

# Towards Accurate and Fast Object Detection

Mateusz Pach

# Rapid Object Detection using a Boosted Cascade of Simple Features

Paul Viola

viola@merl.com

Mitsubishi Electric Research Labs

201 Broadway, 8th FL

Cambridge, MA 02139

Michael Jones

mjones@crl.dec.com

Compaq CRL

One Cambridge Center

Cambridge, MA 02142

## Abstract

*This paper describes a machine learning approach for visual object detection which is capable of processing images extremely rapidly and achieving high detection rates. This*

tested at 15 frames per second on a conventional 700 MHz Intel Pentium III. In other face detection systems, auxiliary information, such as image differences in video sequences, or pixel color in color images, have been used to achieve high frame rates. Our system achieves high frame rates

# Motivation

- Need for rapid and robust detection of faces, e.g in low power devices.
- Current methods rely on color and multiple frames to get high accuracy.



Compaq iPaq H3641 pocket PC

# Method: Rectangle features

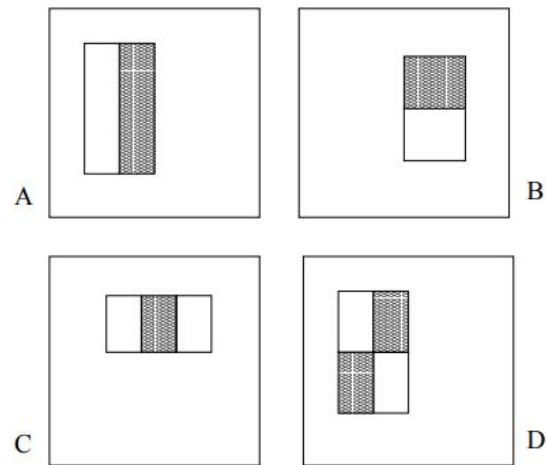
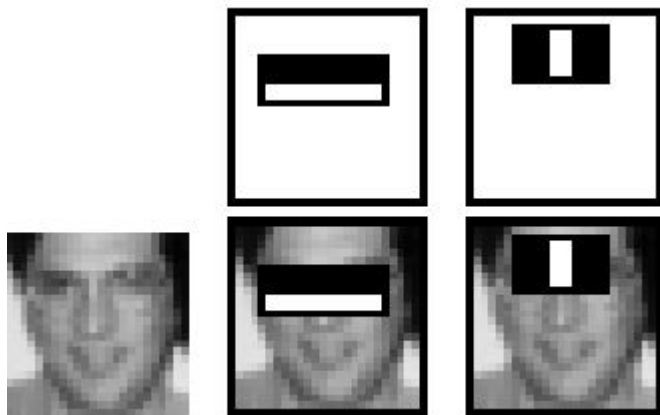


Figure 1: Example rectangle features shown relative to the enclosing detection window. The sum of the pixels which lie within the white rectangles are subtracted from the sum of pixels in the grey rectangles. Two-rectangle features are shown in (A) and (B). Figure (C) shows a three-rectangle feature, and (D) a four-rectangle feature.

# Method: Integral image

aka 2D prefix-sum

$$ii(x, y) = \sum_{x' \leq x, y' \leq y} i(x', y')$$

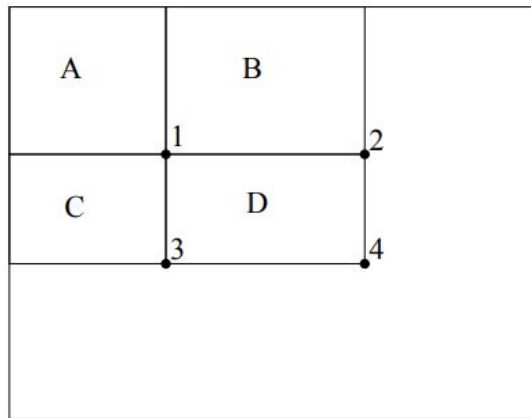
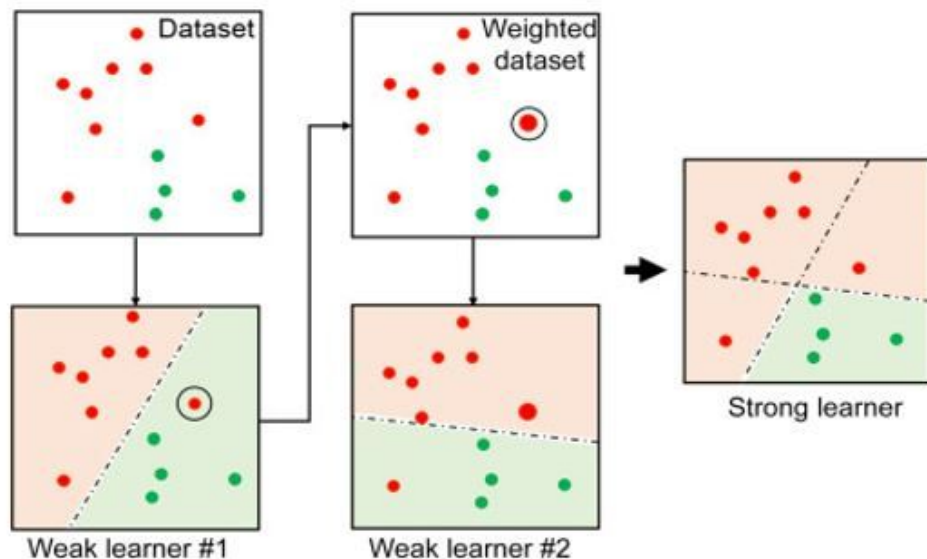


Figure 2: The sum of the pixels within rectangle  $D$  can be computed with four array references. The value of the integral image at location 1 is the sum of the pixels in rectangle  $A$ . The value at location 2 is  $A + B$ , at location 3 is  $A + C$ , and at location 4 is  $A + B + C + D$ . The sum within  $D$  can be computed as  $4 + 1 - (2 + 3)$ .

# Method: AdaBoost



- Given example images  $(x_1, y_1), \dots, (x_n, y_n)$  where  $y_i = 0, 1$  for negative and positive examples respectively.
- Initialize weights  $w_{1,i} = \frac{1}{2m}, \frac{1}{2l}$  for  $y_i = 0, 1$  respectively, where  $m$  and  $l$  are the number of negatives and positives respectively.
- For  $t = 1, \dots, T$ :

1. Normalize the weights,

$$w_{t,i} \leftarrow \frac{w_{t,i}}{\sum_{j=1}^n w_{t,j}}$$

so that  $w_t$  is a probability distribution.

2. For each feature,  $j$ , train a classifier  $h_j$  which is restricted to using a single feature. The error is evaluated with respect to  $w_t$ ,  $e_j = \sum_i w_i |h_j(x_i) - y_i|$ .
3. Choose the classifier,  $h_t$ , with the lowest error  $e_t$ .
4. Update the weights:

$$w_{t+1,i} = w_{t,i} \beta_t^{1-e_i}$$

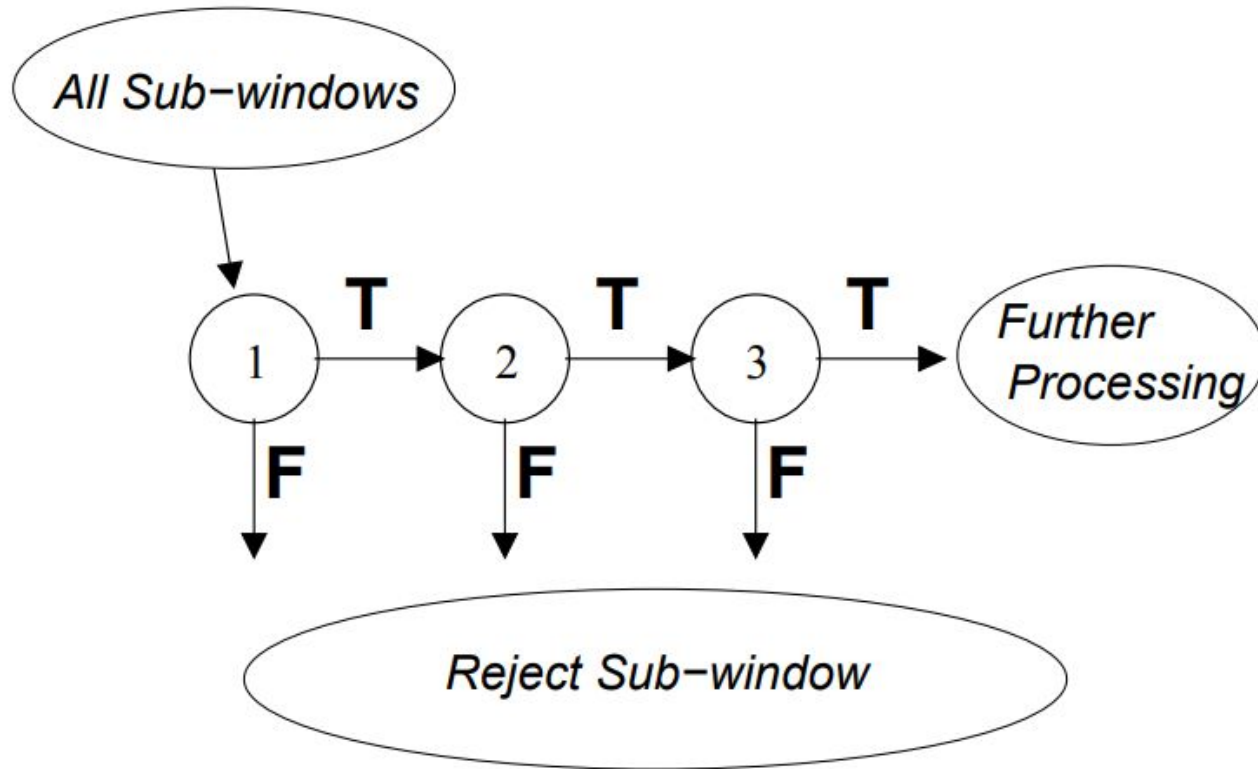
where  $e_i = 0$  if example  $x_i$  is classified correctly,  $e_i = 1$  otherwise, and  $\beta_t = \frac{e_t}{1-e_t}$ .

- The final strong classifier is:

$$h(x) = \begin{cases} 1 & \sum_{t=1}^T \alpha_t h_t(x) \geq \frac{1}{2} \sum_{t=1}^T \alpha_t \\ 0 & \text{otherwise} \end{cases}$$

where  $\alpha_t = \log \frac{1}{\beta_t}$

## Method: Attentional cascade

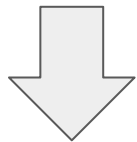
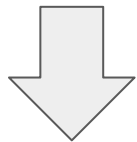
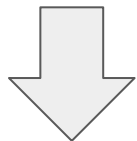


# Method

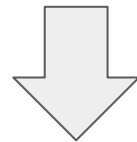
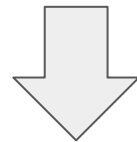
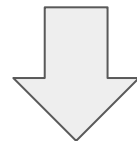
Sequentially eliminate sub-windows with strong classifiers.

Each strong classifier has bounded number of false positives and decrease in detection.

Stages  
Strong classifiers



AdaBoost  
Weak classifiers



Find sufficiently many weak classifiers by boosting.

Each weak classifier uses a single feature and has simple form

$$h_j(x) = \begin{cases} 1 & \text{if } p_j f_j(x) < p_j \theta_j \\ 0 & \text{otherwise} \end{cases}$$

38 stages with over 6000 features, average of 10 feature evaluations per subwindow



# Databases

- Training
  - web crawl
  - 9832 faces from 4916 hand labeled faces
  - 10000 non-faces from 9544 non-face images
  - 24 by 24 pixels
- Evaluation
  - MIT-CMU



Figure 5: Example of frontal upright face images used for training.

# Results

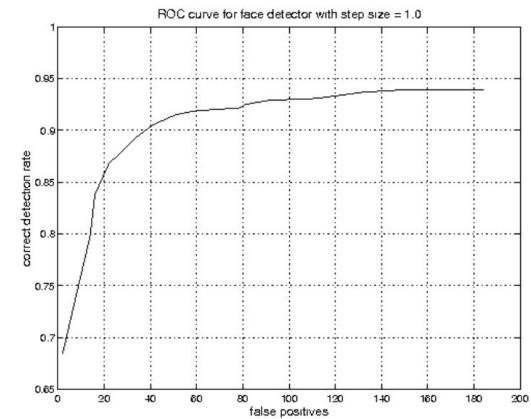
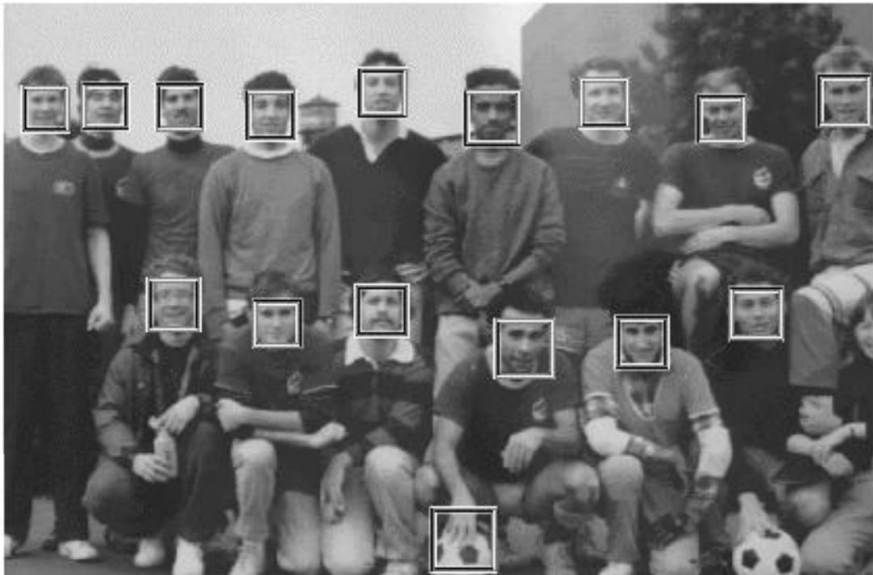


Figure 6: ROC curve for our face detector on the MIT+CMU test set. The detector was run using a step size of 1.0 and starting scale of 1.0 (75,081,800 sub-windows scanned).



Detector	False detections						
	10	31	50	65	78	95	167
Viola-Jones	76.1%	88.4%	91.4%	92.0%	92.1%	92.9%	93.9%
Viola-Jones (voting)	81.1%	89.7%	92.1%	93.1%	93.1%	93.2 %	93.7%
Rowley-Baluja-Kanade	83.2%	86.0%	-	-	-	89.2%	90.1%
Schneiderman-Kanade	-	-	-	94.4%	-	-	-
Roth-Yang-Ahuja	-	-	-	-	(94.8%)	-	-

Time

15x

600x

Table 2: Detection rates for various numbers of false positives on the MIT+CMU test set containing 130 images and 507 faces.

# Conclusions

- The approach achieves high detection accuracy with minimal computation time, making it faster than previous approaches.
- The algorithms, representations, and insights presented have broader applications in computer vision and image processing.
- The experiments conducted on a difficult face detection dataset demonstrate the conclusions drawn are unlikely to be experimental artifacts.
  
- The method may be prone to racial biases.

# Focal Loss for Dense Object Detection

Tsung-Yi Lin Priya Goyal Ross Girshick Kaiming He Piotr Dollár

Facebook AI Research (FAIR)

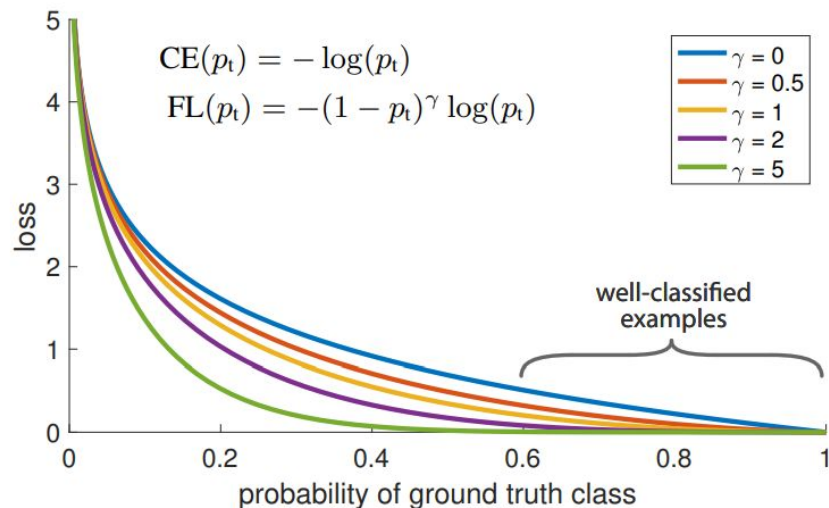


Figure 1. We propose a novel loss we term the *Focal Loss* that adds a factor  $(1 - p_t)^\gamma$  to the standard cross entropy criterion.

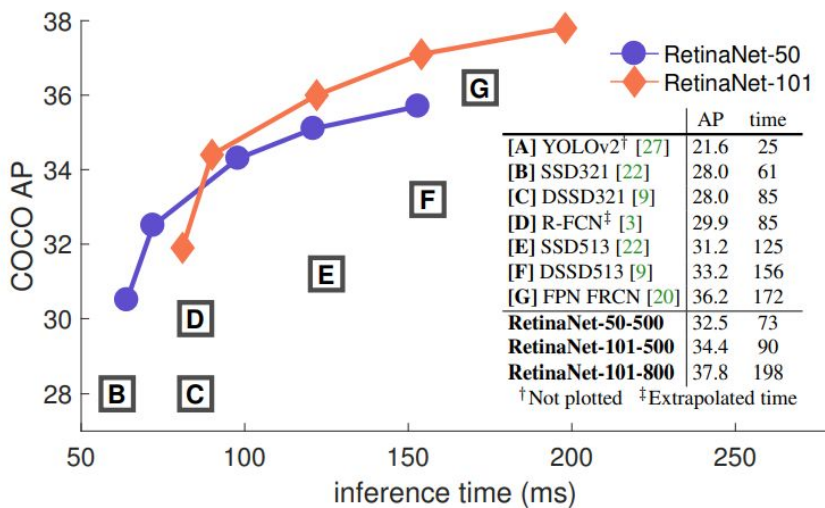
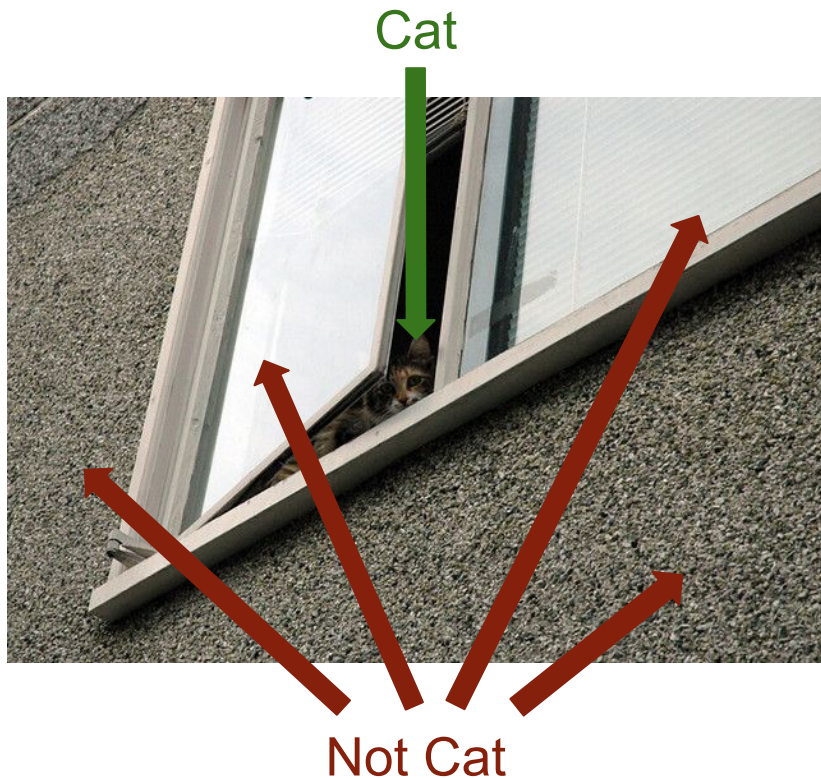


Figure 2. Speed (ms) versus accuracy (AP) on COCO test-dev. Enabled by the focal loss, our simple one-stage *RetinaNet* detec-

# Motivation

- Current SOTA detectors require two stages (propose and classify).
- Single stage detectors perform worse.
- The reason may be the foreground-background class imbalance.
- Let's build a new single stage detector outperforming SOTA by introducing new family of loss functions.

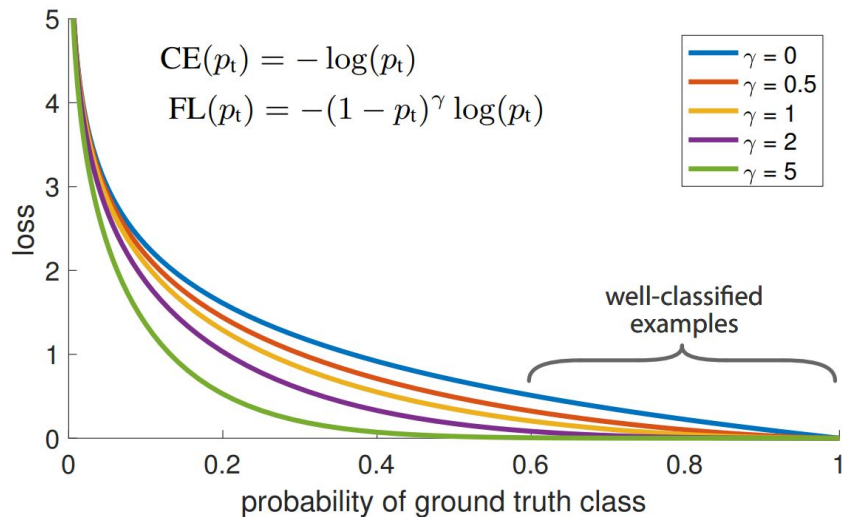


## Method: Focal Loss

$$\text{CE}(p, y) = \begin{cases} -\log(p) & \text{if } y = 1 \\ -\log(1 - p) & \text{otherwise.} \end{cases}$$

$$p_t = \begin{cases} p & \text{if } y = 1 \\ 1 - p & \text{otherwise,} \end{cases}$$

$$\text{CE}(p, y) = \text{CE}(p_t) = -\log(p_t).$$



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



# Method: RetinaNet

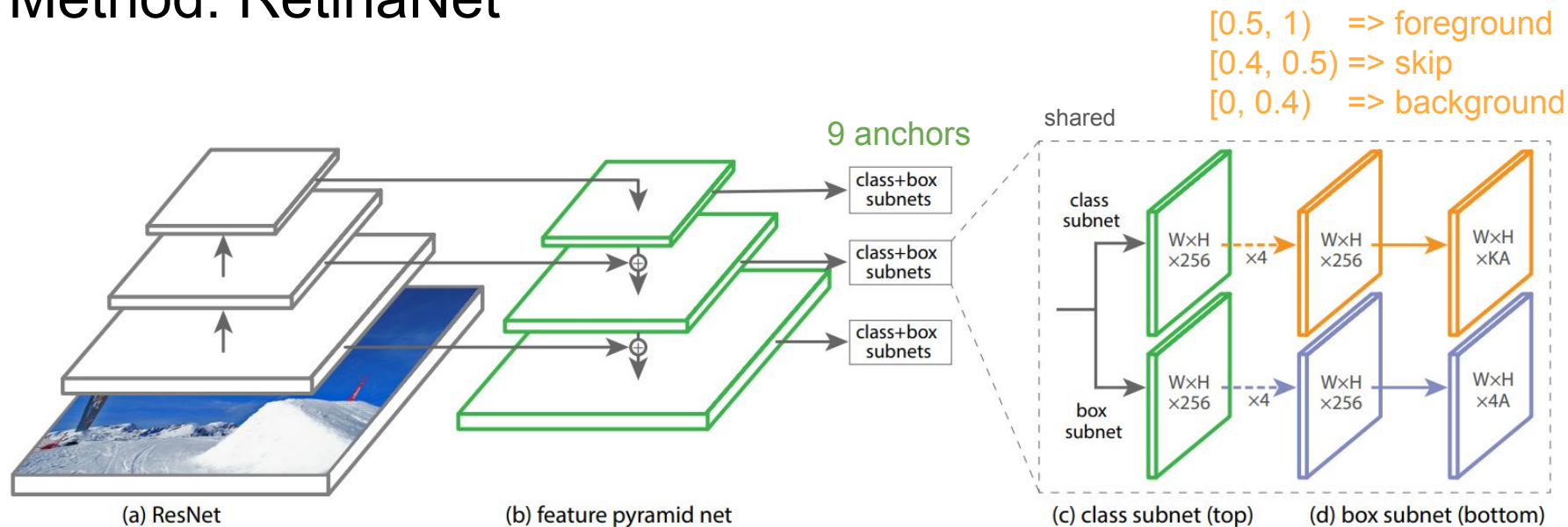


Figure 3. The one-stage **RetinaNet** network architecture uses a Feature Pyramid Network (FPN) [20] backbone on top of a feedforward ResNet architecture [16] (a) to generate a rich, multi-scale convolutional feature pyramid (b). To this backbone RetinaNet attaches two subnetworks, one for classifying anchor boxes (c) and one for regressing from anchor boxes to ground-truth object boxes (d). The network design is intentionally simple, which enables this work to focus on a novel focal loss function that eliminates the accuracy gap between our one-stage detector and state-of-the-art two-stage detectors like Faster R-CNN with FPN [20] while running at faster speeds.

# Databases

## COCO

- Training
  - trainval35k
- Evaluation
  - minival
  - test-dev (hidden from public)





# Results

	backbone	AP	AP <sub>50</sub>	AP <sub>75</sub>	AP <sub>S</sub>	AP <sub>M</sub>	AP <sub>L</sub>
<i>Two-stage methods</i>							
Faster R-CNN+++ [16]	ResNet-101-C4	34.9	55.7	37.4	15.6	38.7	50.9
Faster R-CNN w FPN [20]	ResNet-101-FPN	36.2	59.1	39.0	18.2	39.0	48.2
Faster R-CNN by G-RMI [17]	Inception-ResNet-v2 [34]	34.7	55.5	36.7	13.5	38.1	52.0
Faster R-CNN w TDM [32]	Inception-ResNet-v2-TDM	36.8	57.7	39.2	16.2	39.8	<b>52.1</b>
<i>One-stage methods</i>							
YOLOv2 [27]	DarkNet-19 [27]	21.6	44.0	19.2	5.0	22.4	35.5
SSD513 [22, 9]	ResNet-101-SSD	31.2	50.4	33.3	10.2	34.5	49.8
DSSD513 [9]	ResNet-101-DSSD	33.2	53.3	35.2	13.0	35.4	51.1
<b>RetinaNet</b> (ours)	ResNet-101-FPN	39.1	59.1	42.3	21.8	42.7	50.2
<b>RetinaNet</b> (ours)	ResNeXt-101-FPN	<b>40.8</b>	<b>61.1</b>	<b>44.1</b>	<b>24.1</b>	<b>44.2</b>	51.2

Table 2. **Object detection** *single-model* results (bounding box AP), vs. state-of-the-art on COCO test-dev. We show results for our RetinaNet-101-800 model, trained with scale jitter and for  $1.5\times$  longer than the same model from Table 1e. Our model achieves top results, outperforming both one-stage and two-stage models. For a detailed breakdown of speed versus accuracy see Table 1e and Figure 2.

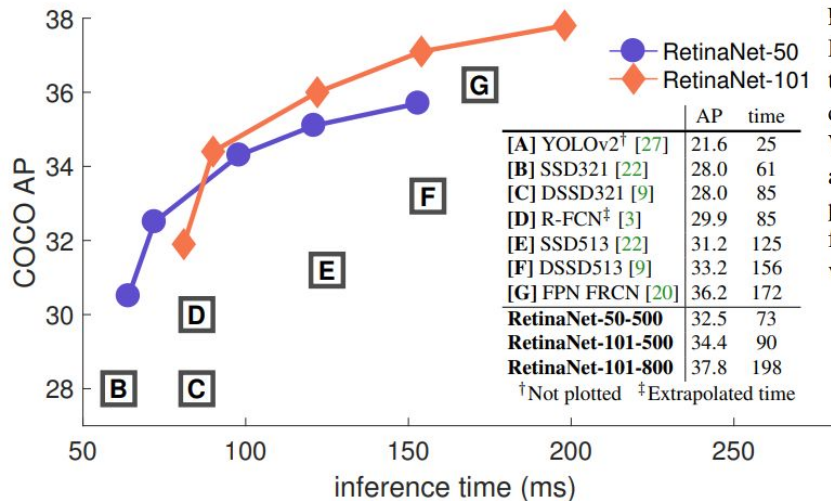


Figure 2. Speed (ms) versus accuracy (AP) on COCO test-dev. Enabled by the focal loss, our simple one-stage *RetinaNet* detector outperforms all previous one-stage and two-stage detectors, including the best reported Faster R-CNN [28] system from [20]. We show variants of RetinaNet with ResNet-50-FPN (blue circles) and ResNet-101-FPN (orange diamonds) at five scales (400-800 pixels). Ignoring the low-accuracy regime ( $AP < 25$ ), RetinaNet forms an upper envelope of all current detectors, and an improved variant (not shown) achieves 40.8 AP. Details are given in §5.

## Average Precision (AP):

AP % AP at IoU=.50:.05:.95 (primary challenge metric)  
 AP<sub>IoU=.50</sub> % AP at IoU=.50 (PASCAL VOC metric)  
 AP<sub>IoU=.75</sub> % AP at IoU=.75 (strict metric)

## AP Across Scales:

AP<sub>small</sub> % AP for small objects: area < 32<sup>2</sup>  
 AP<sub>medium</sub> % AP for medium objects: 32<sup>2</sup> < area < 96<sup>2</sup>  
 AP<sub>large</sub> % AP for large objects: area > 96<sup>2</sup>

# Conclusions

- Class imbalance is the primary obstacle for one-stage object detectors to outperform two-stage methods.
- The focal loss addresses class imbalance by focusing learning on hard negative examples.
- The proposed approach achieves state-of-the-art accuracy and speed, making it a promising solution for overcoming class imbalance in one-stage object detectors.