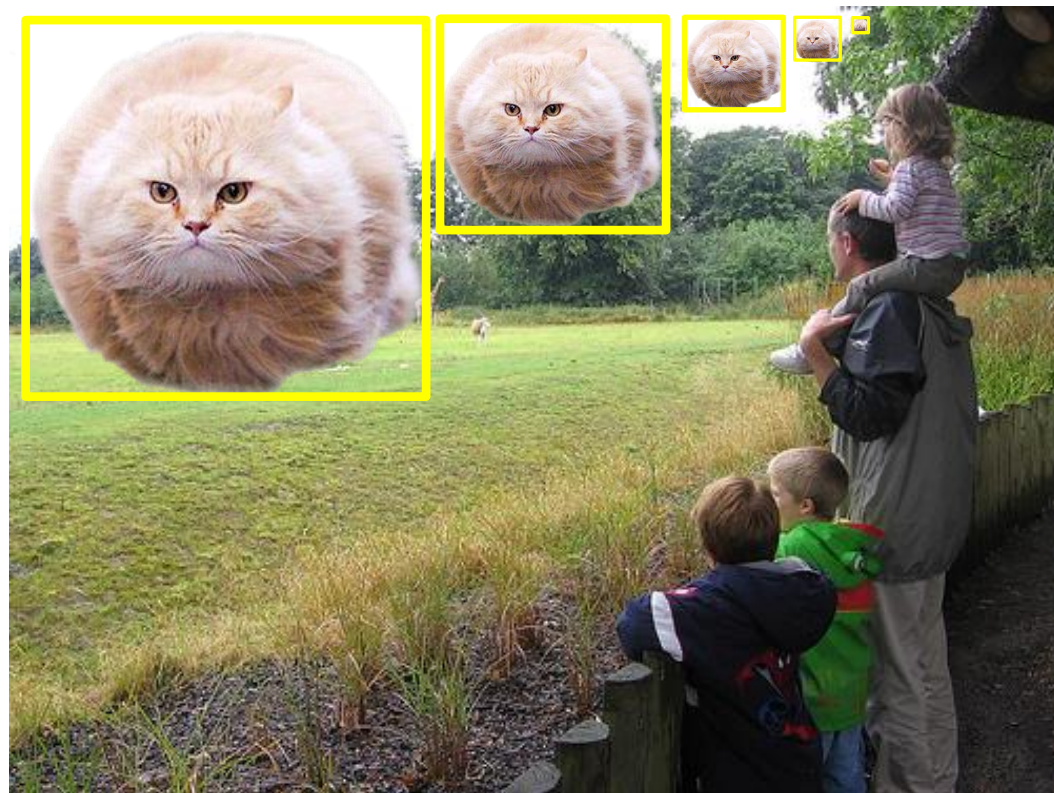


(a) Featurized image pyramid

Standard solution – *slow!*

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

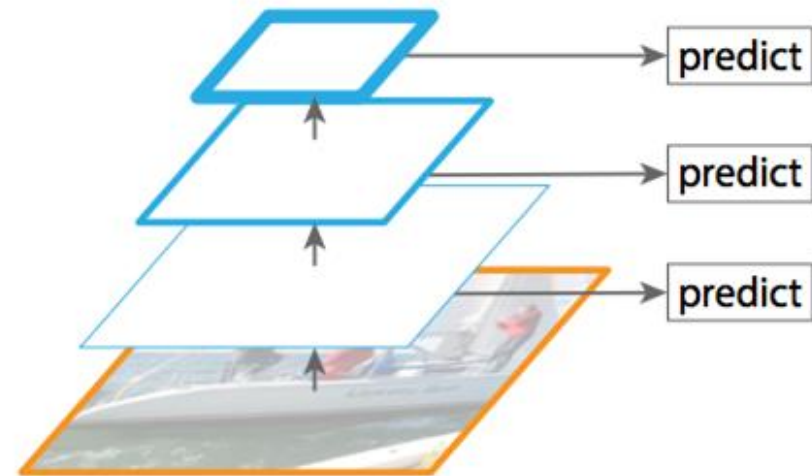
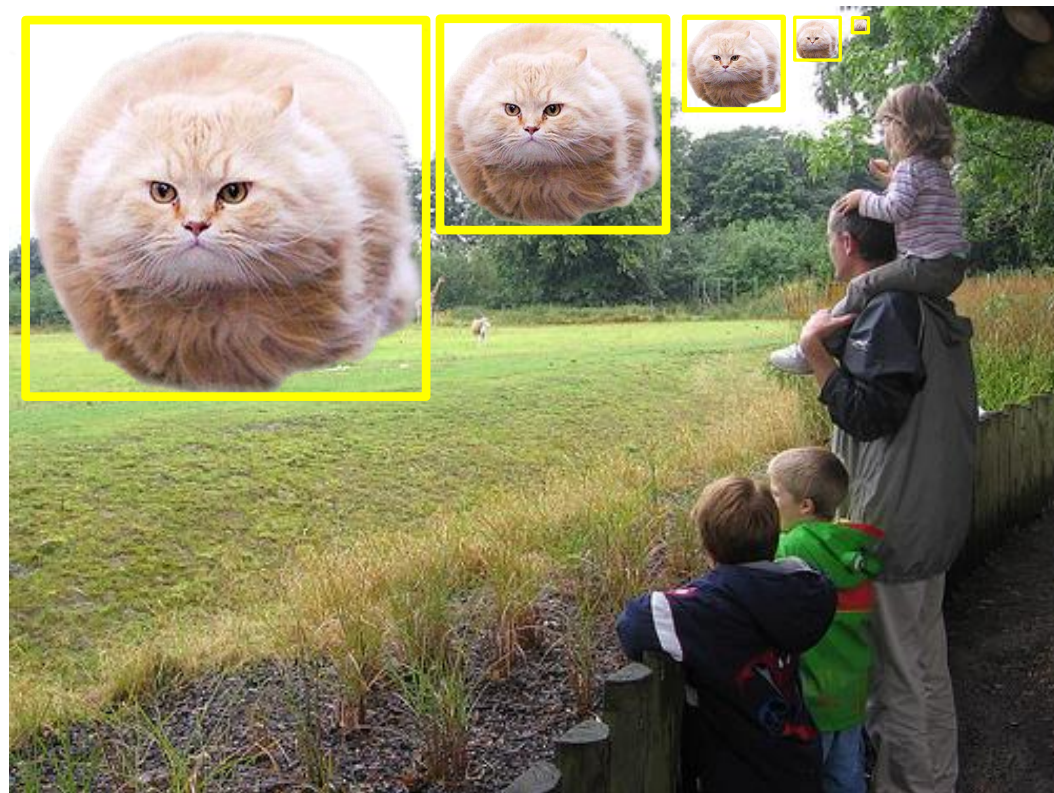
Strategy 2: Multi-scale Features (Single-scale Map)



(b) Single feature map

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

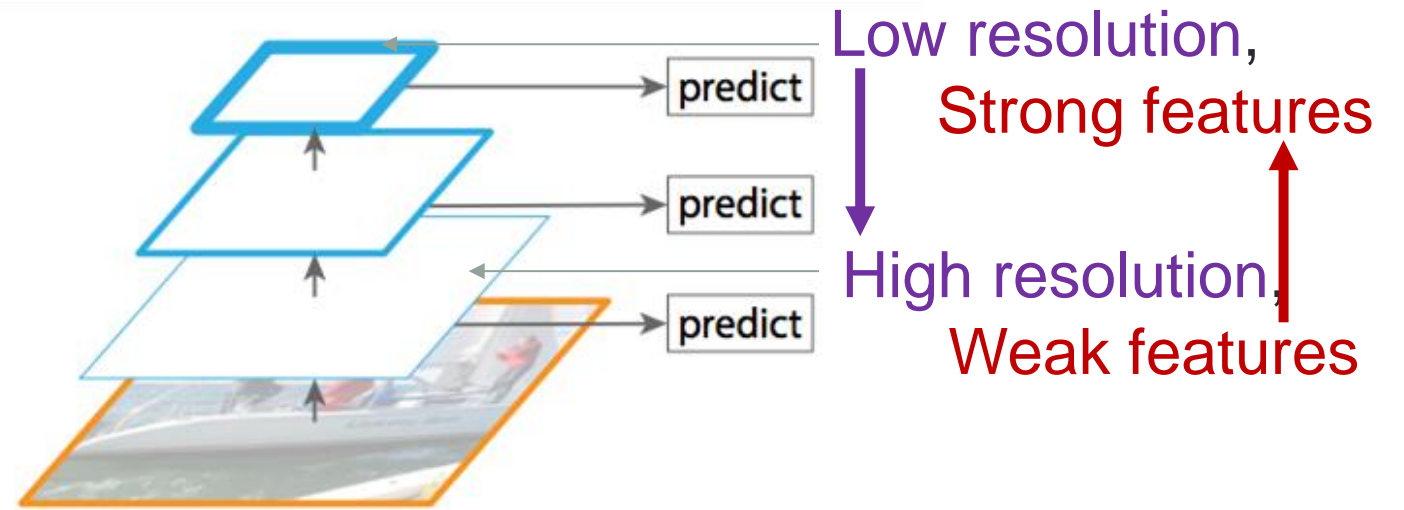
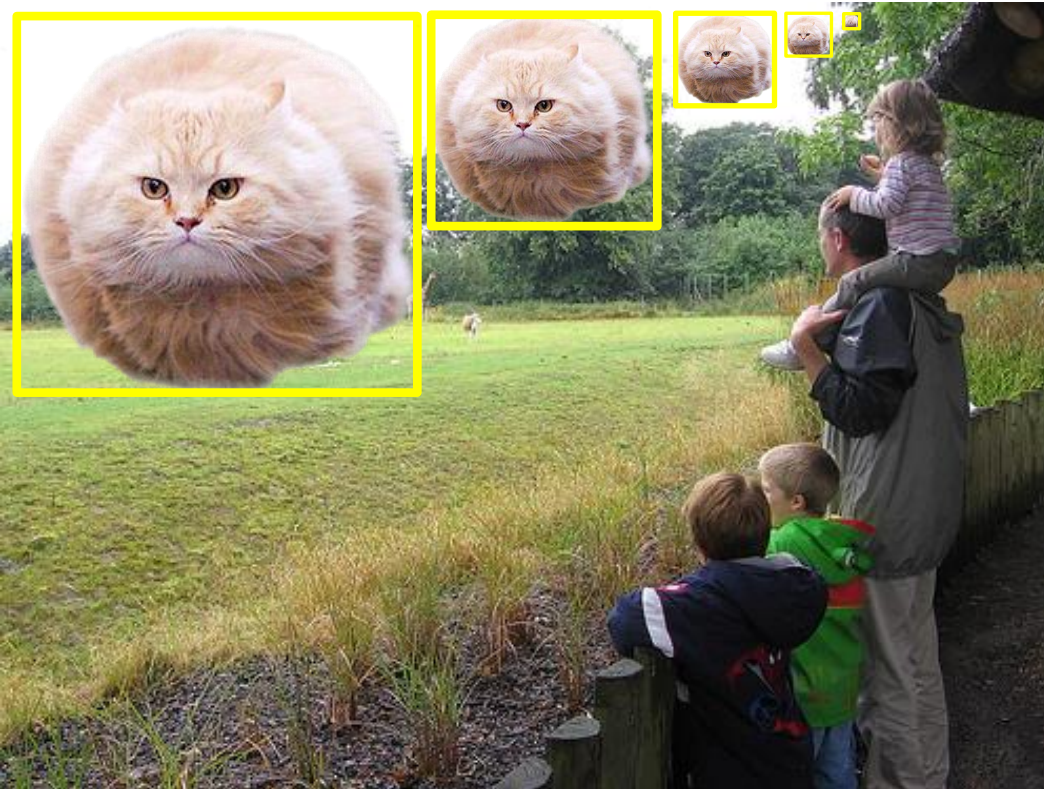
Strategy 3: Naïve In-network Pyramid



(c) Pyramidal feature hierarchy

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

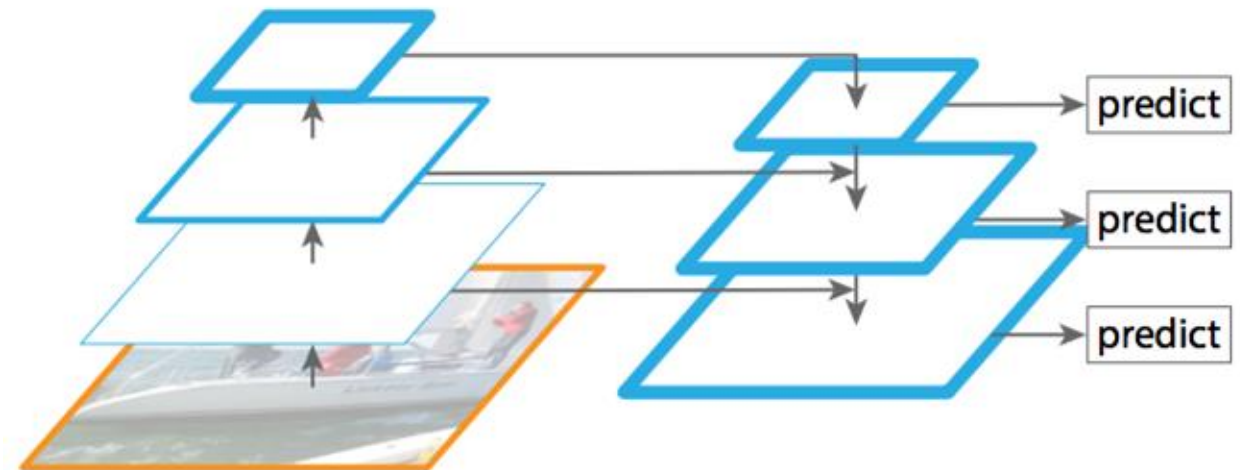
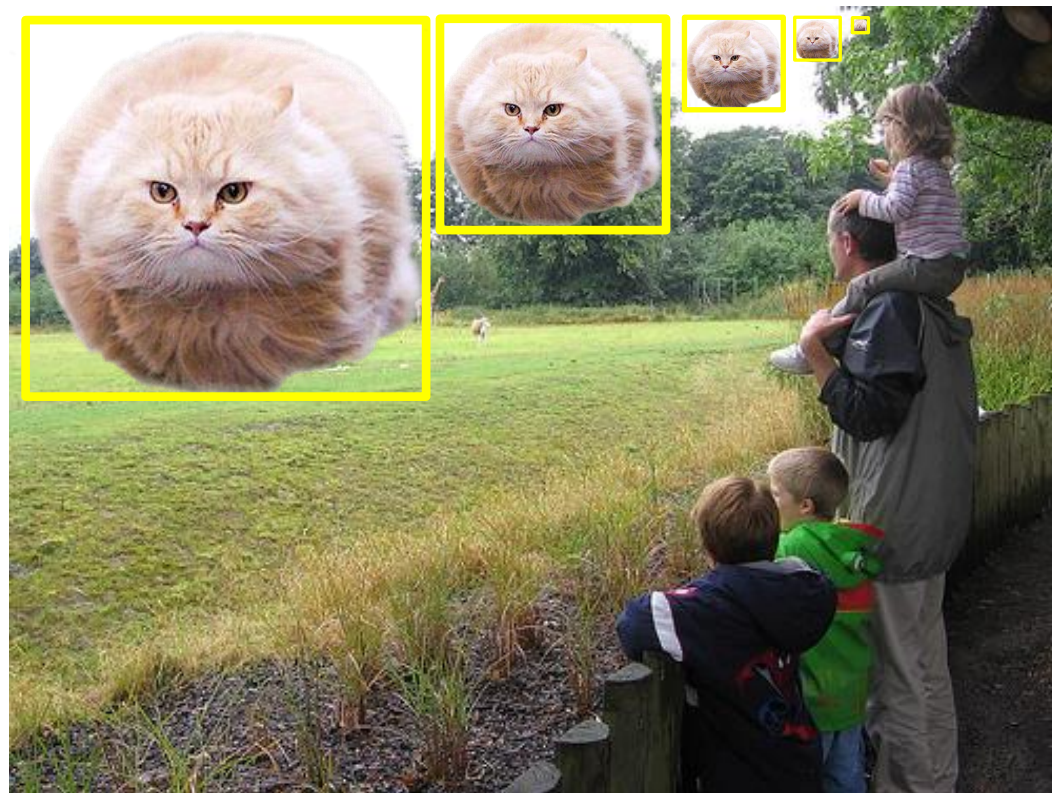
Strategy 3: Naïve In-network Pyramid



(c) Pyramidal feature hierarchy

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

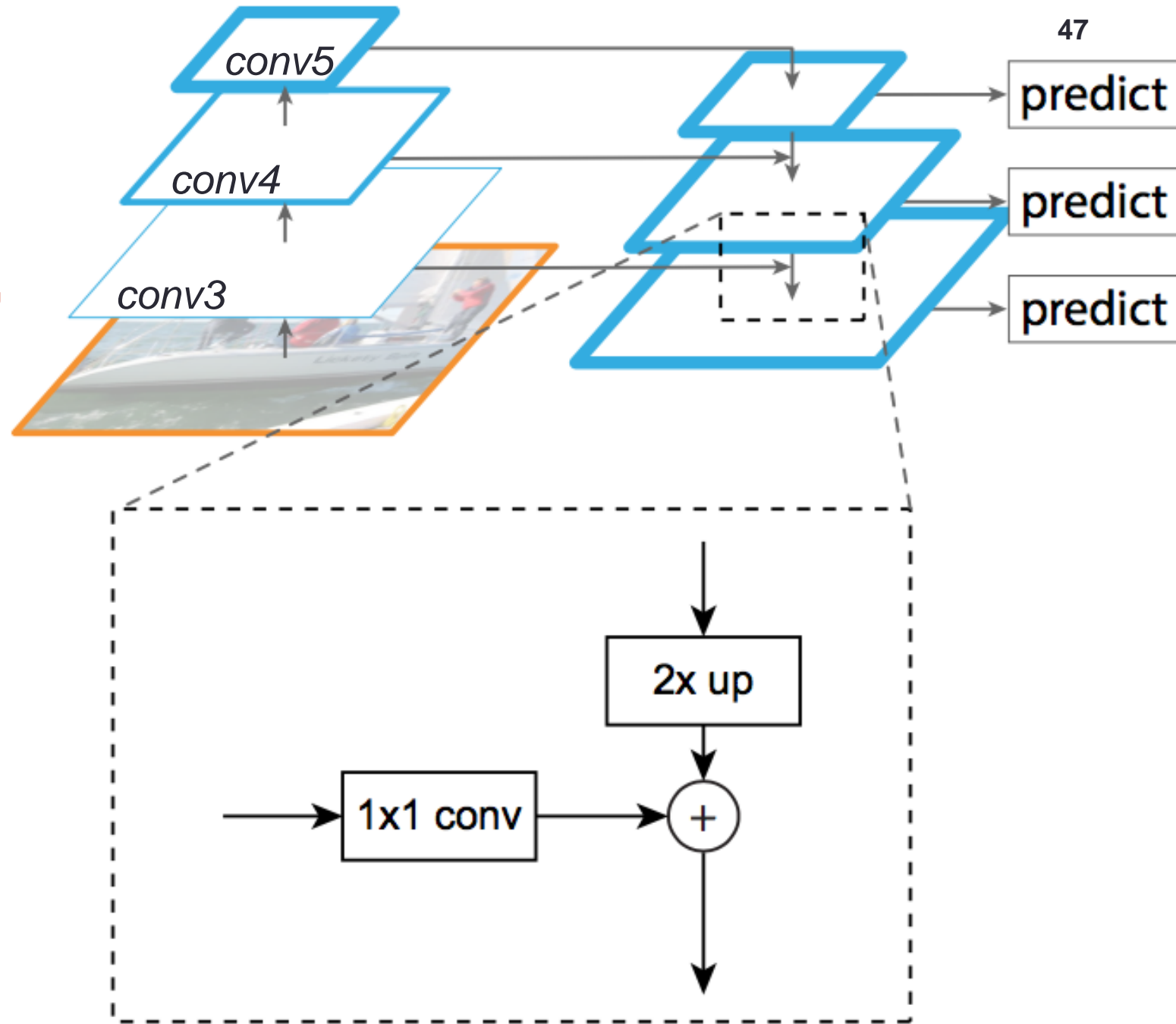
Strategy 4: Feature Pyramid Network



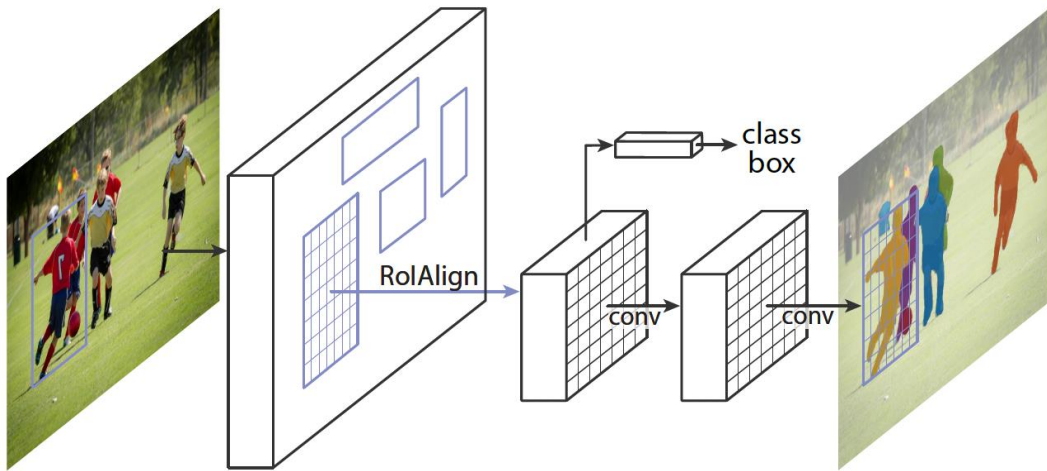
(d) Feature Pyramid Network

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

FPN: Light-weight, Top-down Refinement Module

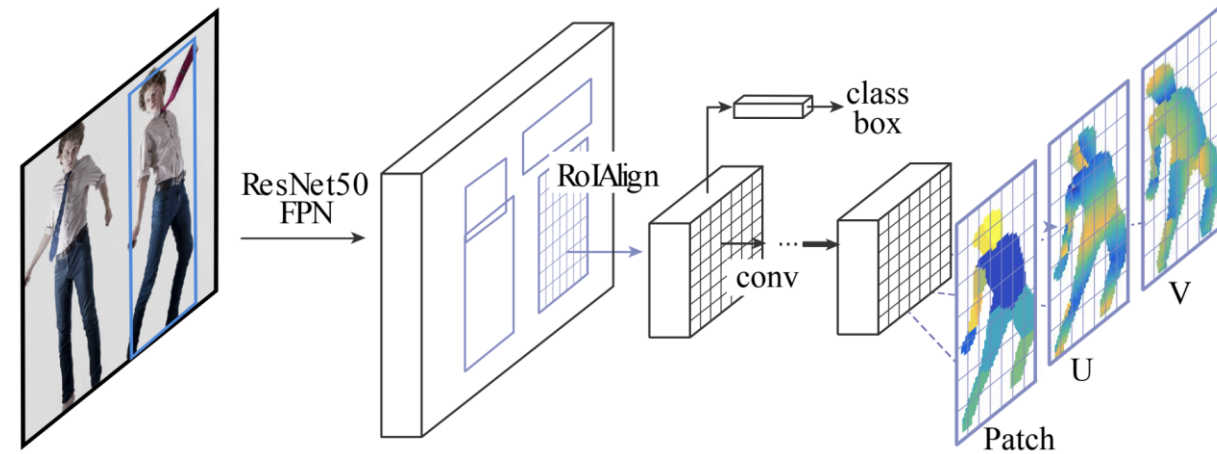


5. Generalized R-CNN: Adding More Heads



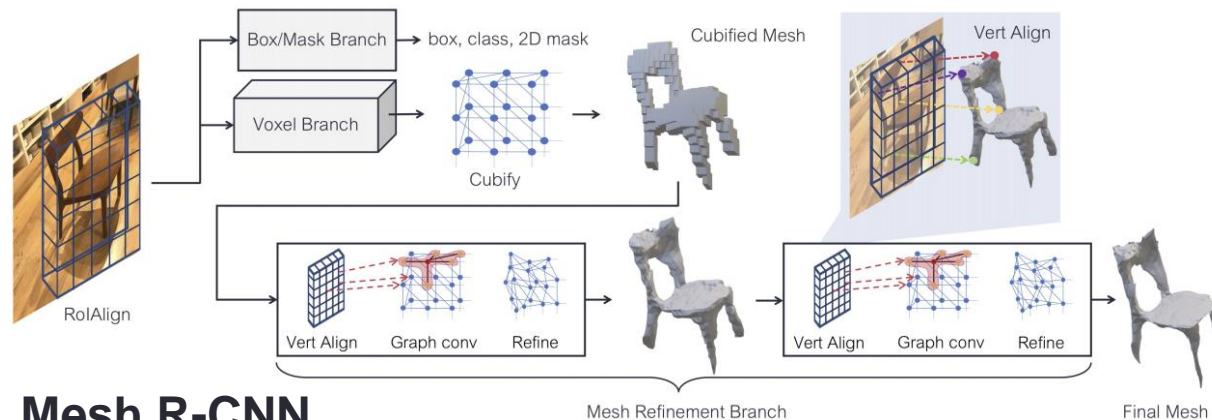
Mask R-CNN

[He, Gkioxari, Dollár, Girshick]



DensePose

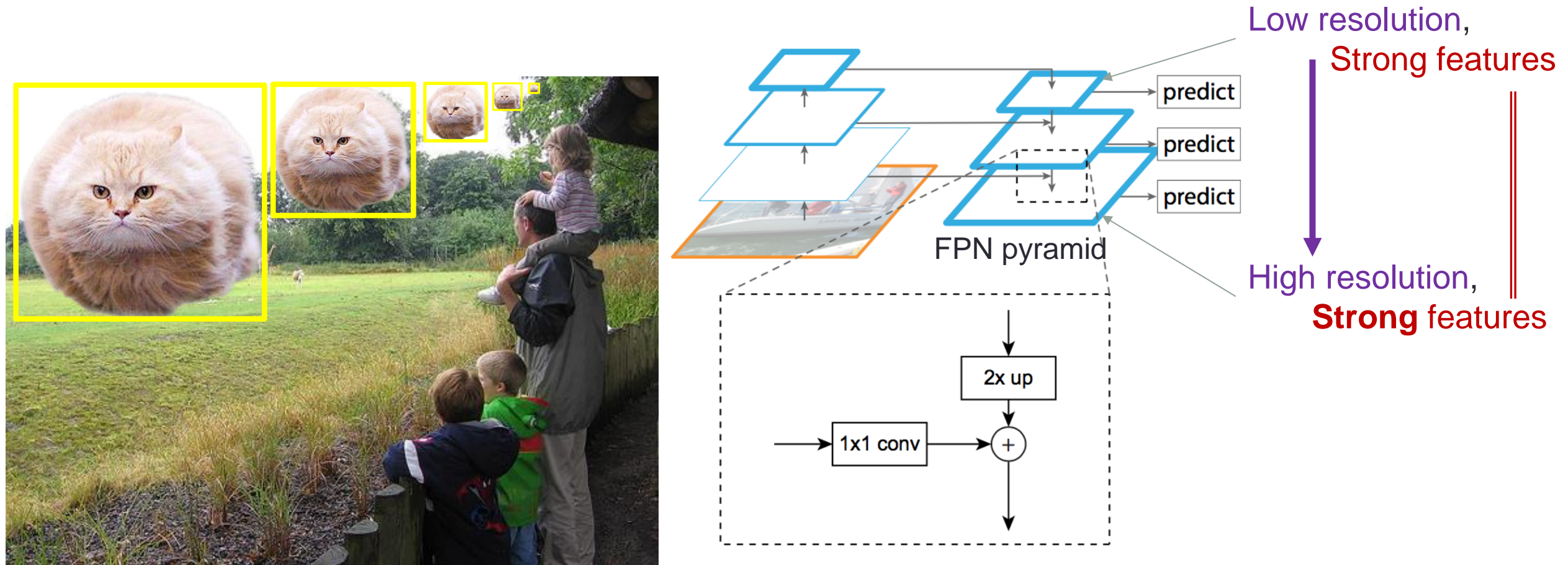
[Güler, Neverova, Kokkinos]



Mesh R-CNN

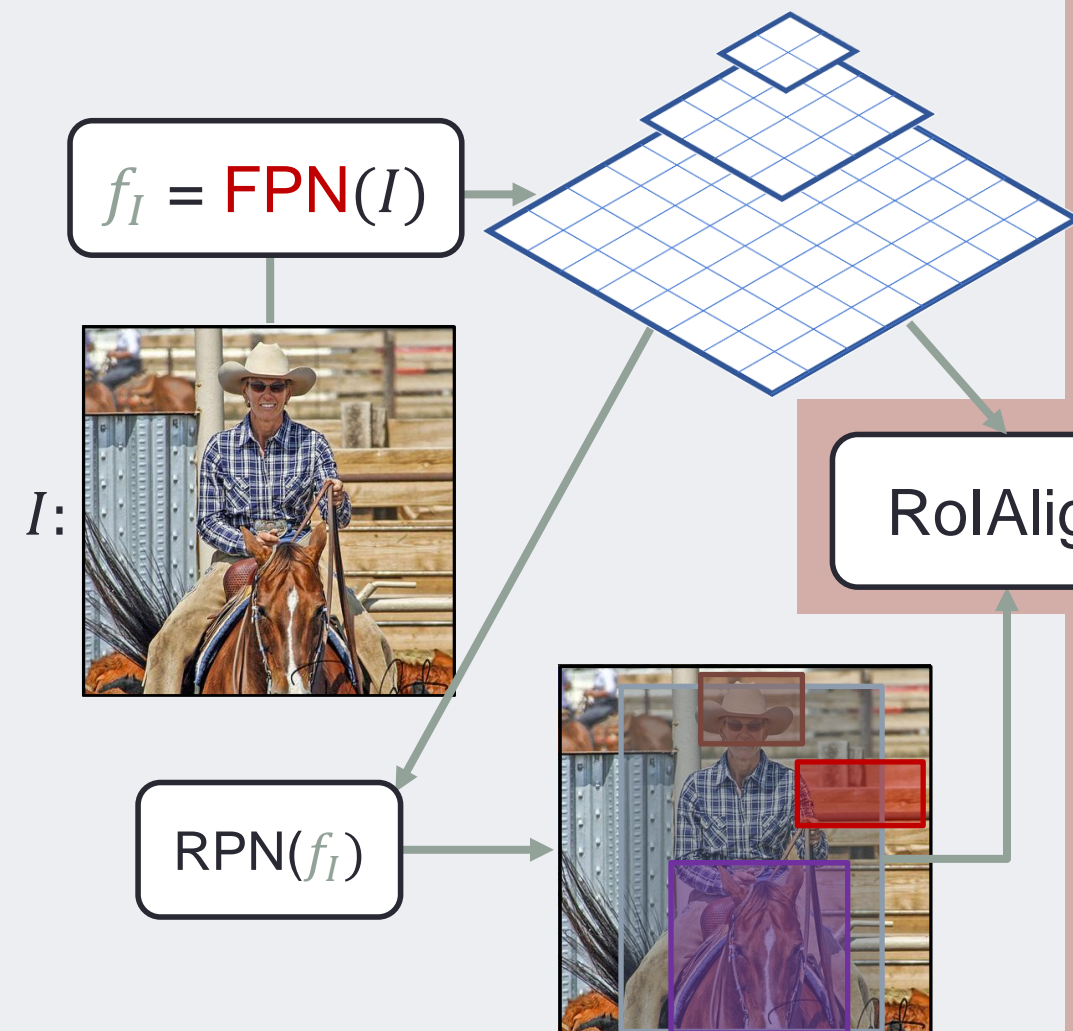
[Gkioxari, Malik, Johnson arXiv 2019]

No Compromise on Feature Quality, Still Fast

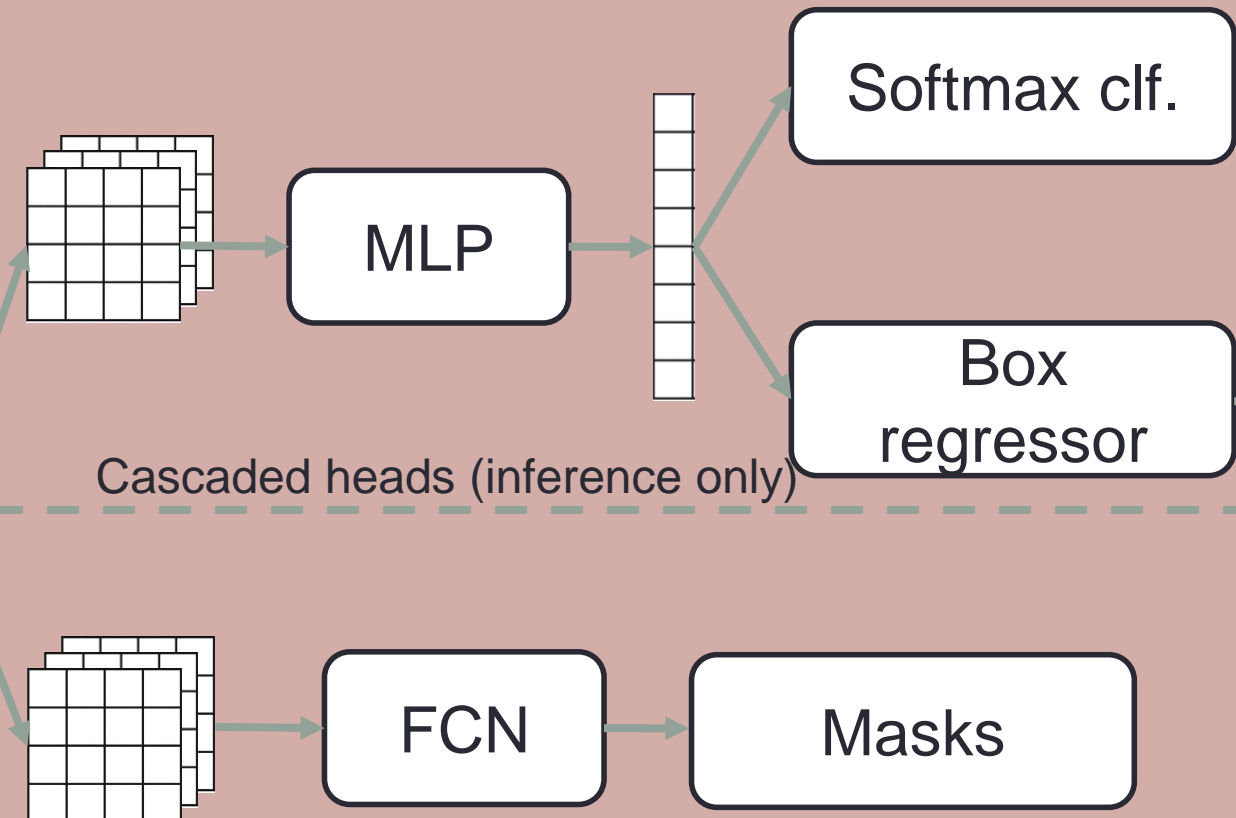


5. Mask R-CNN

Per-image computation



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



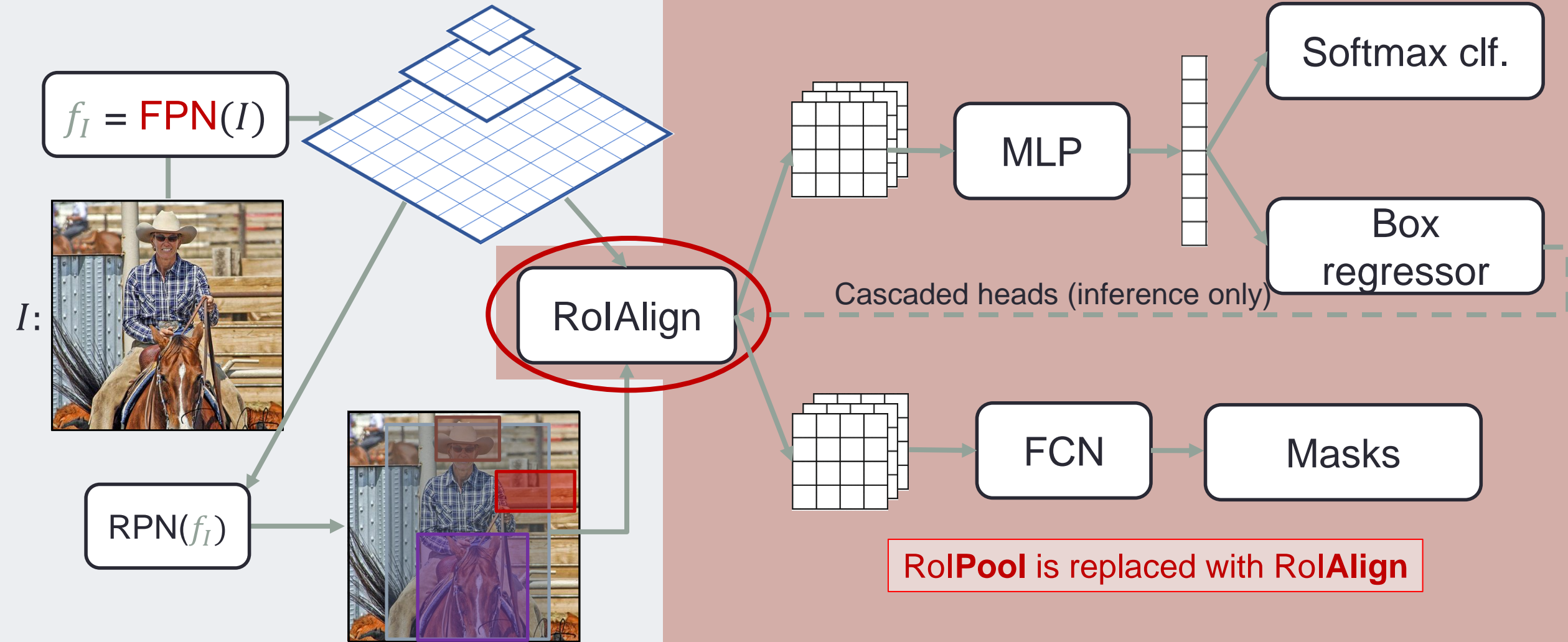
Cascaded heads (inference only)

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

5. Mask R-CNN

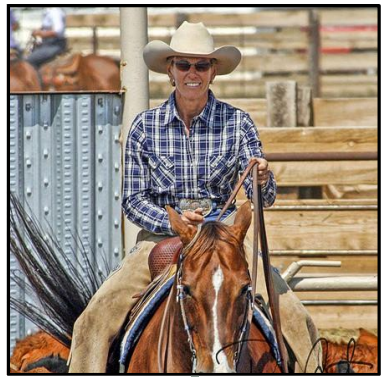
Per-image computation

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

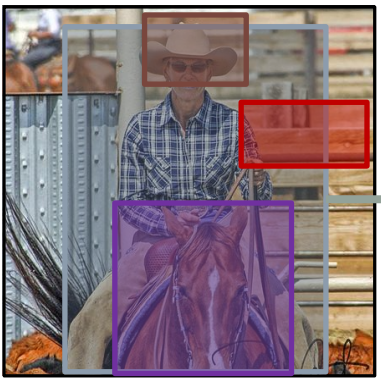


5. RoIAlign Operation (on each Proposal)

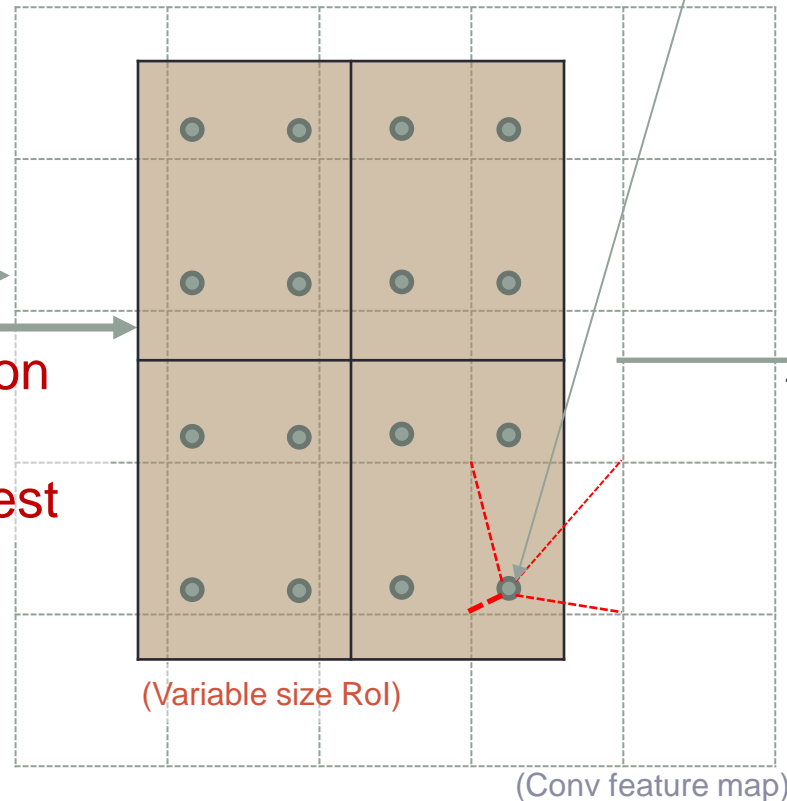
Transform **arbitrary size proposal** into a **fixed-dimensional** representation (e.g., 2x2)



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



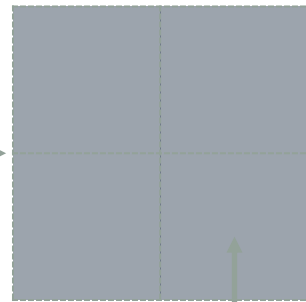
Region
of
Interest
(RoI)



Grid of bilinear
interpolation points

RoIAlign
transform

(Fixed dimensional
representation)

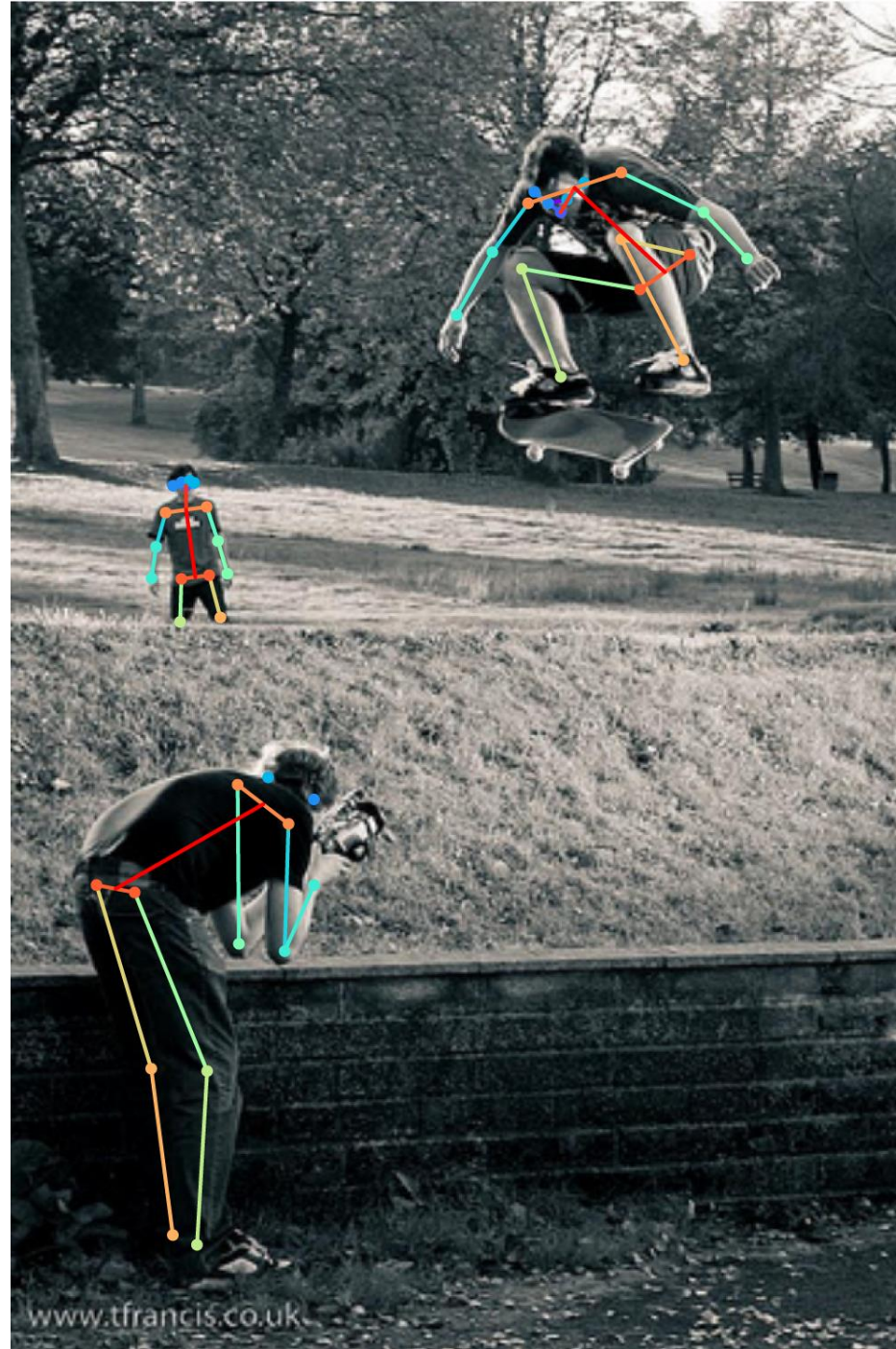


MLP/FCN

Feature value
is **average** of
interpolated values



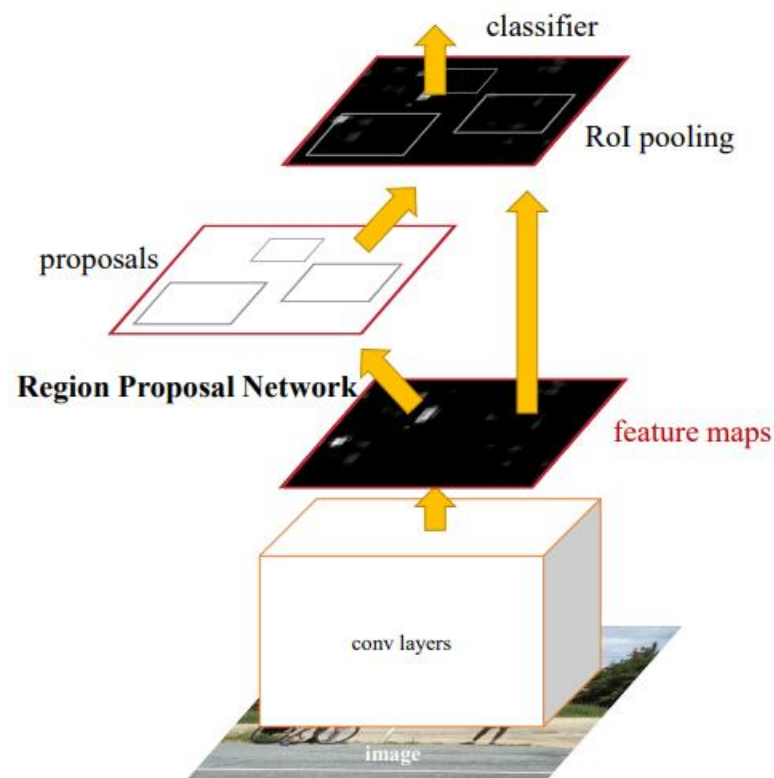
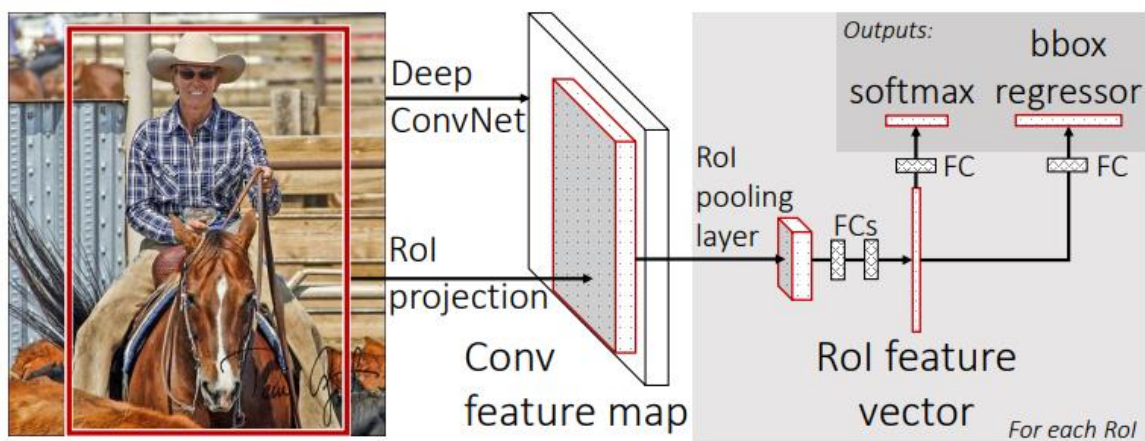
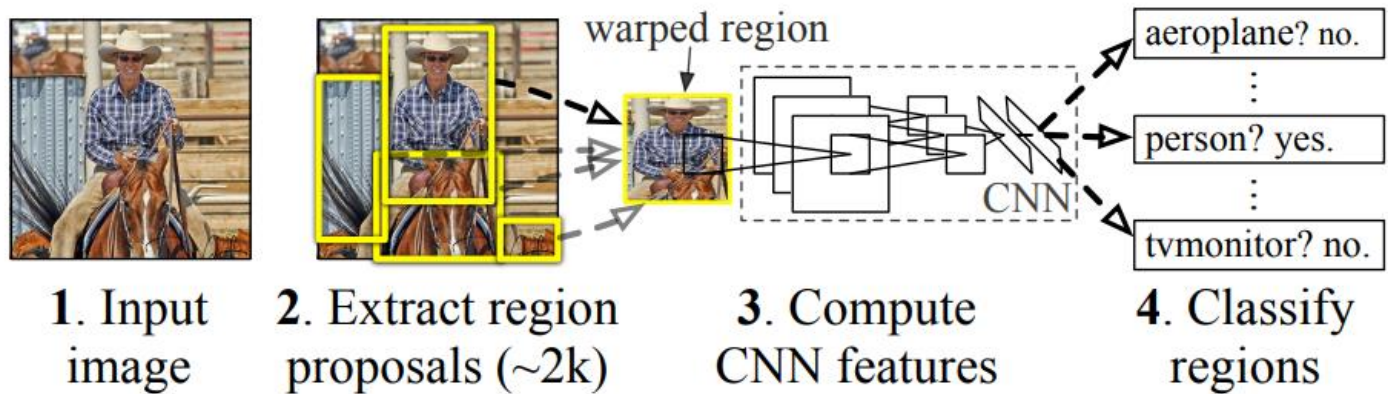






Summary of 1-5

R-CNN: *Regions with CNN features*



Summary of 1-5

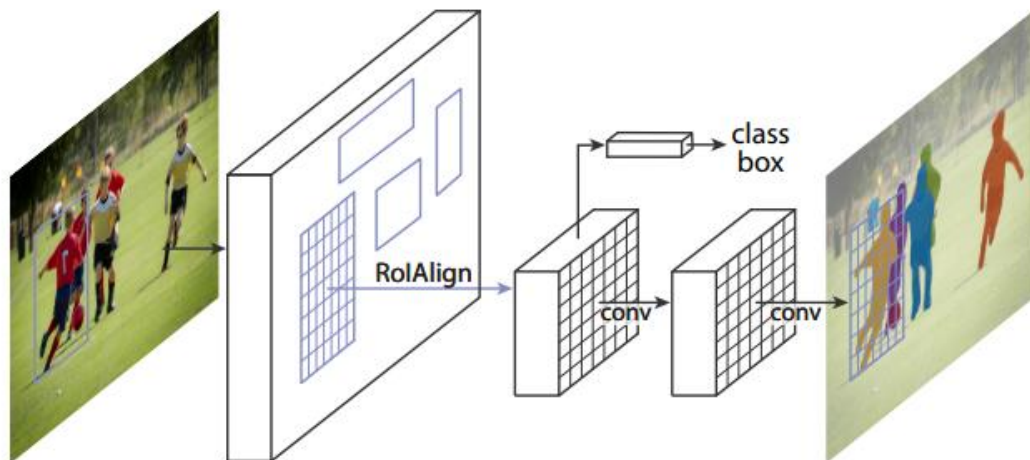


Figure 1. The **Mask R-CNN** framework for instance segmentation.

• Pixel-to-pixel aligned

The diagram shows the process of generating a final mask for a specific instance. It starts with a photograph of a person sitting on a bench by a river, with a red bounding box around a duck in the water. This bounding box is used to extract a Region of Interest (RoI) from the feature map. The RoI is then processed by a 28x28 FCN prediction network to produce a soft prediction. This soft prediction is then resized to match the original image dimensions, resulting in a final mask that is pixel-to-pixel aligned with the original image.

RoI

28x28 FCN prediction

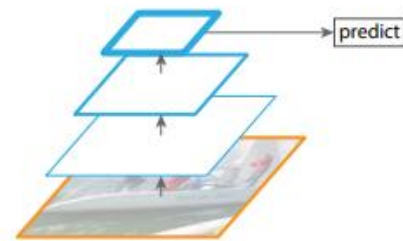
resized soft prediction

final mask

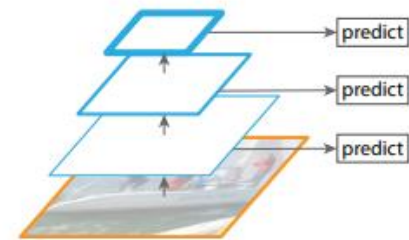
Summary of 1-5



(a) Featurized image pyramid



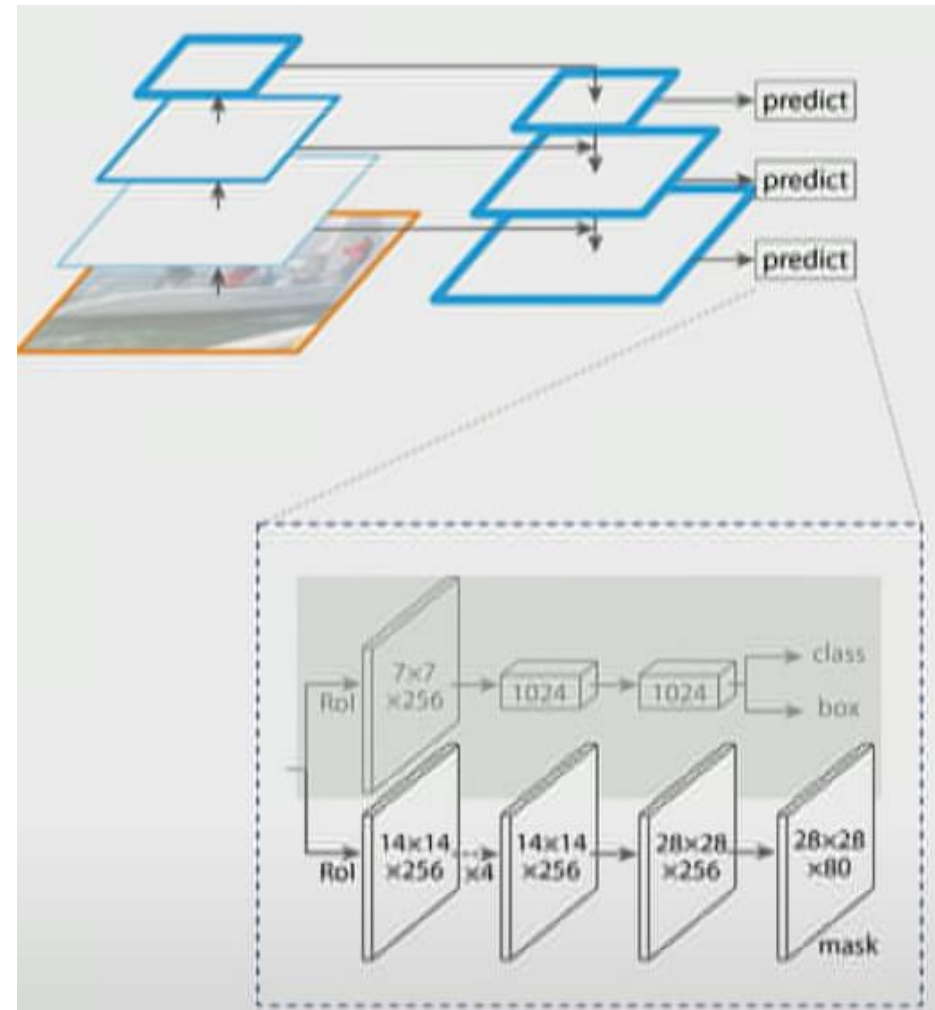
(b) Single feature map



(c) Pyramidal feature hierarchy



(d) Feature Pyramid Network



6. RetinaNet

Lin, Tsung-Yi, Priya Goyal, Ross Girshick, Kaiming He, and Piotr Dollár. "Focal loss for dense object detection." In *IEEE International Conference on Computer Vision*, pp. 2980-2988. 2017.

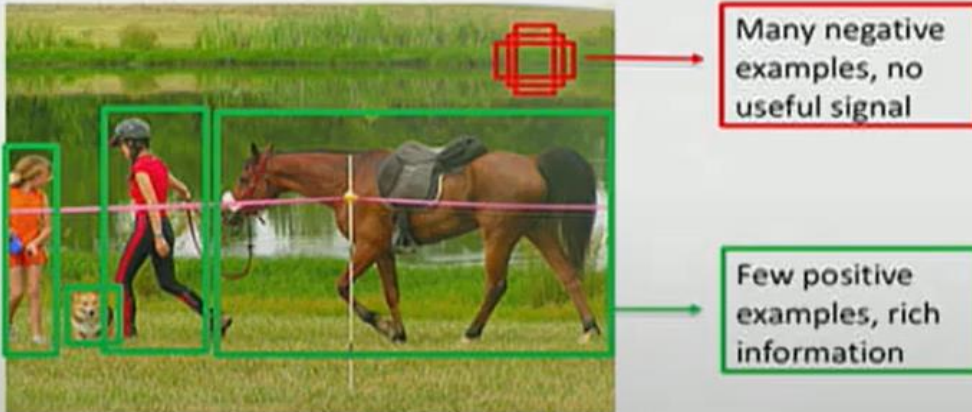
Best Student Paper Award

- Identify **class imbalance** is the major issue for training one-stage dense detector
- Propose **Focal Loss** to address class imbalance
- Achieve state-of-the-art **accuracy** and **speed**

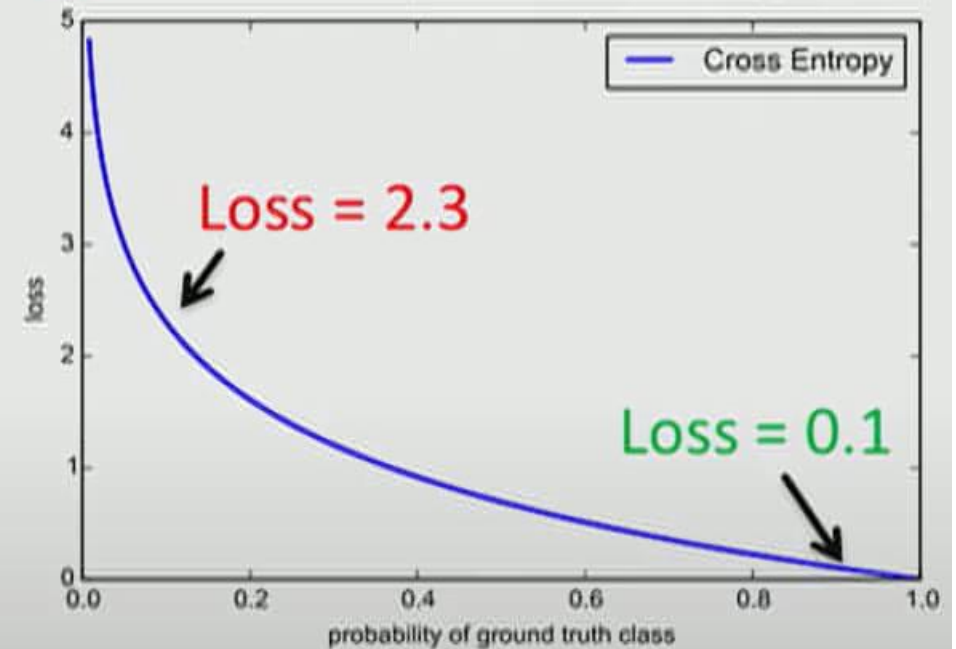
6. RetinaNet

Class Imbalance

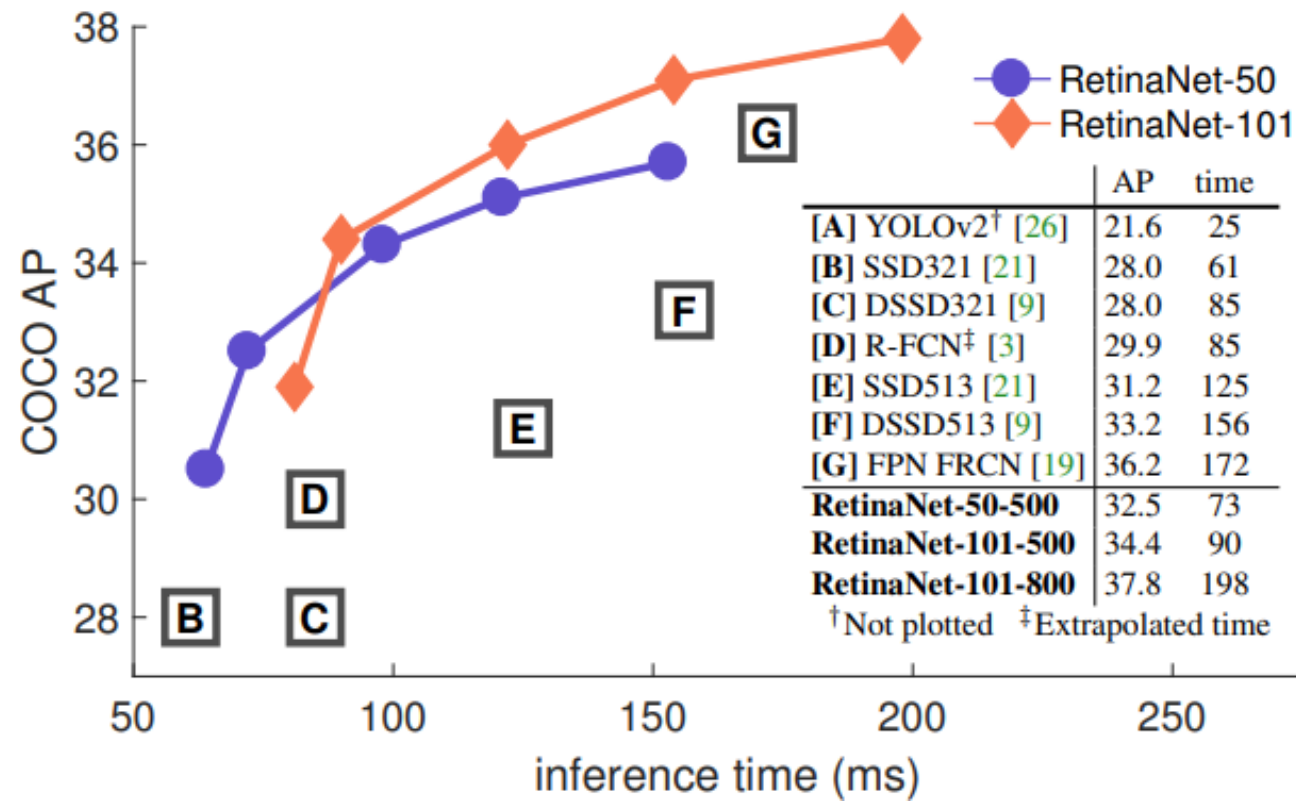
- Few training examples from foreground
- Most examples from background
 - Easy and uninformative
 - Distracting



- 100000 easy : 100 hard examples
- 40x bigger loss from easy examples

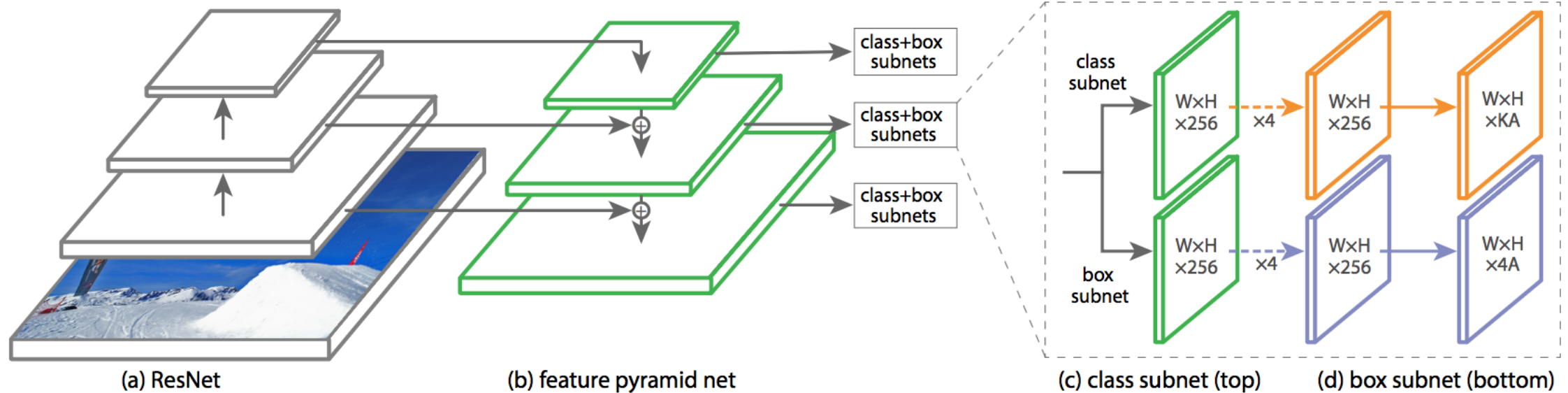


6. RetinaNet



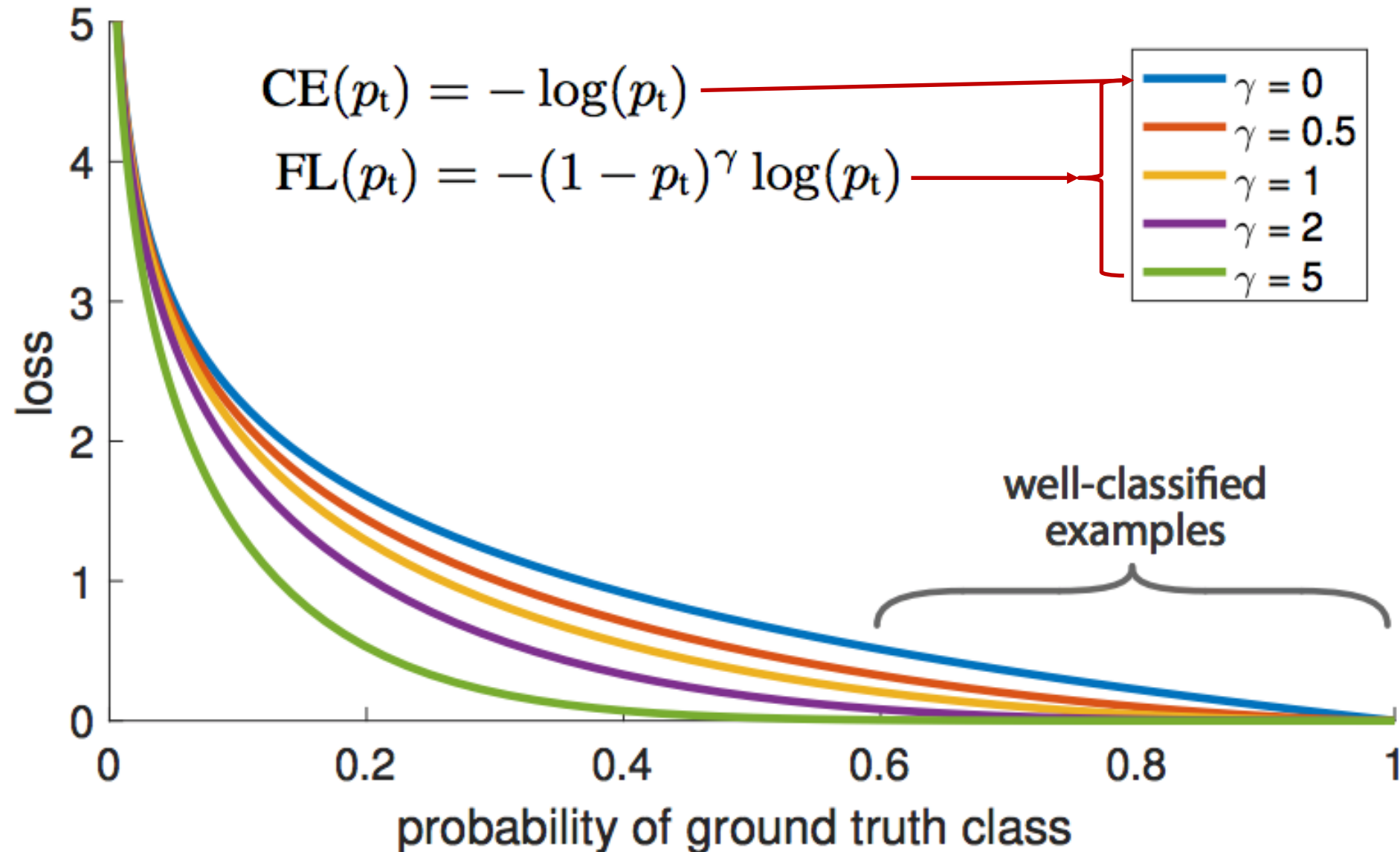
The one-stage RetinaNet detector outperforms all previous one-stage and two-stage detectors

6. RetinaNet Model Description



- Backbone with FPN
- 6 anchors per location (2 scales \times 3 aspect ratios)
- 100 – 200k anchor boxes to classify per image \rightarrow “dense” detection

6. Focal Loss: larger Gamma focuses more on Hard-Example



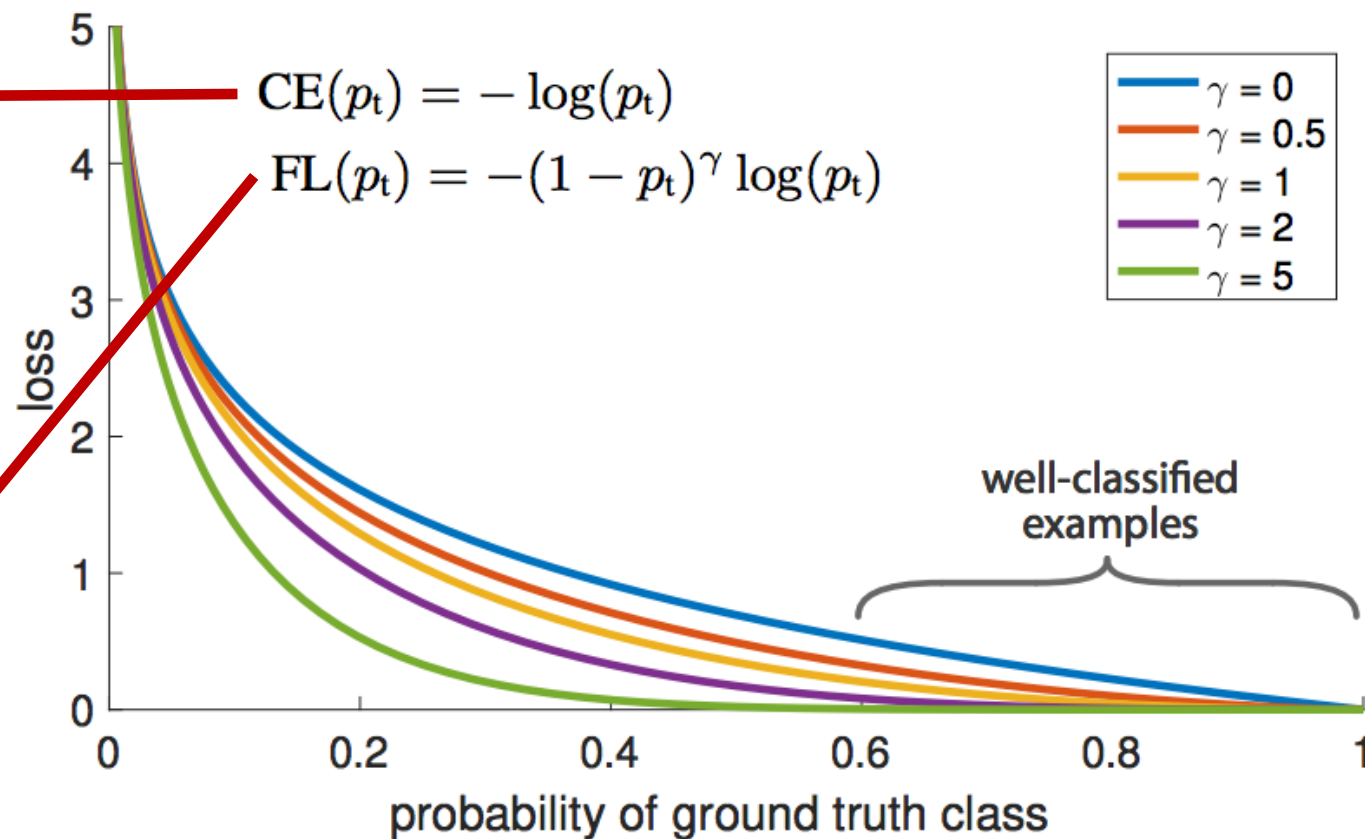
6. Impact of Focal Loss (FL)

$$\text{FL}(p_t) = -\alpha_t(1 - p_t)^\gamma \log(p_t).$$

γ	α	AP	AP ₅₀	AP ₇₅
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	34.0	52.5	36.5
5.0	.25	32.2	49.6	34.8

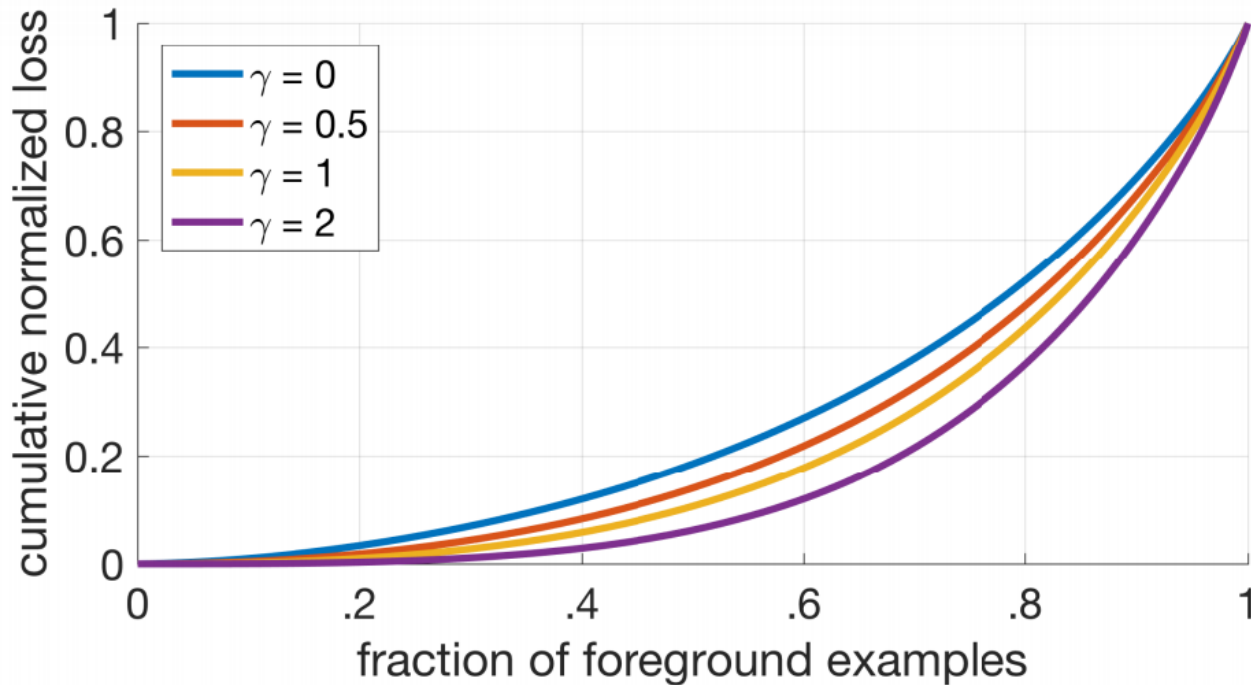
(b) **Varying γ for FL** (w. optimal α)

(ResNet-50-FPN 600px input image)

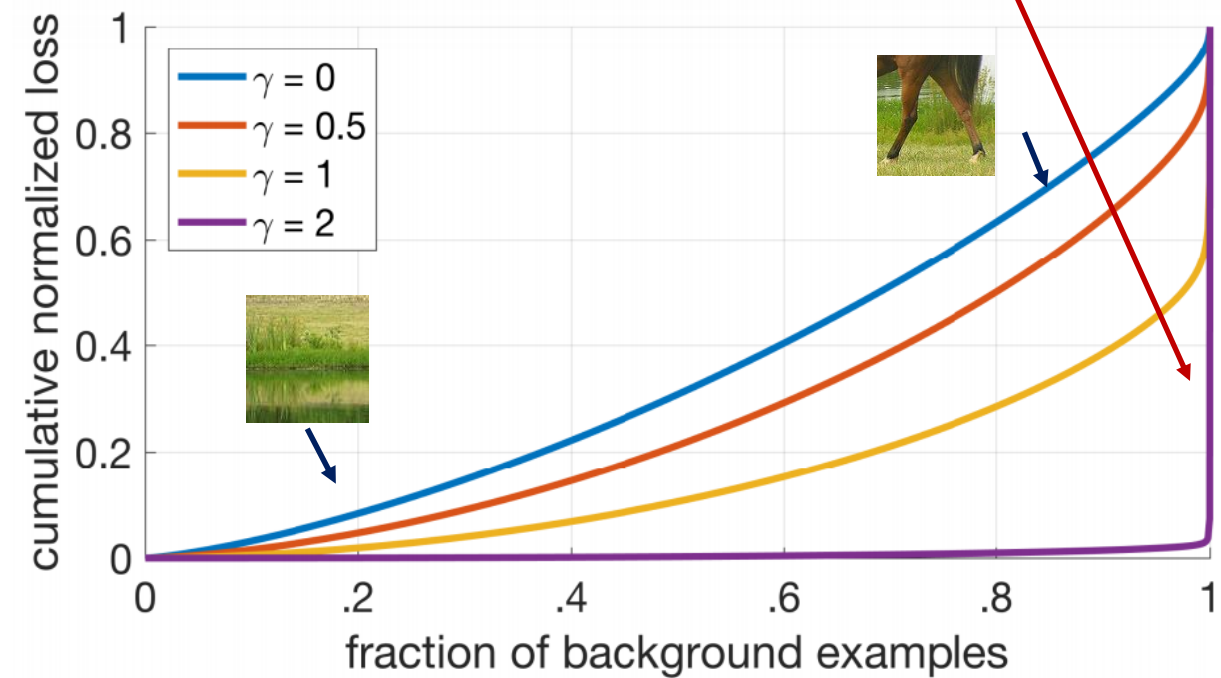


6. Loss Distribution under Focal Loss

Foreground Boxes



Background Boxes



Reference:

1. Dr. Ross Girshick, Facebook AI Research, CVPR 2019 Tutorial on Visual Recognition and Beyond
https://www.dropbox.com/s/rakgo35h6b8p7uv/cvpr2019_tutorial_ross_girshick.pptx?dl=0
2. Dr. Kaiming He, ICCV 2017 Oral Talk for the Mask R-CNN paper.
https://www.youtube.com/watch?v=g7z4mkfRjI4&ab_channel=ComputerVisionFoundationVideos
3. Dr. Tsung-Yi Lin, ICCV 2017 Oral Talk for the Focal Loss paper.
https://www.youtube.com/watch?v=44tlnmmt3h0&ab_channel=ComputerVisionFoundationVideos