## **Object Detection**

Lecture 5

## **Object Detection**



#### Image Classification (what?)



Object Detection (what + where?)

#### **Detection with ConvNets**

• So far, all about classification

• What about localizing objects within the scene?



Groundtruth: tv or monitor tv or monitor (2) tv or monitor (3) person remote control remote control (2)

### **Two General Approaches**

- 1. Examine very position / scale
  - E.g. Overfeat: Integrated recognition, localization and detection using convolutional networks, Sermanet et al., ICLR 2014

- Use some kind of proposal mechanism to attend to a set of possible regions
  - E.g. Region-CNN [Rich feature hierarchies for accurate object detection and semantic segmentation, Girshick et al., CVPR 2014]

#### **Sliding Window with ConvNet**



#### **Sliding Window with ConvNet**



<u>22</u>4



Input Window

#### **Sliding Window with ConvNet**





Input Window

No need to compute two separate windows --- Just one big input window

#### Multi-Scale Sliding Window ConvNet











#### Multi-Scale Sliding Window ConvNet











#### **OverFeat – Output before NMS**



#### **Overfeat Detection Results**

#### [Sermanet et al. ICLR 2014]



**Top predictions:** trombone (confidence 26.8) oboe (confidence 17.5) oboe (confidence 11.5) ILSVRC2012\_val\_00000614.JPEG



Groundtruth:

person hat with a wide brim hat with a wide brim (2) hat with a wide brim (3) oboe oboe (2) saxophone trombone person (2) person (3) person (4)



**Top predictions:** watercraft (confidence 72.2) watercraft (confidence 2.1)





Groundtruth: watercraft watercraft (2)







. Quem come doce caga azedo

#### **Top predictions:** microwave (confidence 5.6) refrigerator (confidence 2.5)

ILSVRC2012 val 00000519.IPEG

a. Quem come doce caga azedi

Groundtruth: bowl microwave

**Top predictions:** tennis ball (confidence 3.5) banana (confidence 2.4) banana (confidence 2.1) hotdog (confidence 2.0) banana (confidence 1.9)

ILSVRC2012\_val\_00000320.JPEG

### **Two General Approaches**

- 1. Examine very position / scale
  - E.g. Overfeat: Integrated recognition, localization and detection using convolutional networks, Sermanet et al., ICLR 2014

- Use some kind of proposal mechanism to attend to a set of possible regions
  - E.g. Region-CNN [Rich feature hierarchies for accurate object detection and semantic segmentation, Girshick et al., CVPR 2014]

Reproducible research – get the code!





## Fast R-CNN

**Ross Girshick** 

Facebook AI Research (FAIR)

Work done at Microsoft Research

## Fast Region-based ConvNets (R-CNNs) for Object Detection



Figure adapted from Kaiming He

#### Object detection renaissance (2013-present)



#### Object detection renaissance (2013-present)



#### Object detection renaissance (2013-present)



#### Region-based convnets (R-CNNs)

- R-CNN (aka "slow R-CNN") [Girshick et al. CVPR14]
- SPP-net [He et al. ECCV14]





Regions of Interest (Rol) from a proposal method (~2k)







Girshick et al. CVPR14.



Girshick et al. CVPR14.

Post hoc component

Apply bounding-box regressors

#### • Ad hoc training objectives

- Fine-tune network with softmax classifier (log loss)
- Train post-hoc linear SVMs (hinge loss)
- Train post-hoc bounding-box regressors (squared loss)

- Ad hoc training objectives
  - Fine-tune network with softmax classifier (log loss)
  - Train post-hoc linear SVMs (hinge loss)
  - Train post-hoc bounding-box regressors (squared loss)
- Training is slow (84h), takes a lot of disk space

- Ad hoc training objectives
  - Fine-tune network with softmax classifier (log loss)
  - Train post-hoc linear SVMs (hinge loss)
  - Train post-hoc bounding-box regressions (least squares)
- Training is slow (84h), takes a lot of disk space
- Inference (detection) is slow
  - 47s / image with VGG16 [Simonyan & Zisserman. ICLR15]
  - Fixed by SPP-net [He et al. ECCV14]



~2000 ConvNet forward passes per image













Kaiming He, Xiangyu Zhang, Shaoqing Ren, & Jian Sun. "Spatial Pyramid Pooling in Deep Convolutional Networks for Visual Recognition". ECCV 2014.

Classify regions with SVMs SVMs Fully-connected layers FCs Spatial Pyramid Pooling (SPP) layer "conv5" feature map of image Regions of -Interest (Rols) from a proposal Forward *whole* image through ConvNet method ConvNet Input image

He et al. ECCV14.



#### He et al. ECCV14.

#### What's good about SPP-net?

• Fixes one issue with R-CNN: makes testing fast


#### What's wrong with SPP-net?

- Inherits the rest of R-CNN's problems
  - Ad hoc training objectives
  - Training is slow (25h), takes a lot of disk space

#### What's wrong with SPP-net?

- Inherits the rest of R-CNN's problems
  - Ad hoc training objectives
  - Training is slow (though faster), takes a lot of disk space
- Introduces a new problem: cannot update parameters below SPP layer during training

#### SPP-net: the main limitation



He et al. ECCV14.

Post hoc component

#### Fast R-CNN

• Fast test-time, like SPP-net

#### Fast R-CNN

- Fast test-time, like SPP-net
- One network, trained in one stage

#### Fast R-CNN

- Fast test-time, like SPP-net
- One network, trained in one stage
- Higher mean average precision than slow R-CNN and SPP-net









## Fast R-CNN (training)



### Fast R-CNN (training)



#### Multi-task loss

### Fast R-CNN (training)



#### Obstacle #1: Differentiable Rol pooling

Region of Interest (RoI) pooling must be (sub-) differentiable to train conv layers



Slow R-CNN and SPP-net use region-wise sampling to make mini-batches

- Sample 128 example Rols uniformly at random
- Examples will come from different images with high probability



Note the receptive field for one example Rol is often very large

• Worst case: the receptive field is the entire image





## Worst case cost per mini-batch (crude model of computational complexity)

input size for Fast R-CNN

input size for slow R-CNN

# 128\*600\*1000 / (128\*224 \*224) = 12x more computation than slow R-CNN





#### Solution: use hierarchical sampling to build minibatches







#### Solution: use hierarchical sampling to build minibatches



 Sample a small number of images (2)

#### Solution: use hierarchical sampling to build minibatches



- Sample a small number of images (2)
- Sample many examples from each image (64)

Use the test-time trick from SPP-net during training

• Share computation between overlapping examples from the same image





Cost per mini-batch compared to slow R-CNN (same crude cost model)

• 2\*600\*1000 / (128\*224\*224) = 0.19x less computation than slow R-CNN





#### Main results

	Fast R-CNN	R-CNN [1]	SPP-net [2]
Train time (h)	9.5	84	25
- Speedup	8.8x	1x	3.4x
Test time / image	0.32s	47.0s	2.3s
Test speedup	146x	1x	20x
mAP	66.9%	66.0%	63.1%

Timings exclude object proposal time, which is equal for all methods. All methods use VGG16 from Simonyan and Zisserman.

[1] Girshick et al. CVPR14.[2] He et al. ECCV14.

#### Main results

	Fast R-CNN	R-CNN [1]	SPP-net [2]
Train time (h)	9.5	84	25
- Speedup	8.8x	1x	3.4x
Test time / image	0.32s	47.0s	2.3s
Test speedup	146x	1x	20x
mAP	66.9%	66.0%	63.1%

Timings exclude object proposal time, which is equal for all methods. All methods use VGG16 from Simonyan and Zisserman.

[1] Girshick et al. CVPR14.[2] He et al. ECCV14.

#### Main results

	Fast R-CNN	R-CNN [1]	SPP-net [2]
Train time (h)	9.5	84	25
- Speedup	8.8x	1x	3.4x
Test time / image	0.32s	47.0s	2.3s
Test speedup	146x	1x	20x
mAP	66.9%	66.0%	63.1%

Timings exclude object proposal time, which is equal for all methods. All methods use VGG16 from Simonyan and Zisserman.

[1] Girshick et al. CVPR14.[2] He et al. ECCV14.

#### Further test-time speedups



Fully connected layers take 45% of the forward pass time

#### Further test-time speedups



Compress these layers with truncated SVD

J. Xue, J. Li, and Y. Gong.

Restructuring of deep neural network acoustic models with singular value decomposition. *Interspeech*, 2013.

#### Further test-time speedups



### Other findings

### End-to-end training matters

	Fast R-CNN (VGG16)			
Fine-tune layers	$\geq$ fc6	$\geq$ conv3_1	$\geq$ conv2_1	
VOC07 mAP	61.4%	66.9%	67.2%	
Test time per image	0.32s	0.32s	0.32s	
			1.4x slower training	

	Fast R-CNN (VGG16)			
Multi-task training?		Υ		Υ
Stage-wise training?			Υ	
Test-time bbox reg.			Υ	Y
VOC07 mAP	62.6%	63.4%	64.0%	66.9%

	Fast R-CNN (VGG16)			
Multi-task training?		Υ		Υ
Stage-wise training?			Υ	
Test-time bbox reg.			Υ	Υ
VOC07 mAP	62.6%	63.4%	64.0%	66.9%

Trained without a bbox regressor

	Fast R-CNN (VGG16)			
Multi-task training?		Υ		Υ
Stage-wise training?			Υ	
Test-time bbox reg.			Υ	Υ
VOC07 mAP	62.6%	63.4%	64.0%	66.9%

Trained with a bbox regressor, but it's disabled at test time

	Fast R-CNN (VGG16)			
Multi-task training?		Υ		Y
Stage-wise training?			Υ	
Test-time bbox reg.			Υ	Y
VOC07 mAP	62.6%	63.4%	64.0%	66.9%

Post hoc bbox regressor, used at test time

	Fast R-CNN (VGG16)			
Multi-task training?		Υ		Υ
Stage-wise training?			Υ	
Test-time bbox reg.			Υ	Υ
VOC07 mAP	62.6%	63.4%	64.0%	66.9%

Multi-task objective, using bbox regressors at test time
## More proposals is harmful



## What's still wrong?

- Out-of-network region proposals
  - Selective search: 2s / im; EdgeBoxes: 0.2s / im
- Fortunately, we have a solution
  - Our follow-up work was presented last week at NIPS

Shaoqing Ren, Kaiming He, Ross Girshick & Jian Sun. "Faster R-CNN: Towards Real-Time Object Detection with Region Proposal Networks." NIPS 2015.

## **Object Detection: Faster R-CNN**

- Faster R-CNN
  - Solely based on CNN
  - No external modules
  - Each step is end-to-end



Shaoqing Ren, Kaiming He, Ross Girshick, & Jian Sun. "Faster R-CNN: Towards Real-Time Object Detection with Region Proposal Networks". NIPS 2015.



## **Region Proposal Network**

- Slide a small window on the feature map
- Build a small network for:
  - classifying object or not-object, and
  - regressing bbox locations
- Position of the sliding window provides localization information with reference to the image
- Box regression provides finer localization information with reference to this sliding window







## Anchors as references

- Anchors: pre-defined reference boxes
  - Box regression is with reference to anchors: regressing an anchor box to a ground-truth box
  - Object probability is with reference to anchors, e.g.:
    - anchors as positive samples: if IoU > 0.7 or IoU is max
    - anchors as negative samples: if IoU < 0.3





Shaoqing Ren, Kaiming He, Ross Girshick, & Jian Sun. "Faster R-CNN: Towards Real-Time Object Detection with Region Proposal Networks". NIPS 2015.



## Anchors as references

- Anchors: pre-defined reference boxes
- Multi-scale/size anchors:
  - multiple anchors are used at each position: e.g., 3 scales (128<sup>2</sup>, 256<sup>2</sup>, 512<sup>2</sup>) and 3 aspect ratios (2:1, 1:1, 1:2) yield 9 anchors
  - each anchor has its own prediction function
  - single-scale features, multi-scale predictions







## Region Proposal Network

• RPN is fully convolutional [Long et al. 2015] *n* anchors *4n* coordinates *n* scores RPN is trained end-to-end 256-d • RPN shares convolutional feature maps with the detection network (covered in Ross's section)



Shaoqing Ren, Kaiming He, Ross Girshick, & Jian Sun. "Faster R-CNN: Towards Real-Time Object Detection with Region Proposal Networks". NIPS 2015.



Shaoqing Ren, Kaiming He, Ross Girshick, & Jian Sun. "Faster R-CNN: Towards Real-Time Object Detection with Region Proposal Networks". NIPS 2015.

#### Fast R-CNN take-aways

- End-to-end training of deep ConvNets for detection
- Fast training times
- Open source for easy experimentation
   "I think [the Fast R-CNN] code is average-somewhat above average for what it is."
   — sporkles on r/MachineLearning
- A large number of ImageNet detection and COCO detection methods are built on Fast R-CNN Checkout the ImageNet / COCO Challenge workshop on Thursday!

# Focal Loss for Dense Object Detection

Tsung-Yi Lin, Google Brain

Work done at Facebook AI Research with Priya Goyal, Ross Girshick, Kaiming He, Piotr Dollár

## Viola and Jones (2001)



Image from OpenCV 3.3 website

#### Fast R-CNN



#### One-stage vs. Two-stage

- One-stage
  - Fast
  - Simple
- Two-stage
  - 10 40% better accuracy



https://arxiv.org/pdf/1611.10012.pdf

#### One-stage vs. Two-stage



Speed/accuracy trade-offs for modern convolutional object detectors, Huang et al., CVPR 2017

#### Toward dense detection

- YOLOv1 98 boxes
- YOLOv2 ~1k
- OverFeat ~1-2k
- SSD ~8-26k

• This work – ~100k

## **Class Imbalance**

- Few training examples from foreground
- Most examples from background
  - Easy and uninformative
  - Distracting



#### **Cross Entropy**



#### **Cross Entropy**



## Cross Entropy with Imbalance Data

- 100000 easy : 100 hard examples
- 40x bigger loss from easy examples



### Focal Loss

$$CE(p_t) = -\log(p_t)$$
$$FL(p_t) = -(1 - p_t)^{\gamma}\log(p_t)$$



### Focal Loss



# Prior

• α-balanced Cross entropy

$$\operatorname{CE}(p_{\mathsf{t}}) = -\alpha_t \log(p_{\mathsf{t}})$$

• α-balanced Focal Loss

$$FL(p_t) = -\alpha_t (1 - p_t)^{\gamma} \log(p_t)$$

- γ: focus more on hard examples
- α: offset class imbalance of number of examples

## Feature Pyramid Network

- Multiscale
- Semantically strong at all scales
- Fast to compute



Feature Pyramid Network for Object Detection, Lin et al., CVPR 2017

### Architecture

- RetinaNet
  - FPN + 100k boxes
  - Focal loss



#### Loss Distribution under Focal Loss

**Background Boxes** 



#### Loss Distribution under Focal Loss

**Foreground Boxes** 



## vs. Cross Entropy

• + 2.9 AP to  $\alpha$ -balanced cross entropy



## vs. OHEM

• +3.2 AP to best OHEM (ResNet-101 FPN)

	method	batch size	nms thr	AP	
	OHEM	128	.7	31.1	
	OHEM	256	.7	31.8	
	OHEM	512	.7	30.6	
C	OHEM	128	.5	32.8	-> Best OHEM
	OHEM	256	.5	31.0	
	OHEM	512	.5	27.6	
(	OHEM 1:3	128	.5	31.1	-
(	OHEM 1:3	256	.5	28.3	
_	OHEM 1:3	512	.5	24.0	
C	FL	n/a	n/a	36.0	🗲 Best Focal Loss

Online Hard Example Mining, Shrivastava et al., 2016

#### **RetinaNet performance**



## Summary

- Identify class imbalance is the major issue for training onestage dense detector
- Propose Focal Loss to address class imbalance
- Achieve state-of-the-art accuracy and speed



ICCV 2017 Tutorial, Venice, Italy

Kaiming He

in collaboration with: Georgia Gkioxari, Piotr Dollár, and Ross Girshick Facebook AI Research (FAIR)

# Introduction

### Visual Perception Problems



## A Challenging Problem...



## **Object Detection**

Fast/Faster R-CNN
 ✓ Good speed
 ✓ Good accuracy
 ✓ Intuitive
 ✓ Easy to use



Ross Girshick. "Fast R-CNN". ICCV 2015.

Shaoqing Ren, Kaiming He, Ross Girshick, & Jian Sun. "Faster R-CNN: Towards Real-Time Object Detection with Region Proposal Networks". NIPS 2015.

## Semantic Segmentation

Fully Convolutional Net (FCN)
✓ Good speed
✓ Good accuracy
✓ Intuitive
✓ Easy to use



Jonathan Long, Evan Shelhamer, & Trevor Darrell. "Fully Convolutional Networks for Semantic Segmentation". CVPR 2015.
### Instance Segmentation

• Goals of Mask R-CNN

✓ Good speed
 ✓ Good accuracy
 ✓ Intuitive
 ✓ Easy to use



#### **Instance Segmentation Methods R-CNN driven**







#### **FCN driven**









#### **Instance Segmentation Methods**



**RCNN-driven** 

- SDS [Hariharan et al, ECCV'14]
- HyperCol [Hariharan et al, CVPR'15]
  - CFM [Dai et al, CVPR'15]
  - MNC [Dai et al, CVPR'16]

• **PFN** [Liang et al, arXiv'15]

#### **FCN-driven**

- InstanceCut [Kirillov et al, CVPR'17]
- Watershed [Bai & Urtasun, CVPR'17]
- FCIS [Li et al, CVPR'17]
- DIN [Arnab & Torr, CVPR'17]

### Mask R-CNN

• Mask R-CNN = Faster R-CNN with FCN on Rols



### Parallel Heads

• Easy, fast to implement and train



### Invariance vs. Equivariance

- Convolutions are translation-equivariant
- *Fully*-ConvNet (FCN) is translation-equivariant
- ConvNet becomes translation-invariant due to fully-connected or global pool layers

#### Equivariance in Mask R-CNN



1. Fully-Conv Features:

equivariant to global (image) translation

#### Equivariance in Mask R-CNN



2. Fully-Conv on Rol: equivariant to translation within Rol

## Fully-Conv on Rol



#### target masks on Rols



Translation of object in Rol => Same translation of mask in Rol

- Equivariant to small translation of Rols
- More robust to Rol's localization imperfection

#### Equivariance in Mask R-CNN



3. RolAlign:

3a. maintain translation-equivariance before/after Rol

## RolAlign

FAQs: how to sample grid points within a cell?

- 4 regular points in 2x2 sub-cells
- other implementation could work



25

## RolAlign vs. RolPool

RoIPool breaks pixel-to-pixel translation-equivariance



#### Equivariance in Mask R-CNN



3. RolAlign:

**3b.** Scale-equivariant (and aspect-ratio-equivariant)

## Equivariance in Mask R-CNN: Summary

- Translation-equivariant
  - FCN features
  - FCN mask head
  - RolAlign (pixel-to-pixel behavior)
- Scale-equivariant (and aspect-ratio-equivariant)
  - RolAlign (warping and normalization behavior) + paste-back
  - FPN features



#### Mask R-CNN results on COCO

# Result Analysis

## Ablation: RolPool vs. RolAlign

in case of big stride (32)

baseline: ResNet-50-Conv5 backbone, stride=32

		mask AP			box AP	
	AP	$AP_{50}$	AP <sub>75</sub>	AP <sup>bb</sup>	$AP_{50}^{bb}$	$AP_{75}^{bb}$
RoIPool	23.6	46.5	21.6	28.2	52.7	26.9
RoIAlign	30.9	51.8	32.1	34.0	55.3	36.4
	+7.3	+ 5.3	+10.5	+5.8	+2.6	+9.5
• hu	ige gain a	it high Iol	J,			

### Ablation: RolPool vs. RolAlign

baseline: ResNet-50-Conv5 backbone, stride=32

		mask AP			box AP	
	AP	$AP_{50}$	AP <sub>75</sub>	AP <sup>bb</sup>	$AP_{50}^{bb}$	$AP_{75}^{bb}$
RoIPool	23.6	46.5	21.6	28.2	52.7	26.9
RoIAlign	30.9	51.8	32.1	34.0	55.3	36.4
	+7.3	+ 5.3	+10.5	+5.8	+2.6	+9.5

• nice box AP without dilation/upsampling

### Instance Segmentation Results on COCO

	backbone	AP	$AP_{50}$	AP <sub>75</sub>	$AP_S$	$AP_M$	$AP_L$
MNC [7]	ResNet-101-C4	24.6	44.3	24.8	4.7	25.9	43.6
FCIS [20] +OHEM	ResNet-101-C5-dilated	29.2	49.5	-	7.1	31.3	50.0
FCIS+++ [20] +OHEM	ResNet-101-C5-dilated	33.6	54.5	-	-	-	-
Mask R-CNN	ResNet-101-C4	33.1	54.9	34.8	12.1	35.6	51.1
Mask R-CNN	ResNet-101-FPN	35.7	58.0	37.8	15.5	38.1	52.4
Mask R-CNN	ResNeXt-101-FPN	37.1	60.0	39.4	16.9	39.9	53.5

- 2 AP better than SOTA w/ R101, without bells and whistles
- 200ms / img

### **Instance Segmentation** Results on COCO

	backbone	AP	$AP_{50}$	AP <sub>75</sub>	$AP_S$	$AP_M$	$AP_L$
MNC [7]	ResNet-101-C4	24.6	44.3	24.8	4.7	25.9	43.6
FCIS [20] +OHEM	ResNet-101-C5-dilated	29.2	49.5	-	7.1	31.3	50.0
FCIS+++ [20] +OHEM	ResNet-101-C5-dilated	33.6	54.5	-	-	-	-
Mask R-CNN	ResNet-101-C4	33.1	54.9	34.8	12.1	35.6	51.1
Mask R-CNN	ResNet-101-FPN	35.7	58.0	37.8	15.5	38.1	52.4
Mask R-CNN	ResNeXt-101-FPN	37.1	60.0	<b>39.4</b>	16.9	<b>39.9</b>	53.5

• benefit from better features (ResNeXt [Xie et al. CVPR'17])

### **Object Detection** Results on COCO

	backbone	AP <sup>bb</sup>	$AP_{50}^{bb}$	$AP_{75}^{bb}$	$\mathrm{AP}^{\mathrm{bb}}_S$	$\mathrm{AP}^{\mathrm{bb}}_M$	$AP_L^{bb}$
Faster R-CNN+++ [15]	ResNet-101-C4	34.9	55.7	37.4	15.6	38.7	50.9
Faster R-CNN w FPN [22]	ResNet-101-FPN	36.2	59.1	39.0	18.2	39.0	48.2
Faster R-CNN by G-RMI [17]	Inception-ResNet-v2 [32]	34.7	55.5	36.7	13.5	38.1	52.0
Faster R-CNN w TDM [31]	Inception-ResNet-v2-TDM	36.8	57.7	39.2	16.2	39.8	52.1
Faster R-CNN, RoIAlign	ResNet-101-FPN	37.3	59.6	40.3	19.8	40.2	48.8
Mask R-CNN	ResNet-101-FPN	38.2	60.3	41.7	20.1	41.1	50.2
Mask R-CNN	ResNeXt-101-FPN	39.8	62.3	43.4	22.1	43.2	51.2

bbox detection improved by:

RolAlign

### **Object Detection** Results on COCO

	backbone	AP <sup>bb</sup>	$AP_{50}^{bb}$	$AP_{75}^{bb}$	$\mathrm{AP}^{\mathrm{bb}}_S$	$\mathrm{AP}^{\mathrm{bb}}_M$	$AP_L^{bb}$
Faster R-CNN+++ [15]	ResNet-101-C4	34.9	55.7	37.4	15.6	38.7	50.9
Faster R-CNN w FPN [22]	ResNet-101-FPN	36.2	59.1	39.0	18.2	39.0	48.2
Faster R-CNN by G-RMI [17]	Inception-ResNet-v2 [32]	34.7	55.5	36.7	13.5	38.1	52.0
Faster R-CNN w TDM [31]	Inception-ResNet-v2-TDM	36.8	57.7	39.2	16.2	39.8	52.1
Faster R-CNN, RoIAlign	ResNet-101-FPN	37.3	59.6	40.3	19.8	40.2	48.8
Mask R-CNN	ResNet-101-FPN	38.2	60.3	41.7	20.1	41.1	50.2
Mask R-CNN	ResNeXt-101-FPN	39.8	62.3	43.4	22.1	43.2	51.2

bbox detection improved by:

- RolAlign
- Multi-task training w/ mask



#### Mask R-CNN results on COCO



Mask R-CNN results on COCO

small objects



#### Mask R-CNN results on CityScapes

#### Failure case: detection/segmentation



#### Mask R-CNN results on COCO

#### Failure case: recognition



Mask R-CNN results on COCO

not a kite



Validation image with box detection shown in red

#### 28x28 soft prediction from Mask R-CNN (enlarged)



#### Soft prediction resampled to image coordinates

(bilinear and bicubic interpolation work equally well)



#### Final prediction (threshold at 0.5)





#### 28x28 soft prediction



#### Resized Soft prediction



Final mask



Validation image with box detection shown in red

### Mask R-CNN: for Human Keypoint Detection

- 1 keypoint = 1-hot "mask"
- Human pose = 17 masks
- Softmax over spatial locations
  - e.g. 56<sup>2</sup>-way softmax on 56x56
- Desire the same equivariances
  - translation, scale, aspect ratio





## Mask R-CNN frame-by-frame

## Conclusion



#### Mask R-CNN

- ✓ Good speed
  ✓ Good accuracy
  ✓ Intuitive
- ✓ Easy to use
- ✓ Equivariance matters

Code will be open-sourced as Facebook AI Research's **Detectron** platform