

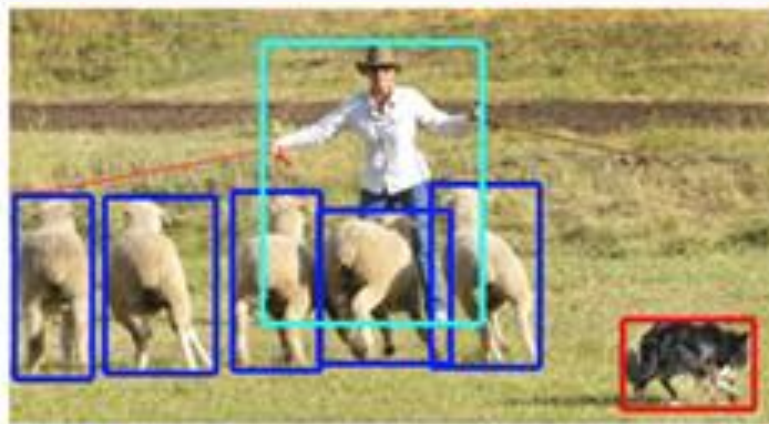
The background of the slide is a light gray gradient. It is decorated with numerous realistic water droplets of various sizes. Some droplets are large and prominent, while others are small and scattered. They are primarily located in the top-left and bottom-right corners, with a few smaller ones in the center and top-right areas. The droplets have highlights and shadows, giving them a three-dimensional appearance.

# **RETINANET: FOCAL LOSS FOR OBJECT DETECTION**

ZHIYUAN WEN



(a) classification



(b) detection



(c) segmentation

# OBJECT DETECTION

## R-CNN: *Regions with CNN features*

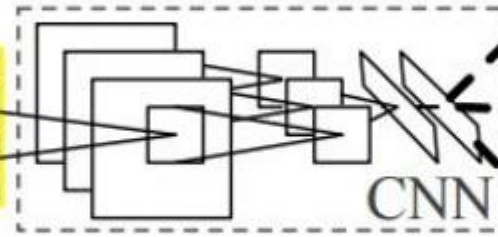


1. Input image



2. Extract region proposals (~2k)

warped region



3. Compute CNN features

aeroplane? no.

⋮

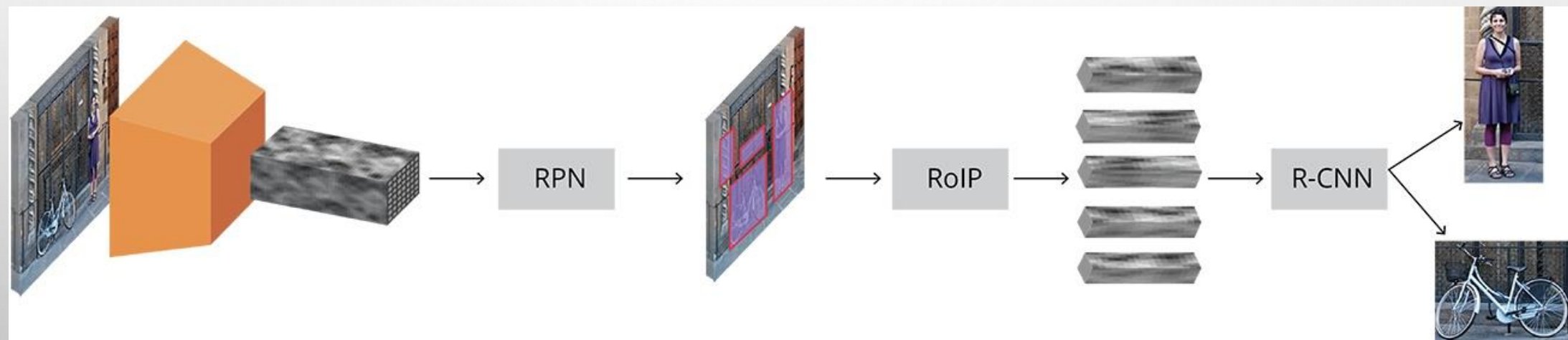
person? yes.

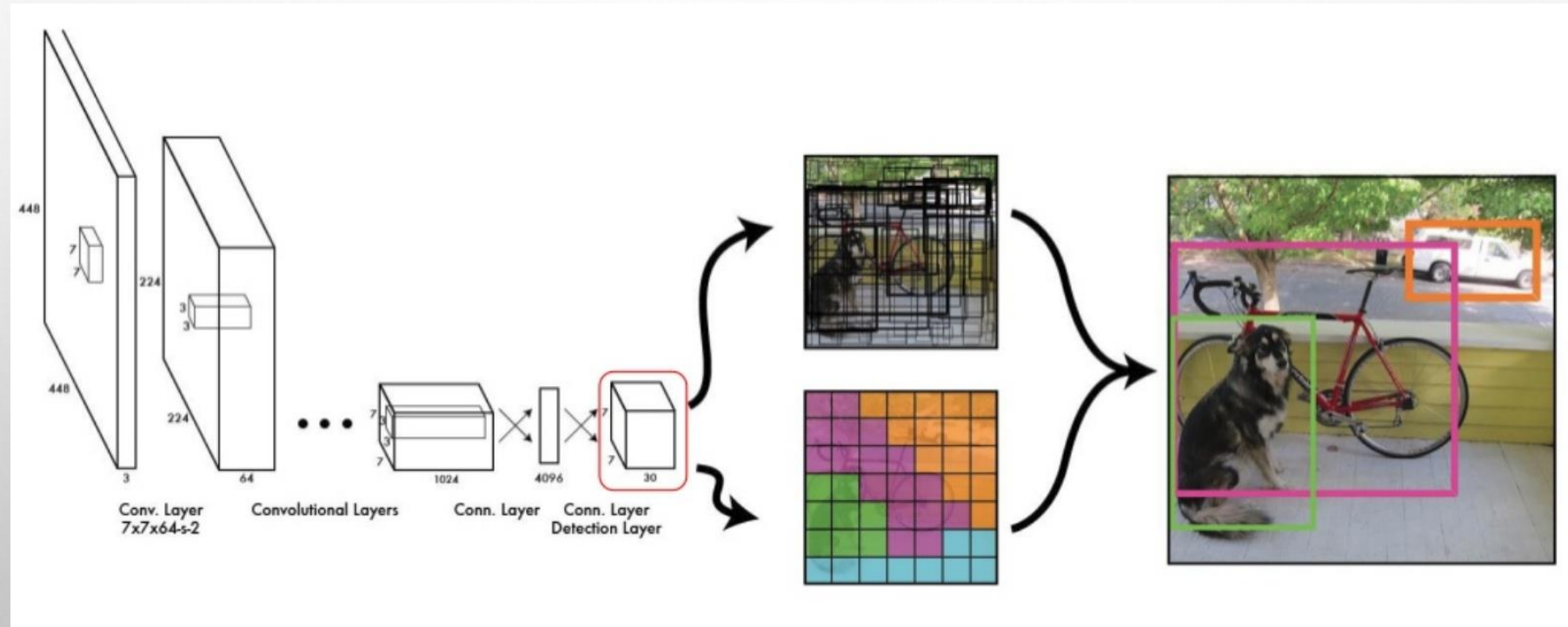
⋮

tvmonitor? no.

4. Classify regions







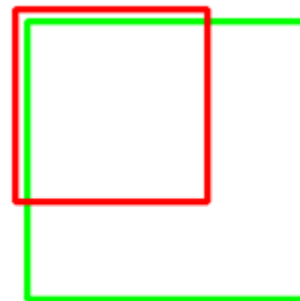
# IOU

$$\text{IoU} = \frac{\text{Area of Overlap}}{\text{Area of Union}}$$



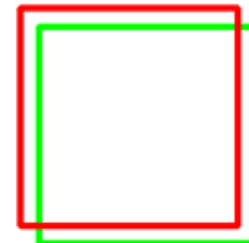
[https://blog.csdn.net/weixin\\_41278720](https://blog.csdn.net/weixin_41278720)

IoU: 0.4034



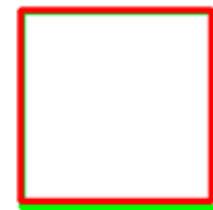
**Poor**

IoU: 0.7330



**Good**

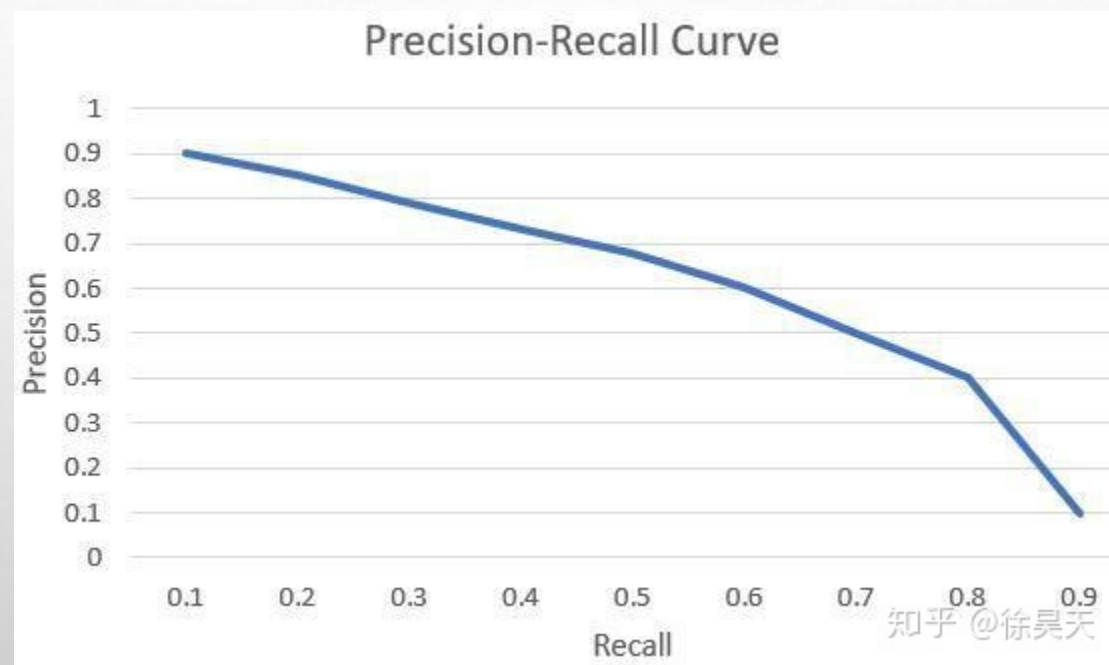
IoU: 0.9264

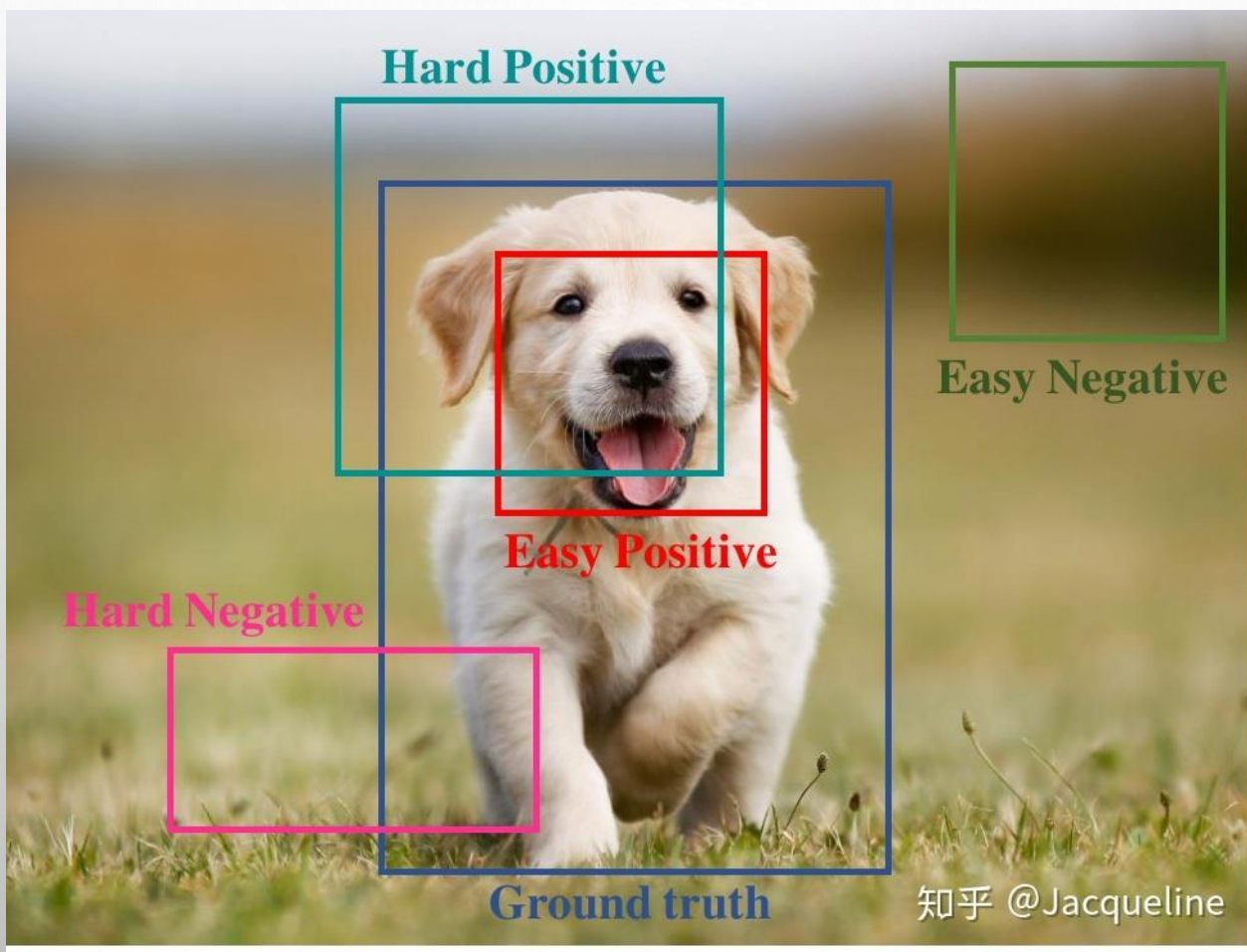


**Excellent**

[https://blog.csdn.net/weixin\\_41278720](https://blog.csdn.net/weixin_41278720)

# MAP





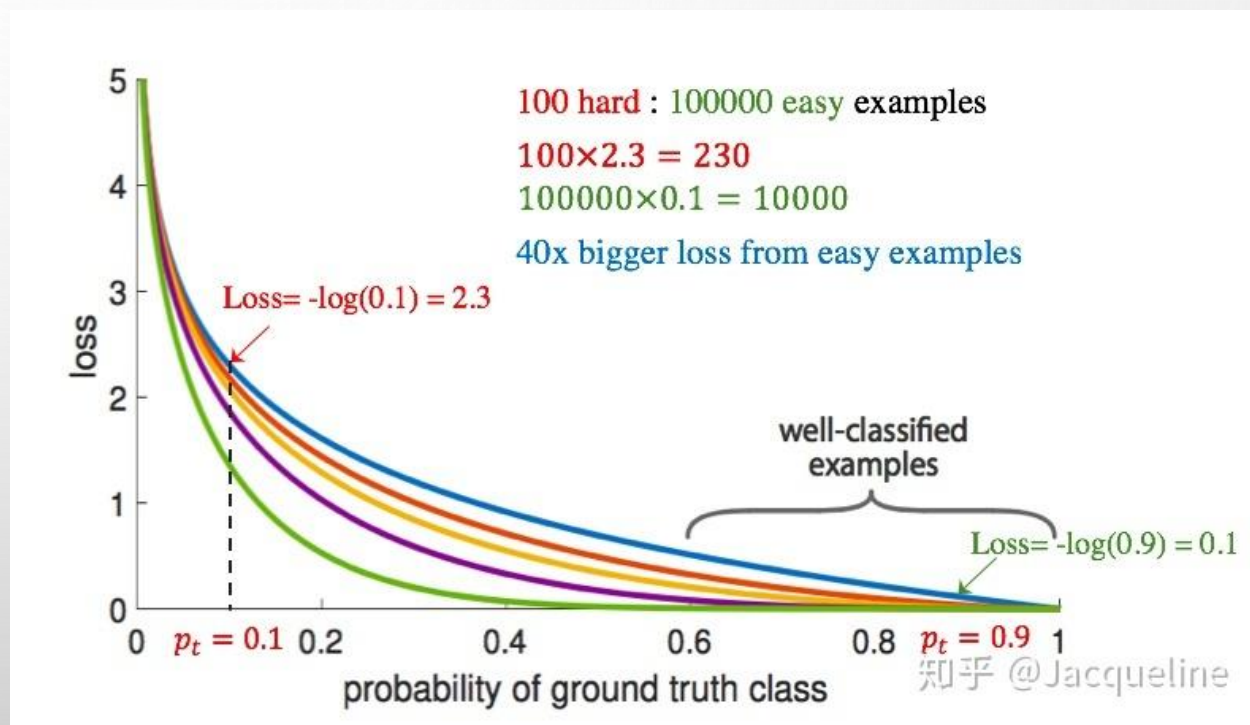


# CROSS ENTROPY (CE) LOSS

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

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

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



# BALANCED CROSS ENTROPY

$$\text{CE}(p_t) = -\alpha_t \log(p_t)$$

$$\alpha_t \begin{cases} \alpha & \text{if } y = 1 \\ 1 - \alpha & \text{otherwise} \end{cases}$$

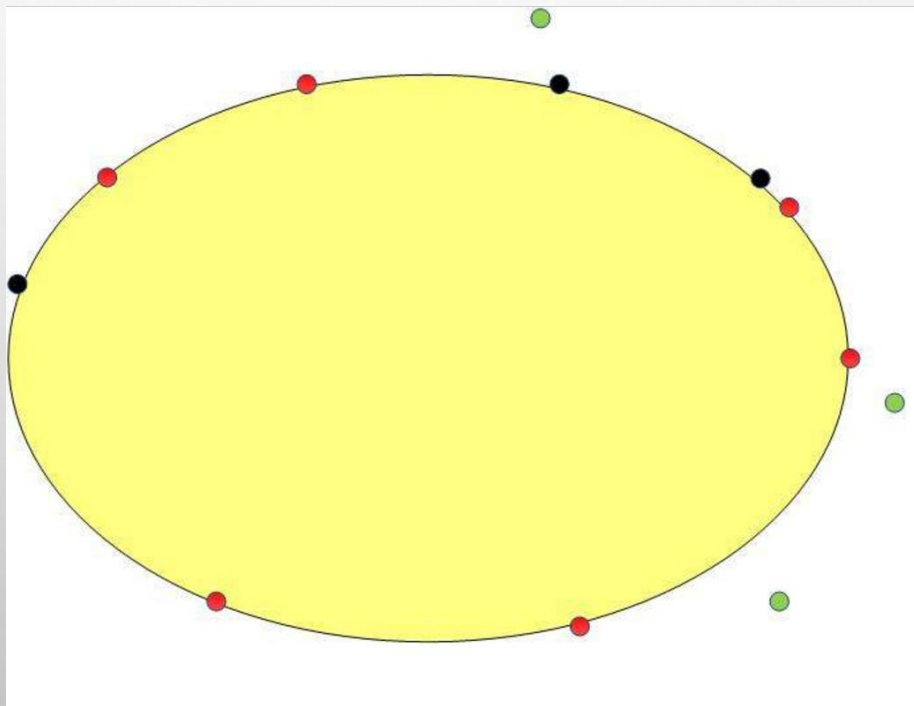
# FOCAL LOSS

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

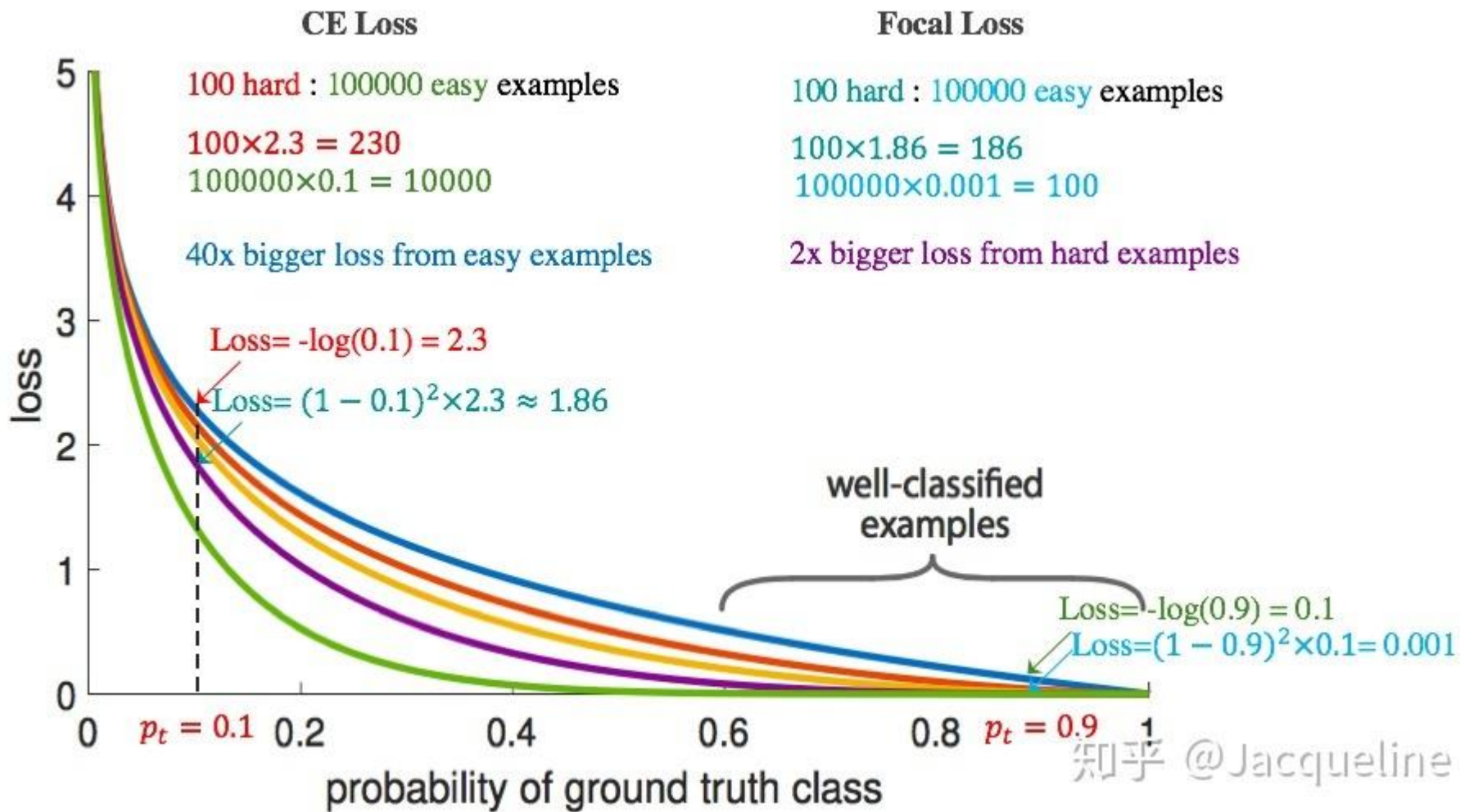
$$\text{FL}(p_t) = -\alpha_t(1 - p_t)^\gamma \log(p_t)$$

# INLIERS

| Loss        | 数量多的类别(如: Background) | 数量少的类别 |
|-------------|-----------------------|--------|
| 被正确分类时的loss | 大幅下降↓                 | 稍微下降↓  |
| 被错误分类时的loss | 适当下降↓                 | 几乎不变   |







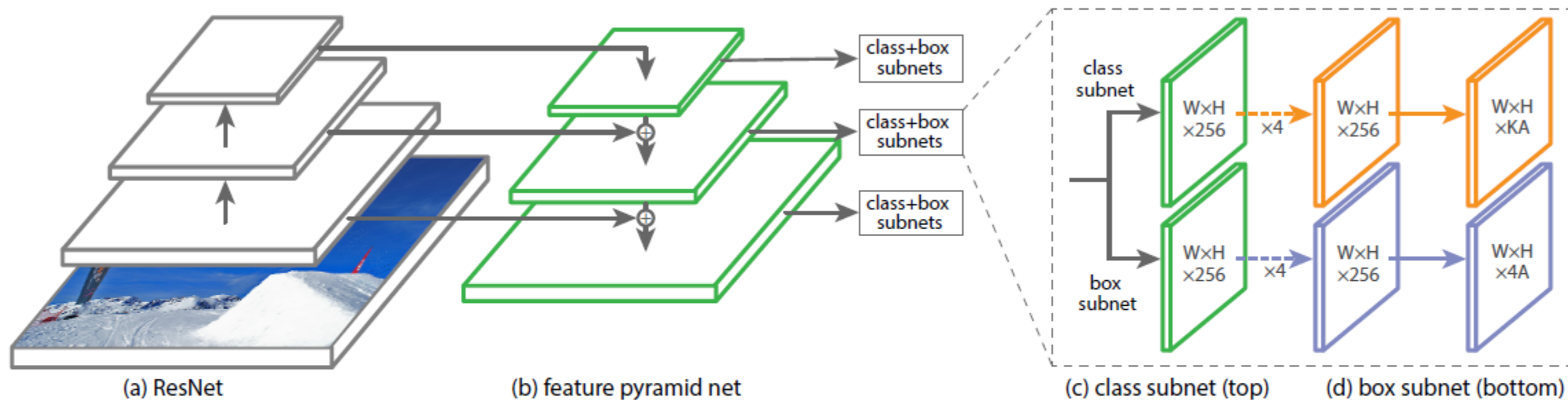


Figure 3. The one-stage **RetinaNet** network architecture uses a Feature Pyramid Network (FPN) [20] backbone on top of a feedforward ResNet architecture [16] (a) to generate a rich, multi-scale convolutional feature pyramid (b). To this backbone RetinaNet attaches two subnetworks, one for classifying 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.

| $\alpha$ | AP   | AP <sub>50</sub> | AP <sub>75</sub> |
|----------|------|------------------|------------------|
| .10      | 0.0  | 0.0              | 0.0              |
| .25      | 10.8 | 16.0             | 11.7             |
| .50      | 30.2 | 46.7             | 32.8             |
| .75      | 31.1 | 49.4             | 33.0             |
| .90      | 30.8 | 49.7             | 32.3             |
| .99      | 28.7 | 47.4             | 29.9             |
| .999     | 25.1 | 41.7             | 26.1             |

(a) Varying  $\alpha$  for CE loss ( $\gamma = 0$ )

| $\gamma$ | $\alpha$ | 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  $\alpha$ )

| #sc | #ar | AP          | AP <sub>50</sub> | AP <sub>75</sub> |
|-----|-----|-------------|------------------|------------------|
| 1   | 1   | 30.3        | 49.0             | 31.8             |
| 2   | 1   | 31.9        | 50.0             | 34.0             |
| 3   | 1   | 31.8        | 49.4             | 33.7             |
| 1   | 3   | 32.4        | 52.3             | 33.9             |
| 2   | 3   | <b>34.2</b> | <b>53.1</b>      | <b>36.5</b>      |
| 3   | 3   | 34.0        | 52.5             | <b>36.5</b>      |
| 4   | 3   | 33.8        | 52.1             | 36.2             |

(c) Varying anchor scales and aspects

| method    | batch size | nms thr | AP          | AP <sub>50</sub> | AP <sub>75</sub> |
|-----------|------------|---------|-------------|------------------|------------------|
| OHEM      | 128        | .7      | 31.1        | 47.2             | 33.2             |
| OHEM      | 256        | .7      | 31.8        | 48.8             | 33.9             |
| OHEM      | 512        | .7      | 30.6        | 47.0             | 32.6             |
| OHEM      | 128        | .5      | 32.8        | 50.3             | 35.1             |
| OHEM      | 256        | .5      | 31.0        | 47.4             | 33.0             |
| OHEM      | 512        | .5      | 27.6        | 42.0             | 29.2             |
| OHEM 1:3  | 128        | .5      | 31.1        | 47.2             | 33.2             |
| OHEM 1:3  | 256        | .5      | 28.3        | 42.4             | 30.3             |
| OHEM 1:3  | 512        | .5      | 24.0        | 35.5             | 25.8             |
| <b>FL</b> | n/a        | n/a     | <b>36.0</b> | <b>54.9</b>      | <b>38.7</b>      |

(d) FL vs. OHEM baselines (with ResNet-101-FPN)

| depth | scale | AP   | AP <sub>50</sub> | AP <sub>75</sub> | AP <sub>S</sub> | AP <sub>M</sub> | AP <sub>L</sub> | time |
|-------|-------|------|------------------|------------------|-----------------|-----------------|-----------------|------|
| 50    | 400   | 30.5 | 47.8             | 32.7             | 11.2            | 33.8            | 46.1            | 64   |
| 50    | 500   | 32.5 | 50.9             | 34.8             | 13.9            | 35.8            | 46.7            | 72   |
| 50    | 600   | 34.3 | 53.2             | 36.9             | 16.2            | 37.4            | 47.4            | 98   |
| 50    | 700   | 35.1 | 54.2             | 37.7             | 18.0            | 39.3            | 46.4            | 121  |
| 50    | 800   | 35.7 | 55.0             | 38.5             | 18.9            | 38.9            | 46.3            | 153  |
| 101   | 400   | 31.9 | 49.5             | 34.1             | 11.6            | 35.8            | 48.5            | 81   |
| 101   | 500   | 34.4 | 53.1             | 36.8             | 14.7            | 38.5            | 49.1            | 90   |
| 101   | 600   | 36.0 | 55.2             | 38.7             | 17.4            | 39.6            | 49.7            | 122  |
| 101   | 700   | 37.1 | 56.6             | 39.8             | 19.1            | 40.6            | 49.4            | 154  |
| 101   | 800   | 37.8 | 57.5             | 40.8             | 20.2            | 41.1            | 49.2            | 198  |

(e) Accuracy/speed trade-off RetinaNet (on test-dev)



# CONCLUSION

- DATA AUGMENTATION / TRUNCATION
- LOSS

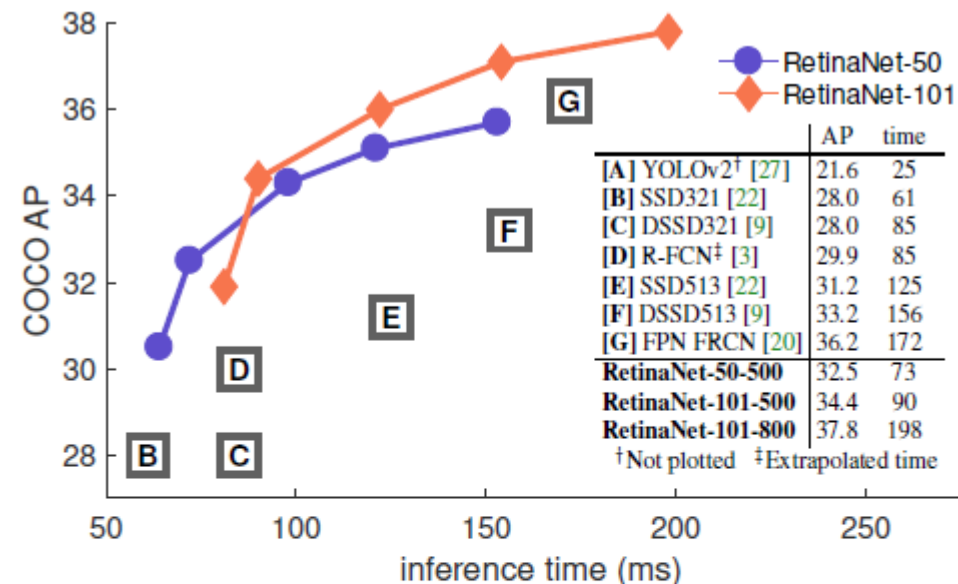


Figure 2. Speed (ms) versus accuracy (AP) on COCO test-dev. Enabled by the focal loss, our simple one-stage *RetinaNet* detector outperforms all previous one-stage and two-stage detectors, including the best reported Faster R-CNN [28] system from [20]. We show variants of RetinaNet with ResNet-50-FPN (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.