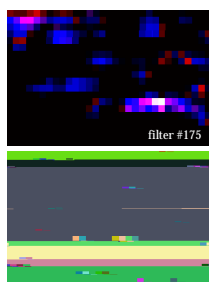


# Spatial Pyramid Pooling in Deep Convolutional Networks for Visual Recognition

Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun





be of any sizes. This not only allows arbitrary aspect ratios, but also allows arbitrary scales. We can resize the input image to any scale (e.g.,  $\min(w; h)=180, 224, \dots$ ) and apply the same deep network. When the input image is at different scales, the network (with

model	conv <sub>1</sub>	conv <sub>2</sub>	conv <sub>3</sub>	conv <sub>4</sub>	conv <sub>5</sub>	conv <sub>6</sub>	conv
-------	-------------------	-------------------	-------------------	-------------------	-------------------	-------------------	------









method	VOC 2007	Caltech101
VQ [15] <sup>y</sup>		

the feature of this window. If the pre-defined scales are dense enough and the window is approximately square, our method is roughly equivalent to resizing the window to  $224 \times 224$  and then extracting features from it. Nevertheless, our method only requires computing the feature maps once (at each scale) from the entire image, regardless of the number of candidate windows.

We also fine-tune our pre-trained network, following [7]. Since our features are pooled from the  $\text{conv}_5$  feature maps from windows of any sizes, for simplicity we only fine-tune the fully-connected layers. In this case, the data layer accepts the fixed-length pooled features after  $\text{conv}_5$ , and the  $\text{fc}_{6,7}$  layers and a new 21-way (one extra negative category)  $\text{fc}_8$  layer follow. The  $\text{fc}_8$

fdistribtino





Figure 6: Example detection results of “SPP-net ftfc<sub>7</sub> bb” on the Pascal VOC 2007 testing set (59.2% mAP). All windows with scores >

## REFERENCES

- [1] Y. LeCun, B. Boser, J. S. Denker, D. Henderson, R. E. Howard,