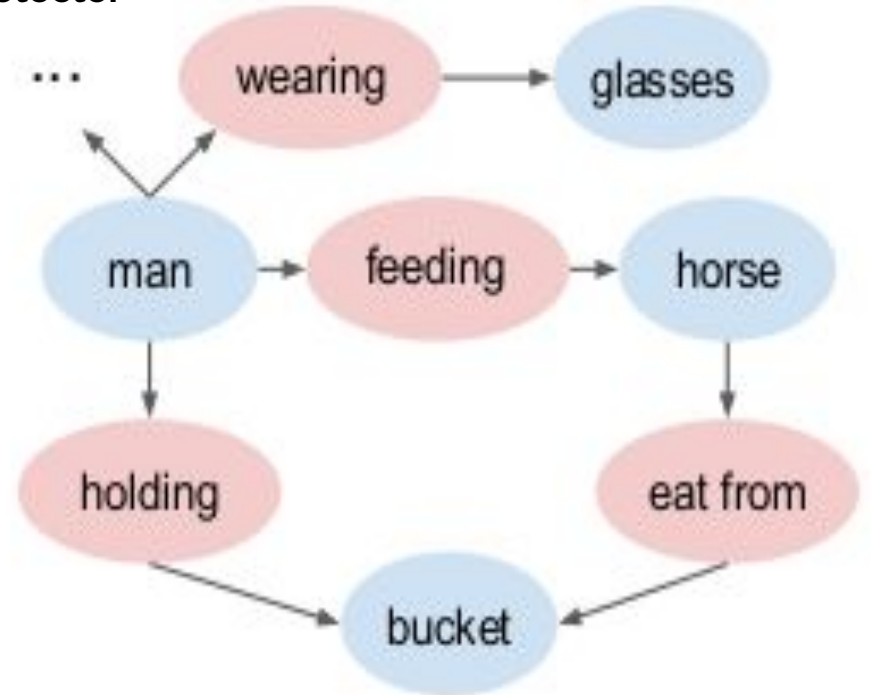

Scene Graph Generation by Iterative Message Passing

Danfei Xu, Yuke Zhu, Christopher B. Choy, Li Fei-Fei (2017)

<https://arxiv.org/abs/1701.02426>

Scene Graph

In every image, there's more than meets the eye





Problem statement

Given an image I and a set of boxes B

from pretrained Region Proposal Network,
we want to identify:

→ **Object classes**

For each box, the object class

→ **Box offsets**

For each box, the offset w.r.t. the
proposed box coordinates

→ **Pairwise relationships**

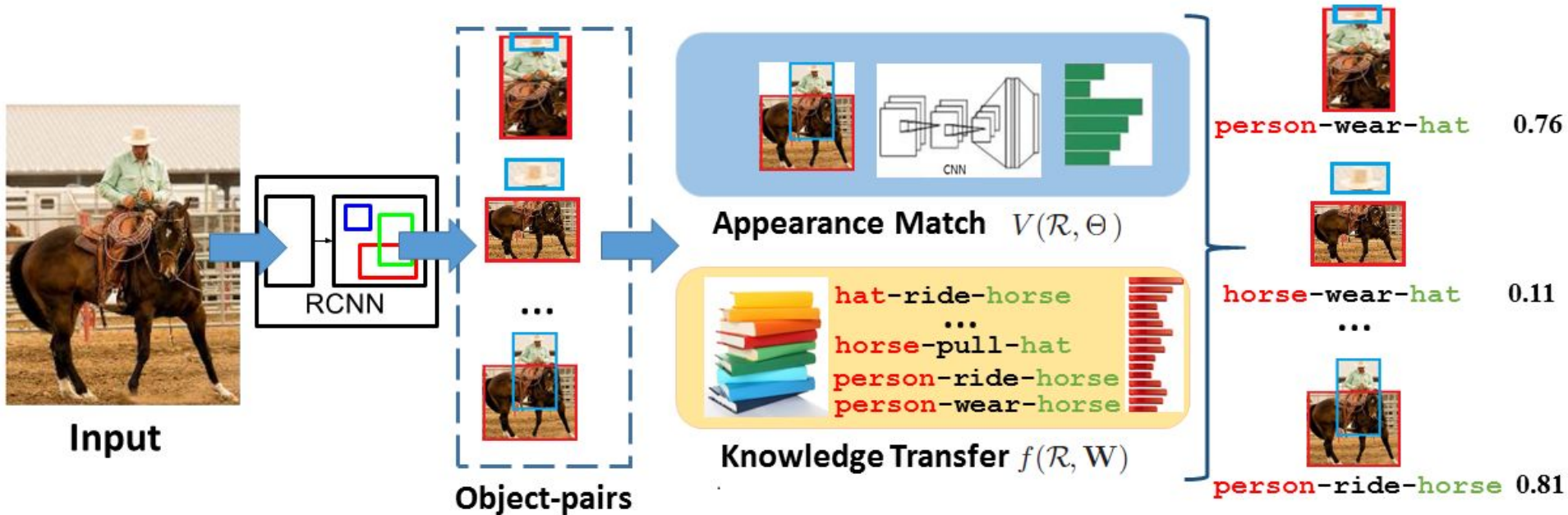
For each pair of boxes, the most likely
relationship between their objects

The baseline

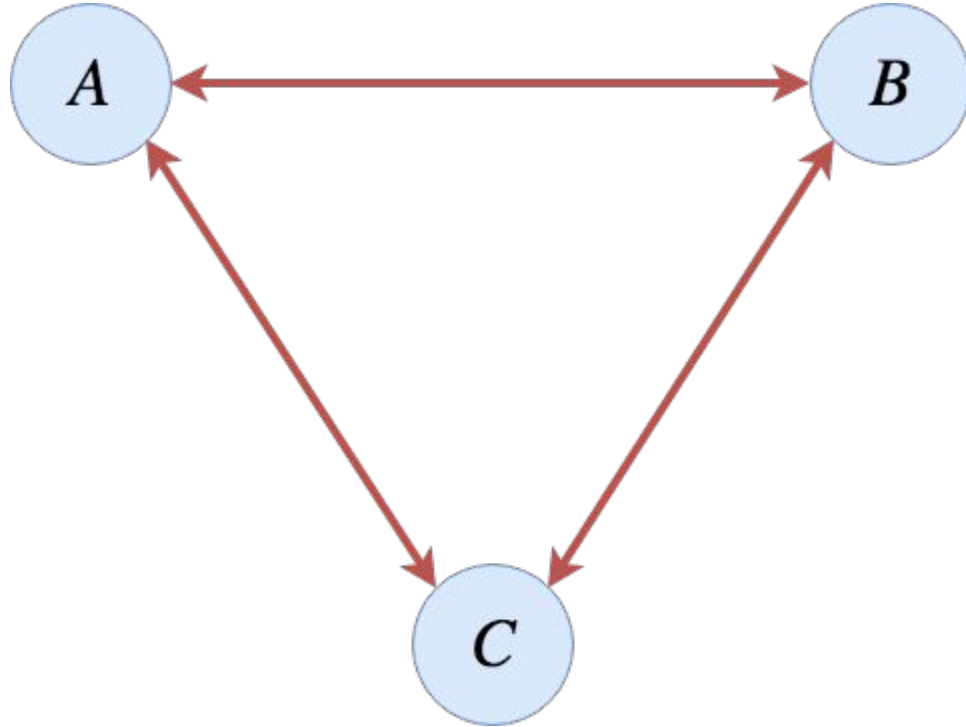
Visual Relationship Detection with Language Priors

Cewu Lu, Ranjay Krishna, Michael Bernstein, Li Fei-Fei (2016)

- Uses visual features from the region containing 2 objects
- Uses language priors to cluster relationships together
- Similar to the approach of this paper, without message passing

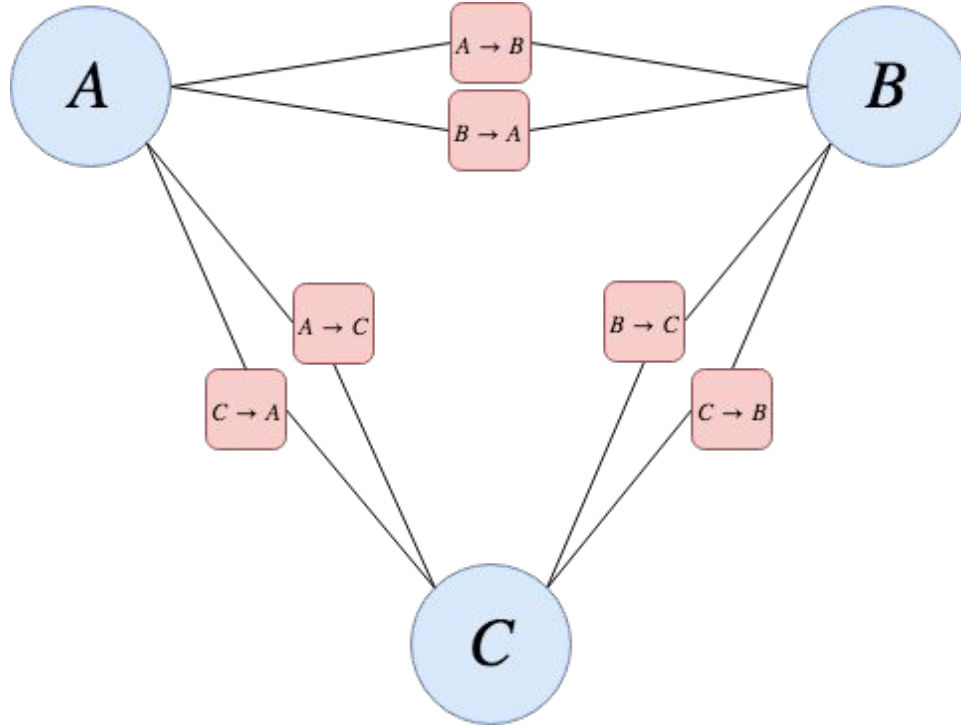


Primal Dual Graph

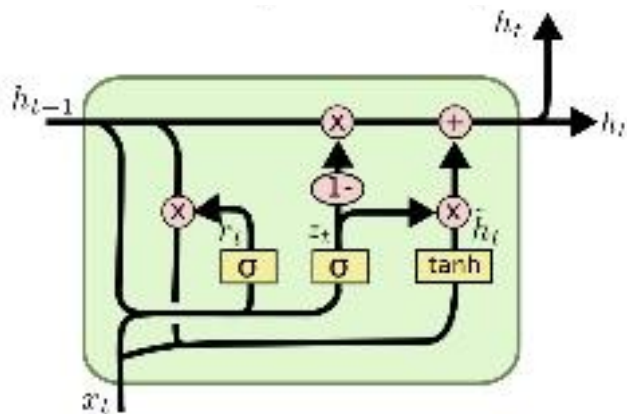


- **Nodes**
Represent objects in the scene
- **Edges**
Represent object relationships

Primal Dual Graph



- **Object nodes**
Represent objects in the scene
 - **Relationship nodes**
Represent object relationships
 - **Edges**
Represent messages exchanged between object and relationship nodes
- Object nodes and relationship nodes form a bipartite graph



Gated Recurrent Unit

Learning Phrase Representations using RNN
Encoder-Decoder for Statistical Machine Translation

Kyunghyun Cho, Bart van Merriënboer, Caglar Gulcehre, Dzmitry Bahdanau, Fethi Bougares, Holger Schwenk, Yoshua Bengio (2014)

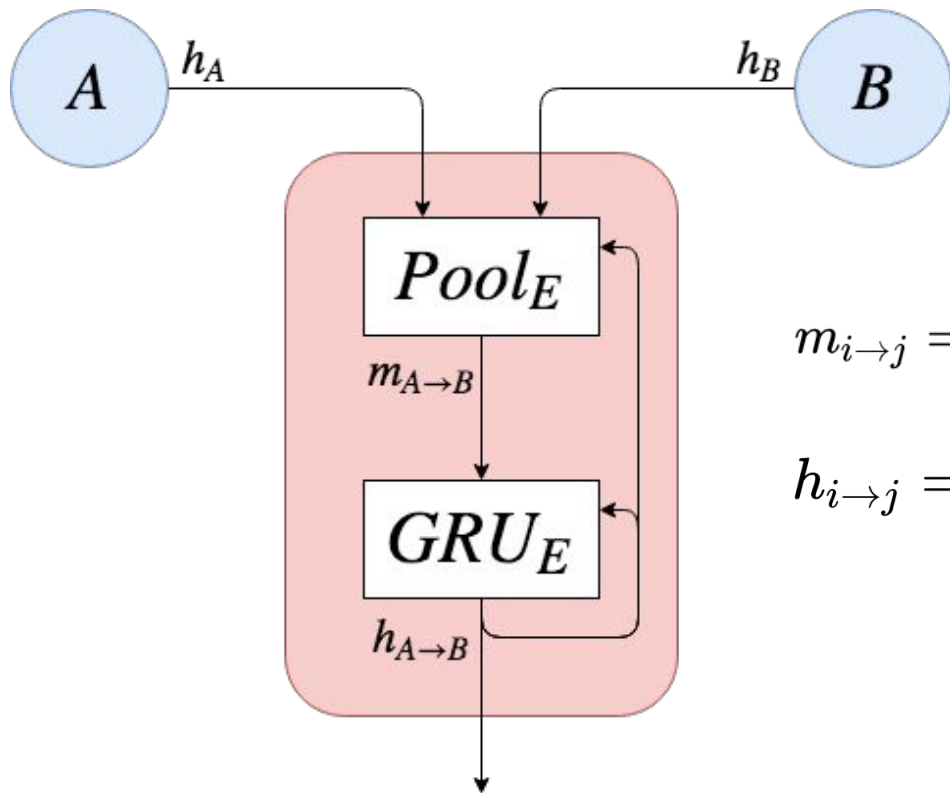
$$z_t = \sigma(W_z \cdot [h_{t-1}, x_t])$$

$$r_t = \sigma(W_r \cdot [h_{t-1}, x_t])$$

$$\tilde{h}_t = \tanh(W \cdot [r_t * h_{t-1}, x_t])$$

$$h_t = (1 - z_t) * h_{t-1} + z_t * \tilde{h}_t$$

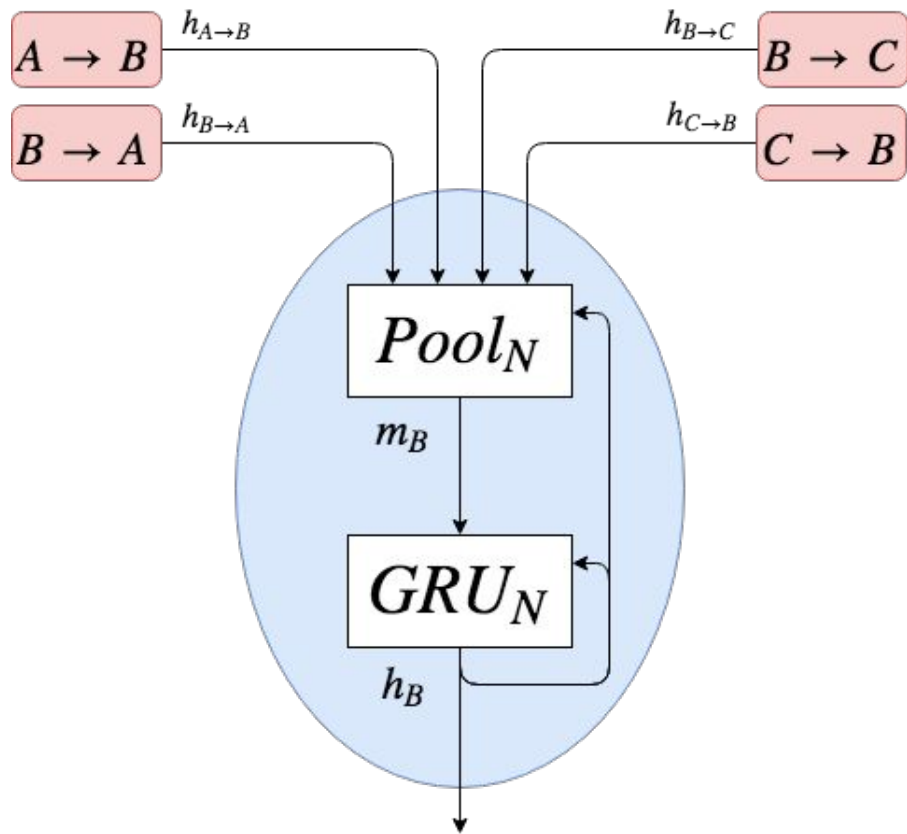
Relationship node $A \rightarrow B$



$$m_{i \rightarrow j} = \sigma(\mathbf{w}_1^T [h_i, h_{i \rightarrow j}])h_i + \sigma(\mathbf{w}_2^T [h_j, h_{i \rightarrow j}])h_j$$

$$h_{i \rightarrow j} = GRU(m_{i \rightarrow j}, h_{i \rightarrow j})$$

Object node B



$$m_i = \sum_{j:i \rightarrow j} \sigma(\mathbf{v}_1^T [h_i, h_{i \rightarrow j}]) h_{i \rightarrow j} + \sum_{j:j \rightarrow i} \sigma(\mathbf{v}_2^T [h_i, h_{j \rightarrow i}]) h_{j \rightarrow i}$$

$$h_i = GRU(m_i, h_i)$$



Training procedure

- Pretrained VGG-16 for region proposals and visual feature
- 512-dimensional vectors for state and messages
- For each image, 128 boxes are randomly selected from the top 2.000 proposed boxes
- For each image, 128 labeled relationships are randomly selected from the 8.128 possible object pairs
- For inference, only the top 50 boxes and all their pairs are considered

Visual Genome

- 100k images
- Top 150 object classes
(avg. 25 per image)
- Top 50 relationships
(avg. 6.2 per image)



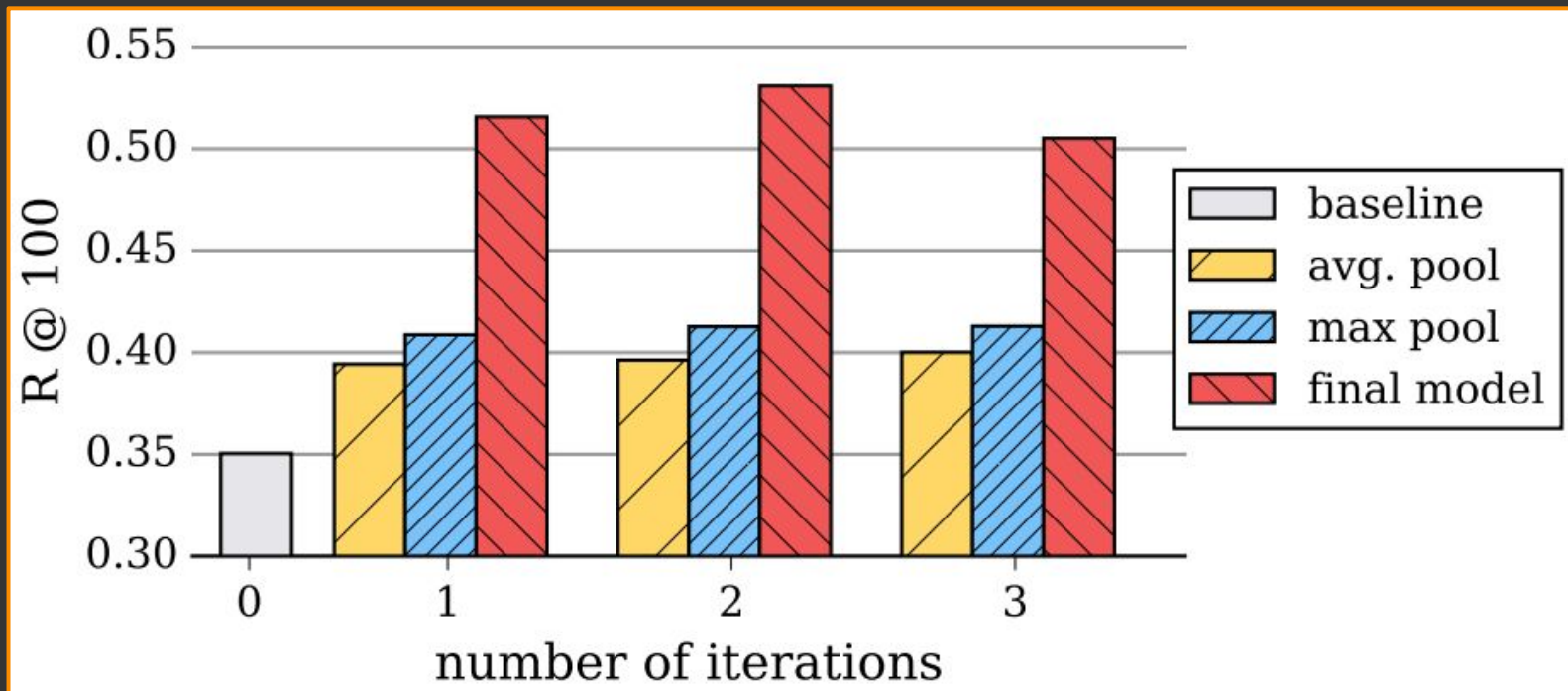
Note

This is a cleaned version of the [VG dataset](#), because the original annotations were of poor quality.

Visual Genome

		[26]	avg. pool	max pool	final
PREDCLS	R@50	27.88	32.39	34.33	44.75
	R@100	35.04	39.63	41.99	53.08
SGCLS	R@50	11.79	15.65	16.31	21.72
	R@100	14.11	18.27	18.70	24.38
SGGEN	R@50	0.32	2.70	3.03	3.44
	R@100	0.47	3.42	3.71	4.24

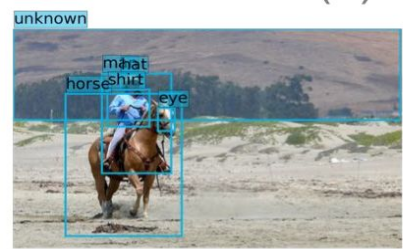
Visual Genome



Visual Genome

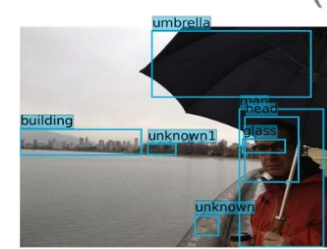
predicate	[26]	ours	predicate	[26]	ours
on	99.71	99.25	under	28.64	52.73
has	98.03	97.25	sitting on	31.74	50.17
in	80.38	88.30	standing on	44.44	61.90
of	82.47	96.75	in front of	26.09	59.63
wearing	98.47	98.23	attached to	8.45	29.58
near	85.16	96.81	at	54.08	70.41
with	31.85	88.10	hanging from	0.00	0.00
above	49.19	79.73	over	9.26	0.00
holding	61.50	80.67	for	12.20	31.71
behind	79.35	92.32	riding	72.43	89.72

unknown




```

    graph LR
      man -- wearing --> hat
      man -- wearing --> shirt
      man -- riding --> horse
      eye -- riding --> horse
      unknown -- on --> horse
  
```



```

    graph LR
      unknown1 -- on --> building
      umbrella -- holding --> man
      unknown -- wearing --> man
      glass -- wearing --> man
      head -- wearing --> man
      man -- holding --> holding
  
```

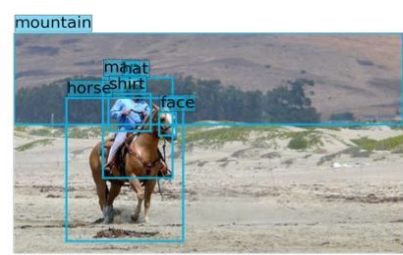


```

    graph LR
      vase -- on --> counter
      vase -- in --> flower
      counter -- on --> bear
      counter -- on --> vase
      bear -- on --> counter
      flower -- in --> vase
  
```

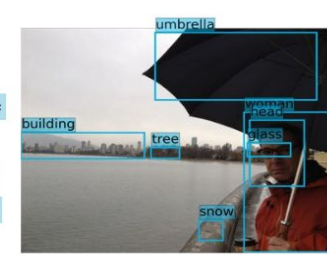
N=0 (baseline)

mountain




```

    graph LR
      face -- of --> horse
      mountain -- behind --> horse
      man -- riding --> horse
      man -- wearing --> hat
      man -- wearing --> shirt
  
```



```

    graph LR
      tree -- behind --> building
      umbrella -- on --> woman
      glass -- of --> woman
      head -- of --> woman
      snow -- on --> woman
      woman -- holding --> holding
  
```

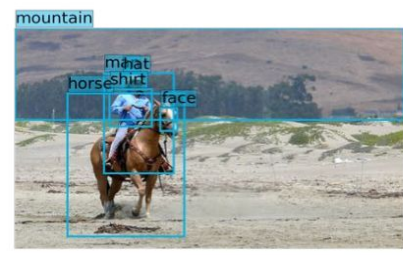


```

    graph LR
      vase -- on --> table
      vase -- in --> flower
      table -- at --> bear
      table -- in --> vase
      bear -- on --> table
      flower -- in --> vase
  
```


N=1

mountain




```

    graph LR
      face -- of --> horse
      mountain -- behind --> horse
      man -- riding --> horse
      man -- wearing --> hat
      man -- wearing --> shirt
  
```



```

    graph LR
      tree -- near --> building
      umbrella -- behind --> man
      head -- of --> man
      window -- on --> man
      glass -- on --> man
      man -- holding --> holding
  
```

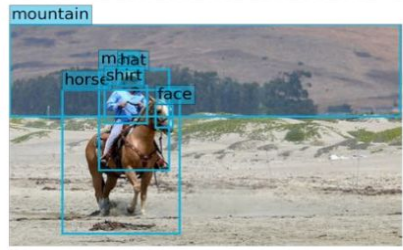


```

    graph LR
      vase -- on --> table
      vase -- with --> flower
      table -- under --> bear
      table -- under --> vase
      bear -- on --> table
      flower -- in --> vase
  
```

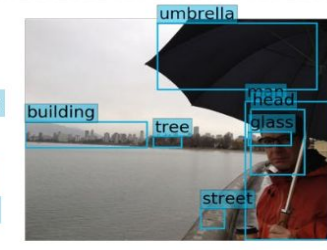
N=2

mountain




```

    graph LR
      face -- of --> horse
      mountain -- behind --> horse
      man -- on --> horse
      man -- has --> hat
      man -- has --> shirt
  
```



```

    graph LR
      tree -- in front of --> building
      umbrella -- over --> man
      head -- of --> man
      street -- on --> man
      glass -- of --> man
      man -- holding --> holding
  
```



```

    graph LR
      vase -- on --> table
      vase -- has --> flower
      table -- has --> bear
      table -- has --> vase
      bear -- on --> table
      flower -- in --> vase
  
```

ground truth

NYU Depth v2

- 1.449 RGB-D images
- 4 object classes
(floor, structure, furniture, prop)
- 3 support relationships
(behind, below, hidden)

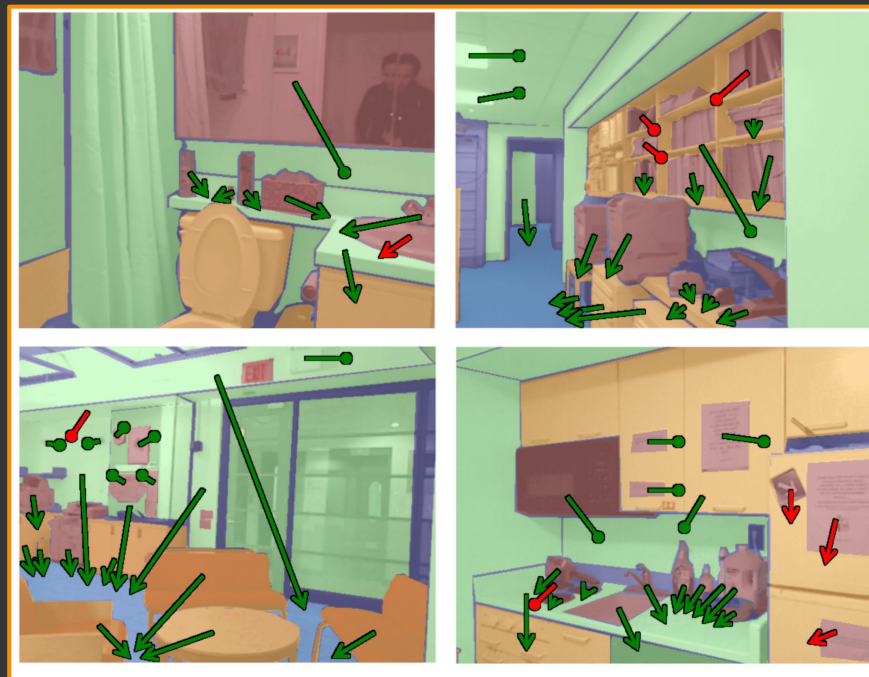


Notes

The depth channel is not used during the experiments.

Ground-truth object locations are provided as inputs, not predicted.

NYU Depth v2



NYU Depth v2

	Support Accuracy		PREDCLS	
	t-ag	t-aw	R@50	R@100
Silberman <i>et al.</i> [28]	75.9	72.6	-	-
Liao <i>et al.</i> [24]	88.4	82.1	-	-
Baseline [26]	87.7	85.3	34.1	50.3
Final model (ours)	91.2	89.0	41.8	55.5



Discussion