

Object Detection – Part 2

Danna Gurari

University of Colorado Boulder
Fall 2023



Review

- Last lecture:
 - Object Detection Problem
 - Object Detection Applications
 - Object Detection Datasets
 - Object Detection Evaluation Metric
 - Overview of object detection algorithms and baseline (R-CNN)
- Assignments (Canvas)
 - Reading assignment was due earlier today
 - Next reading assignment due next Wednesday (guest lecture for Monday)
 - Project proposal due in two weeks
- Questions?

Object Detection: Today's Topics

- Overview of object detection algorithms
- Fast R-CNN
- Faster R-CNN
- YOLO
- Discussion

Recall Motivation: Go Faster While Getting Good Accuracy

Person?

Person?

Person?

Person?

Person?

Person?

Person?

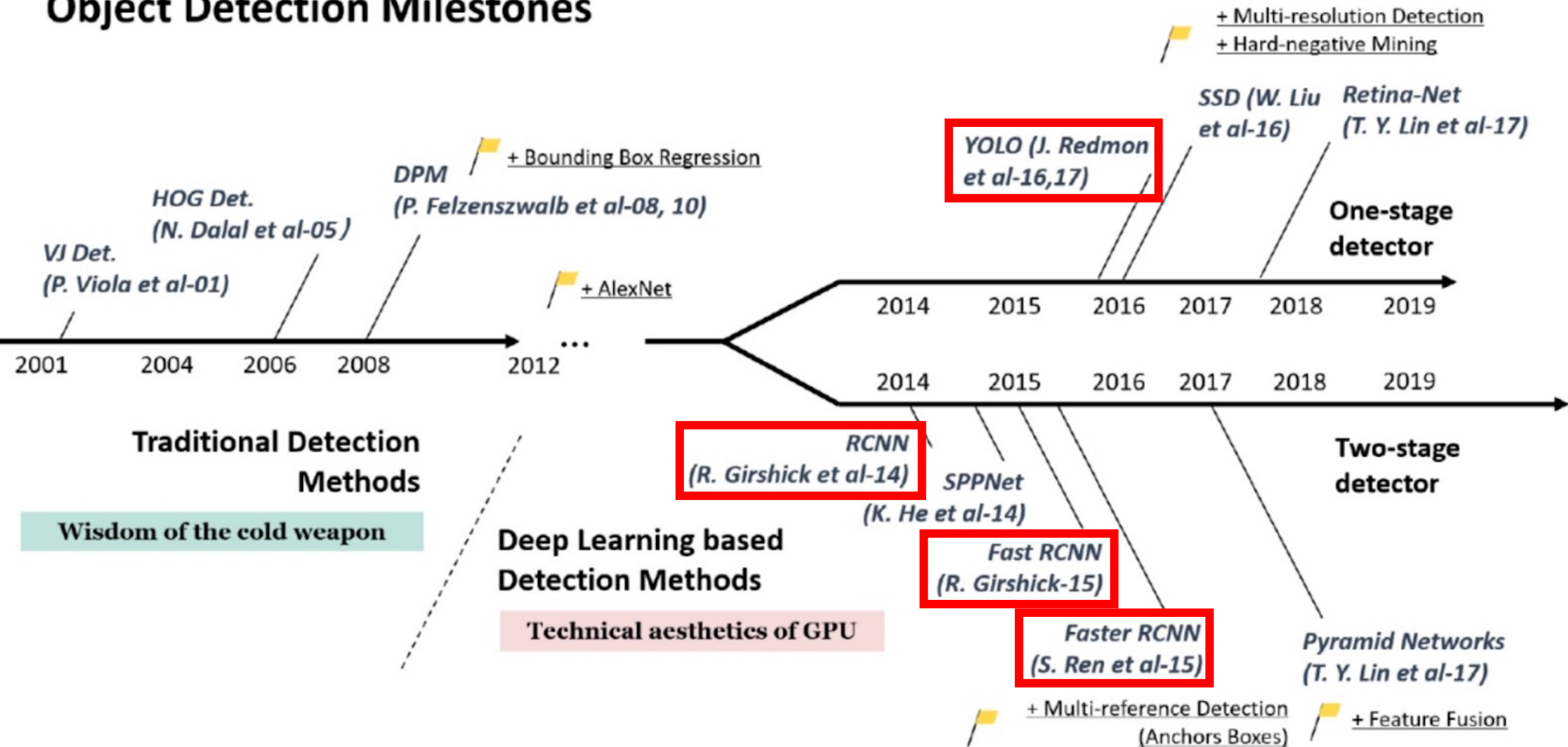
Person?

Person?



Image Source: <https://yourboulder.com/boulder-neighborhood-downtown/>

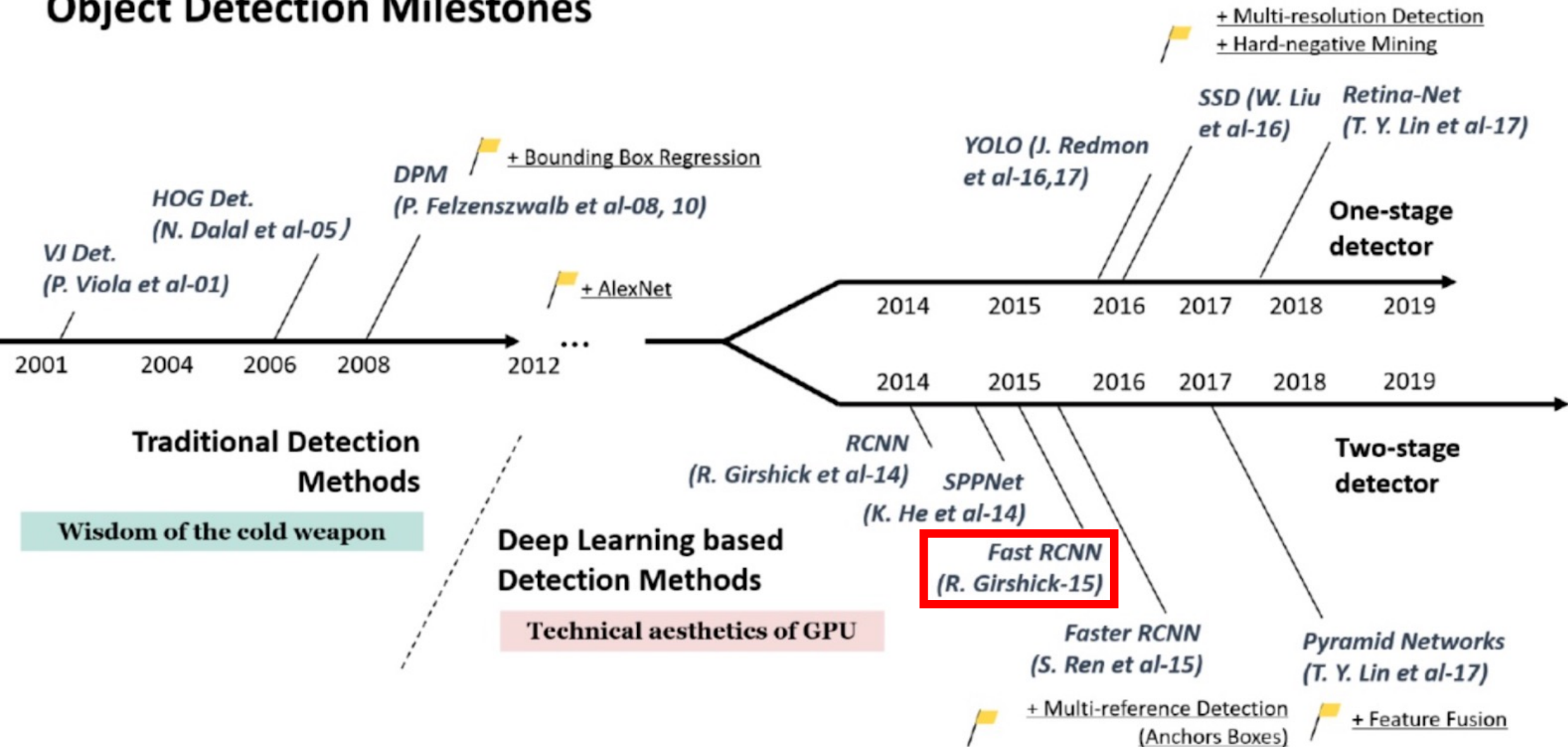
Object Detection Milestones



Object Detection: Today's Topics

- Overview of object detection algorithms
- **Fast R-CNN**
- Faster R-CNN
- YOLO
- Discussion

Object Detection Milestones



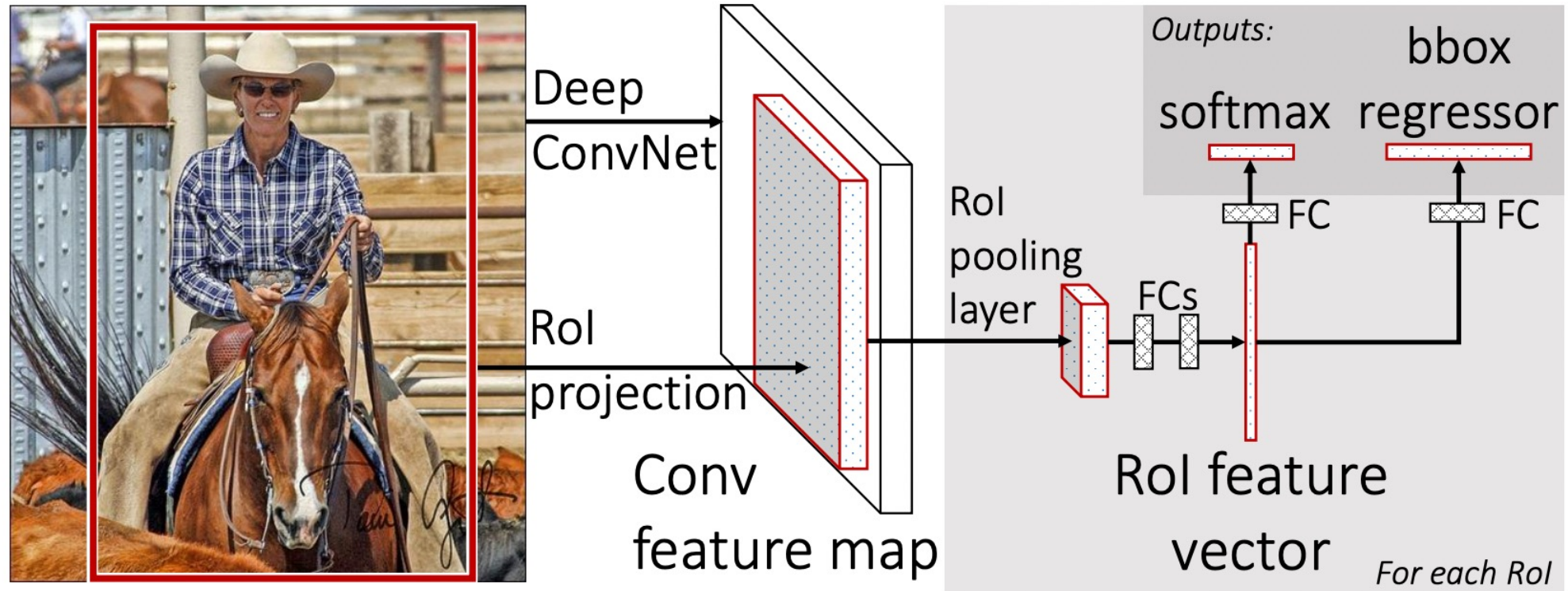
Key Contributions of Fast R-CNN

1. State of art object detection model in terms of accuracy and speed
 1. 9x faster than R-CNN
 2. mAP of 66% vs 62% for R-CNN on VOC2012
2. Reduced storage requirements by not requiring features to be stored for each region proposal
3. A training algorithm that learns in a single stage (rather than the three stages required by R-CNN)

Architecture

(Also, like R-CNN, apply non-maximum suppression to remove redundant predictions for the same object)

Like R-CNN, assign each region proposal to a class and refine it

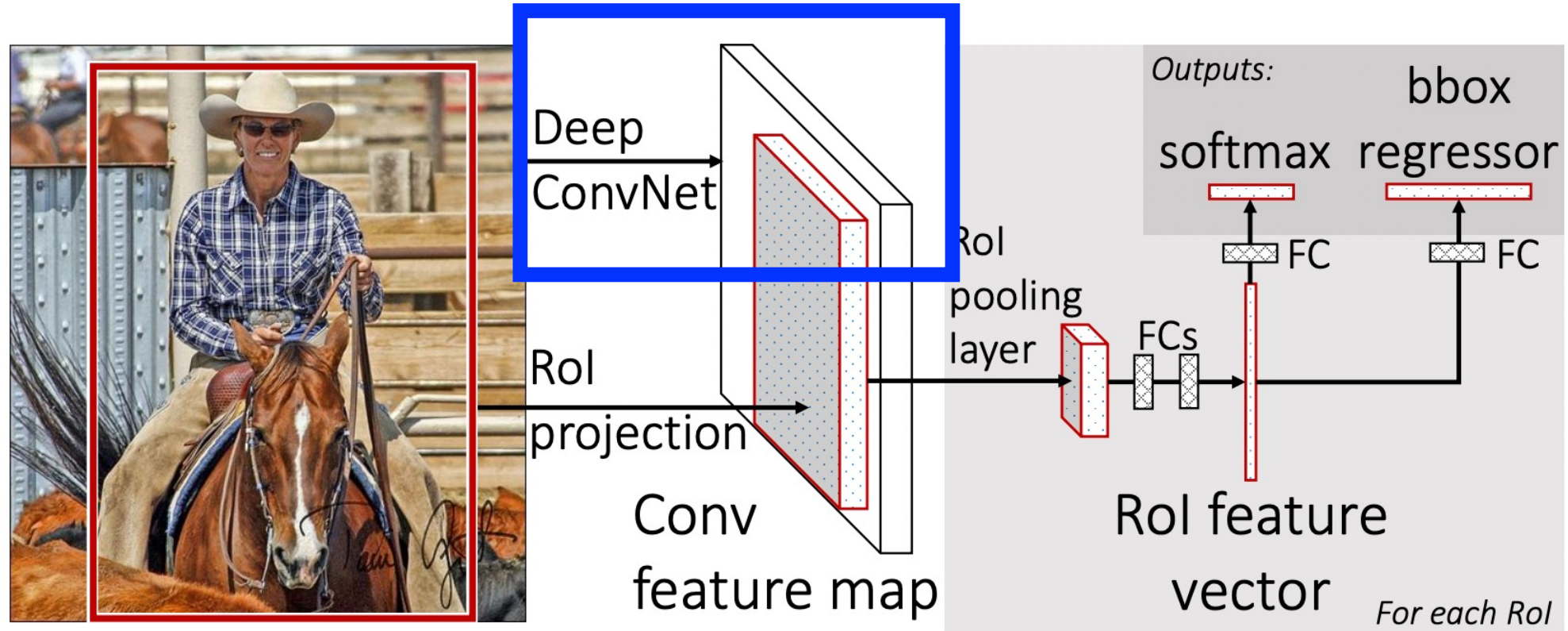


Input: any image size + “object”-like regions found using selective search (same as R-CNN)

Architecture

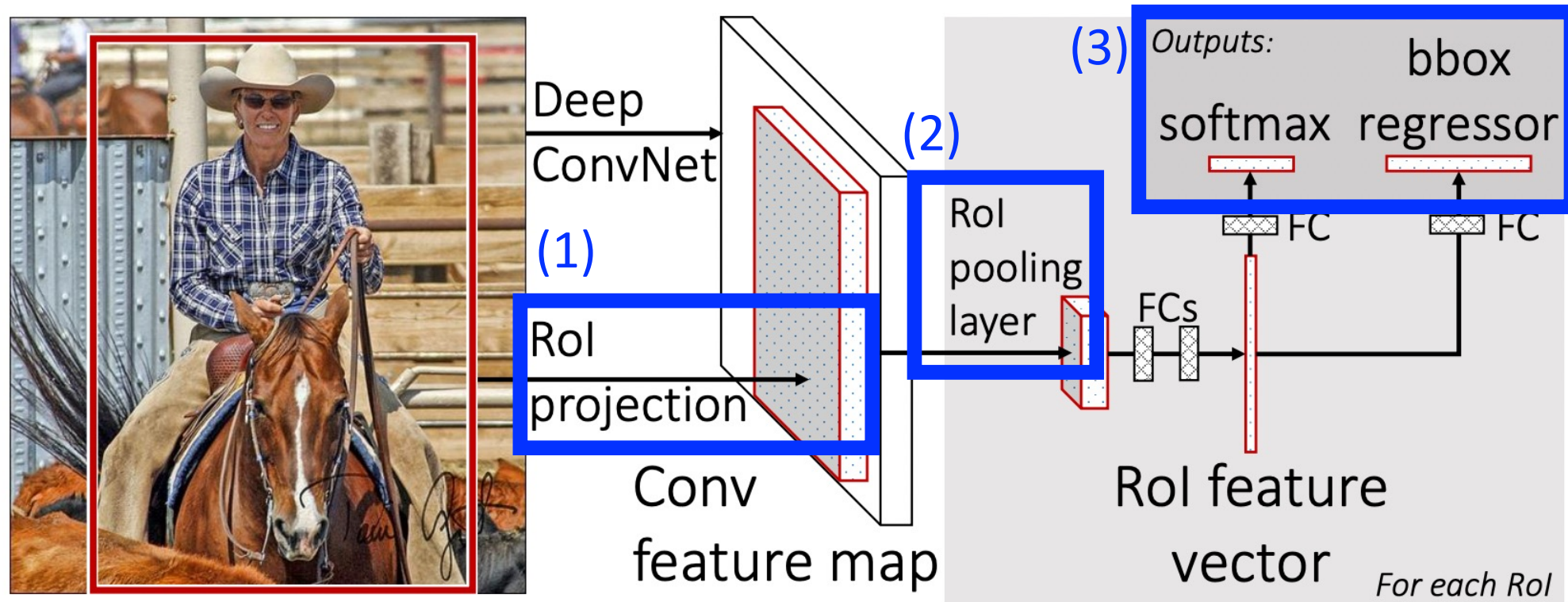
Image is fed through a pretrained ImageNet model to generate a feature map for the **entire image**. Which models could we use?

- Was tested with AlexNet (small), VGG16 (large), and another VGG variant

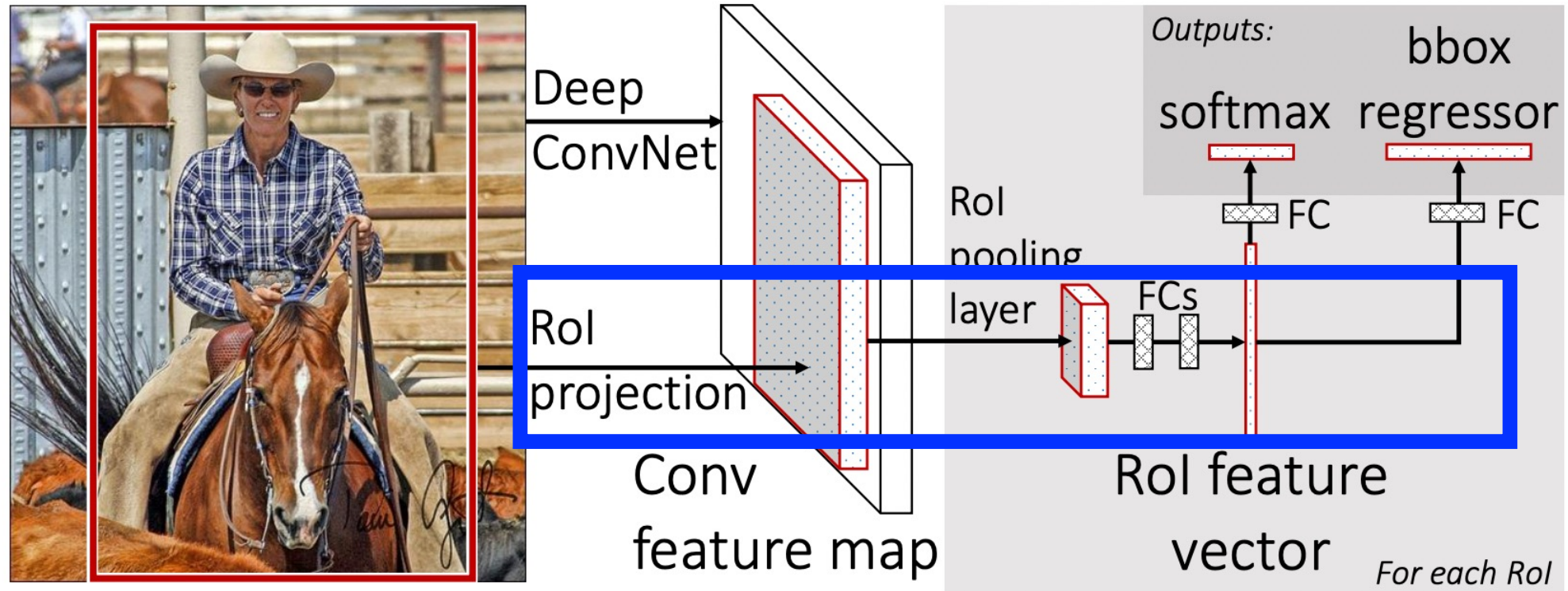


Architecture

Each pretrained model is modified to: (1) update input to accept list of region proposals, (2) use ROI pooling layer to support arbitrary image sizes, and (3) use output layer for object detection instead of image classification

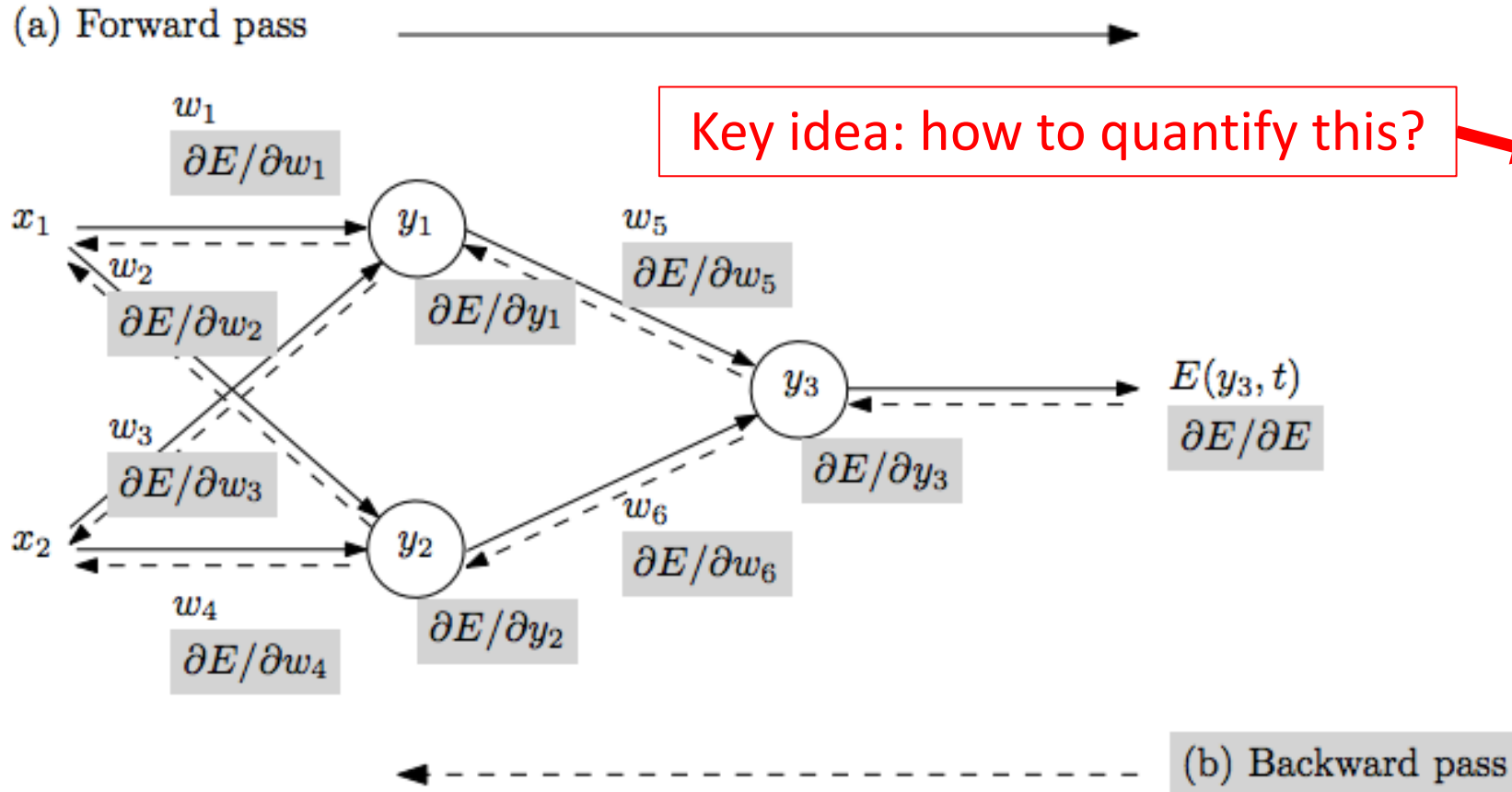


Architecture



The modified input enables passing only the portion of the feature map inside each region proposal for prediction

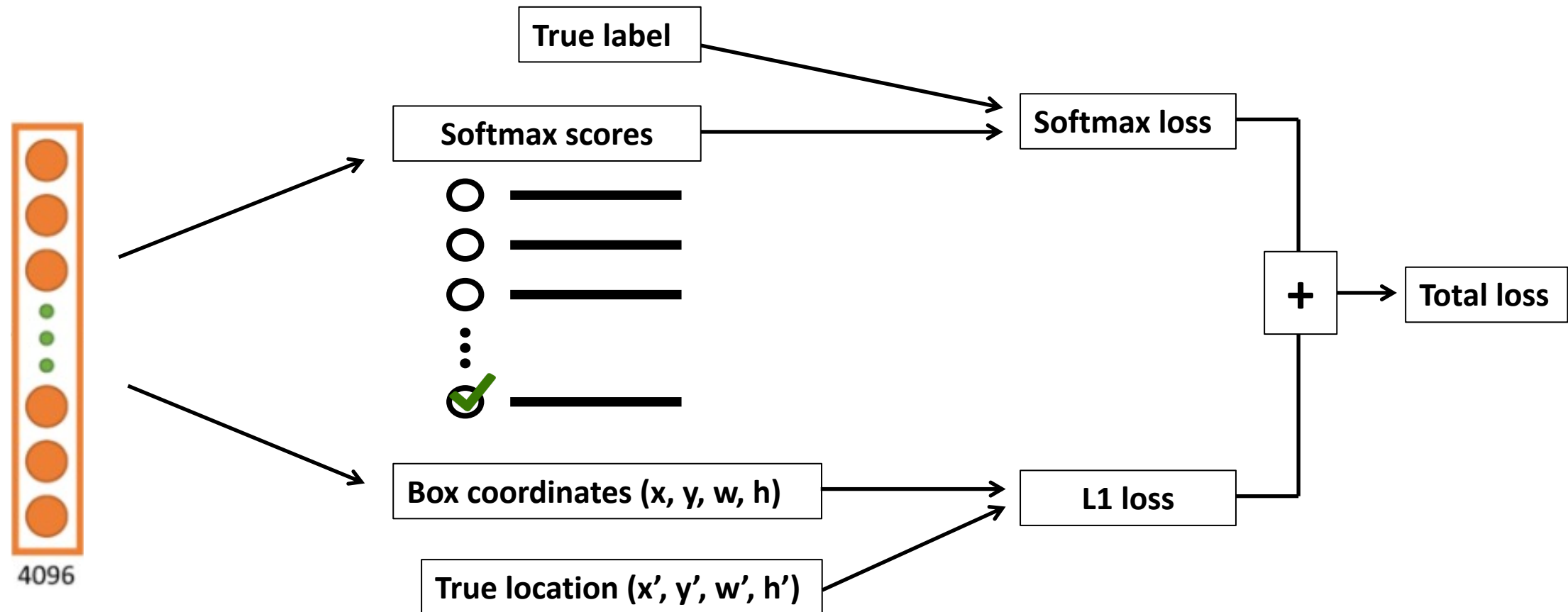
Training in One Stage



- Repeat until stopping criterion met:
 1. **Forward pass:** propagate training data through model to make predictions
 2. **Error quantification:** measure dissatisfaction with a model's predictions on training data
 3. **Backward pass:** using predicted output, calculate gradients backward to assign blame to each model parameter; **account for weight sharing by using average of all connections for a parameter**
 4. Update each parameter using calculated gradients

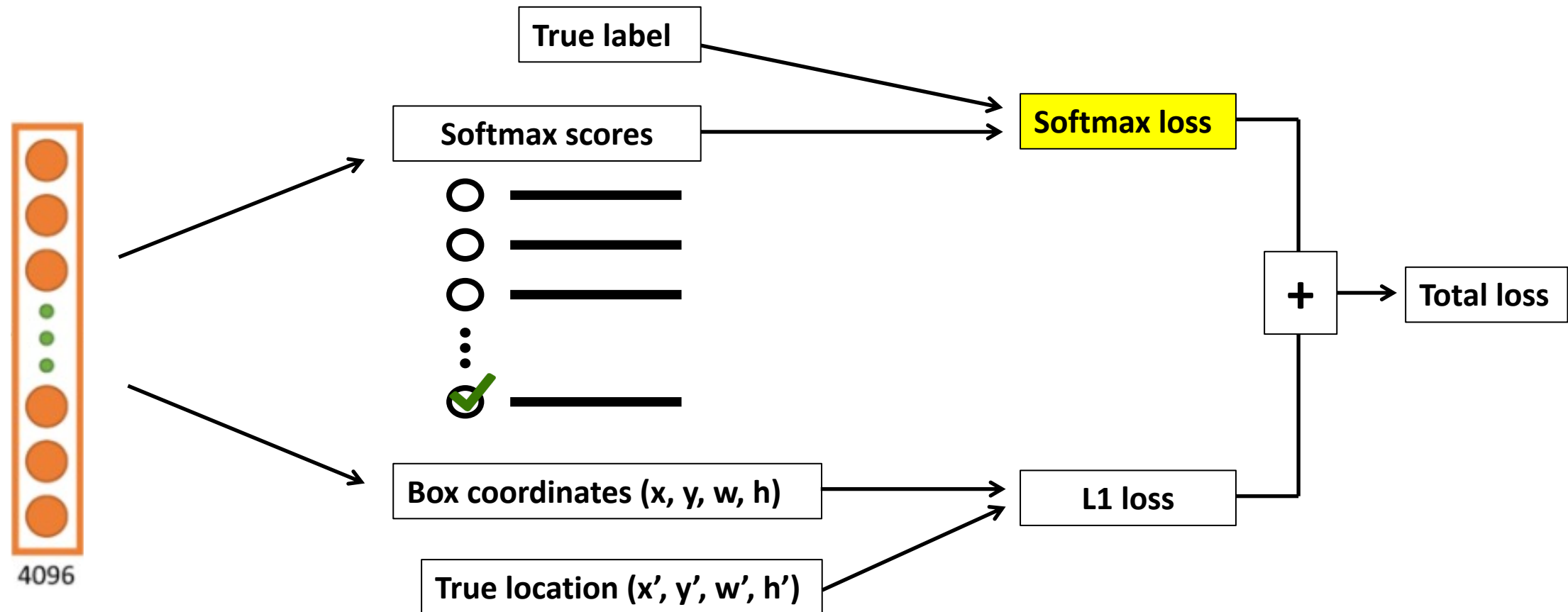
Training in One Stage: Multi-task Loss

Objective function sums classification and localization losses for each region proposal

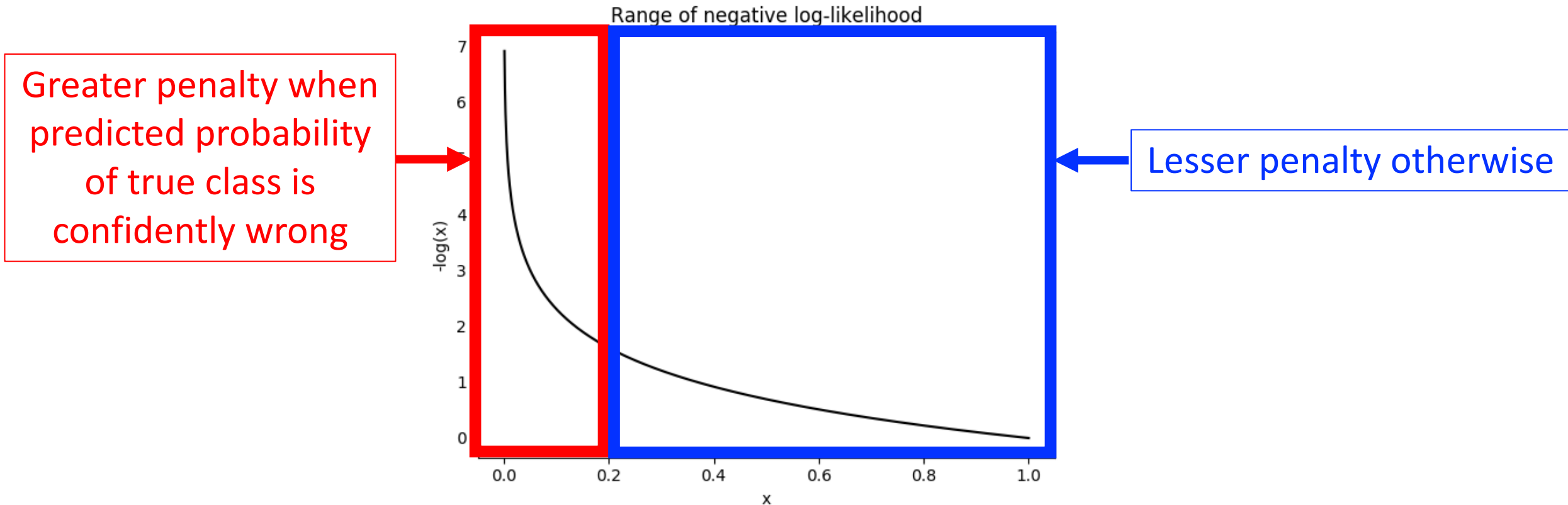


Training in One Stage: Multi-task Loss

Objective function sums classification and localization losses for each region proposal



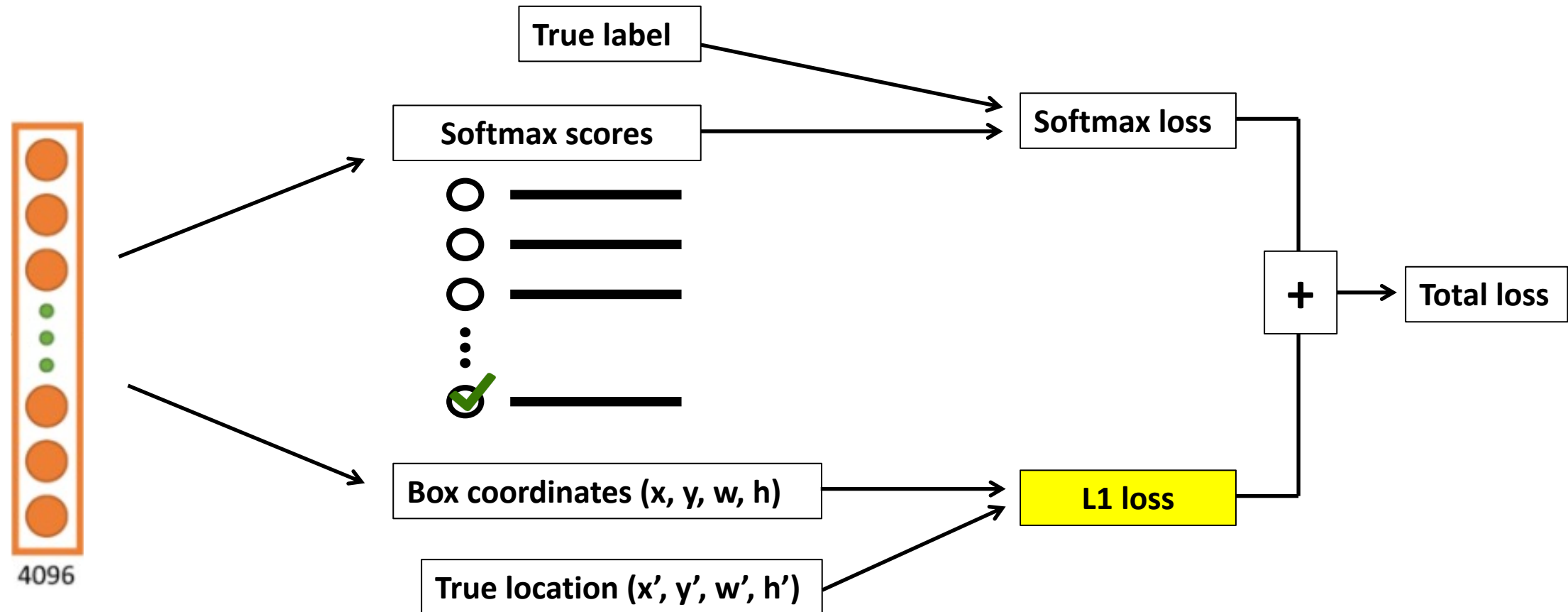
Classification Loss (Recall Cross Entropy Loss)



$$-\log \frac{\exp(w_k \cdot x + b_k)}{\sum_{j=1}^K \exp(w_j \cdot x + b_j)}$$

Training in One Stage: Multi-task Loss

Objective function sums classification and localization losses for each region proposal



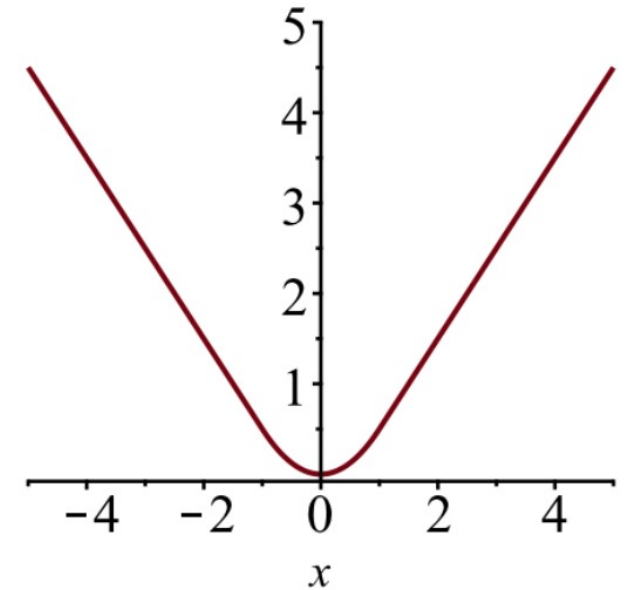
Localization Loss

$$\mathcal{L}_{\text{box}}(t^u, v) = \sum_{i \in \{x, y, w, h\}} L_1^{\text{smooth}}(t_i^u - v_i) \rightarrow L_1^{\text{smooth}}(x) = \begin{cases} 0.5x^2 & \text{if } |x| < 1 \\ |x| - 0.5 & \text{otherwise} \end{cases}$$

True location for true class "u"

Predicted location for class u

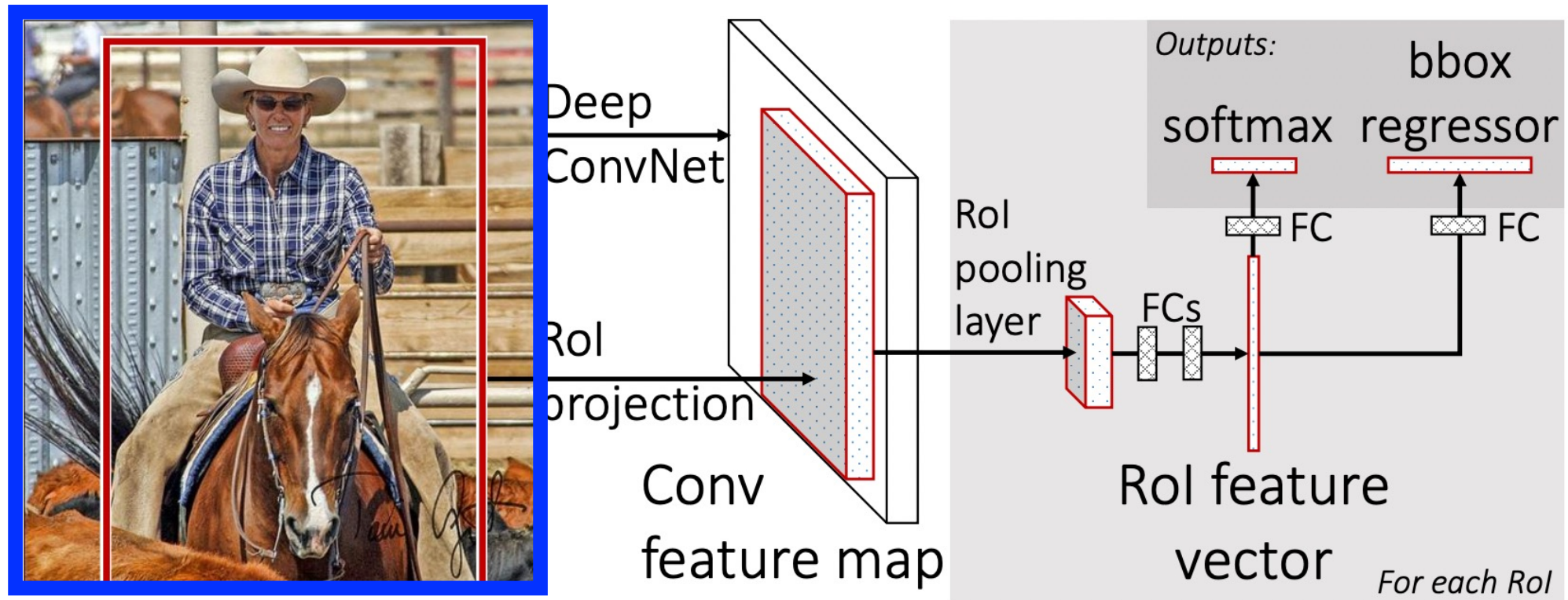
Less sensitive to outliers than SSE



Interesting Experimental Analysis

- Which layers to fine-tune?
 - Fine-tune layers after the first layer since including it doesn't boost performance, possibly because it captures task-independent features (e.g., lines, colors)
- Does multi-task training help?
 - Yes; removing either the classification loss or inserting regression loss independently worsens results
- Does more training data help?
 - Yes!
- Does more object proposals help?
 - The evidence suggests it can actually hurt performance

Key Limitation: Still Slow and Complex

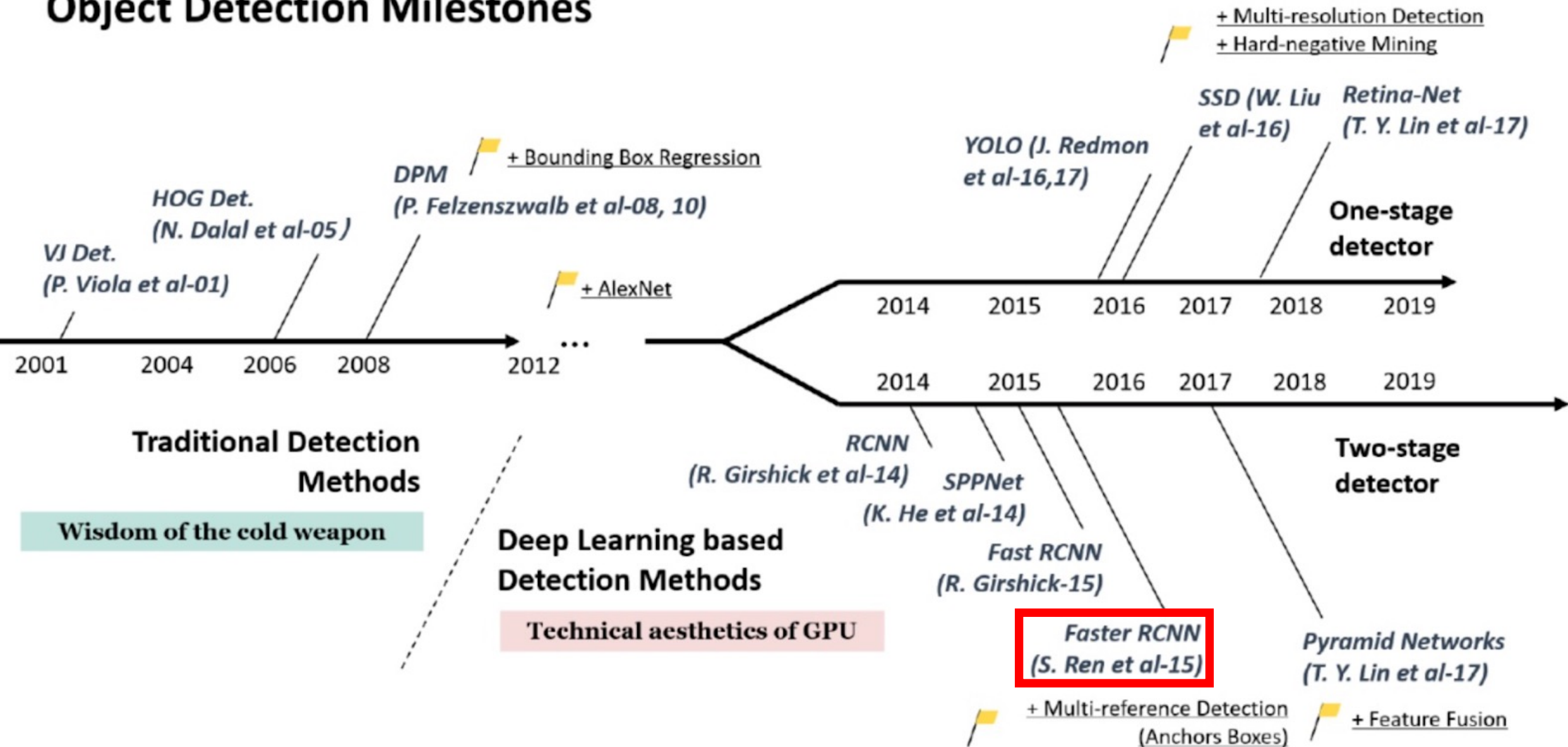


Still, requires generating “object”-like region proposals to provide as input (estimated then to take 2 sec on a CPU) and treated as a separate step.

Object Detection: Today's Topics

- Overview of object detection algorithms
- Fast R-CNN
- **Faster R-CNN**
- YOLO
- Discussion

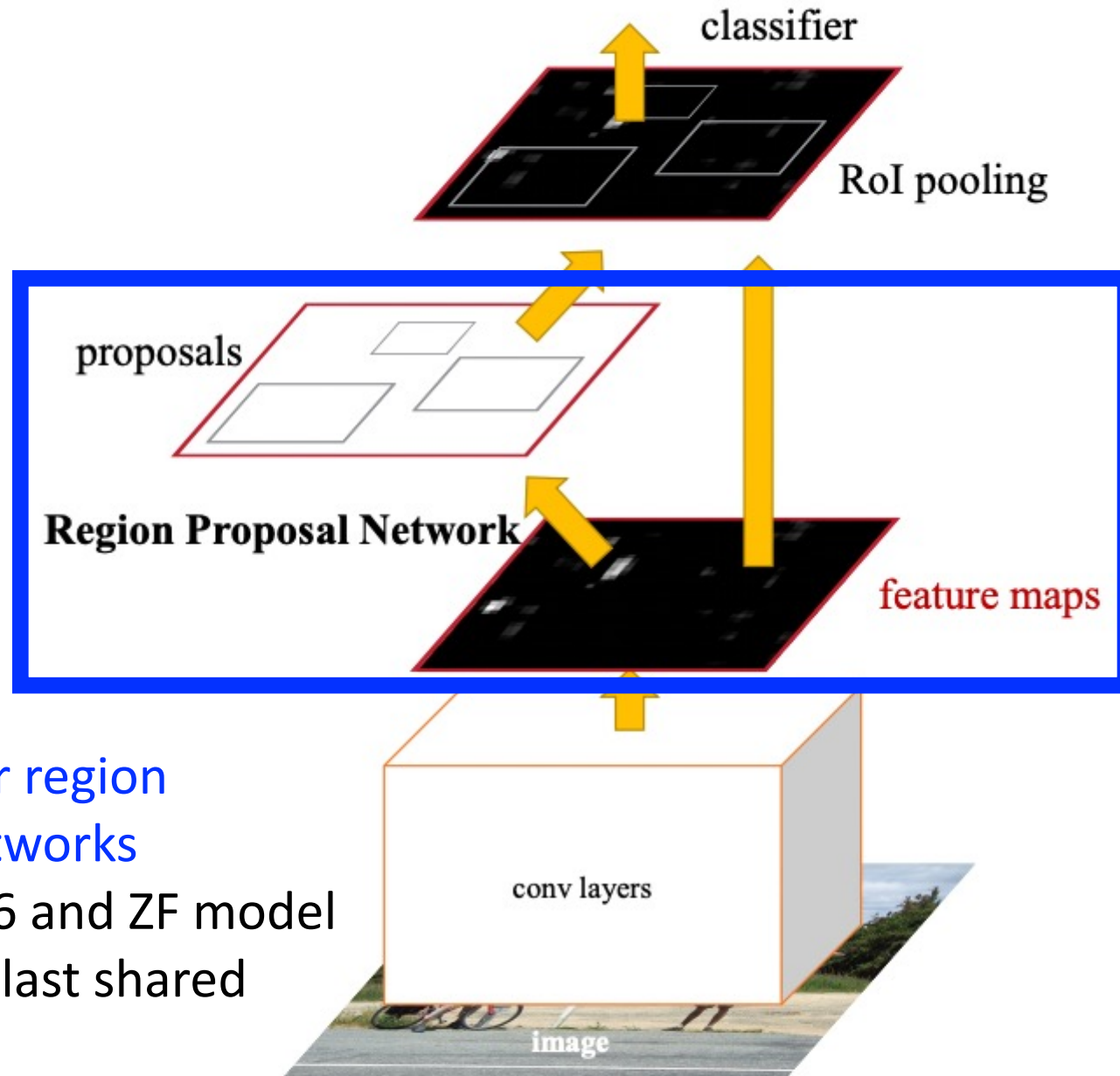
Object Detection Milestones



Key Contributions of Faster R-CNN

1. State of art object detection model in terms of accuracy and speed
2. An end-to-end trained model that learns all parts of the pipeline, particularly with the addition of region proposals

Architecture



Key novelty: adds finding region proposals to network so that all parts of model are learned in end-to-end fashion

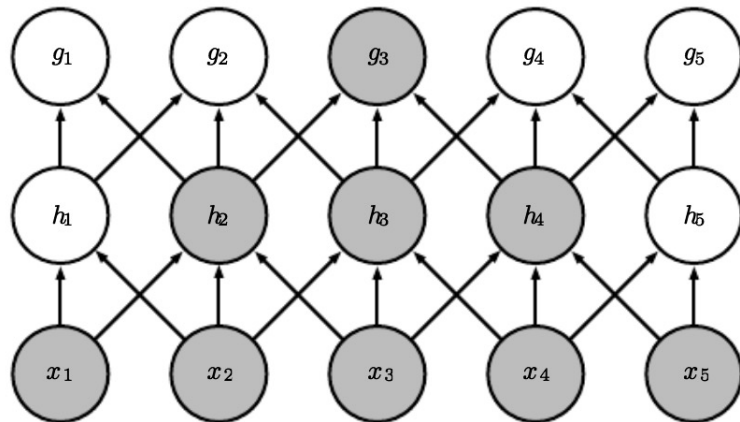
Convolutional layers shared for region proposal and detection subnetworks

- 2 architectures tested: VGG16 and ZF model
- input to both subnetworks is last shared convolutional layer

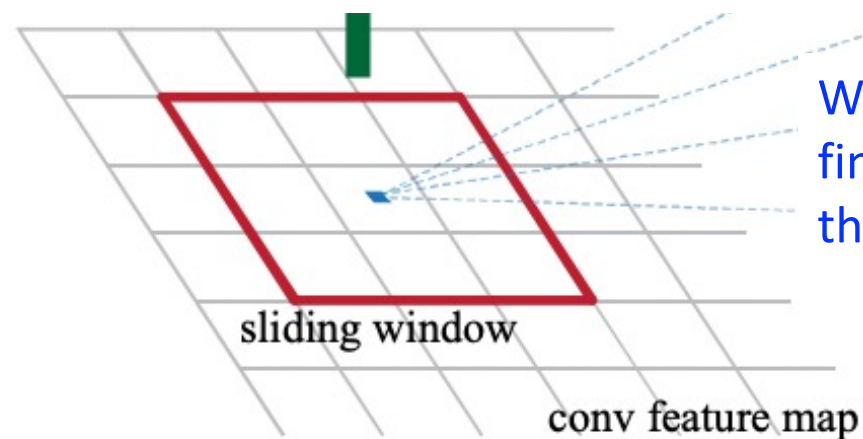
Architecture: Region Proposal Network

Input: convolutional feature map from last layer shared with object detector

Step 1: 3 x 3 convolutional filter applied to identify candidate proposals (recall, filter in the middle of an architecture maps to a larger input space, aka receptive field)



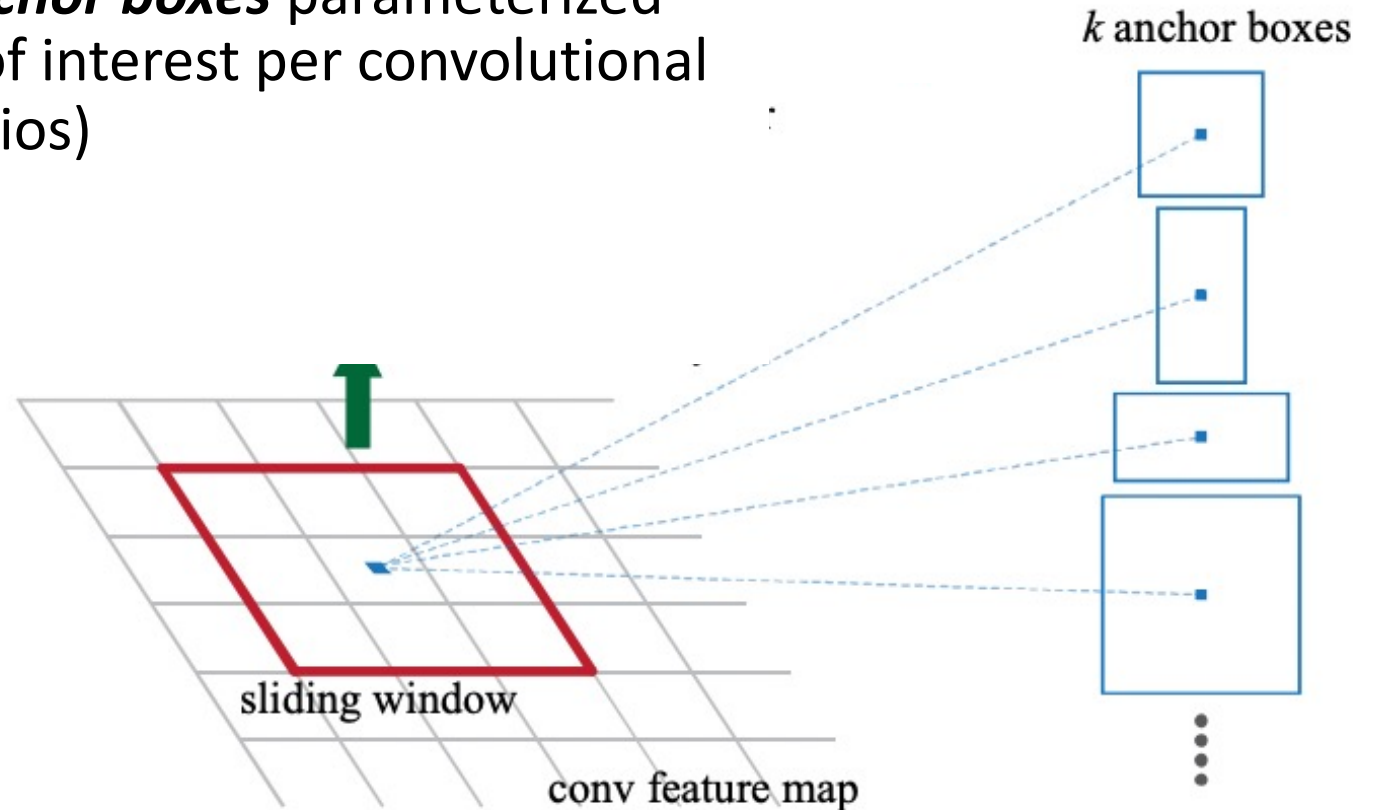
<https://www.deeplearningbook.org/contents/convnets.html>



Why can such an approach find objects larger than the receptive field?

Architecture: Region Proposal Network

Step 2: candidate regions of multiple scales and aspect ratios are supported efficiently using $k = 9$ **anchor boxes** parameterized relative to the feature map region of interest per convolutional operation (3 scales and 3 aspect ratios)

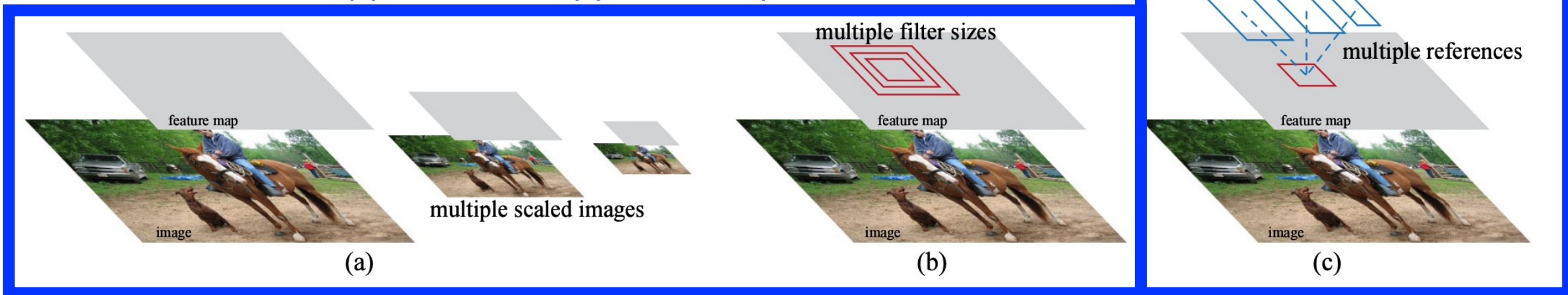


Architecture: Region Proposal Network

Step 2: candidate regions of multiple scales and aspect ratios are supported efficiently using $k = 9$ **anchor boxes** parameterized relative to the feature map region of interest per convolutional operation (3 scales and 3 aspect ratios)

Idea: images, feature maps, and filters all are a single size by using a pyramid of anchors (more efficient!)

Prior work's approach to support multiple scales



Architecture: Region Proposal Network

(k independent regressors learned to support k anchor box dimensions)

For each region, predict:

(1) Probability of object/not object

(2) Parameters to regress anchor box to GT box (center, width, and height)

$2k$ scores

cls layer

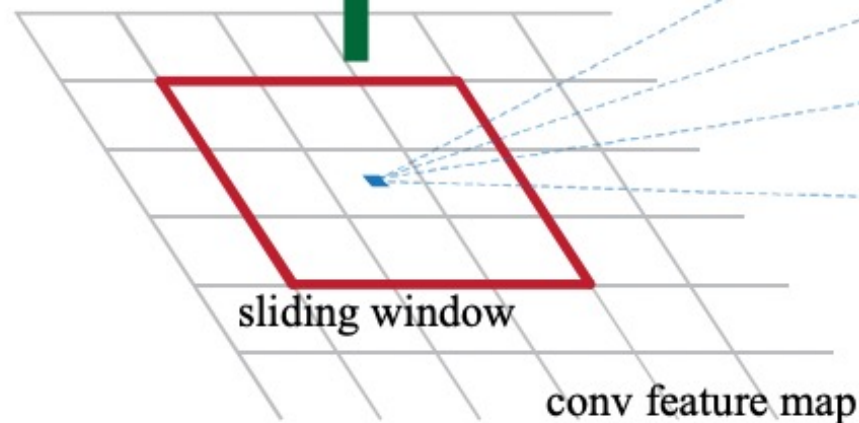
$4k$ coordinates

reg layer

256-d

intermediate layer

k anchor boxes



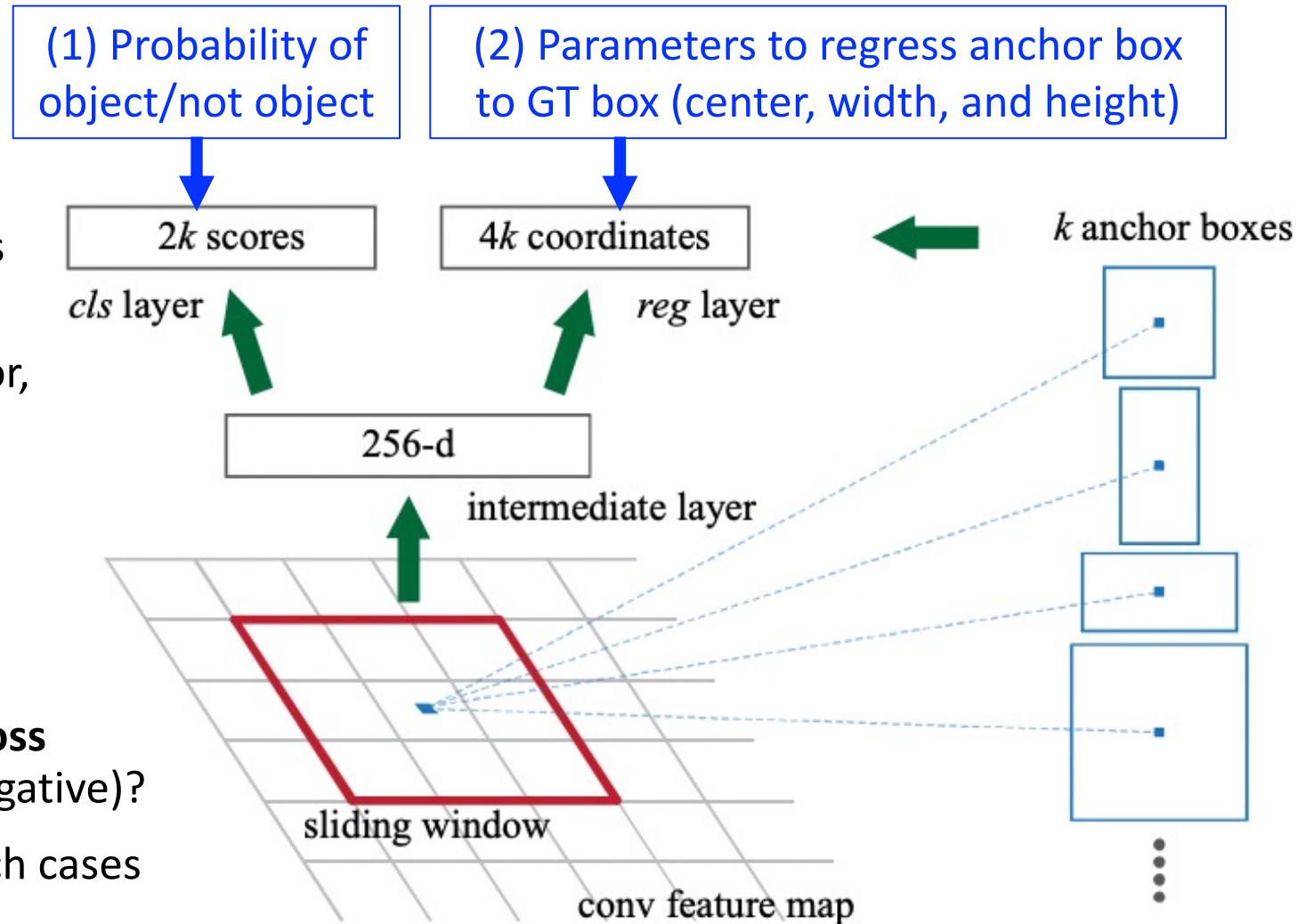
Training: Region Proposal Network

Training

- Loss: for each region proposal, sum classification and localization losses
- GT positive: anchors with IoU > 0.7 with GT (can be multiple anchors) or, when none, highest scoring one
- GT negative: non-positive anchors with IoU < 0.3 with GT
- Any non-assigned anchors ignored

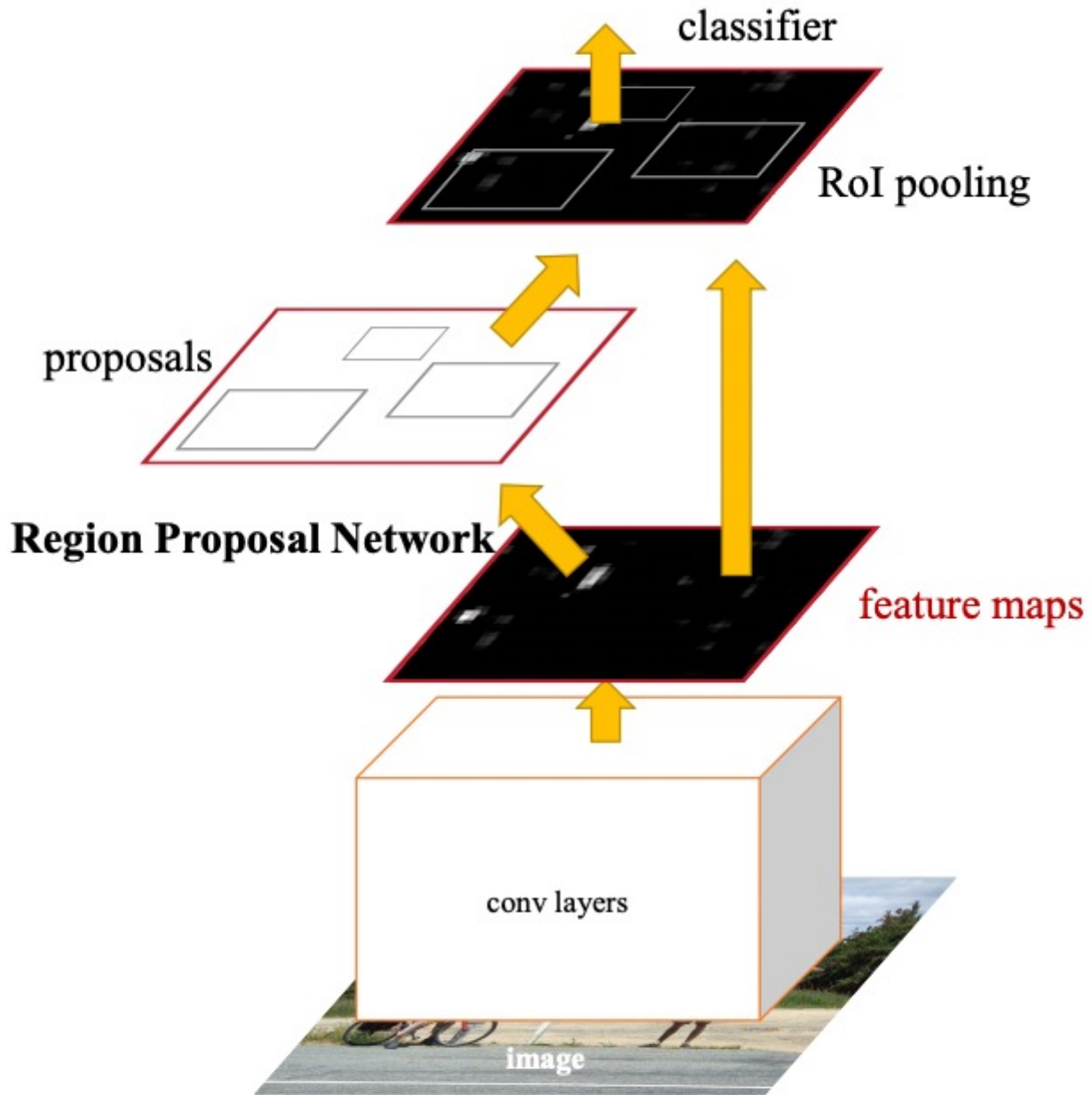
What is relevance of the **regression loss** when no object is present (i.e., GT negative)?

- none; regression loss disabled in such cases



Training: Overall

1. Train RPN
2. Train Fast R-CNN using proposals from pretrained RPN
3. Fine-tune layers unique to RPN
4. Fine-tune the fully connected layers of Fast R-CNN



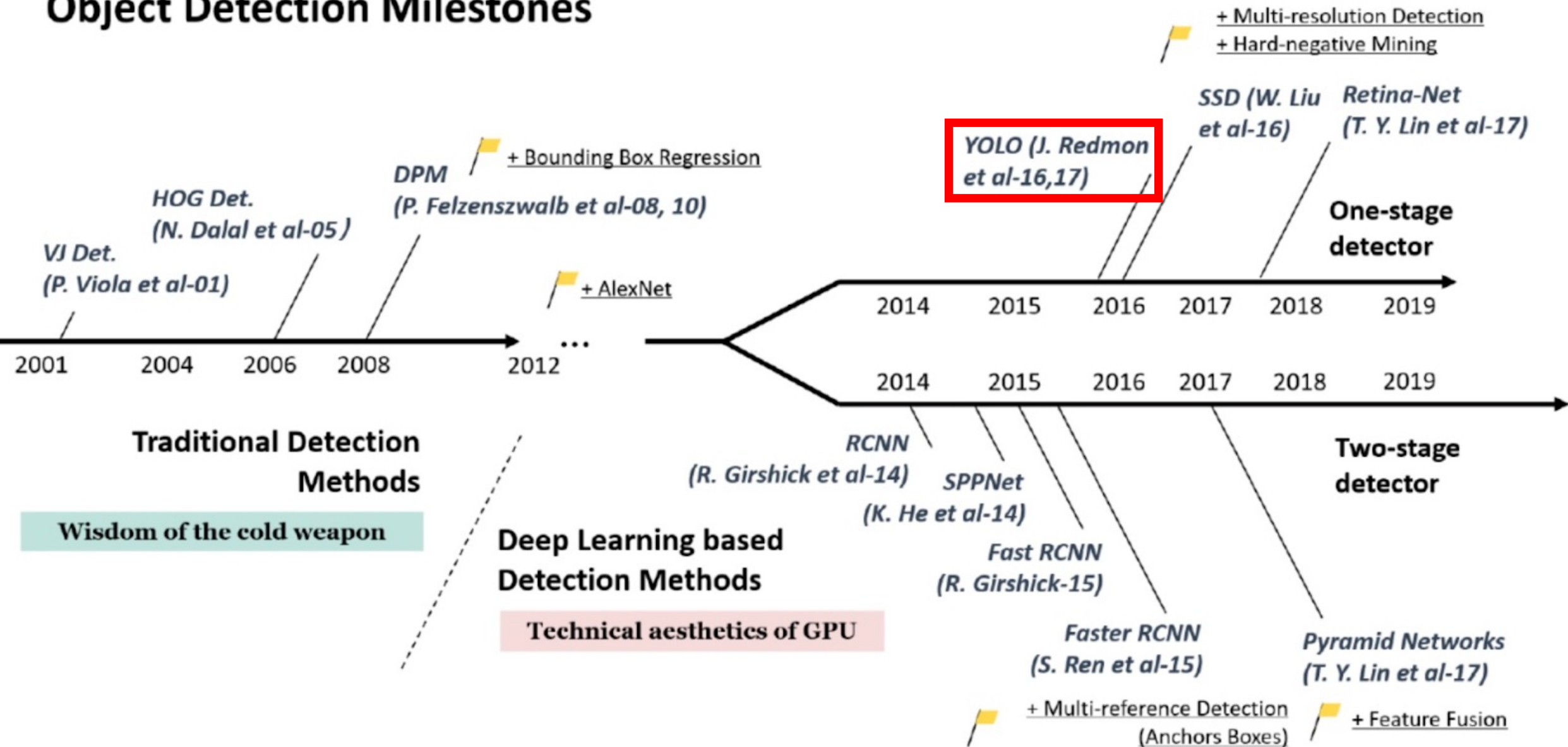
Limitations

- Still relatively slow; i.e., does not support real-time performance

Object Detection: Today's Topics

- Overview of object detection algorithms
- Fast R-CNN
- Faster R-CNN
- YOLO
- Discussion

Object Detection Milestones



Why YOLO?

Named after the proposed technique: **You Only Look Once**

Joseph Redmon, Santosh Divvala, Ross Girshick, and Ali Farhadi. “You Only Look Once: Unified, Real-Time Object Detection.” CVPR 2016.

Key Contributions

- A CNN architecture that detects objects by looking at the entire image **once**, treating the problem as a regression problem
 - i.e., a paradigm shift away from using classifiers to label image regions
- Most accurate **real-time** object detection system (i.e., 30+ fps)
- Generalizes better than existing approaches to out-of-domain data

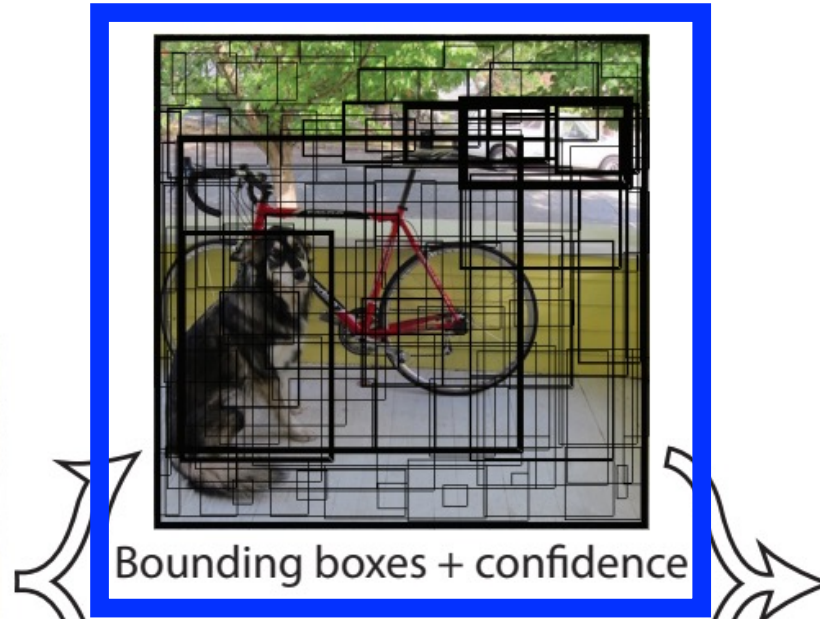
Approach

2. For each grid cell, (1) locate (potentially multiple) objects and (2) predict a probability distribution for class labels (assuming an object is present)

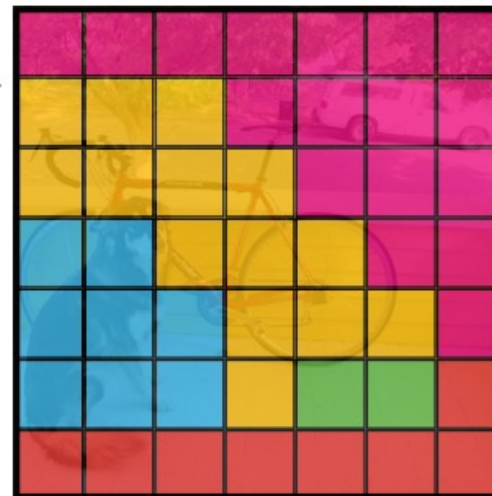
1. Divide image into grid



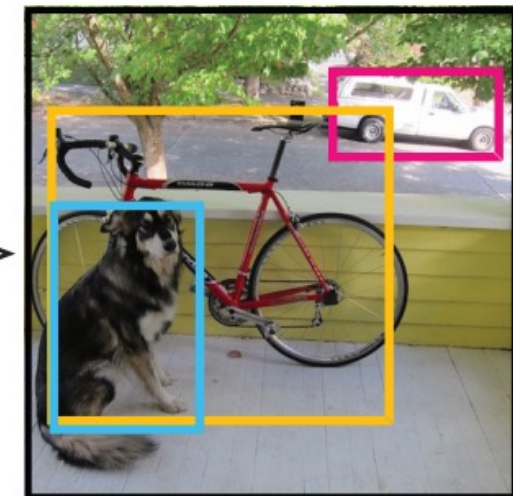
$S \times S$ grid on input



Bounding boxes + confidence



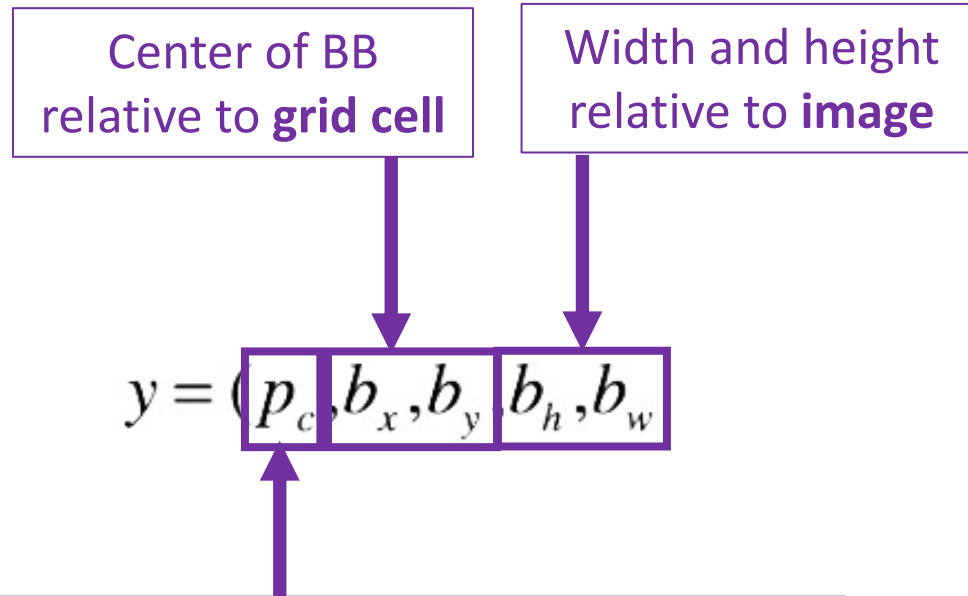
Class probability map



Final detections

3. For each grid cell, (1) locate (potentially multiple) objects

Approach: BB Prediction Per Grid Cell

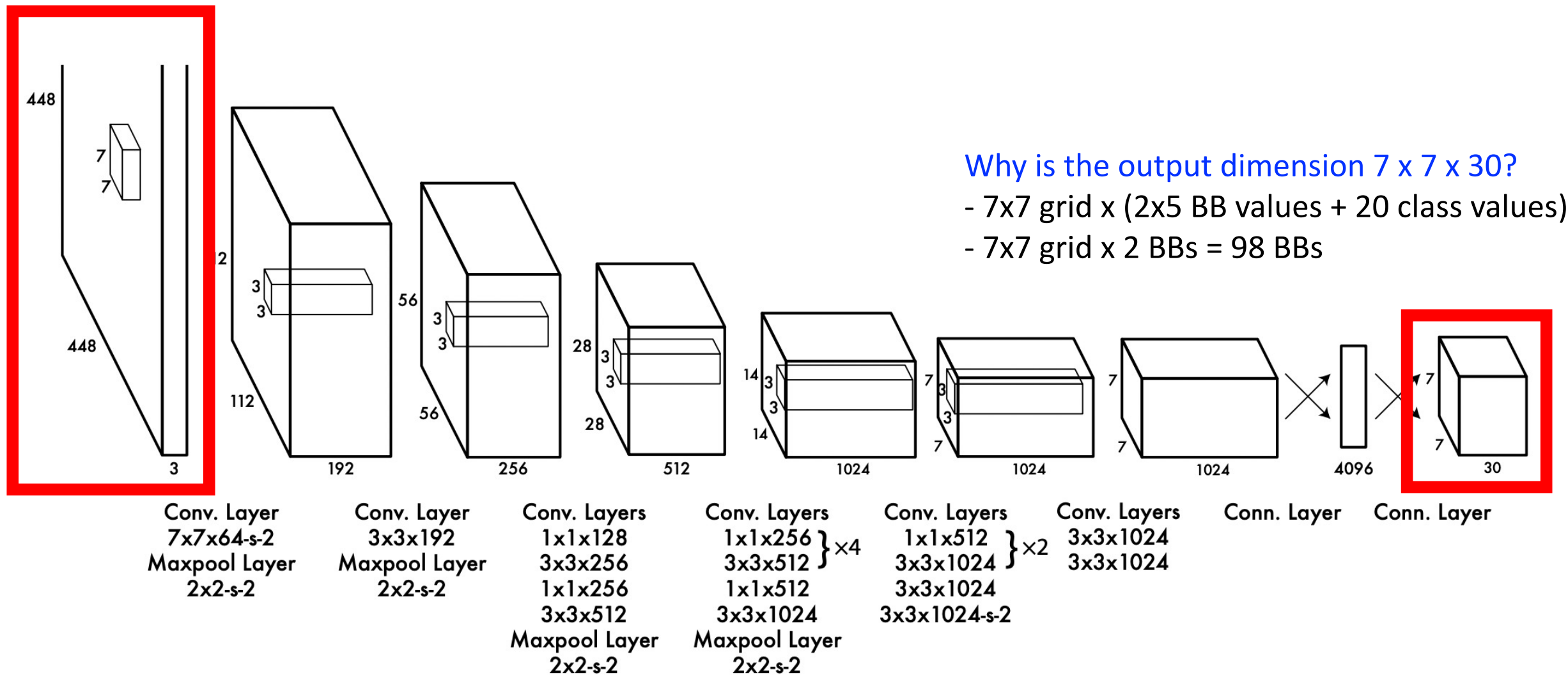


1. What should p_c equal if no object is present?
- 0
2. What should p_c equal if an object is present?
- IoU between predicted and ground truth boxes
3. Although multiple BBs are predicted per grid cell, only the BB with the highest IoU to the GT is used. So why have multiple BBs per grid cell?
- Encourage each BB predictor to specialize to different BB properties (e.g., sizes, aspect ratios, object category types)

Architecture

Input: RGB image resized to fixed input size

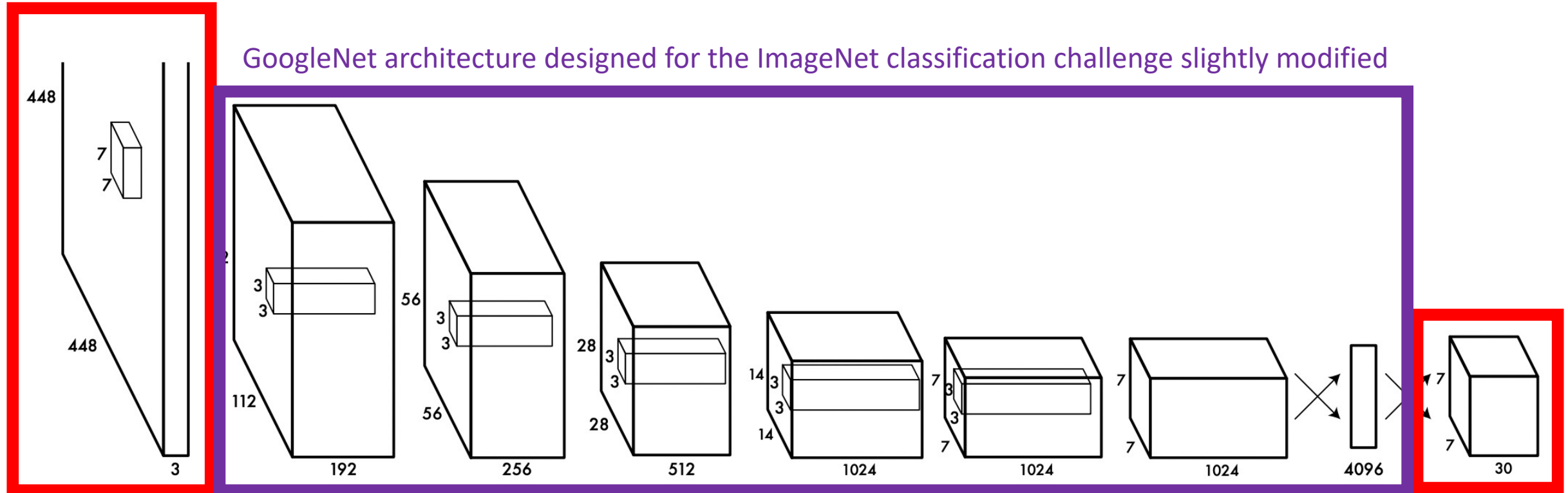
Output: 98 BB per image w/ class probabilities



Architecture

Input: RGB image resized to fixed input size

Output: 98 BB per image w/ class probabilities
(i.e., 7x7 grid x 2 BB per grid cell = 98 BB)



Conv. Layer 7x7x64-s-2	Conv. Layer 3x3x192	Conv. Layers 1x1x128	Conv. Layers 1x1x256	Conv. Layers 1x1x512	Conv. Layers 1x1x512	Conv. Layers 3x3x1024	Conn. Layer	Conn. Layer
Maxpool Layer 2x2-s-2	Maxpool Layer 2x2-s-2	3x3x256	3x3x512	3x3x1024	3x3x1024	3x3x1024		
		1x1x256	1x1x512	3x3x1024-s-2				
		Maxpool Layer 2x2-s-2	Maxpool Layer 2x2-s-2					

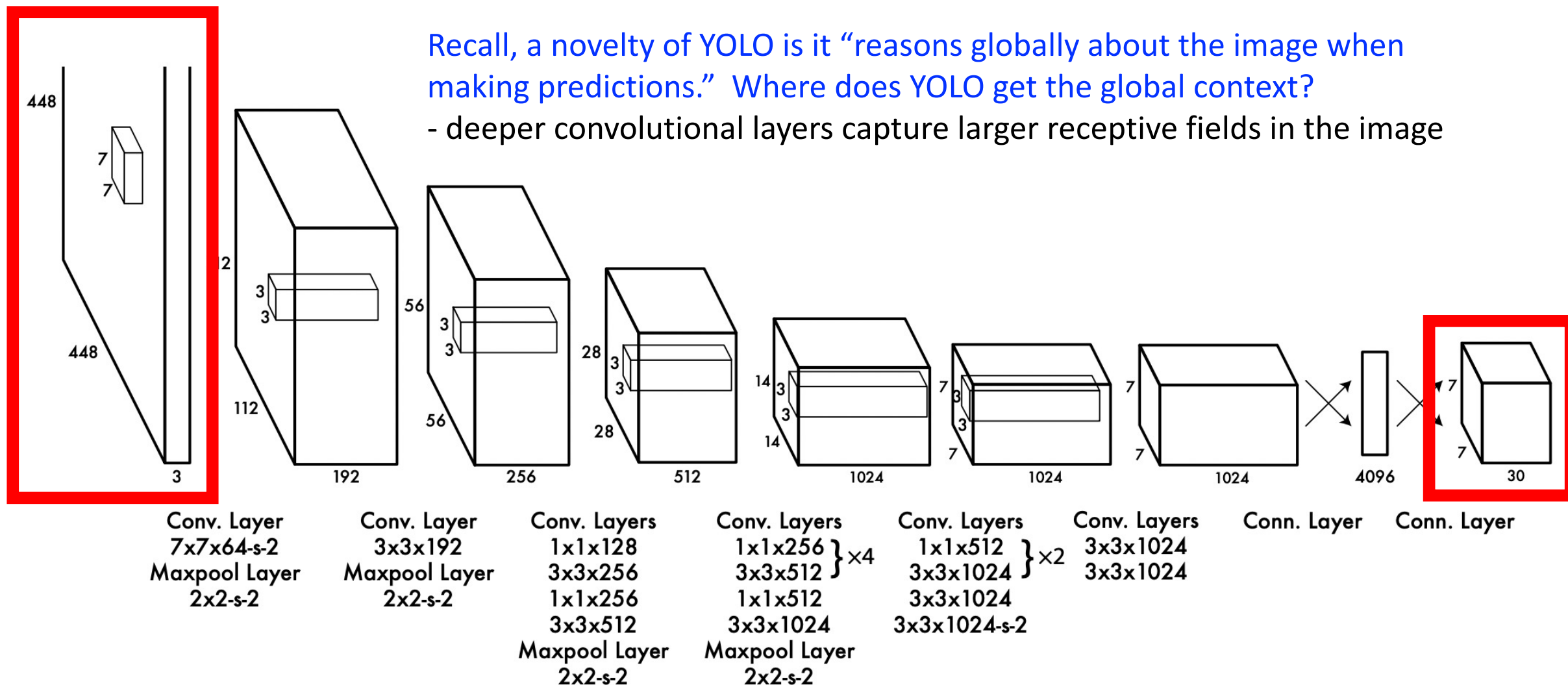
Architecture

Input: RGB image resized to fixed input size

Output: 98 BB per image w/ class probabilities
(i.e., 7x7 grid x 2 BB per grid cell = 98 BB)

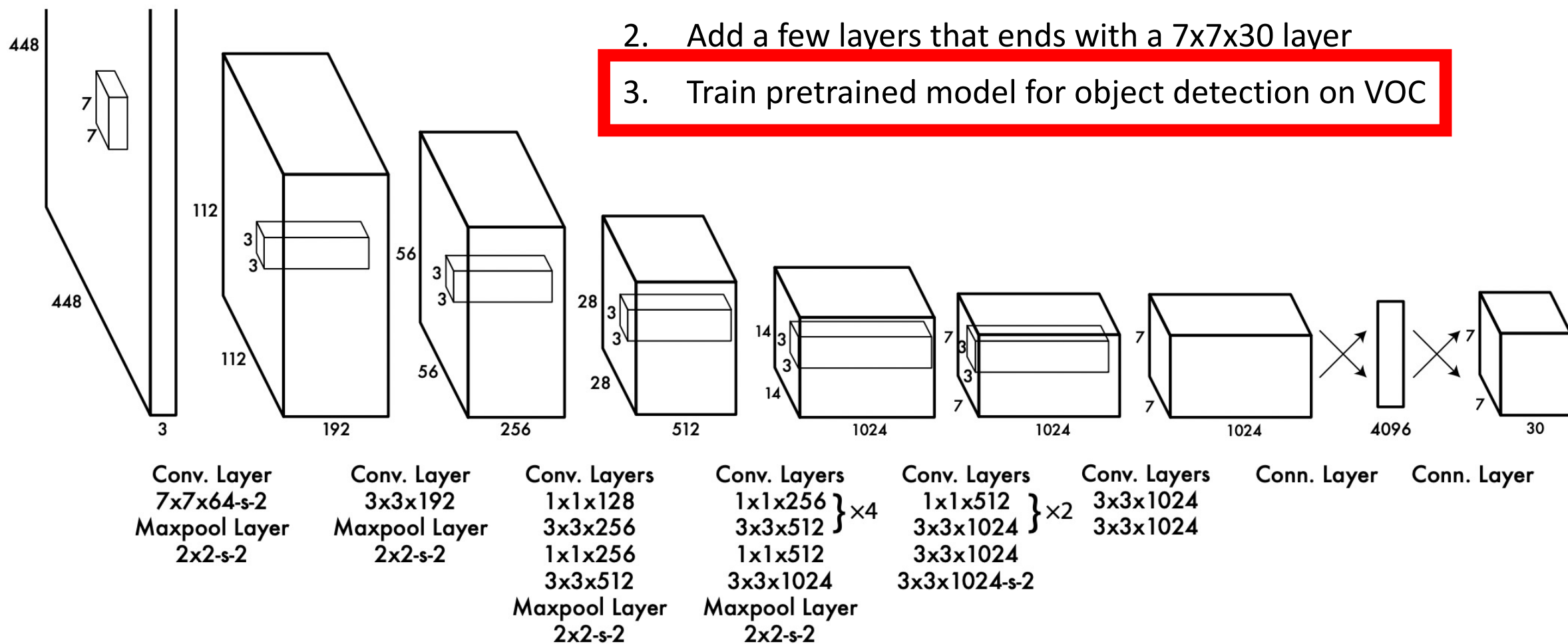
Recall, a novelty of YOLO is it “reasons globally about the image when making predictions.” Where does YOLO get the global context?

- deeper convolutional layers capture larger receptive fields in the image

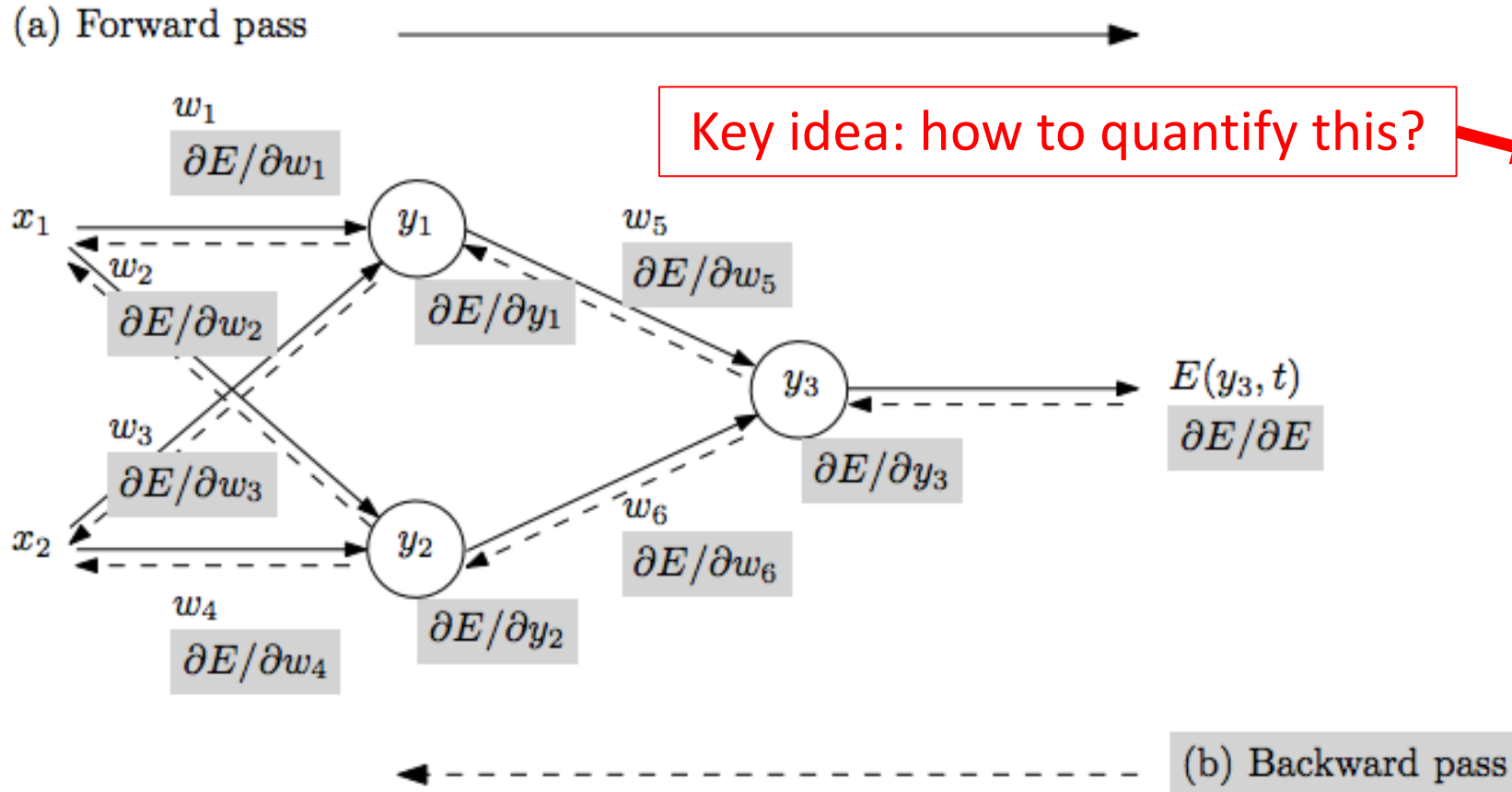


Training

1. Pretrain first 20 convolutional layers for ImageNet classification
2. Add a few layers that ends with a 7x7x30 layer
3. Train pretrained model for object detection on VOC



Training



- Repeat until stopping criterion met:
 1. **Forward pass:** propagate training data through model to make predictions
 2. **Error quantification:** measure dissatisfaction with a model's predictions on training data
 3. **Backward pass:** using predicted output, calculate gradients backward to assign blame to each model parameter; **account for weight sharing by using average of all connections for a parameter**
 4. Update each parameter using calculated gradients

Training: Multi-Part Loss Function

Number of grid cells; what does YOLO use?

Number of bounding box predictors per cell; what does YOLO use?

Penalty only for most confident BB predictor from j predictors when an object is present

Penalizes incorrect classifications (for each grid cell)

Penalty for classification occurs ONLY if an object is actually present in the grid cell

$$\lambda_{\text{coord}} \sum_{i=0}^{S^2} \sum_{j=0}^B \mathbb{1}_{ij}^{\text{obj}} \left[(x_i - \hat{x}_i)^2 + (y_i - \hat{y}_i)^2 \right]$$

Penalizes imperfect positioning of object centers (for each BB)

$$+ \lambda_{\text{coord}} \sum_{i=0}^{S^2} \sum_{j=0}^B \mathbb{1}_{ij}^{\text{obj}} \left[\left(\sqrt{w_i} - \sqrt{\hat{w}_i} \right)^2 + \left(\sqrt{h_i} - \sqrt{\hat{h}_i} \right)^2 \right]$$

Penalizes imperfect widths and heights (for each BB)

$$+ \sum_{i=0}^{S^2} \sum_{j=0}^B \mathbb{1}_{ij}^{\text{obj}} \left(C_i - \hat{C}_i \right)^2$$

Pushes the BB confidence score to match the IoU between the prediction and GT

$$+ \lambda_{\text{noobj}} \sum_{i=0}^{S^2} \sum_{j=0}^B \mathbb{1}_{ij}^{\text{noobj}} \left(C_i - \hat{C}_i \right)^2$$

Punishes the network for predicting an object is present when no object is in the grid cell (i.e., pushes confidence to 0)

$$+ \sum_{i=0}^{S^2} \mathbb{1}_i^{\text{obj}} \sum_{c \in \text{classes}} \left(p_i(c) - \hat{p}_i(c) \right)^2$$

Training: Loss Function

$$\begin{aligned} & \lambda_{\text{coord}} \sum_{i=0}^{S^2} \sum_{j=0}^B \mathbb{1}_{ij}^{\text{obj}} \left[(x_i - \hat{x}_i)^2 + (y_i - \hat{y}_i)^2 \right] \\ & + \lambda_{\text{coord}} \sum_{i=0}^{S^2} \sum_{j=0}^B \mathbb{1}_{ij}^{\text{obj}} \left[\left(\sqrt{w_i} - \sqrt{\hat{w}_i} \right)^2 + \left(\sqrt{h_i} - \sqrt{\hat{h}_i} \right)^2 \right] \\ & + \sum_{i=0}^{S^2} \sum_{j=0}^B \mathbb{1}_{ij}^{\text{obj}} (C_i - \hat{C}_i)^2 \\ & + \lambda_{\text{noobj}} \sum_{i=0}^{S^2} \sum_{j=0}^B \mathbb{1}_{ij}^{\text{noobj}} (C_i - \hat{C}_i)^2 \\ & + \sum_{i=0}^{S^2} \mathbb{1}_i^{\text{obj}} \sum_{c \in \text{classes}} (p_i(c) - \hat{p}_i(c))^2 \end{aligned}$$

Set coefficient of 5 to prioritize loss from BB shape predictions over misclassification predictions

Set coefficient of 0.5 to reduce loss from boxes that don't contain objects to mitigate learning to push the confidence to 0 since most cells don't contain objects

Square root used to weight a small deviation in large boxes less than for small boxes

Object Detection: Today's Topics

- Overview of object detection algorithms
- Fast R-CNN
- Faster R-CNN
- YOLO
- Discussion

A dark gray background with a central circular glow. The glow is a gradient from light gray in the center to dark gray at the edges. The text "The End" is centered within this glow. The entire scene is framed by a white film strip border with sprocket holes on the left and right sides.

The End