

Introduction	Method	Experiments	Limitations	Conclusion	Future work
○○○○○	○○○○○○○ ○○○	○○○○○○○ ○○			

R-FCN: Object Detection via Region-based Fully Convolutional Networks

Kris Korrel Diederik Rusticus

Master Artificial Intelligence
University of Amsterdam

Computer Vision 2, 2017

Kris Korrel, Diederik Rusticus
R-FCN: Object Detection via Region-based Fully Convolutional Networks

Introduction	Method	Experiments	Limitations	Conclusion	Future work
○○○○○	○○○○○○○ ○○○	○○○○○○○ ○○			

Outline

- 1 Introduction
 - Authors
 - Previous work
- 2 Method
 - Architecture
 - Training
 - Inference
- 3 Experiments
 - Data: PASCAL VOC
 - Data: MC COCO
- 4 Limitations
- 5 Conclusion
- 6 Future work

Kris Korrel, Diederik Rusticus
R-FCN: Object Detection via Region-based Fully Convolutional Networks

Introduction	Method	Experiments	Limitations	Conclusion	Future work
●○○○○	○○○○○○○ ○○○	○○○○○○○ ○○			

Authors

Microsoft Researchers - Paper: June 2016



(a) Jian Sun



(b) Jifeng Dai



(c) Kaiming He



(d) Yi Li

Kris Korrel, Diederik Rusticus
R-FCN: Object Detection via Region-based Fully Convolutional Networks

Introduction	Method	Experiments	Limitations	Conclusion	Future work
●○○○○	○○○○○○○ ○○○	○○○○○○○ ○○			

Authors

R-FCN

Introducing R-FCN:

Region-based Fully Convolutional Network

- Shared, fully convolutional (like FCN)
- Translation variance incorporated

Kris Korrel, Diederik Rusticus
R-FCN: Object Detection via Region-based Fully Convolutional Networks

Introduction	Method	Experiments	Limitations	Conclusion	Future work
●○○○○	○○○○○○○ ○○○	○○○○○○○ ○○			

Previous work

AlexNet (2012), VGGNet (2015), ResNet (2016), etc.

- Designed for image classification
- Convolutional and pooling layers, followed by fully connected layers
- Used convolutional layers for translation invariance

Kris Korrel, Diederik Rusticus
R-FCN: Object Detection via Region-based Fully Convolutional Networks

Introduction	Method	Experiments	Limitations	Conclusion	Future work
●○○○○	○○○○○○○ ○○○	○○○○○○○ ○○			

Previous work

R-CNN (2014)

- Designed for object detection
- Rols proposed by region proposal algorithm (like Selective Search)
- Similar design, but Rol-wise

Kris Korrel, Diederik Rusticus
R-FCN: Object Detection via Region-based Fully Convolutional Networks

Introduction	Method	Experiments	Limitations	Conclusion	Future work
○○●○○	○○○○○○○ ○○○	○○○○○○○ ○○			

Previous work

Fast R-CNN (2015)

- Extends on R-CNN
- One forward pass through convolutional layers followed by Rol-wise fully connected layers
- Introduces Rol pooling layer
- 0.3s test time excluding generation of region proposals

Kris Korrel, Diederik Rusticus
R-FCN: Object Detection via Region-based Fully Convolutional Networks

Introduction	Method	Experiments	Limitations	Conclusion	Future work
○○○○●	○○○○○○○ ○○○	○○○○○○○ ○○			

Previous work

Faster R-CNN (2015)

- Extends on Fast R-CNN
- Introduces Region Proposal Network
- Trained alternatively with object detection network

Kris Korrel, Diederik Rusticus
R-FCN: Object Detection via Region-based Fully Convolutional Networks

Introduction	Method	Experiments	Limitations	Conclusion	Future work
○○○○●	○○○○○○○ ○○○	○○○○○○○ ○○			

Previous work

Subdivided networks

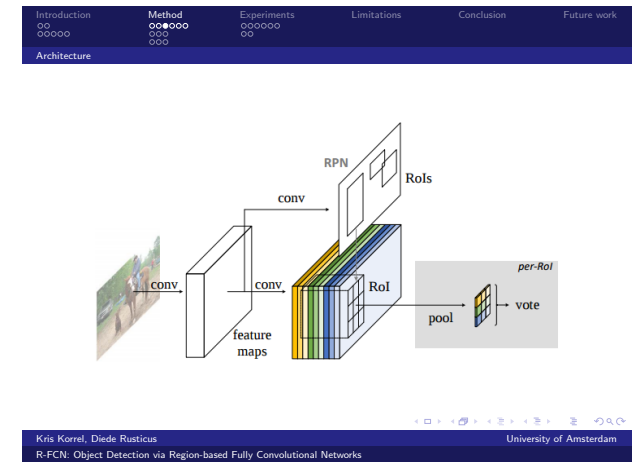
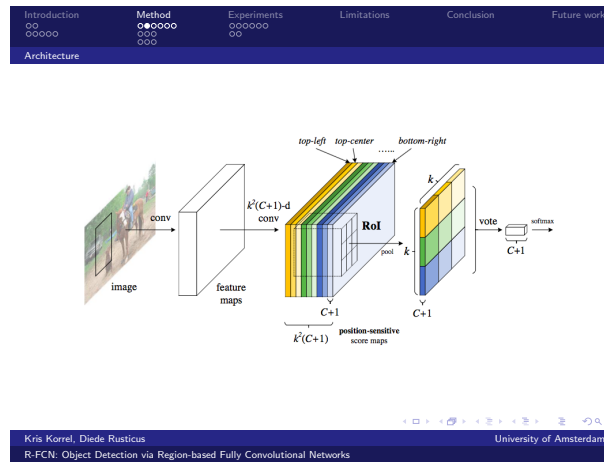
- Translation invariance needed for image classification
- Translation variance needed for object detection
- Fast and Faster R-CNN solve this by Rol pooling
- No shared computation in second subnetwork
- Can not take advantage of properties of FCNs like ResNet and GoogLeNet

Kris Korrel, Diederik Rusticus
R-FCN: Object Detection via Region-based Fully Convolutional Networks

Introduction	Method	Experiments	Limitations	Conclusion	Future work
00 00000	000000 000 000	000000 00			
Architecture					
Method					

Overview:

- All learnable weights are fully convolutional
- Last layer produces a bank of position-sensitive score maps
- Ends with position-sensitive RoI pooling and voting, which is not learned
- Negligible RoI-wise computation
- RPN from Faster R-CNN is used for proposing Rols



Introduction	Method	Experiments	Limitations	Conclusion	Future work
00 00000	000000 000 000	000000 00			
Architecture					
Network design					

- ResNet-101 as backbone architecture
- Dimensionality reduction from 2048-d to 1024-d
- Last layer produces bank of $k^2(C+1)$ position-sensitive score maps

Introduction	Method	Experiments	Limitations	Conclusion	Future work
00 00000	000000 000 000	000000 00			
Architecture					
Position-sensitive pooling					

For each RoI:

- Position-sensitive RoI pooling
 $K^2(C+1)$ score maps $\rightarrow K^2(C+1)$ scores
- Voting
 $K^2(C+1)$ scores $\rightarrow C+1$ scores
- Softmax responses

Introduction	Method	Experiments	Limitations	Conclusion	Future work
00 00000	000000 000 000	000000 00			
Architecture					
Bounding boxes					

- Bounding boxes are proposed by RPN as introduced by Faster R-CNN
- R-FCN and RPN are trained alternatively
- Bounding box regression is done similarly to object detection
- bank of $4k^2$ score maps that predict (t_x, t_y, t_w, t_h)

Introduction	Method	Experiments	Limitations	Conclusion	Future work
00 00000	000000 000 000	000000 00			
Training					
Training					

The loss function defined on each RoI is the summation of the **cross-entropy loss** and the **box regression loss**:

$$L(s, t_{x,y,w,h}) = \underbrace{L_{cls}(s_{c^*})}_{\text{classif. loss}} + \lambda \underbrace{[c^* > 0]}_{\text{indicator (t/f)}} \underbrace{L_{reg}(t, t^*)}_{\text{bounding box regr. loss}}$$

Positive example = IoU > 0.5

Introduction	Method	Experiments	Limitations	Conclusion	Future work
00 00000	000000 000 000	000000 00			
Training					
Online hard example mining					

OHEM is easy to adopt:

- The Rols are sorted by their losses, X Rols with the **highest loss** are used for **backpropagation**
- Per-Rol computation is negligible so forward time is not affected by N
- In Fast R-CNN this doubles training time!

Introduction	Method	Experiments	Limitations	Conclusion	Future work
00 00000	000000 000 000	000000 00			
Training					
Default Settings					

weight decay = 0.0005
momentum = 0.9
single-scale training = images are resized: shorter side is 600 px.
OHEM = 128 Rols for backprop.
GPU = 1 RoI per GPU - 8 GPU's used
fine-tune learning rate = 0.001 / 0.0001

Introduction
○○○○○○○
Inference

Method
○○○○○○○
○○○

Experiments
○○○○○○○
○○

Limitations

Conclusion

Future work

Inference

- Feature maps are computed
- RPN proposes Rols
- R-FCN evaluates category scores and regresses bounding boxes
- To compare with Faster R-CNN; also 300 Rols are evaluated
- Non-maximum suppression as post-processing

Introduction
○○○○○○○
Inference

Method
○○○○○○○
○○○

Experiments
○○○○○○○
○○

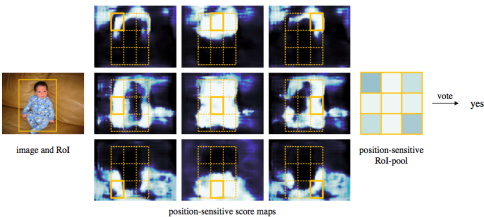
Limitations

Conclusion

Future work

Visualisation

Person category: $(k \times k = 3 \times 3)$



Introduction
○○○○○○○
Inference

Method
○○○○○○○
○○○

Experiments
○○○○○○○
○○

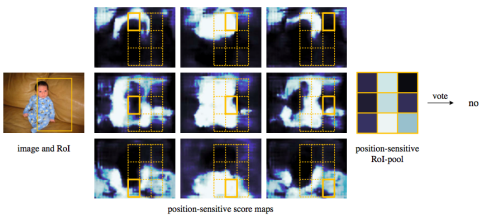
Limitations

Conclusion

Future work

Visualisation

Person category: $(k \times k = 3 \times 3)$



Introduction
○○○○○○○
Data: PASCAL VOC

Method
○○○○○○○
○○○

Experiments
○○○○○○○
○○

Limitations

Conclusion

Future work

Dataset

PASCAL VOC

- 20 object categories
- trained on VOC 2007 + 2012 trainval set, evaluated on VOC 2007 test set
- Performance measure: mAP

Introduction
○○○○○○○
Data: PASCAL VOC

Method
○○○○○○○
○○○

Experiments
○○○○○○○
○○

Limitations

Conclusion

Future work

Results 1/4

Comparison: other Fully Convolutional Strategies

method	RoI output size $(k \times k)$	mAP on VOC 07 (%)
naïve Faster R-CNN	1×1 7×7	61.7 68.9
class-specific RPN	-	67.6
R-FCN (w/o position-sensitivity)	1×1	fail
R-FCN	3×3 7×7	75.5 76.6

Introduction
○○○○○○○
Data: PASCAL VOC

Method
○○○○○○○
○○○

Experiments
○○○○○○○
○○

Limitations

Conclusion

Future work

Results 2/4

Comparison: Faster R-CNN using ResNet-101

	depth of per-RoI subnetwork	training w/ COCO17	train time (sec/img)	test time (sec/img)	mAP (%) on VOC07
Faster R-CNN	10		1.2	0.42	76.4
R-FCN	0		0.45	0.17	76.6
Faster R-CNN	10	✓ (300 Rols)	1.5	0.42	79.3
R-FCN	0	✓ (300 Rols)	0.45	0.17	79.5
Faster R-CNN	10	✓ (2000 Rols)	2.9	0.42	N/A
R-FCN	0	✓ (2000 Rols)	0.46	0.17	79.3

Introduction
○○○○○○○
Data: PASCAL VOC

Method
○○○○○○○
○○○

Experiments
○○○○○○○
○○

Limitations

Conclusion

Future work

Results 3/4

Comparison: Faster R-CNN using ResNet-101

Trained on MS COCO and finetuned with PASCAL VOC

	training data	mAP (%)	test time (sec/img)
Faster R-CNN [9]	07+12	76.4	0.42
Faster R-CNN +++ [9]	07+12+COCO	85.6	3.36
R-FCN	07+12	79.5	0.17
R-FCN multi-sc train	07+12	80.5	0.17
R-FCN multi-sc train	07+12+COCO	83.6	0.17

Introduction
○○○○○○○
Data: PASCAL VOC

Method
○○○○○○○
○○○

Experiments
○○○○○○○
○○

Limitations

Conclusion

Future work

Results 4/4

Comparison: Faster R-CNN using ResNet-101

Trained on MS COCO and finetuned with PASCAL VOC

	training data	mAP (%)	test time (sec/img)
Faster R-CNN [9]	07+12	73.8	0.42
Faster R-CNN +++ [9]	07+12+COCO	83.8	3.36
R-FCN multi-sc train	07+12	77.6 [†]	0.17
R-FCN multi-sc train	07+12+COCO	82.0 [‡]	0.17

Introduction
○○○○○○○
Data: PASCAL VOC

Method
○○○○○○○
○○○

Experiments
○○○○○○○
○○

Limitations

Conclusion

Future work

Further results

Impact of depth

- Increased accuracy from depth 50 to 101
- Saturated accuracy at a depth of 152

Impact of region proposals

- Good generality of the method because competitive performances with other region proposal methods

	training data	test data	RPN [18]	SS [27]	EB [28]
R-FCN	07+12	07	79.5	77.2	77.8

