



北京邮电大学

BEIJING UNIVERSITY OF POSTS AND TELECOMMUNICATIONS

# REAL-TIME OBJECT DETECTION BY A MULTI-FEATURE FULLY CONVOLUTIONAL NETWORK

Yajing Guo, Xiaoqiang Guo, Zhuqing Jiang, Aidong Men, Yun Zhou

Email: [gyj@bupt.edu.cn](mailto:gyj@bupt.edu.cn)



School of Information and Communication Engineering  
Beijing University of Posts and Telecommunications



# Object Detection Overview



## Related Work

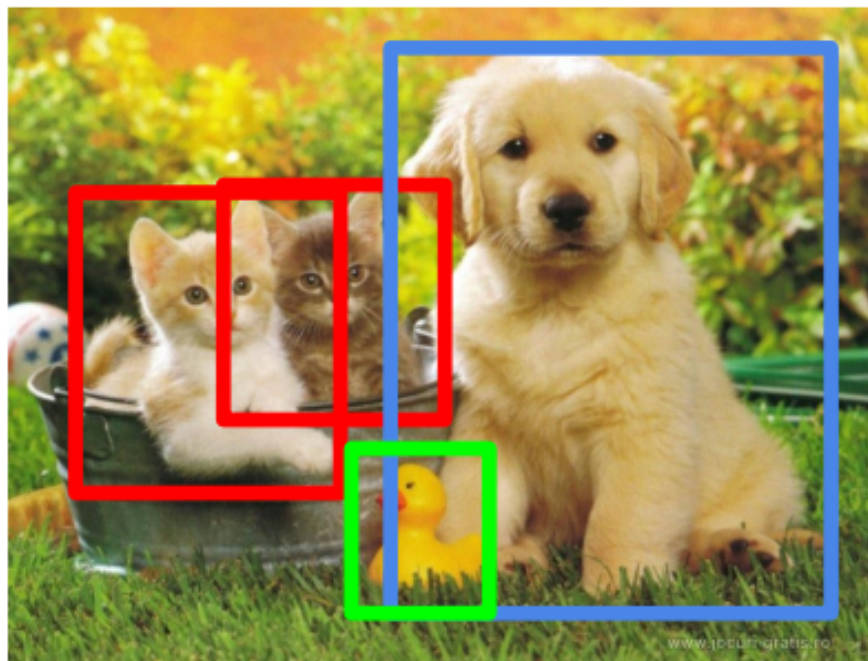


## Proposed Method



## Experiments & Results

- Object detection aims to predict both the category and location of objects in terms of a bounding box.
- Object detection is indispensable for many applications such as video surveillance, autonomous vehicles, augmented reality, and human-computer interaction

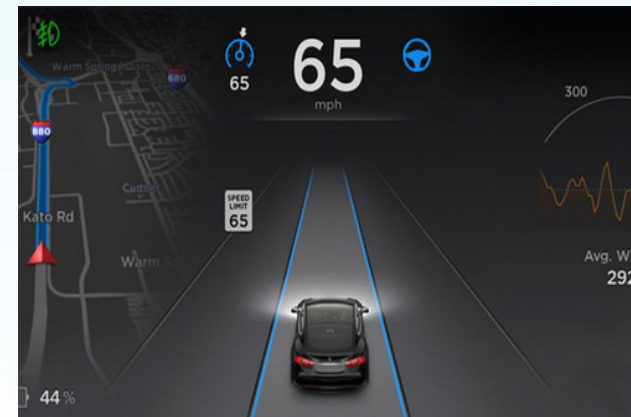


CAT, DOG, DUCK

## 1

# Object Tracking Overview

- Object detection aims to predict both the category and location of objects in terms of a bounding box.
- Object detection is indispensable for many applications such as video surveillance, autonomous vehicles, augmented reality, and human-computer interaction





This paper focus on designing a multi-feature fully convolutional network, aiming to achieve real-time object detection.

## Contributions:

Step1

- A single fully convolutional network treating detection task as a regression problem.

step2

- Multi-feature concatenation fusing shallow and deep information and increasing the detection confidence.

step3

- Anchor boxes mechanism discretizing the space of output box shapes.



**1** Object Detection Overview

**2** Related Work

**3** Proposed Method

**4** Experiments & Results

- Detectors based on region proposals:
  - R-CNN [\*R. Girshick, 2014\*](#)
  - Fast R-CNN [\*R. Girshick, 2015\*](#)
  - Faster R-CNN [\*S. Ren, 2015\*](#)
  - MSCNN [\*S. Gidaris, 2015\*](#)
  
- Detectors free from region proposals:
  - YOLO [\*J. Redmon, 2015\*](#)
  - SSD [\*W. Liu, 2016\*](#)
  - G-CNN [\*M. Najibi, 2016\*](#)



## Object Detection Overview



## Related Work



## Proposed Method



## Experiments & Results

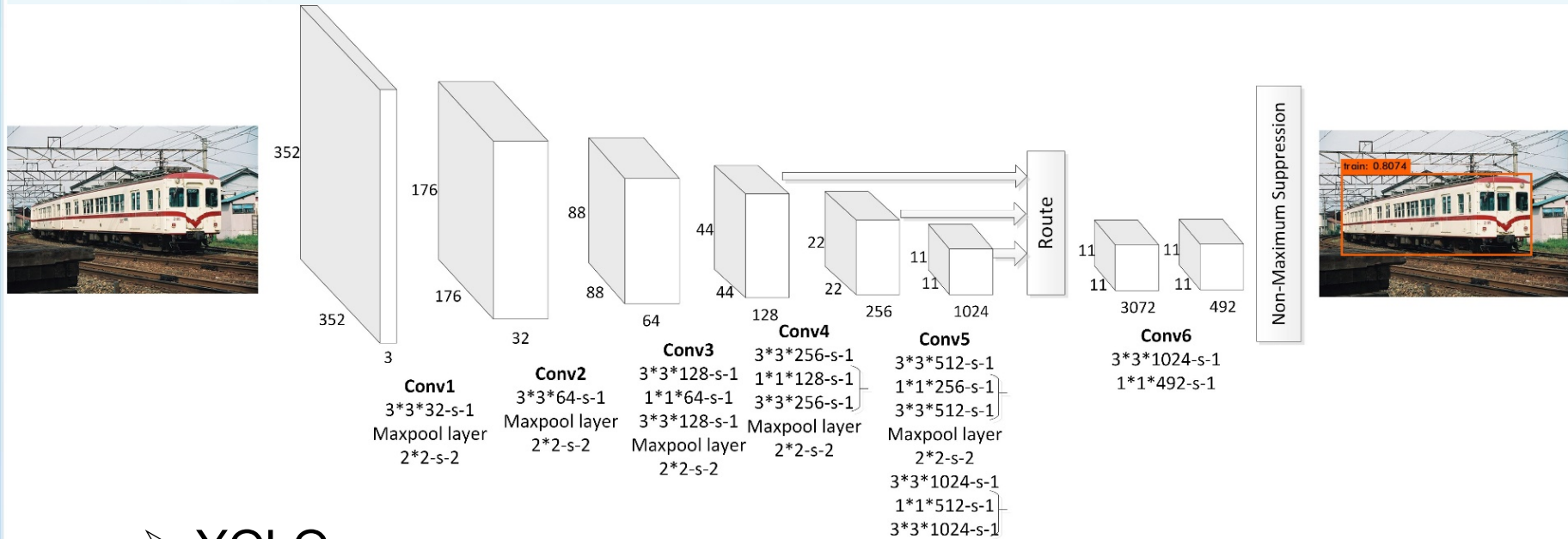


# 3

# Proposed Method

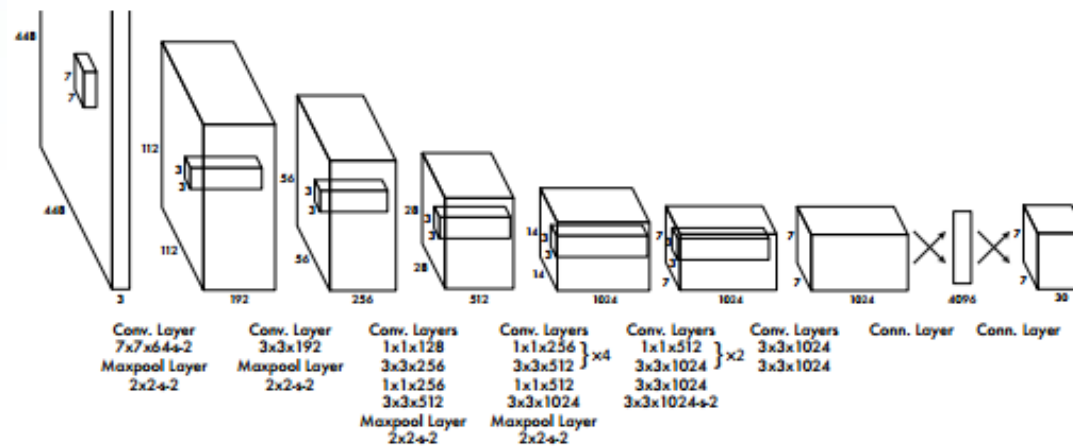
- System Overview

  - The framework of proposed system





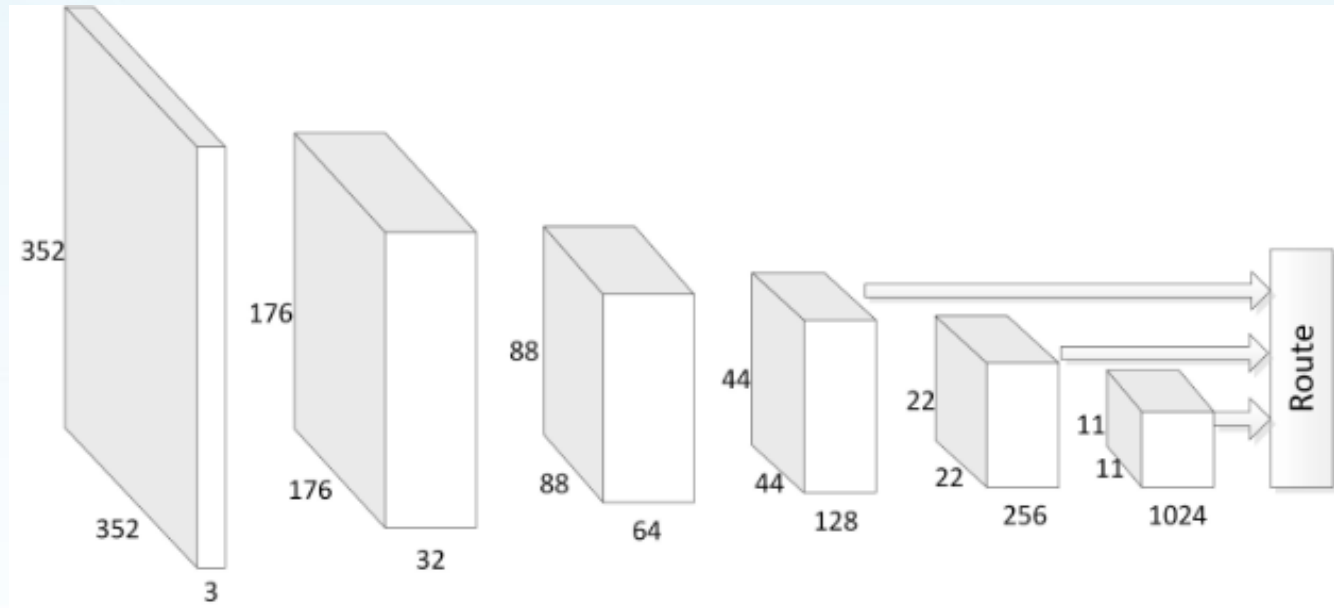
  - YOLO



# 3

## Proposed Method

- Multi-feature concatenation



$$44*44*128 \rightarrow 11*11*2048$$

$$22*22*256 \rightarrow 11*11*1024$$

# 3

## Proposed Method

- Anchor boxes mechanism



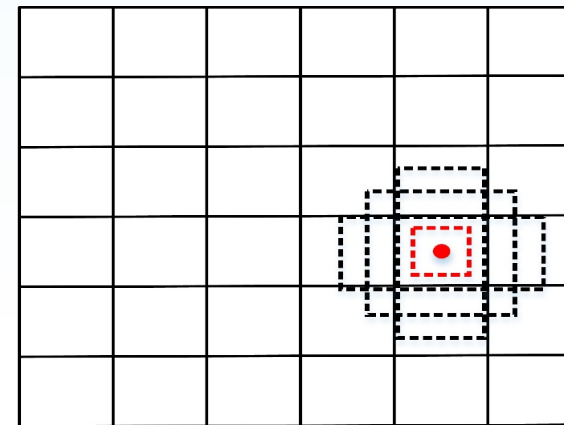
aspect ratio 2:1  
for cats



aspect ratio 1:2  
for cats



aspect ratio 1:1  
for cats



- Anchor boxes mechanism

Each bounding box contains 4+C values

$$b_x = \sigma(t_x) + c_x$$

$$b_y = \sigma(t_y) + c_y$$

$$b_w = p_w e^{t_w}$$

$$b_h = p_h e^{t_h}$$

$$P(\text{object}) * IoU_b^{\text{truth}} = \sigma(t_c)$$





# Object Detection Overview



# Related Work



# Proposed Method



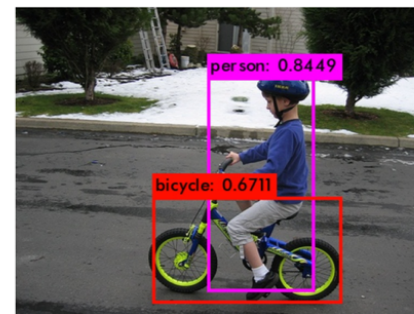
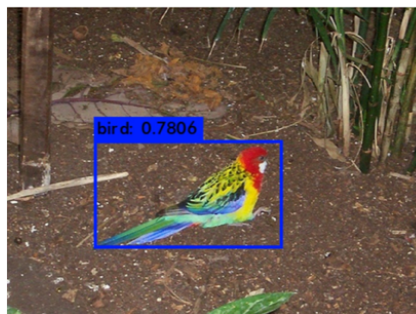
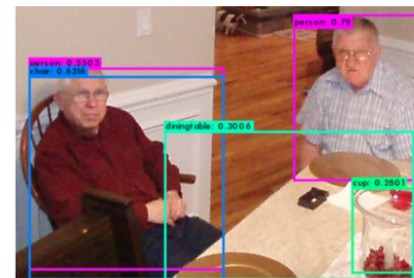
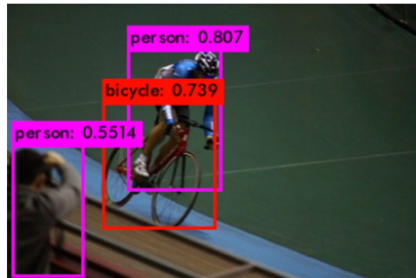
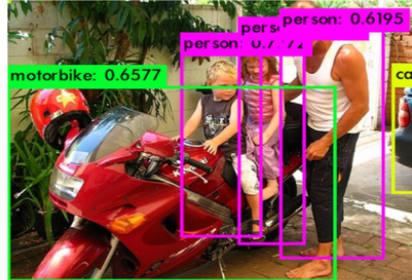
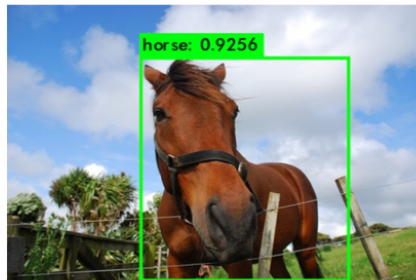
# Experiments & Results

- Evaluating the proposed method on PASCAL VOC dataset with comparisons to other methods.

# 4

# Experiments & Results

- Qualitative Evaluation



- Quantitative Evaluation

- mean average precision (mAP) and frame per second (FPS) are used to quantitatively evaluate the method

**Table 1. PASCAL VOC2012 test detection results.** Each model was trained on PASCAL VOC2012 trainval and VOC2007 trainval and test set. Fast and Faster R-CNN use images with minimum dimension 600, while the image size for YOLO and SSD300 is  $448 \times 448$  and  $300 \times 300$  respectively.

Method	mAP	areo	bike	bird	boat	bottle	bus	car	cat	chair	cow	table	dog	horse	mbike	person	plant	sheep	sofa	train	tv
Fast R-CNN [7]	68.4	82.3	78.4	70.8	52.3	38.7	77.8	71.6	89.3	44.2	73.0	55.0	87.5	80.5	80.8	72.0	35.1	68.3	65.7	80.4	64.2
Faster R-CNN [8]	70.4	84.9	79.8	74.3	53.9	49.8	77.5	75.9	88.5	45.6	77.1	55.3	86.9	81.7	80.9	79.6	40.1	72.6	60.9	81.2	61.5
YOLO [12]	57.9	77.0	67.2	57.7	38.3	22.7	68.3	55.9	81.4	36.2	60.8	48.5	77.2	72.3	71.3	63.5	28.9	52.2	54.8	73.9	50.8
SSD300 [13]	72.4	85.6	80.1	70.5	57.6	46.2	79.4	76.1	89.2	<b>53.0</b>	77.0	<b>60.8</b>	87.0	<b>83.1</b>	82.3	79.4	45.9	75.9	<b>69.5</b>	81.9	67.5
our MFCN	<b>73.2</b>	<b>86.1</b>	<b>82.0</b>	<b>74.4</b>	<b>59.2</b>	<b>50.8</b>	<b>79.6</b>	<b>76.2</b>	<b>90.2</b>	52.1	<b>78.2</b>	58.1	<b>89.0</b>	82.5	<b>83.4</b>	<b>81.1</b>	<b>48.5</b>	<b>77.1</b>	62.4	<b>83.6</b>	<b>68.2</b>

**Table 2.** Detection performance for speed on PASCAL VOC2012 test set. Our MFCN is faster and more accurate than prior detection methods.

Method	mAP (%)	Time (ms)	FPS
Fast R-CNN [7]	68.4	1830	0.5
Faster R-CNN [8]	70.4	142	7
YOLO [12]	57.9	22	45
SSD300 [13]	72.4	21	46
SSD500 [13]	<b>74.9</b>	52	19
our MFCN	73.2	<b>13</b>	<b>75</b>





- **Error Analysis**

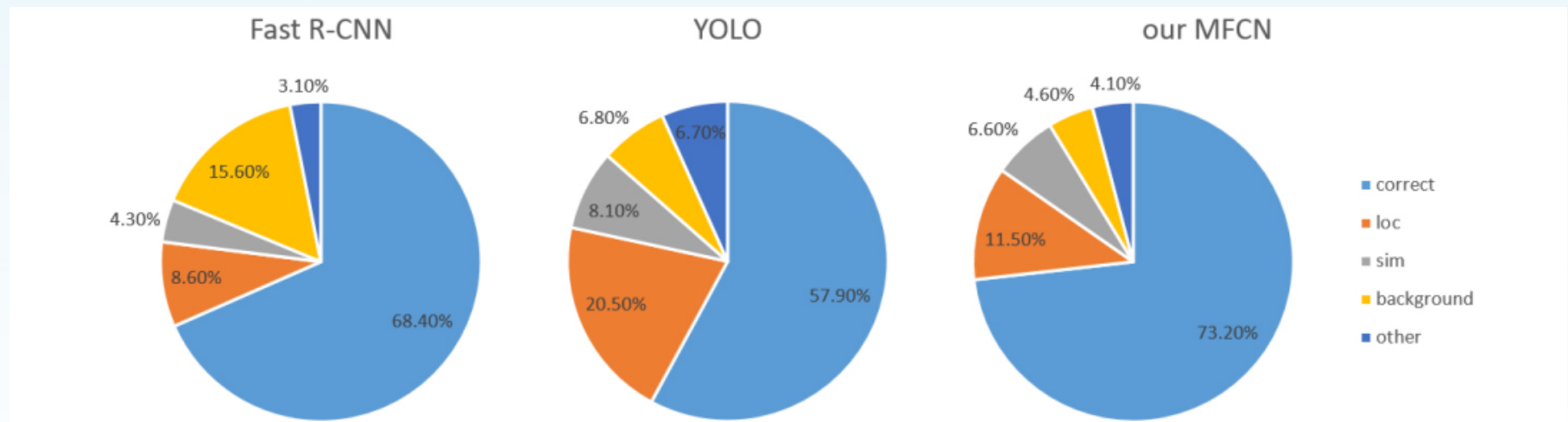


Fig. 3. **Error analysis from PASCAL VOC2012 test.** These charts show the percentage of detections that are correct or false positive due to poor localization (loc), confusion with similar categories (sim), with others, or with background. **Best viewed in color.**

- **Conclusion**
  - framing object detection as a regression problem can simplify detection pipeline and improve the detection speed
  - multi-feature concatenation can efficiently fuse shallow and deep information and increase the detection confidence
  - anchor boxes mechanism helps to discretize the space of output box shapes
- **Future Work**
  - more attention will be paid for tradeoff between accuracy and speed.



北京邮电大学

BEIJING UNIVERSITY OF POSTS AND TELECOMMUNICATIONS

Thank You !

