# Relation Networks for Object Detection

<u>Han Hu</u><sup>1\*</sup>, Jiayuan Gu<sup>2\*</sup>, Zheng Zhang<sup>1\*</sup>, Jifeng Dai<sup>1</sup>, and Yichen Wei<sup>1</sup> <sup>1</sup>Microsoft Research Asia (MSRA) <sup>2</sup>Peking University (\*Equal contribution)

#### pixel-pixel relation



 $\checkmark$  convolution

#### *pixel-pixel* relation



 $\checkmark$  convolution

#### *part-part* relation



✓ *RolPool+FC* 

#### pixel-pixel relation



✓ convolution

#### *part-part* relation



✓ *RolPool+FC* 

#### Effective and Easy to use

- ✓ Parallel
- ✓ Learnable
- ✓ Require no relation supervision
- ✓ Translational invariant
- ✓ Stackable (convolution)

#### pixel-pixel relation



 $\checkmark$  convolution

#### *part-part* relation



✓ RolPool+FC

#### object-object relation



?

## Well Recognized Problem



It is much easier to detect the *glove* if we know there is a *person*.

# Rarely Studied in Deep Learning Era



#### **Irregularities of objects**

- At arbitrary image locations
- Of different scales
- Within different categories
- Of varying number across different images

# Rarely Studied in Deep Learning Era



#### **Irregularities of objects**

- At arbitrary image locations
- Of different scales
- Within different categories
- Of varying number across different images

## Goal



**Goal**: design a simple module to model object-object relation

#### Effective and Easy to use

- ✓ Parallel
- ✓ Learnable
- ✓ Require no relation supervision
- ✓ Translational invariant
- ✓ In-place, stackable

## **Object Relation Module**

#### • **Extension** of *attention module*





**object-object** relation (**2D irregular**)

Left figure credit by A. Vaswani et al.



#### in standard *attention* module

• A novel geometric weight



in standard *attention* module

in object relation module

• A novel geometric weight



#### in standard *attention* module

in object relation module

• A novel geometric weight



in standard *attention* module

in object relation module

## **Relation Aggregation**

object n  $f_{out}(n)$  $f_{in}(m)$ 



## Multi-Branch Relation



## Multi-Branch Relation

branch #1 (person->glove)



branch #2 (playground->glove)



branch #N (duplicate proposals)





## **Object Relation Module**



Kaiming He, Xiangyu Zhang, Shaoqing Ren, Jian Sun. Deep Residual Learning for Image Recognition. CVPR, 2016

## **Object Relation Module**



# Application: Object Detection



• Fast/Faster R-CNN



• Fast/Faster R-CNN



• Fast/Faster R-CNN



• Fast/Faster R-CNN



### Our method: inserting object relation modules (ORMs)



Instance Recognition

**Duplicate Removal** 

### Our method: inserting object relation modules (ORMs)



### Learnable Duplicate Removal



## The First Fully End-to-End Object Detector



back propagation steps

# Results



faster R-CNN		NN	+ object relation modules								
		NIN	<i>W.O</i> .	W.O.	<i>W.O.</i>	$\begin{bmatrix} m_1 & m_2 \end{bmatrix} = \begin{bmatrix} 1 & 1 \end{bmatrix}$			$\int r_1 r_2 = \int \Lambda \Lambda$		
			geometric weight	multi-branch	residual	$\{1, 1, 1, 2\} = \{1, 1\}$		<b>_</b> 1,1 <i>∫</i>	1/1, /2 = 14, 4		
	29.6		30.3	30.5	30.9		31.9		32.8		

\*Faster R-CNN with ResNet-50 model are used

- +2.3 mAP by inserting 2 ORMs
- with +3% FLOPs



faster R_CNN	+ object relation modules									
laster K-CIVIN	<i>W.O</i> .	W.O.	<i>W.O.</i>	W.O. $ \begin{cases} w.o. \\ \int r_1 & r_2 \\ \end{bmatrix} = \int 1 & 1 \\ \end{bmatrix}$		$\int r_1 r_2 = \int A A$				
	geometric weight	multi-branch residual		1/1, /2 = 1, 1	$\{1, 1, 1, 2\} - \{4, 4\}$					
29.6	30.3	30.5	30.9	31.9		32.8				

\*Faster R-CNN with ResNet-50 model are used

• More modules: 8 ORMs



faster R_CNN		+ object relation modules								
laster K-CIVIN	<i>W.O.</i>		W.O.	<i>W.O.</i>	Sm. a	$\begin{bmatrix} m & m_2 \end{bmatrix} = \begin{bmatrix} 1 & 1 \end{bmatrix}$		$\int r_1 r_2 = \int A A$		
	geometric weig		eight	multi-branch	residual	$\{1, 1, 1, 2\} = \{1, 1\}$		_ <b>1</b> , <b>1</b> }	$\{1, 1, 1, 2\} - \{4, 4\}$	
29.6		30.3		30.5	30.9		31.9		32.8	

\*Faster R-CNN with ResNet-50 model are used

• Importance of **relative geometric weight** 



faster R_CNN	+ object relation modules									
	<i>W.O</i> .	<i>w.c</i>		<i>W.O.</i>		$ \int m_1 m_2 = \int 1 1 $		1 1	$\begin{bmatrix} r_1 & r_2 \end{bmatrix} - \begin{bmatrix} 1 & 1 \end{bmatrix}$	
	geometric weight	mι	ılti-brar	nch	residual	$\{1, 1, 1, 2\} = \{1, 1\}$		⊥,⊥∫	1/1, 1/2 = 14, 4	
29.6	30.3		30.5		30.9		31.9		32.8	

\*Faster R-CNN with ResNet-50 model are used

• Importance of **multi-branch relation** 



faster <b>R</b> _CNN	+ object relation modules								
laster K-CIVIN	<i>W.O</i> .	W.O.	<i>W.O</i> .	Sr.	r_l _ ∫	1 1	$\begin{bmatrix} r_1 & r_2 \end{bmatrix} = \begin{bmatrix} 1 & 1 \end{bmatrix}$		
	geometric weight	multi-branch	residual	$\{1, 1, 2\} = \{1, 1\}$			$\left\{ 1, 1, 2 \right\} = \left\{ 4, 4 \right\}$		
29.6	30.3	30.5	30.9		31.9		32.8		

\*Faster R-CNN with ResNet-50 model are used

• Importance of **residual connection** 

## **Duplicate Removal Experiments**



## Duplicate Removal Experiments

	fixed			
method	parameters	mAP	$mAP_{50}$	mAP <sub>75</sub>
NMS	$N_t = 0.3$	29.0	51.4	29.4
NMS	$N_t = 0.4$	29.4	52.1	29.5
NMS	$N_t = 0.5$	29.6	51.9	29.7
NMS	$N_t = 0.6$	29.6	50.9	30.1
NMS	$N_t = 0.7$	28.4	46.6	30.7
SoftNMS	$\sigma = 0.2$	30.0	52.3	30.5
SoftNMS	$\sigma = 0.4$	30.2	51.7	31.3
SoftNMS	$\sigma = 0.6$	30.2	50.9	31.6
SoftNMS	$\sigma = 0.8$	29.9	49.9	31.6
SoftNMS	$\sigma = 1.0$	29.7	49.7	31.6
ours	$\eta = 0.5$	30.3	51.9	31.5
ours	$\eta = 0.75$	30.1	49.0	32.7
ours	$\eta \in [0.5, 0.9]$	30.5	50.2	32.4



- Noticeably better than NMS
- Slightly better than SoftNMS [N. Bodla et al, 2017]

## Fully End-to-End Object Detection



method	parameters	mAP	$mAP_{50}$	mAP <sub>75</sub>
NMS	$N_t = 0.6$	29.6	50.9	30.1
SoftNMS	$\sigma = 0.6$	30.2	50.9	31.6
ours	$\eta = 0.5$	30.3	51.9	31.5
ours	$\eta = 0.75$	30.1	49.0	32.7
ours	$\eta \in [0.5, 0.9]$	30.5	50.2	32.4
ours (e2e)	$\eta \in [0.5, 0.9]$	31.0	51.4	32.8

• Benefit from fully end-to-end training

## Using Stronger Backbones

backbone	setting	mAP	$mAP_{50}$	$mAP_{75}$	#. params	FLOPS	
	2fc+SoftNMS	32.2/32.7	52.9/53.6	34.2/34.7	58.3M	122.2B	
faster RCNN	2fc+RM+SoftNMS	34.7/35.2	55.3/ <b>56.2</b>	37.2/37.8	64.3M	124.6B	+3.0 mAP
	2fc+RM+e2e	35.2/35.4	<b>55.8</b> /56.1	38.2/38.5	64.6M	124.9B	
	2fc+SoftNMS	36.8/37.2	57.8/58.2	40.7/41.4	56.4M	145.8B	
FPN	2fc+RM+SoftNMS	38.1/38.3	59.5/59.9	41.8/42.3	62.4M	157.8B	+2.0 mAP
	2fc+RM+e2e	38.8/38.9	60.3/60.5	42.9/43.3	62.8M	158.2B	
	2fc+SoftNMS	37.5/38.1	57.3/58.1	41.0/41.6	60.5M	125.0B	
DCN	2fc+RM+SoftNMS	38.1/38.8	57.8/ <b>58.7</b>	41.3/42.4	66.5M	127.4B	+1.0 mAP
	2fc+RM+e2e	38.5/39.0	<b>57.8</b> /58.6	42.0/42.9	66.8M	127.7B	

\*Faster R-CNN with ResNet-101 model are used (evaluation on *minival/test-dev* are reported)

less than 10% computation overhead on all backbones

## What is Learnt?

## **Object Pairs with High Relation Weights**

### instance recognition









0.140 570.155

duplicate removal







other objects contributing high weights

## Class Co-Occurrence Information is Learnt





Class Co-occurrence Probability

Learnt Attentional Weights

## Conclusion

• A novel object relation module to model object-object relation

✓ Parallel

✓Learnable

- ✓ Require no relation supervision
- $\checkmark Translational invariant$
- ✓ Stackable
- Application: Object Detection
  - ✓ Improves object detection accuracy✓ The first fully end-to-end object detector



#### code:

https://github.com/msracver/Relation-Networks-for-Object-Detection