

## YOLOv3

- 320x320, 22ms, 28.2 mAP
- It is kind of a "Tech Report"

### Deal

#### BBox Prediction

- ↳ dimension clusters as anchor boxes.
- ↳ predictions for each bbox:  $t_x, t_y, t_w, t_h$

$$\begin{aligned} b_x &= \sigma(t_x) + c_x \\ b_y &= \sigma(t_y) + c_y \\ b_w &= e^{t_w} p_w \\ b_h &= e^{t_h} p_h \end{aligned} \begin{array}{l} \text{cell offset from} \\ \text{top of the img} \\ \text{bbox prior sizes } w, h \\ \text{respectively} \end{array}$$

- sum of squared errors loss.

$$\begin{array}{l} \uparrow \\ \text{↳ } t_x - t_{x'} \end{array} \quad \begin{array}{l} \text{↳ } t_{x'} \text{ calculated by inverting} \\ \text{equations above} \end{array}$$

- Calculate objectness score for each bbox using logistic regression. It is 1 if bbox prior overlaps a ground truth obj by more than any other bbox prior. Even if it's not the best and IoU passes some threshold, it is ignored.

- System assigns only one bbox prior for each gt obj. If a bbox is not assigned, then it incurs no loss for coordinate or class preds.

### Class Prediction

- ↳ each bbox does multilabel classification. NO softmax. USES indep. logistic classifiers.
- Loss: binary cross-entropy

### Across Scales

- ↳ Predicts bboxes at 3 diff. scales.
- ↳ extract features using smtg. like feature pyramid networks
- ↳ Last tensor:  $N \times N \times [3 \cdot (4 + 1 + 80)]$ 
  - ↳ then apply upsampling (p2) for two times

- Use k-means clustering for bbox priors.

### Feature Extractor

	Type	Filters	Size	Output
1x	Convolutional	32	3 × 3	256 × 256
	Convolutional	64	3 × 3 / 2	128 × 128
	Convolutional	32	1 × 1	
	Convolutional	64	3 × 3	128 × 128
2x	Residual			128 × 128
	Convolutional	128	3 × 3 / 2	64 × 64
	Convolutional	64	1 × 1	
	Convolutional	128	3 × 3	64 × 64
8x	Residual			64 × 64
	Convolutional	256	3 × 3 / 2	32 × 32
	Convolutional	128	1 × 1	
	Convolutional	256	3 × 3	32 × 32
8x	Residual			32 × 32
	Convolutional	512	3 × 3 / 2	16 × 16
	Convolutional	256	1 × 1	
	Convolutional	512	3 × 3	16 × 16
4x	Residual			16 × 16
	Convolutional	1024	3 × 3 / 2	8 × 8
	Convolutional	512	1 × 1	
	Convolutional	1024	3 × 3	8 × 8
	Residual			8 × 8
	Avgpool		Global	
	Connected Softmax		1000	

← **Darknet-53**

- Train w/ multi-scale batch normalization
- lots of data augmentation
- ↳ check YOLOv2 for training details

→ COCO's 'weird' meanAP metric?

- logistic activation = sigmoid
- linear (identity) activation  $\Rightarrow f(x) = x$