

### **Feature Pyramid Networks** (FPN) for Object Detection

Tsung-Yi, Piotr Dollar, Ross Girshick, Kiaming He, Bharath Hariharan, and Serge Belongie

IEEE Intl. Conf. on Computer Vision and Pattern Recognition (CVPR), 2017

Speaker: Shih-Shinh Huang

October 29, 2021



# Outline

- Introduction
  - About Object Detection
  - Solution to Multiple Scales
  - Idea of FPN
- Network Architecture
  - Bottom-Up Pathway
  - Lateral Connection
  - Top-Down Pathway

- Region Proposal Network (RPN)
  - About RPN
  - Adopting FPN





- About Object Detection
  - Input:
    - *I*: input image
    - {c<sub>1</sub>, c<sub>2</sub>, ..., c<sub>n</sub>}: object classes to be detected
  - Output:
    - {r<sub>1</sub>, r<sub>2</sub>, ..., r<sub>m</sub>}: bounding boxes
      of *m* detected objects
    - {l<sub>1</sub>, l<sub>2</sub>, ... l<sub>m</sub>}: class labels of all detected objects









- About Object Detection
  - Appearance Variance: object appearances are highly varied.
    - point of view
    - object poses
- Scale Problem: object sizes are significantly different due to perspective phenomena
  - near  $\rightarrow$  large
  - far  $\rightarrow$  small







- Solution to Multiple Scales
  - Featurized Image Pyramid: is a basic approach for addressing scale problem





- Solution to Multiple Scales
  - Single Deep Map: only use <u>a</u> feature map from the last Conv. layer for prediction.





- Solution to Multiple Scales
  - Deep Feature Hierarchy: use the feature maps with enough semantics from the high-up layers.



**drawback:** miss to use high-resolution maps that are important for detecting small objects





- Idea of FPN
  - **propagate** semantics from top layers (low-resolution) to bottom layers (high-resolution)

⇒ all scales have rich semantics





- Overview
  - Input: an image with an arbitrary size
  - Output: a feature pyramid that all scales have rich semantics





- Overview
  - Bottom-Up Pathway: generate multiple-scaled feature maps in a pyramidal shape
  - Lateral Connection: reduce channel numbers of all feature maps to a fixed size for merging.
  - Top-Down Pathway: propagate semantics from topmost feature map to bottommost one.





#### **Feature Pyramid Network (FPN)**





- Bottom-Up Pathway
  - be a backbone Conv. network composed of several stages







- Bottom-Up Pathway
  - perform feed-forward computation to produce a feature hierarchy with a scaling factor of 2





- Bottom-Up Pathway
  - example: <u>residual networks</u> consisting of 5 stages



Kaiming He, et. al., "Deep Residual Learning for Image Classification", CVPR 2016



- Lateral Connection
  - reduce channel numbers of all pyramidal feature maps to a fixed value (i.e. 256)





- Lateral Connection
  - perform channel reduction of all feature maps via convolution layers
    - kernel no.: 256
    - kernel size:  $1 \times 1 \times d$
    - stride size: 1 × 1





- Top-Down Pathway
  - propagate <u>semantics</u> from the topmost feature map to the bottommost one
    - up-sample the upper feature map by a factor of 2 using <u>nearest neighbor strategy</u>
    - merge the up-sampled feature map and lower one by <u>element-wise addition</u>.







 $\mathbf{X}$  x 2 : up-sampling by a factor of 2

**:** element-wise addition



- Top-Down Pathway
  - reduce <u>aliasing effect</u> from up-sampling by a 3 × 3 convolution (padding size = 1)



# Region Proposal Network (RPN)

- About RPN
  - RPN is a network widely used in two-stage object detection.
  - RPN generates object proposals by using <u>dense</u> <u>anchor mechanism</u>.
    - attach 9 anchors (3 aspect ratios × 3 scales) centered at each point of the conv. feature map
    - predict one proposal with 6 parameters ( 4 regression + 2 class probabilities) w.r.t. each anchor

S. Ren, et. al., "Faster R-CNN: Towards Real-Time Object Detection with Region Proposal Network", *NIPS*, 2015









About RPN

Quarter Unit: Region Proposal Network (00:19:09)

Link: https://youtu.be/0tBhRfEzUWs

Web: http://gg.gg/quarter

- predict distant party ithride proposals
- perform convolution operation kernel size:  $1 \times 1 \times d$ 
  - kemelkizaeß:  $30 (= d9 \times 3)$ no. kernels: d



←36→

h

←18→



- Adopting FPN: anchor assignment
  - Basic RPN: attach 9 anchors to each point of a single-scale feature map





- Adopting FPN: anchor assignment
  - FPN+RPN: attach 3 anchors that are different in aspect ratios but with the same size.





- Adopting FPN
  - share the parameters of PRN across all feature pyramid levels.  $12=3\times$



