# Fast R-CNN

Ross Girshick
Microsoft Research
rbg@microsoft.com

## Abstract

*This paper proposes* Fast R-CNN, *a clean and fast framework for object detection. Compared to traditional R-CNN, and its accelerated version SPPnet, Fast R-CNN trains networks using a multi-task loss in a single training stage. The multi-task loss simplifies learning and improves detection accuracy. Unlike SPPnet, all network layers can be updated during fine-tuning. We show that this difference has practical ramifications for very deep networks, such as VGG16, where mAP suffers when only the fully-connected layers are updated. Compared to "slow" R-CNN, Fast R-CNN is 9 faster at training VGG16 for detection, 213 faster at test-time, and achieves a significantly higher mAP*

tures under the proposal's projection onto the feature map

**Back-propagation through RoI pooling layers.** Fast R-CNN mini-batches start from whole images and hence contain all information needed to back-propagate derivatives from the loss function to the image. Back-propagation routes derivatives through the RoI pooling layer, as described below.

During the forward pass, an *input batch* of $N = 2$ images is expanded by the RoI pooling layer into an *output batch* of size $R$

augment the VOC07 trainval set with the VOC12 trainval set, roughly tripling the number of images to 16.5k. The expanded training set improves mAP on VOC07 test from 66.9% to 70.0% (Table 1). When training on this dataset we use 60k mini-batch iterations instead of 40k.

We performed similar experiments for VOC10 and 2012, for which we construct an enlarged dataset of 21.5k images from VOC07 trainval and test union with VOC12 trainval. When training on this dataset we use 100k mini-batch iterations and lower the learning rate by $0.1$ each 40k iterations (instead of each 30k). For VOC10 and 2012, mAP improves from 66.1% to 68.8% and from 65.7% to 68.4%, respecF219(s)-196(and)-197(lo)25(v)15(hly65(1.)4)80(F)15aestR-CNN,t(and)-06(p(reun R-CNN(usst)4458(the)449(softmaxe)449(classi(hrt)449((learst)449(durning219( -11.955 -11.960 Td [(ne-tuining)-290(instead)-290(of)-308( imauct of thischoice2, we(imleimensnd)-211poes-hoceSVMd traicF219(s)-196(and)-197ning(with)-07(harnd)-06(n(e)15g)15a(lo)25(v)15(h)- tio(al)-250R-CNN.t

them with a dense set of "sliding windows" is attractive, since it is essentially free. Yet, these experiments provide the first evidence that sparse proposals do indeed "improve detection quality by reducing spurious false positives" [11