

# Object detection: where then what

D.A. Forsyth

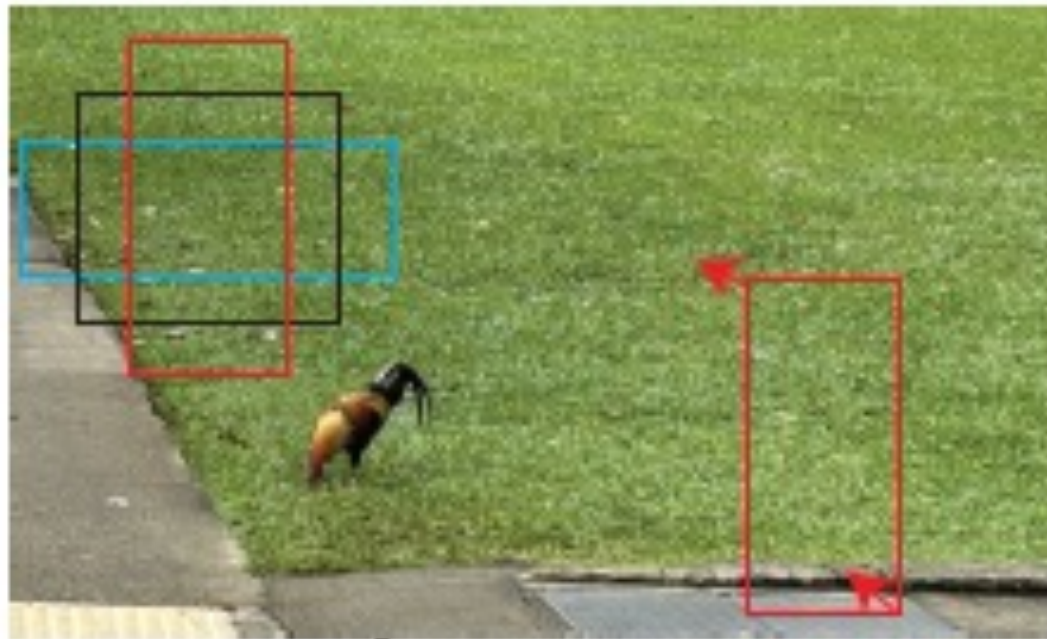
University of Illinois at Urbana Champaign

# Faster RCNN

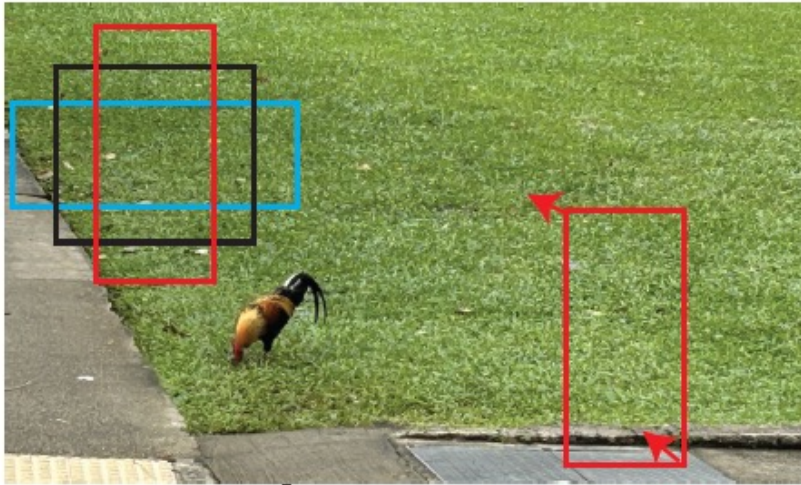
- Ideas:
  - Sample the world of boxes
  - Decide which boxes are worth looking at
  - Score boxes
  - Cleanup
    - Threshold
    - Non-maximum suppression
    - Bounding box regression
- Build all this on top of an image encoder

# Sampling the world of boxes

- 4 dimensional
  - boxes are axis aligned (eg top left, bottom right)
- Samples don't need to be uniform across dims
- Sampling:
  - Even grid of centers across image (S x T)
  - At each center, 3 scales and 3 aspect ratios

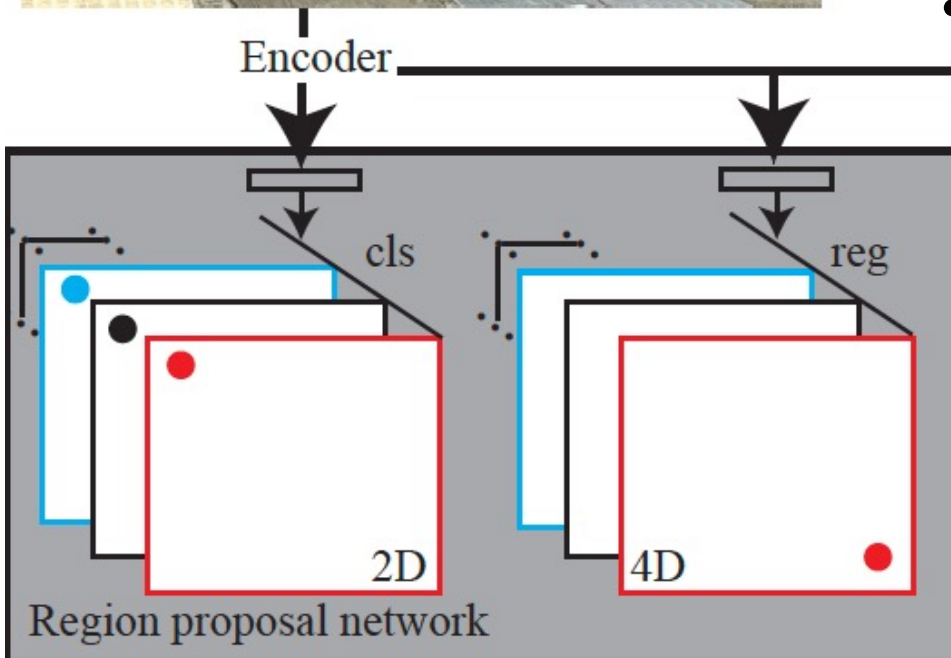


# Which boxes are worth looking at

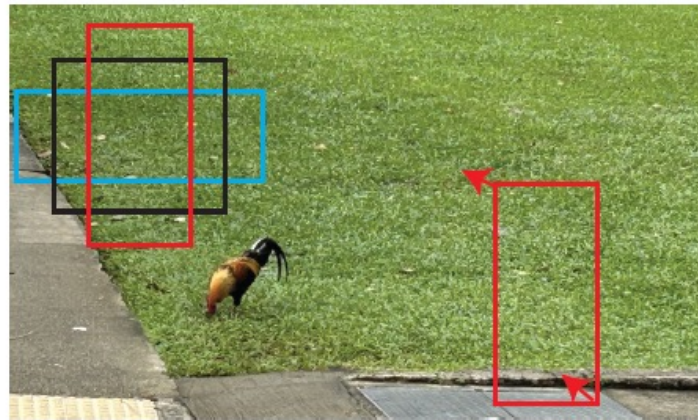


• **Region proposal network** produces:

- $(2 \times 9) \times S \times T$  block of box scores (cls)
  - objectness score
- $(4 \times 9) \times S \times T$  block of offsets (reg)
  - mild improvement to box corners to get score



# Scoring the boxes

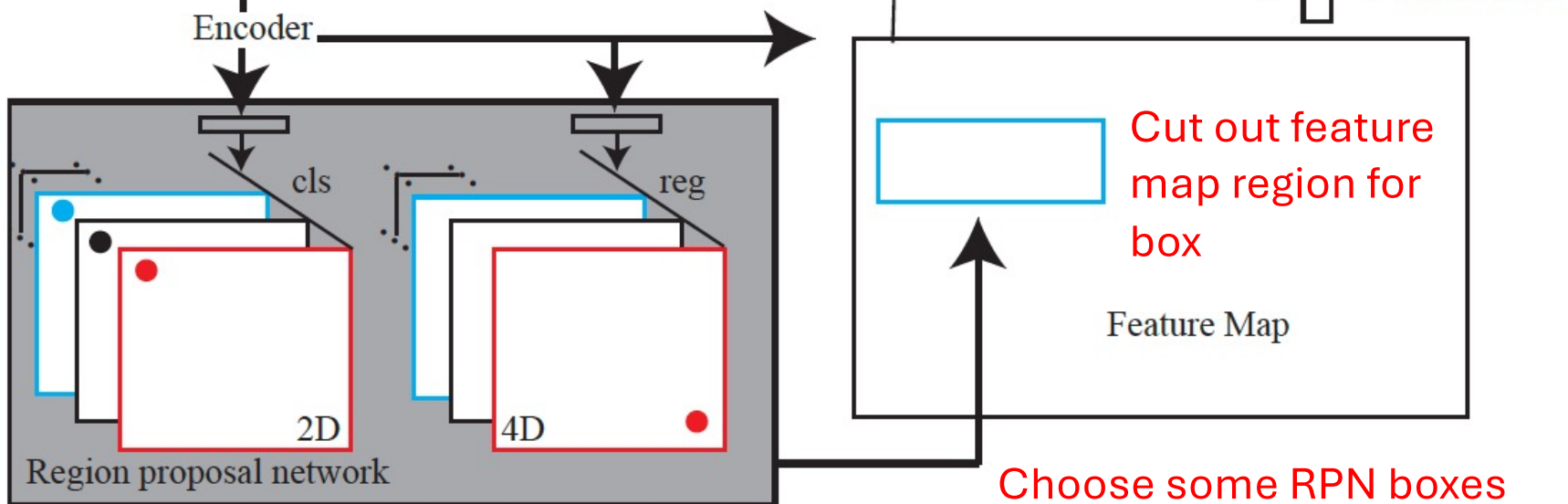


Adjust feature map region to fixed size

ROI pool

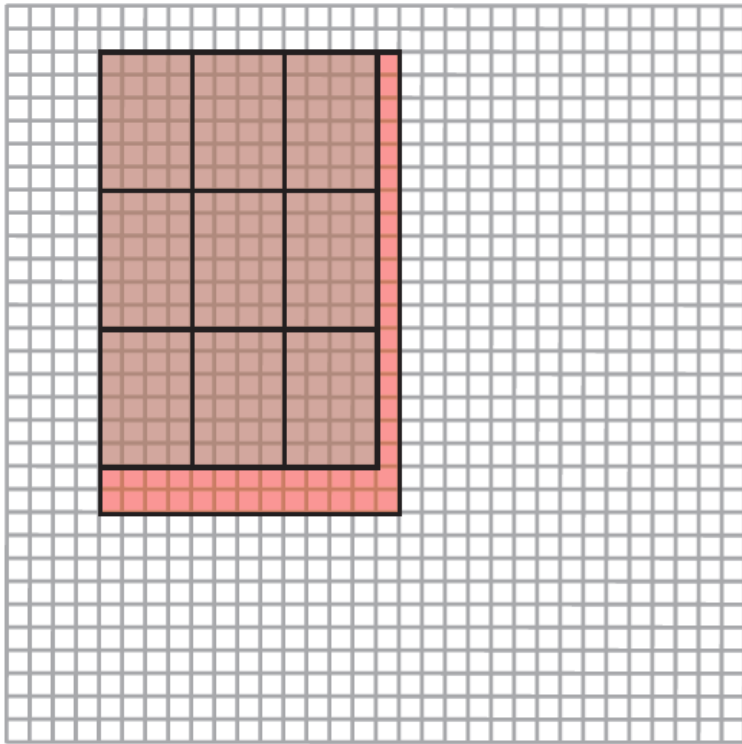
Scores

Box offsets



# Adjusting the feature map region

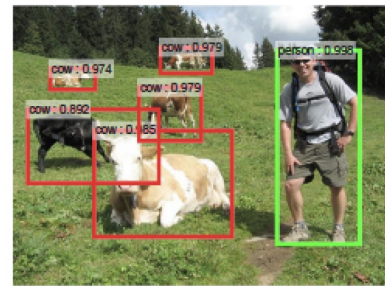
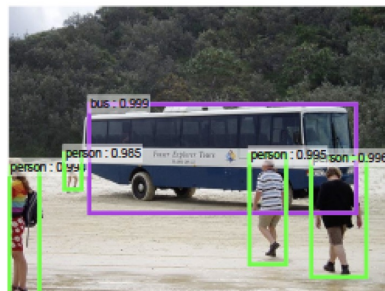
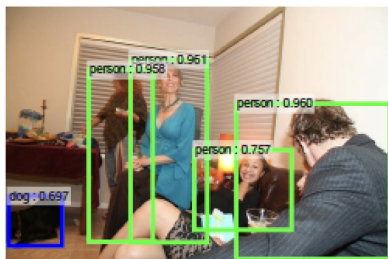
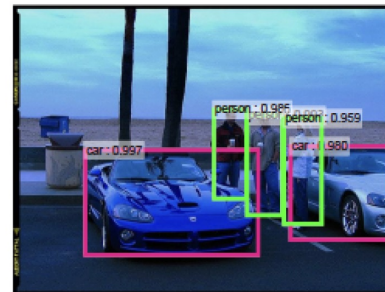
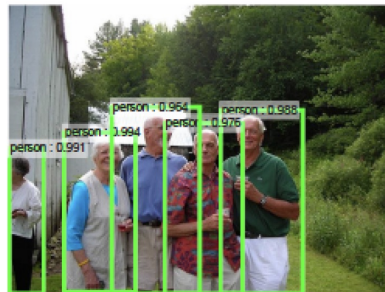
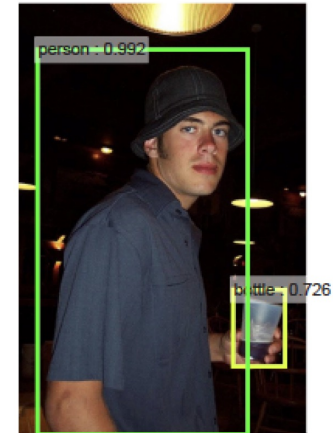
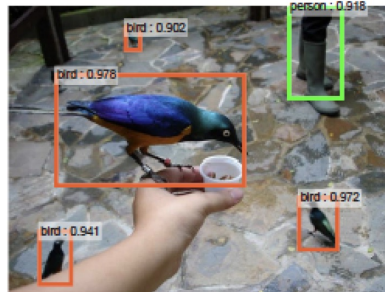
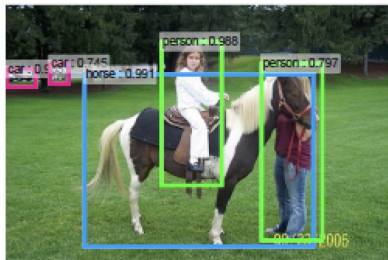
ROI Pool



- Find corresponding window
  - by rounding
- Within subgrid, max-pool
  - subgrid by rounding

Figure 26.2 shows what happens when a  $1024 \times 1024$  image produces a  $32 \times 32$  feature map which is ROI Pooled to a 9 dimensional vector for a given ROI. The top left corner of the ROI is at  $(75, 140)$  in the original image, and the size of the ROI is  $650 \times 430$ . In this case, ROI Pool would round the corner location to  $(2, 4) = (\lfloor 75/32 \rfloor, \lfloor 140/32 \rfloor)$ , the size to  $20 \times 13 = \lfloor 650/32 \rfloor \times \lfloor 430/32 \rfloor$ .

# Detection



# Training requires care

- Data:
  - labelled boxes on images
- I) Train RPN using pretrained encoder
  - Finetune ImageNet encoder
  - Note class imbalance
    - there are more empty than object boxes
- II) Train classifier using pretrained encoder
  - Finetune different ImageNet encoder
  - Use RPN boxes (now two encoders!)
- III) Finetune RPN layers on fixed encoder from II
- IV) Finetune classifier, etc
  - using RPN, encoder from III

# Mask RCNN

- Would like a better localization of object
- Get by:
  - modifying ROI Pool
  - decoding result into mask as well.

# ROIAlign preserves spatial detail

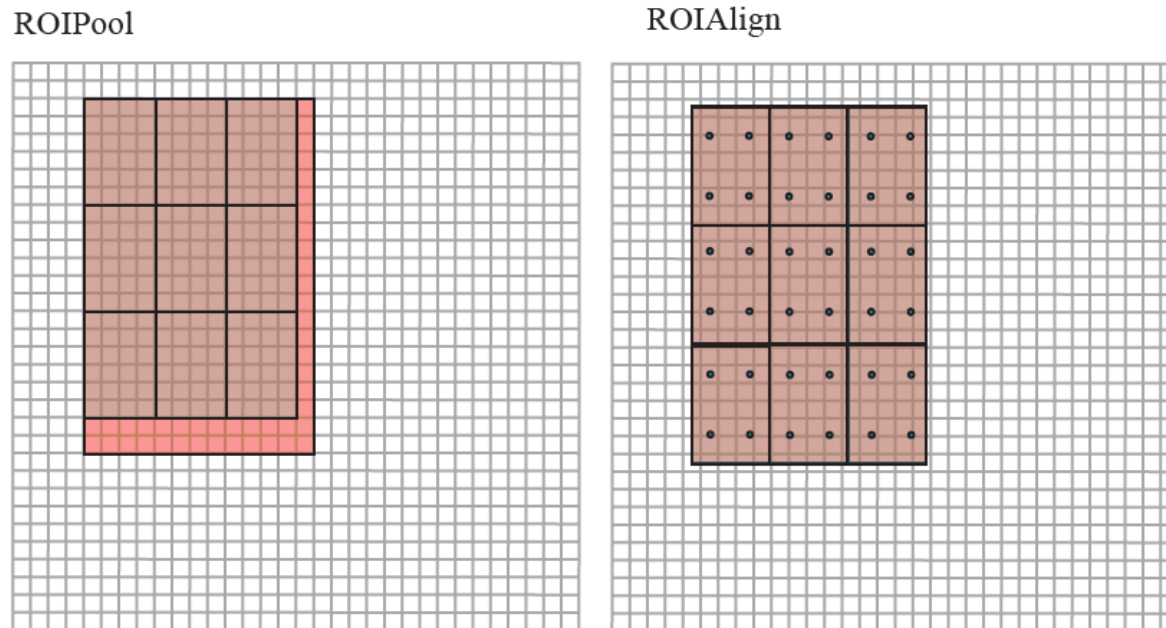
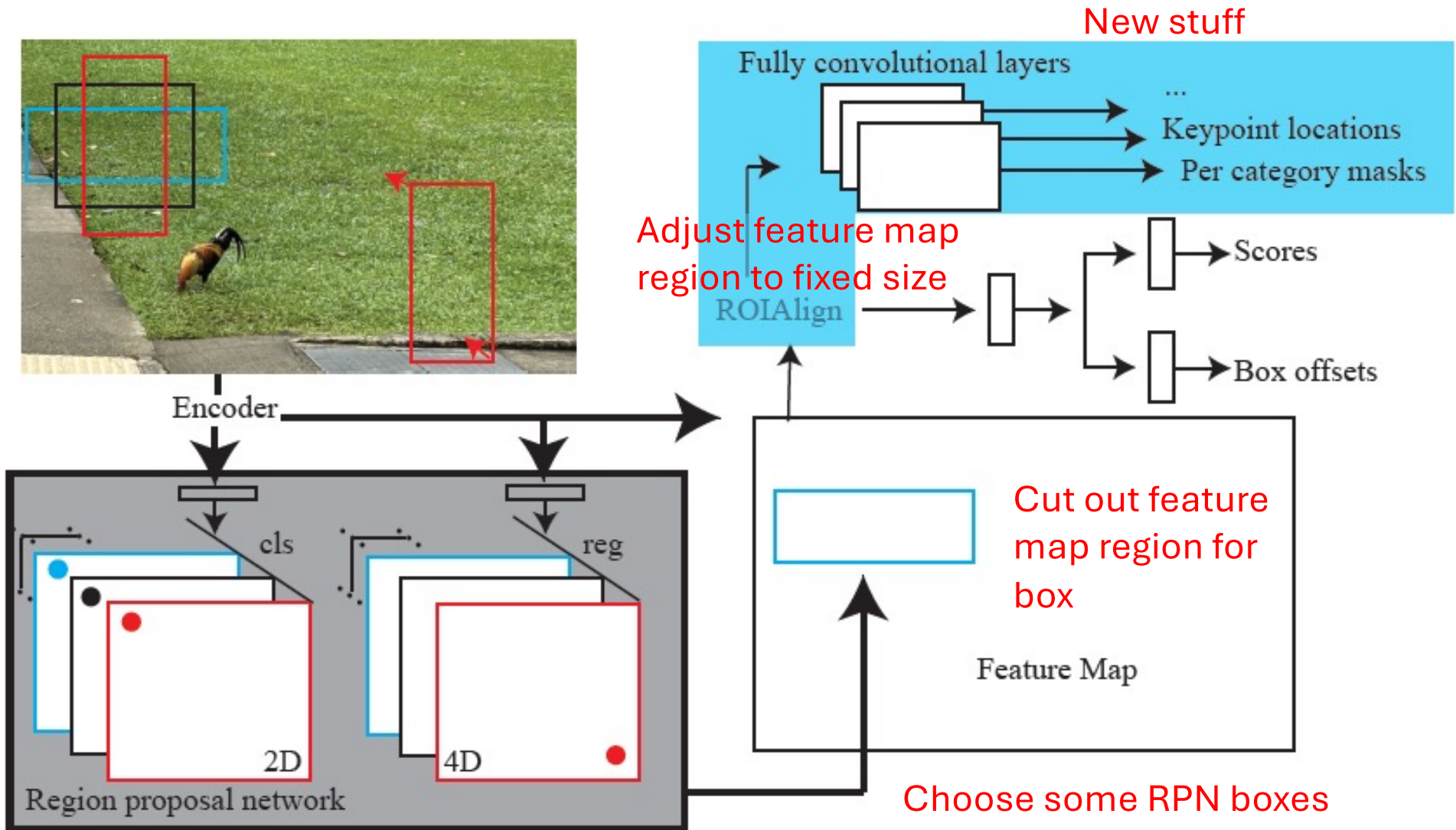
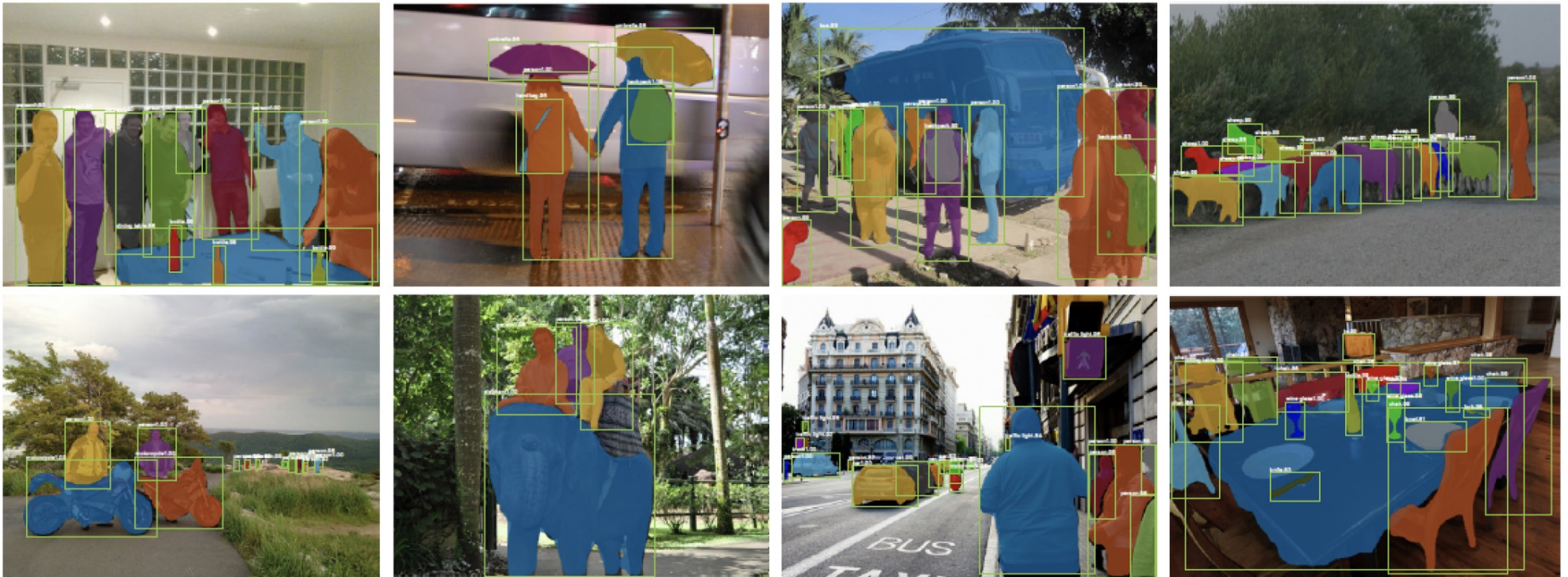


Figure 26.2 shows what happens when a  $1024 \times 1024$  image produces a  $32 \times 32$  feature map which is ROI Pooled to a 9 dimensional vector for a given ROI. The top left corner of the ROI is at  $(75, 140)$  in the original image, and the size of the ROI is  $650 \times 430$ . In this case, ROI Pool would round the corner location to  $(2, 4) = (\lfloor 75/32 \rfloor, \lfloor 140/32 \rfloor)$ , the size to  $20 \times 13 = \lfloor 650/32 \rfloor \times \lfloor 430/32 \rfloor$ .

# MaskRCNN



# Detection with masks



- Evaluation:
  - you can compute IoU for masks, too...

# Variants

- Predict more stuff from ROI
  - eg body keypoints



- Bolt onto a semantic segmenter
  - to get panoptic segmentation

# Became Detectron, Detectron 2

<https://detectron2.readthedocs.io/en/latest/>

# Things to think about...

- What other useful things could you predict with MaskRCNN?
- How do you guarantee performance?