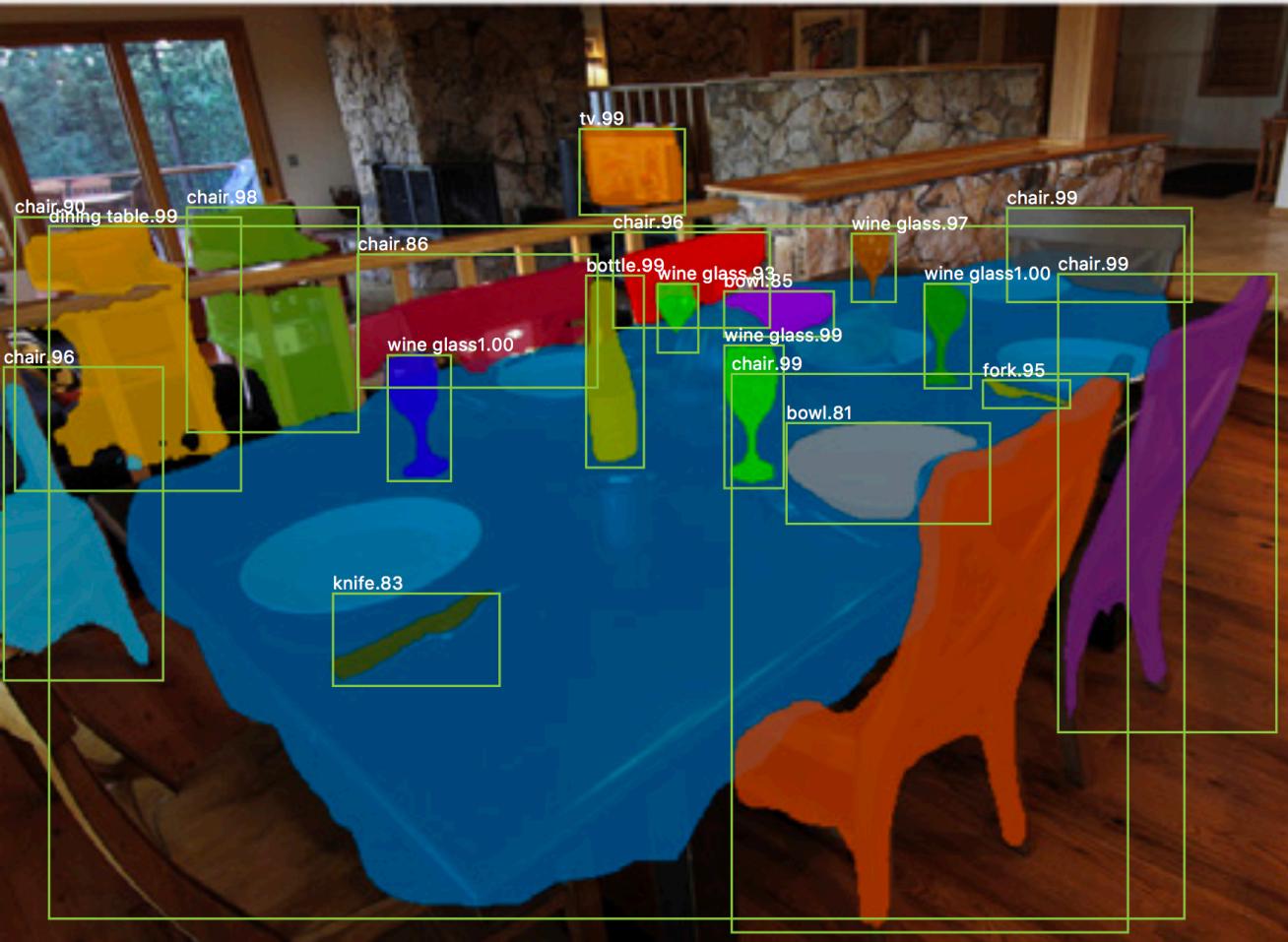


# Scene Graphs for Recognition and Generation

Justin Johnson

6/22/2018

# Object Detection: Image -> Objects



person

llama

person

# *Llama next to person*



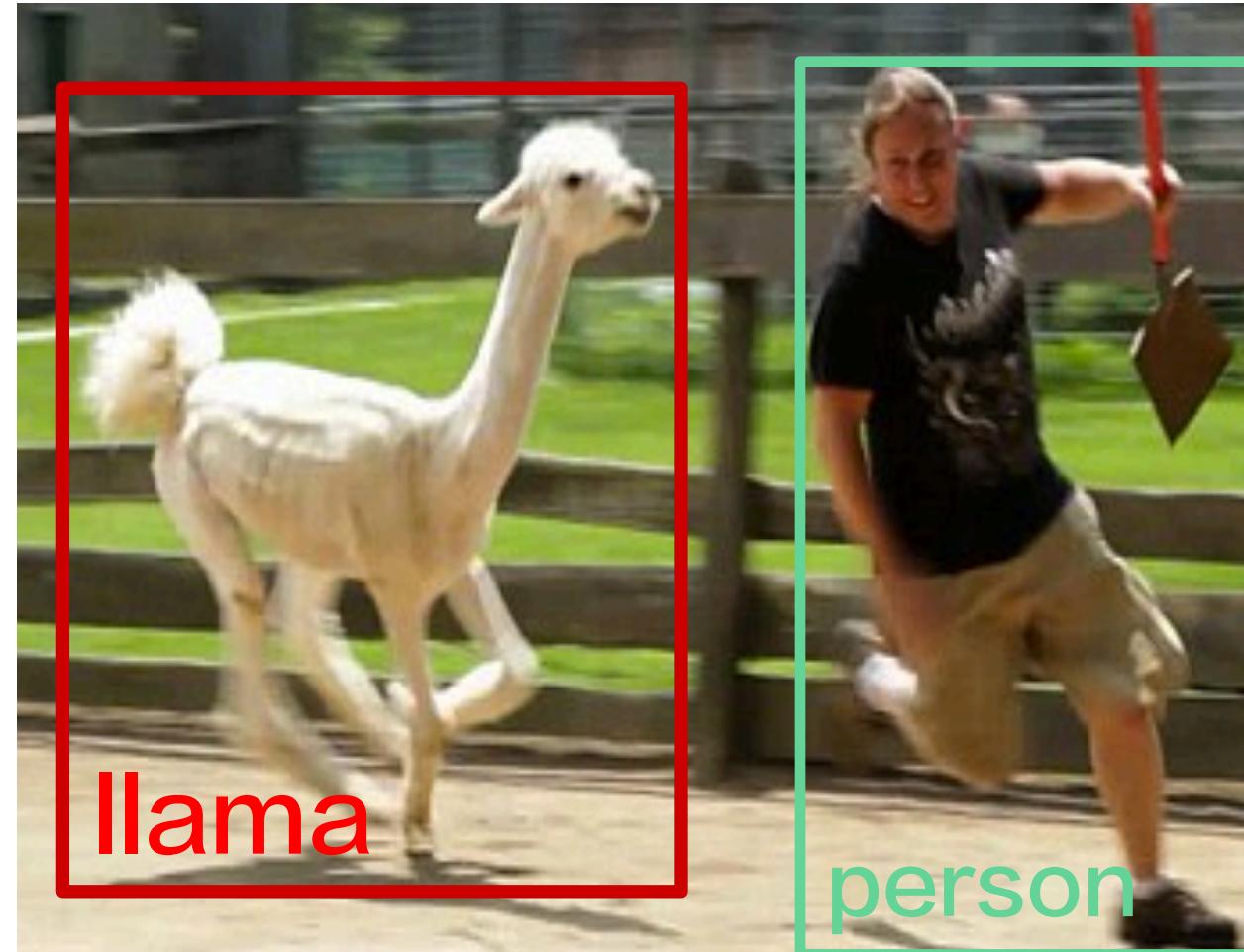
llama

person

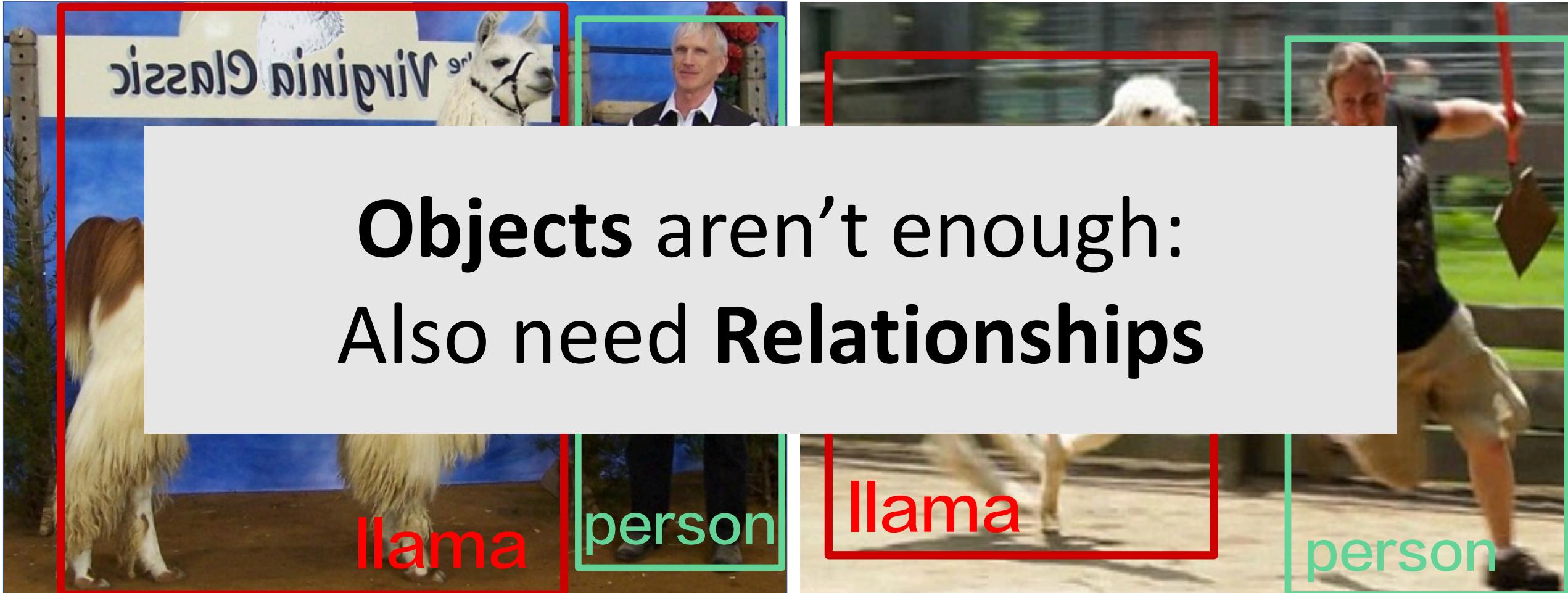
*Llama next to person*



*Llama chasing person*

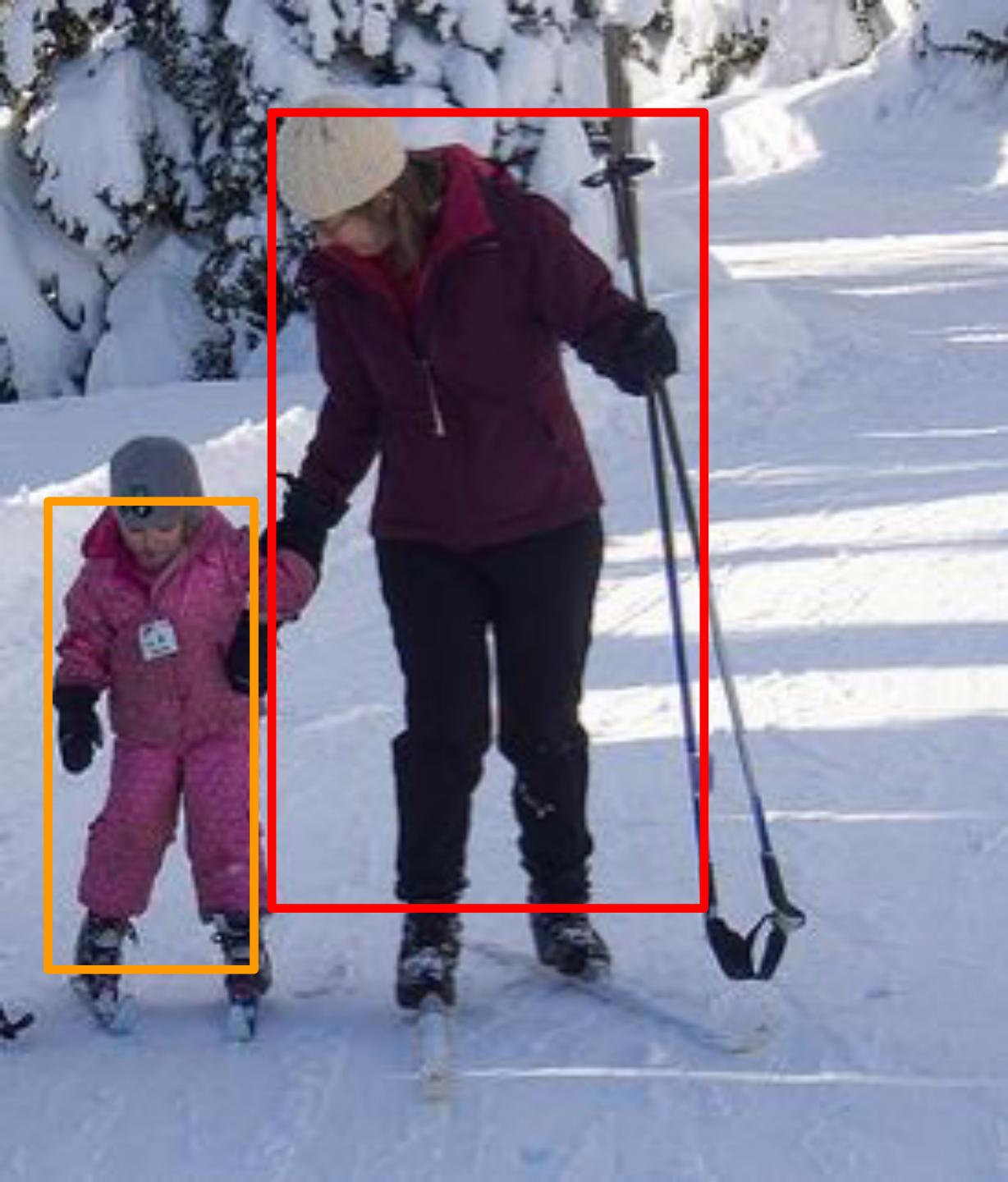


*Llama next to person*



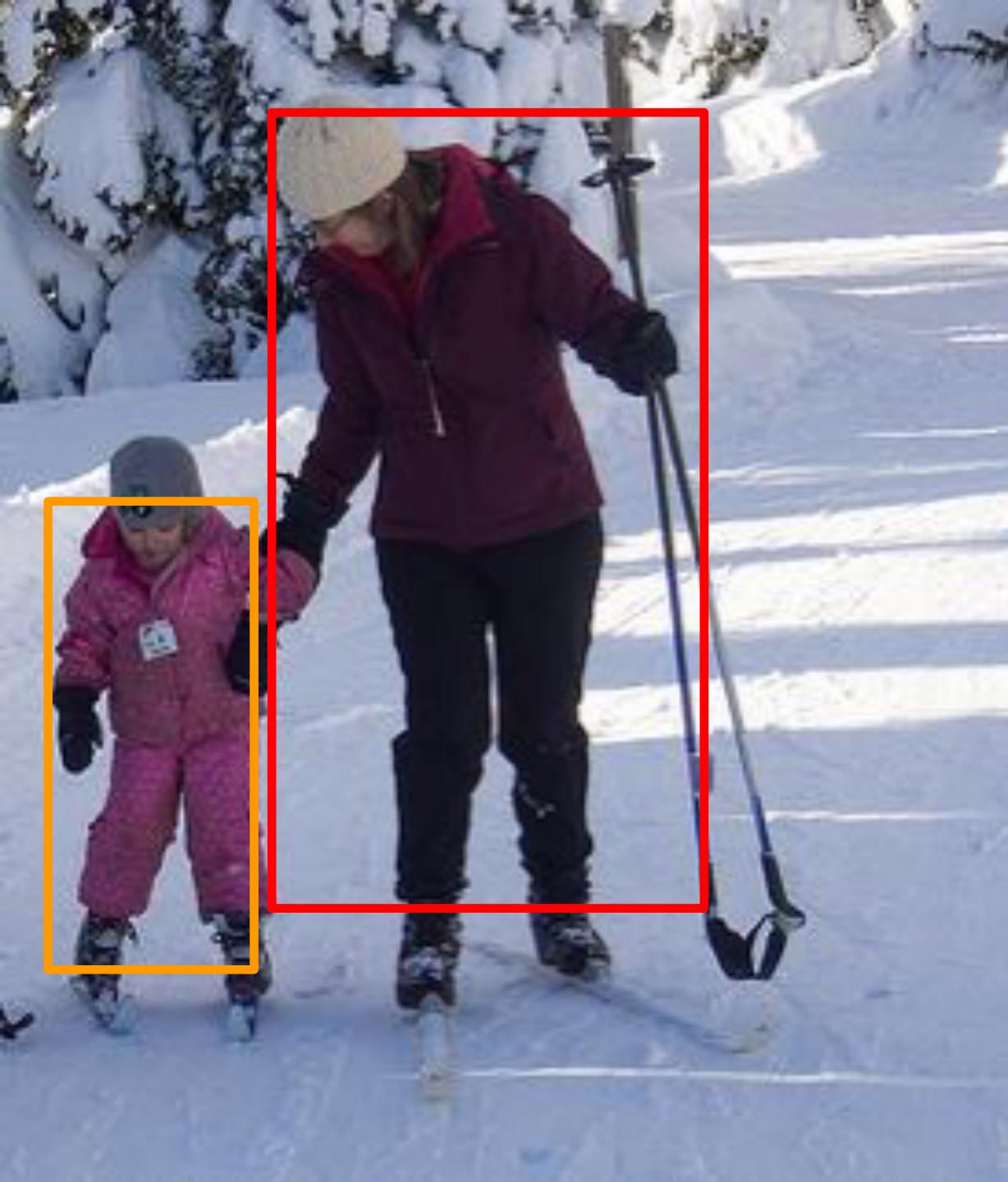


# Relationships



Relationships:  
spatial, comparative,  
verb, prepositional

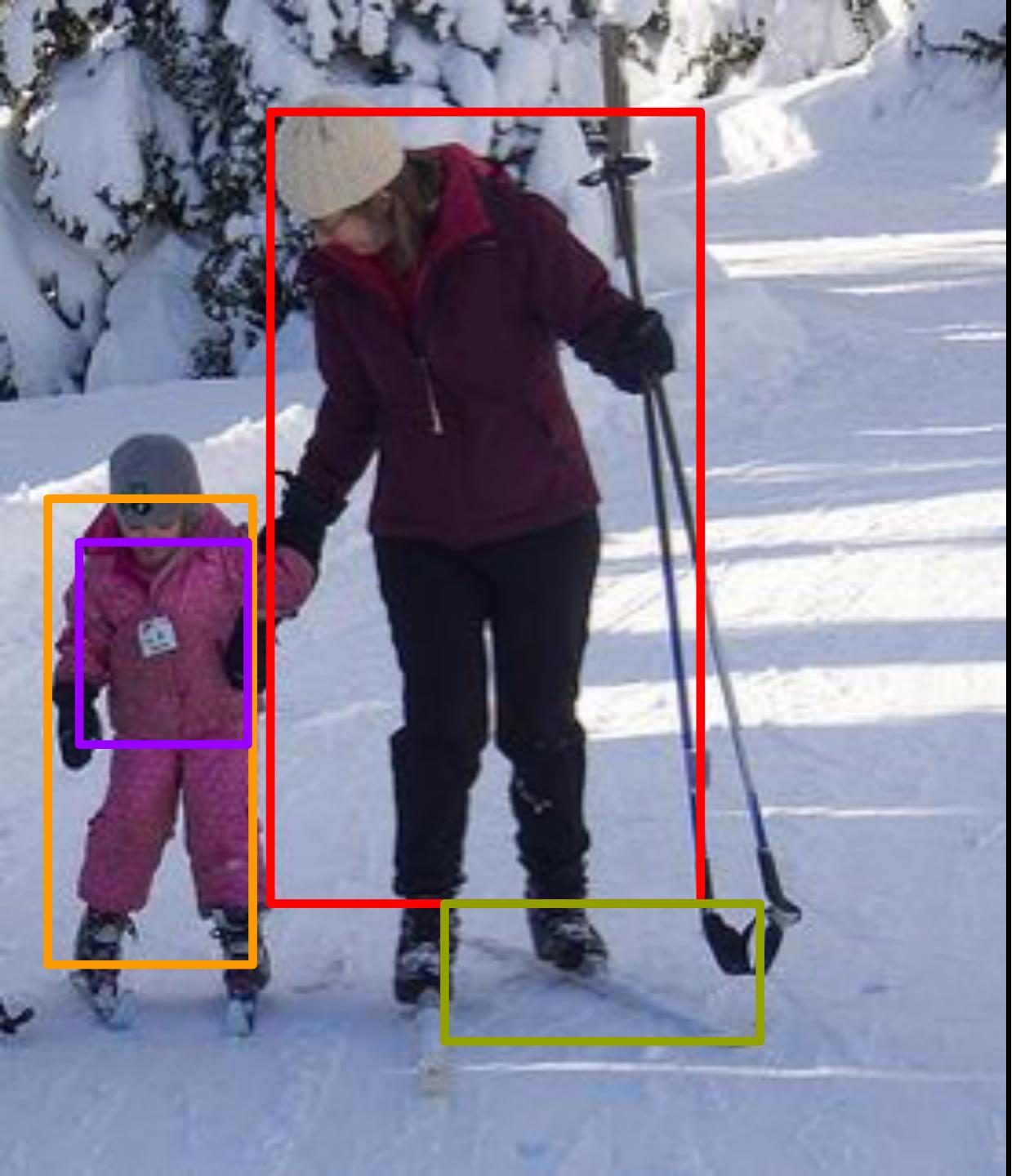
person → left of → person



# Relationships:

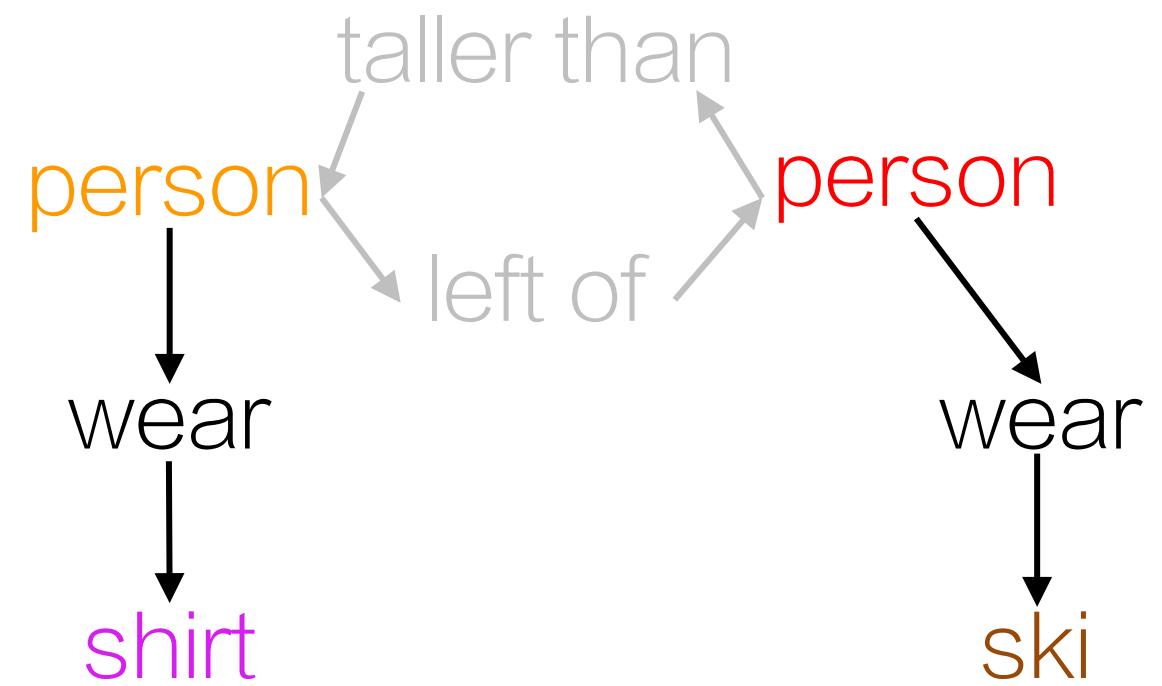
spatial, comparative,  
verb, prepositional

taller than  
person → person  
left of



# Relationships:

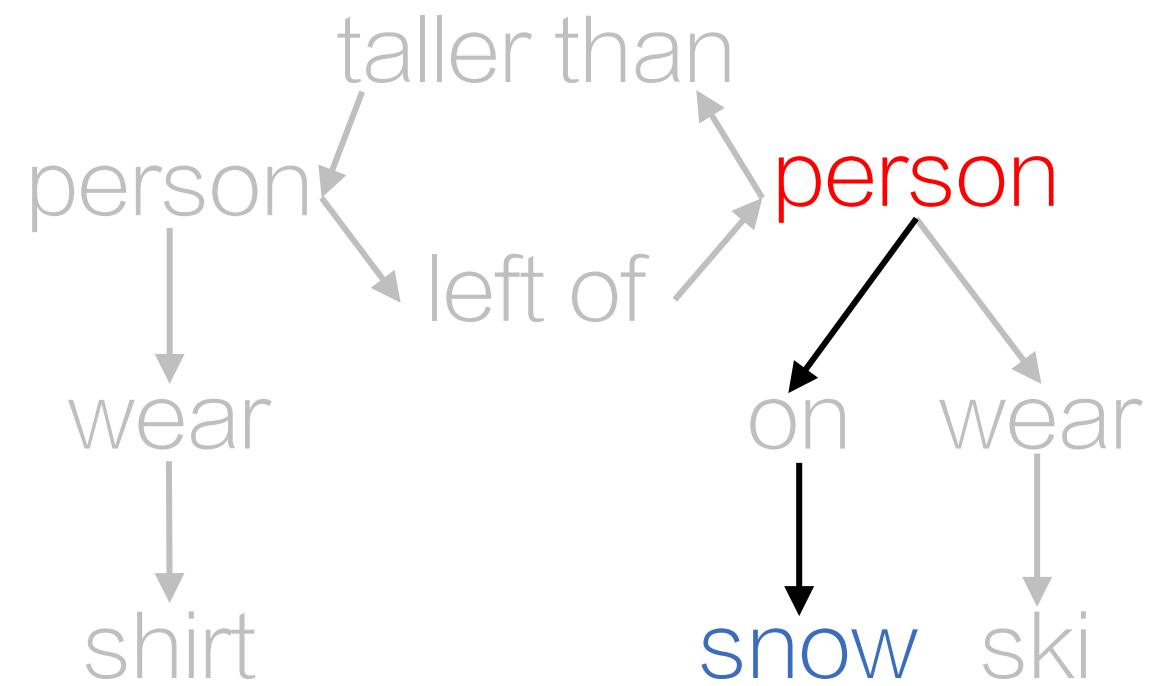
spatial, comparative,  
verb, prepositional

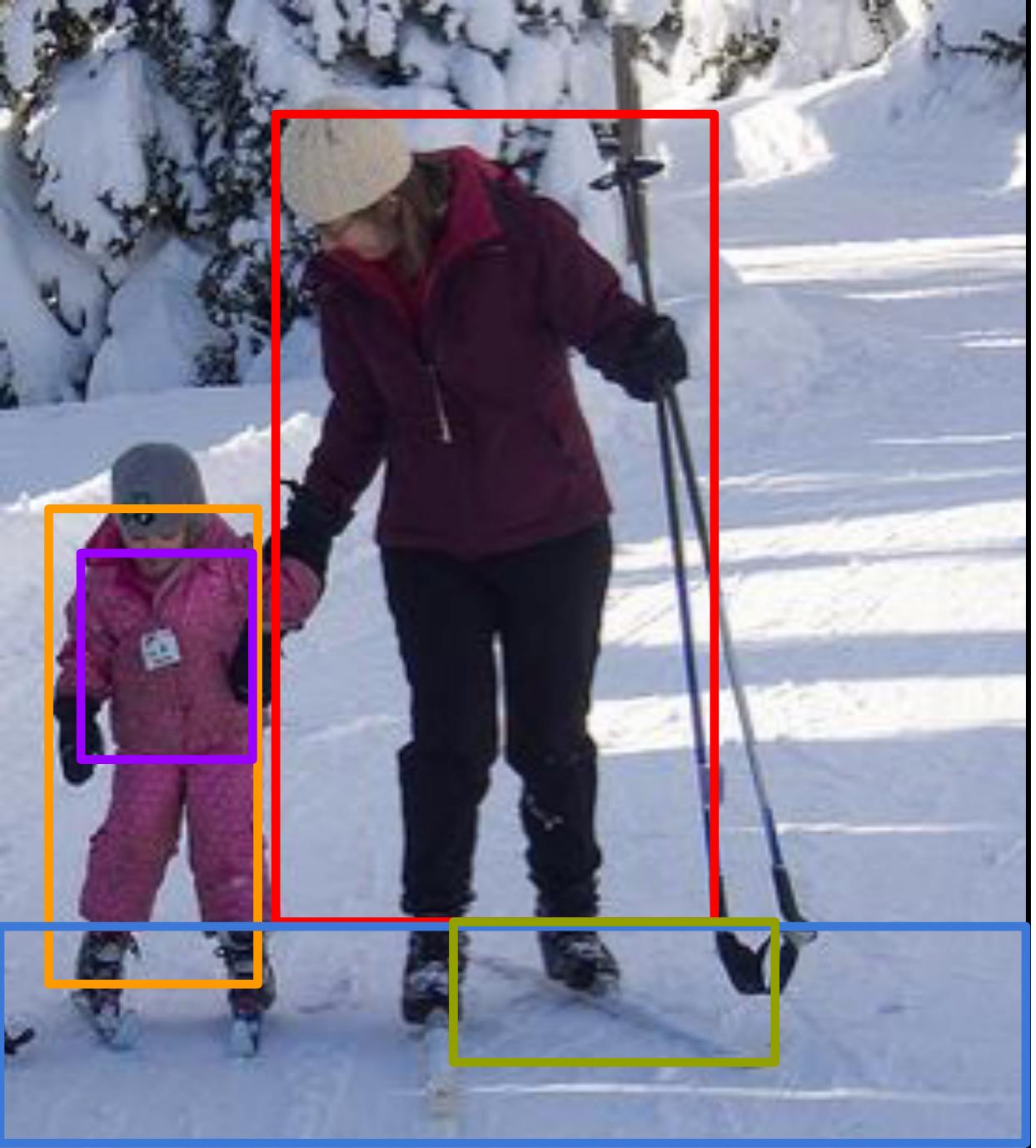




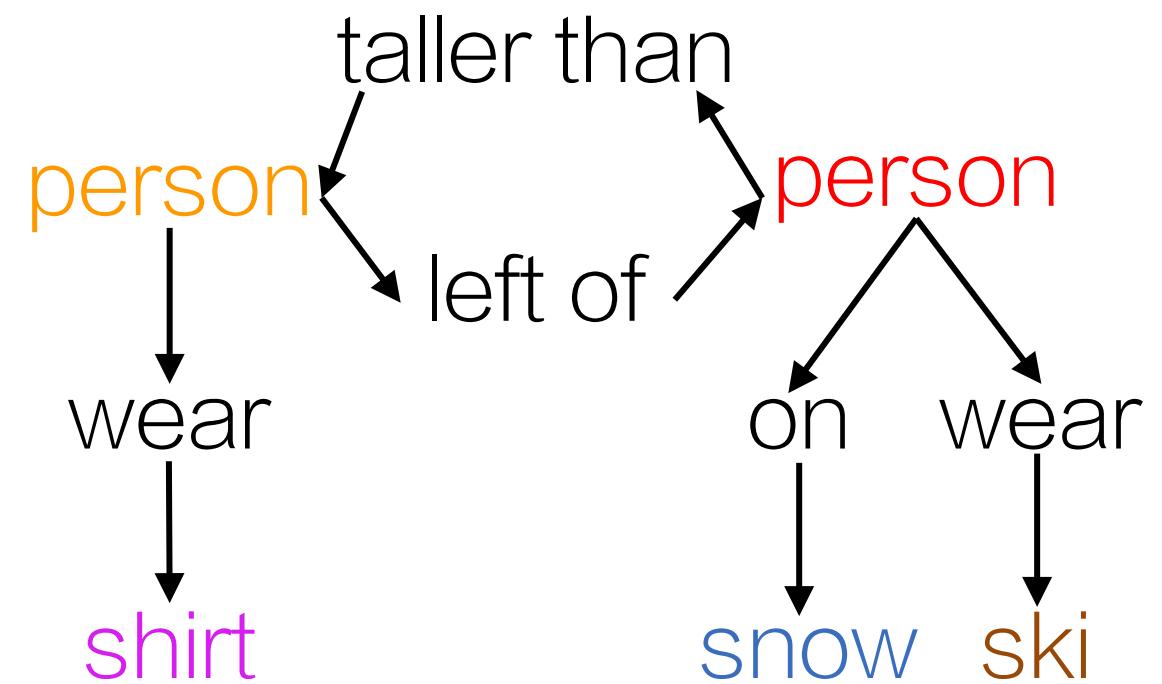
# Relationships:

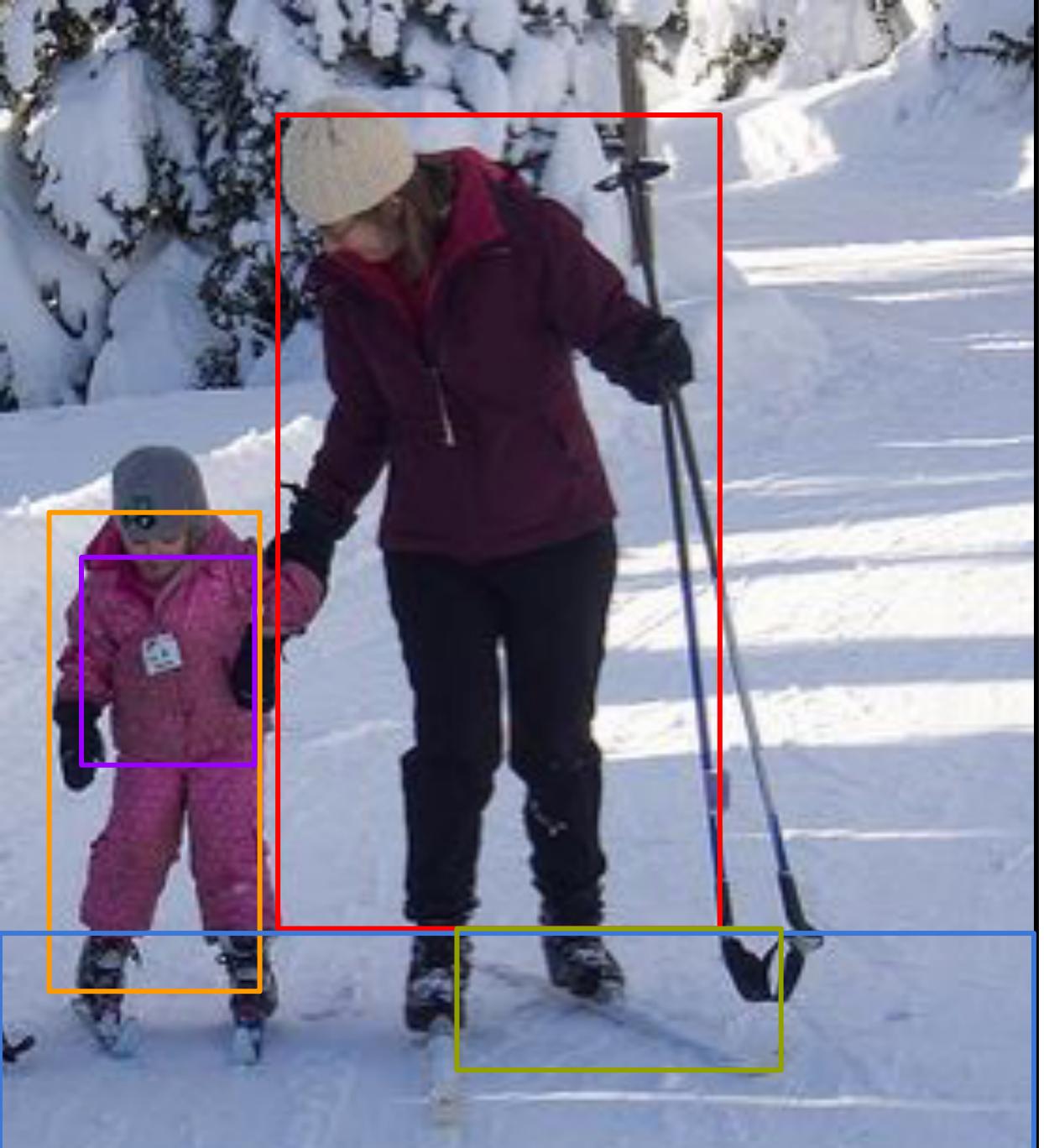
spatial, comparative,  
verb, prepositional



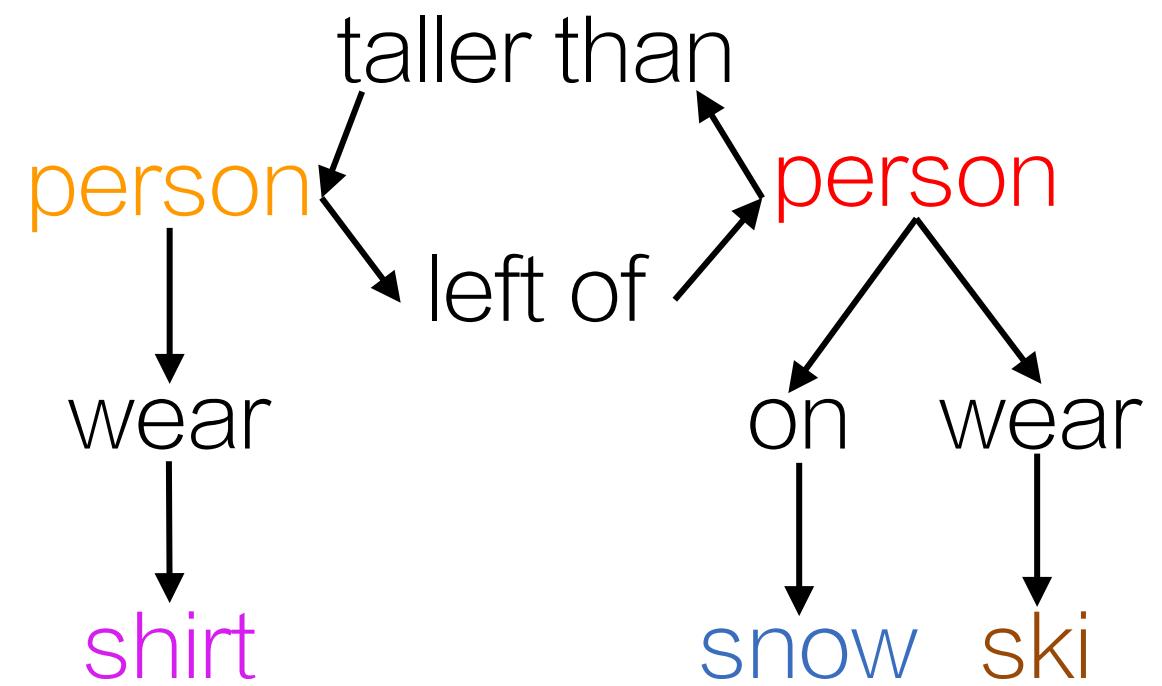


Relationships:  
spatial, comparative,  
verb, prepositional





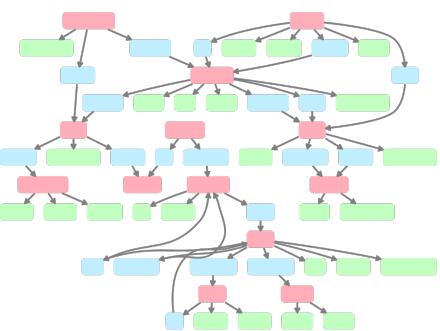
# Objects + Relationships = Scene Graphs



# Scene Graph Datasets

## Scene Graphs 5K

Johnson et al, CVPR 2015



- 5000 images
- 6745 object categories
- 1310 relationship types
- Long-tailed

## Visual Relationships

Lu et al, ECCV 2016



- 5000 images
- 100 object categories
- 70 relationship types
- Fully-annotated

## Visual Genome

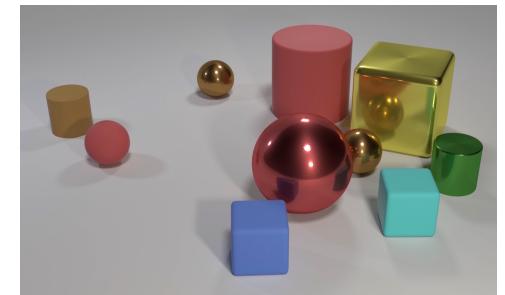
Krishna et al, IJCV 2017



- 108K images
- 33K object categories
- 42K relationship types
- Long-tailed

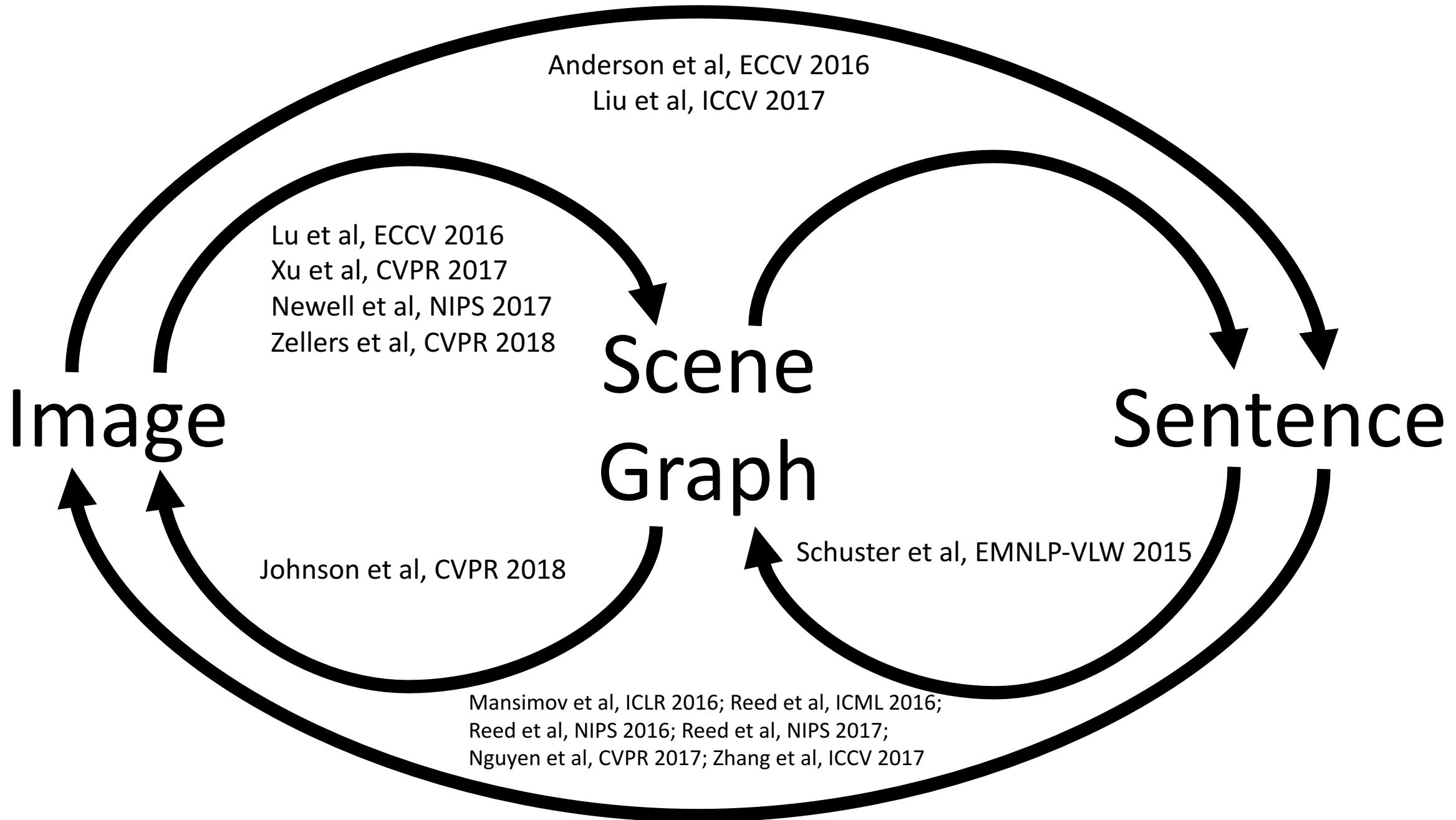
## CLEVR

Johnson et al, CVPR 2017



- 100K images
- 3 object categories
- 8 relationship types
- Fully-annotated

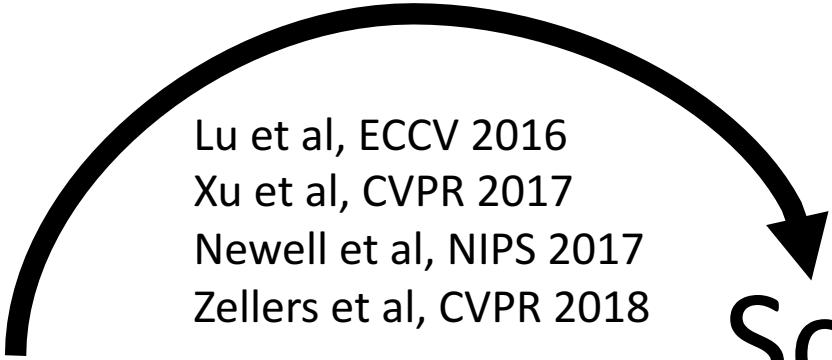
What can we do with scene graphs?



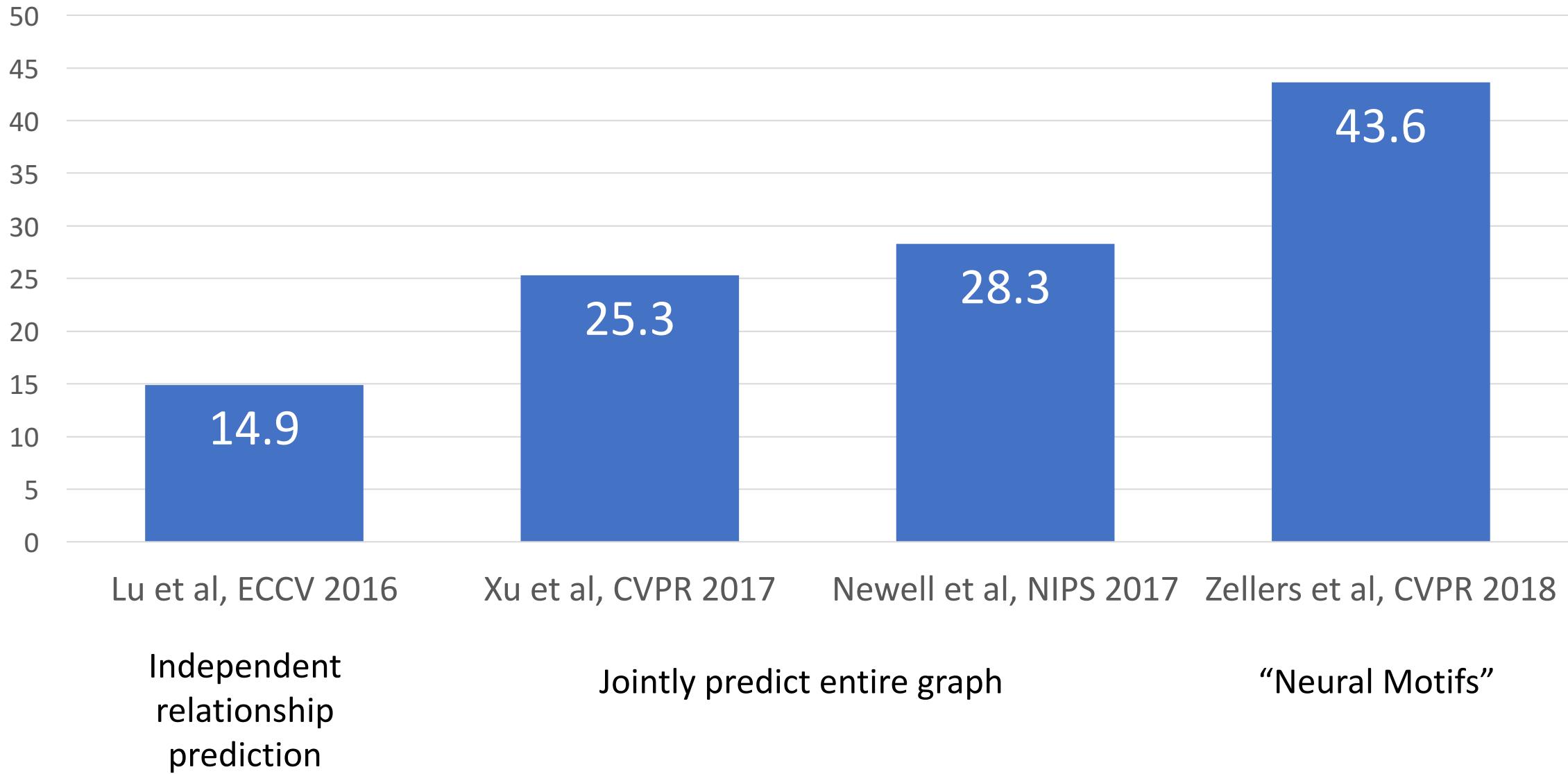
Image

Lu et al, ECCV 2016  
Xu et al, CVPR 2017  
Newell et al, NIPS 2017  
Zellers et al, CVPR 2018

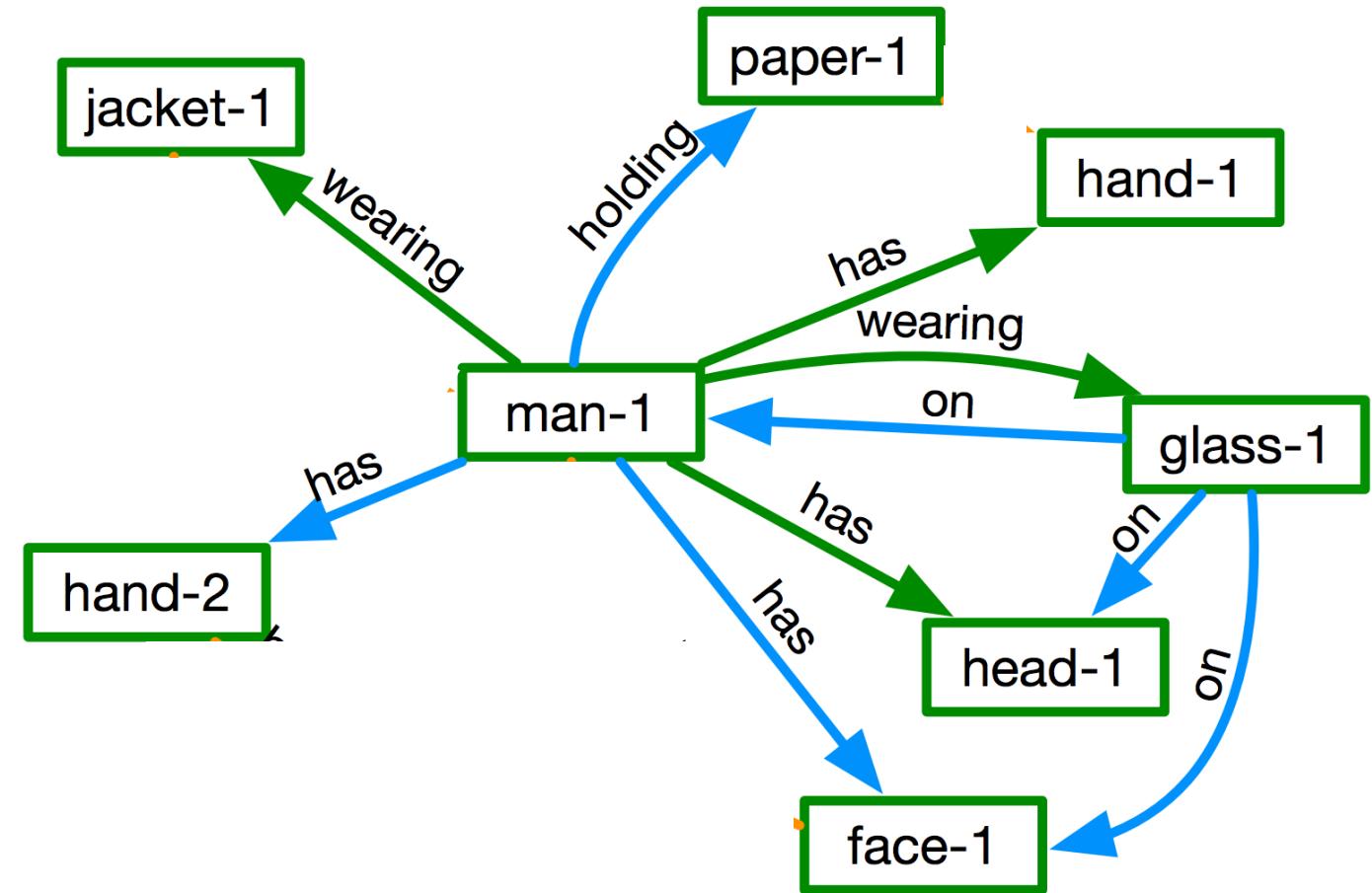
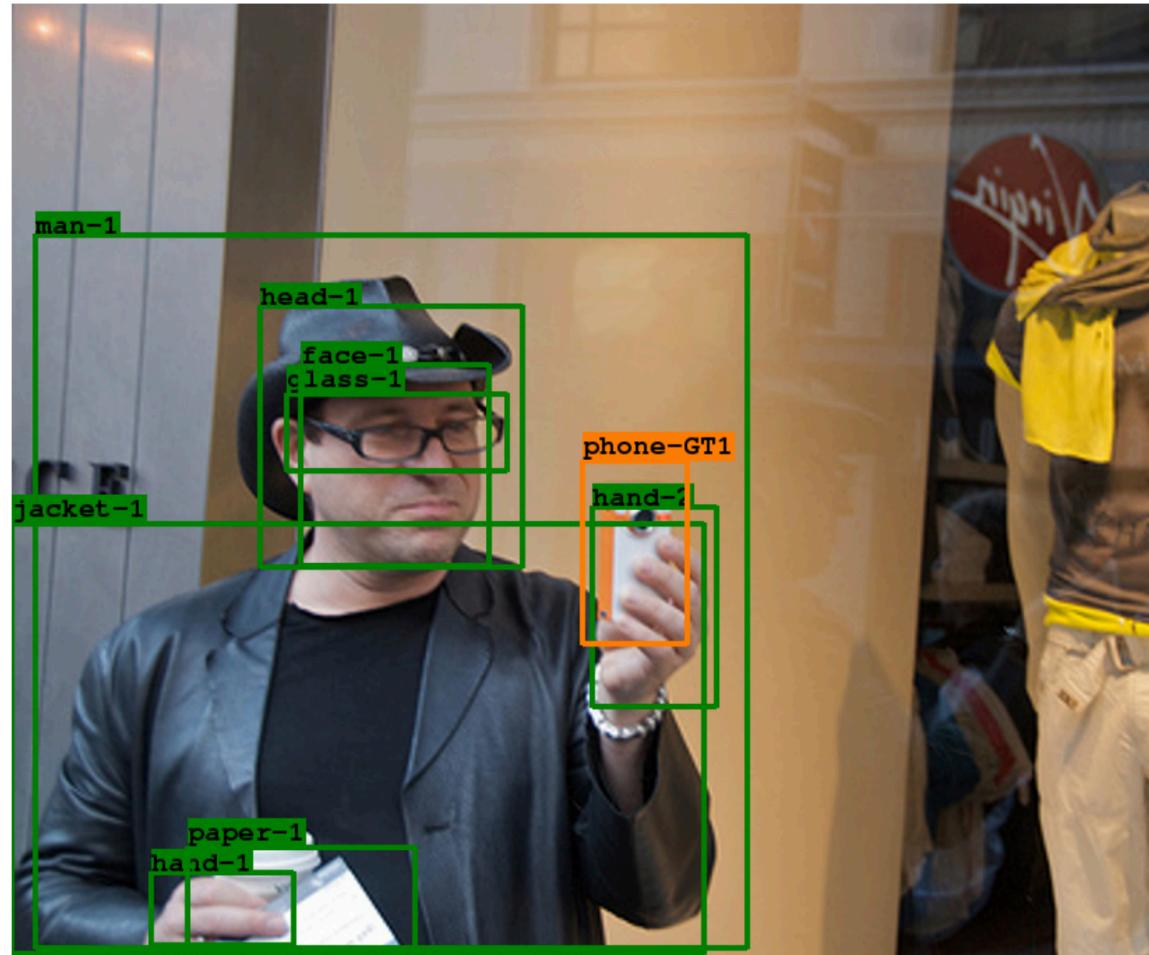
Scene  
Graph



# Image to Scene Graph



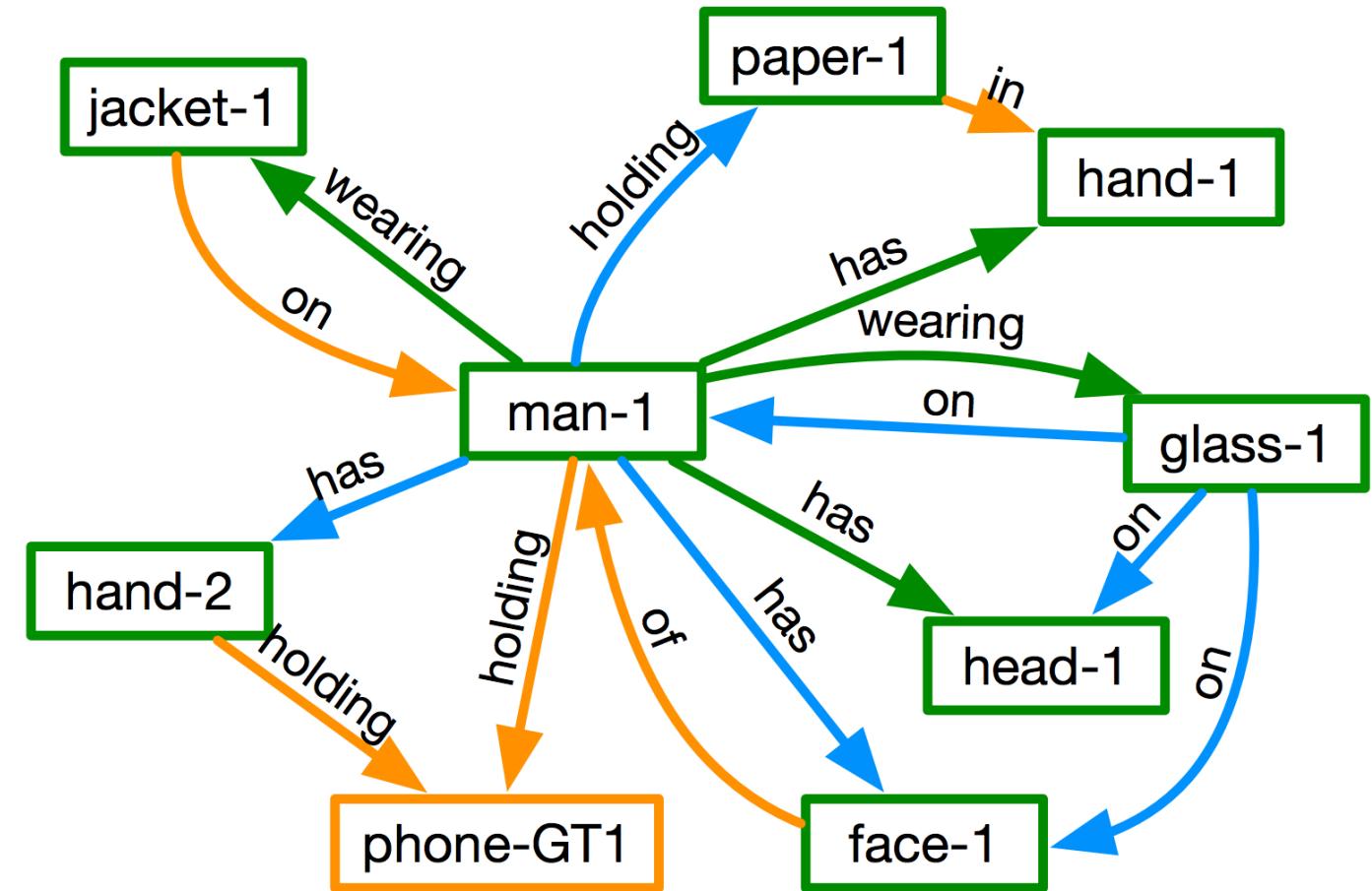
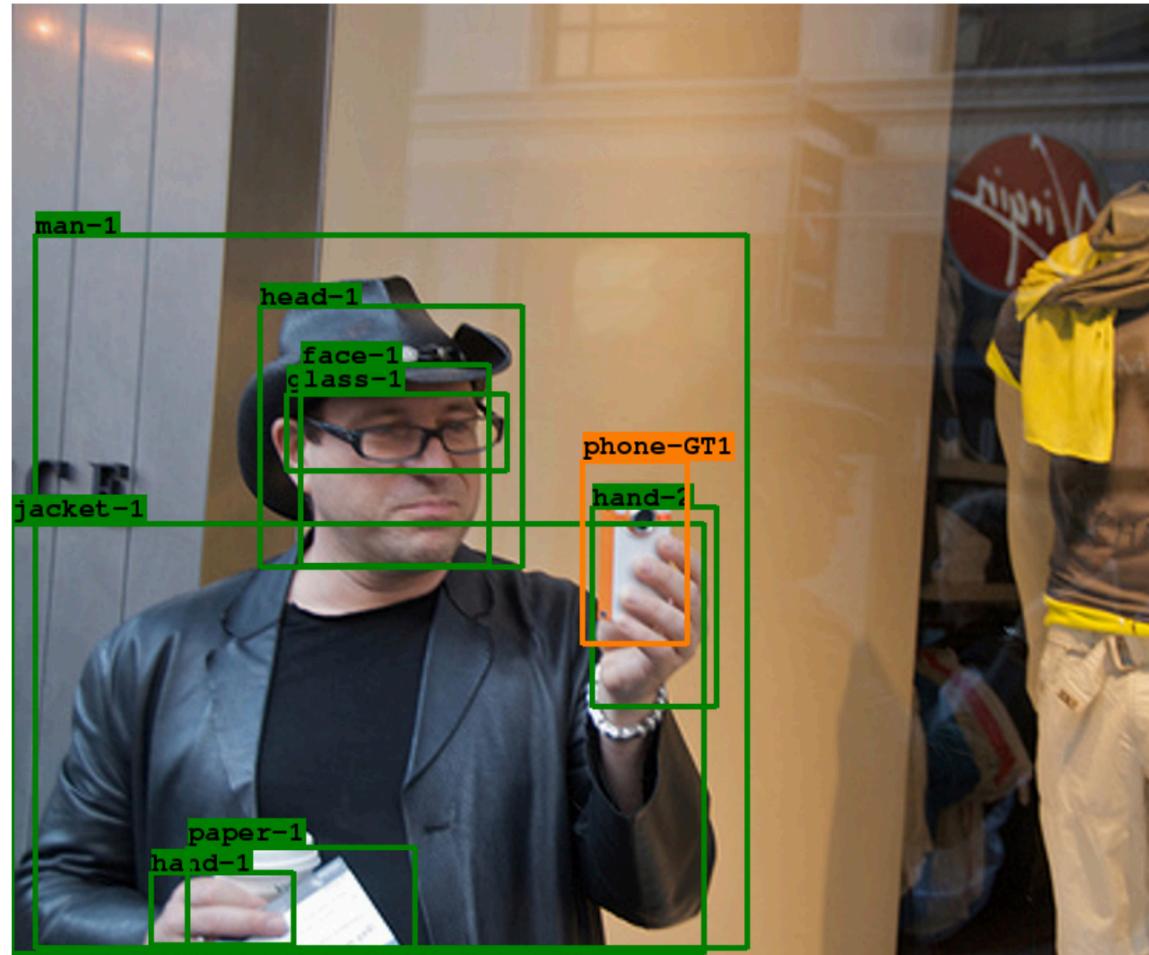
# Image to Scene Graph



True positive

False Positive

# Image to Scene Graph



True positive

False Positive

False Negative

Image

# Scene Graph

Sentence

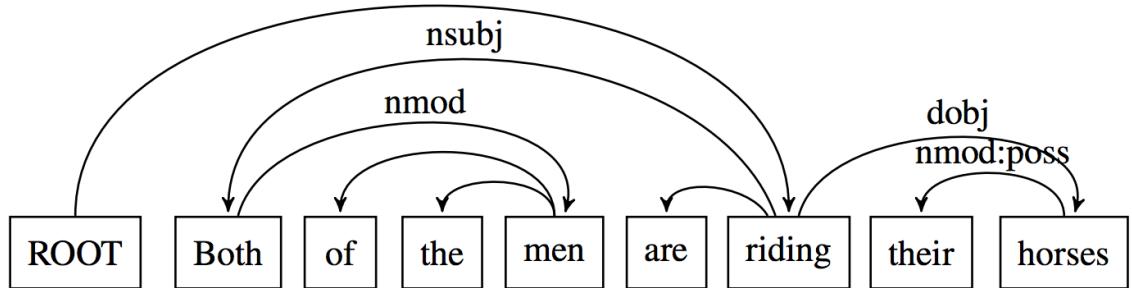
Lu et al, ECCV 2016  
Xu et al, CVPR 2017  
Newell et al, NIPS 2017  
Zellers et al, CVPR 2018

Schuster et al, EMNLP-VLW 2015

# Sentence to Scene Graph

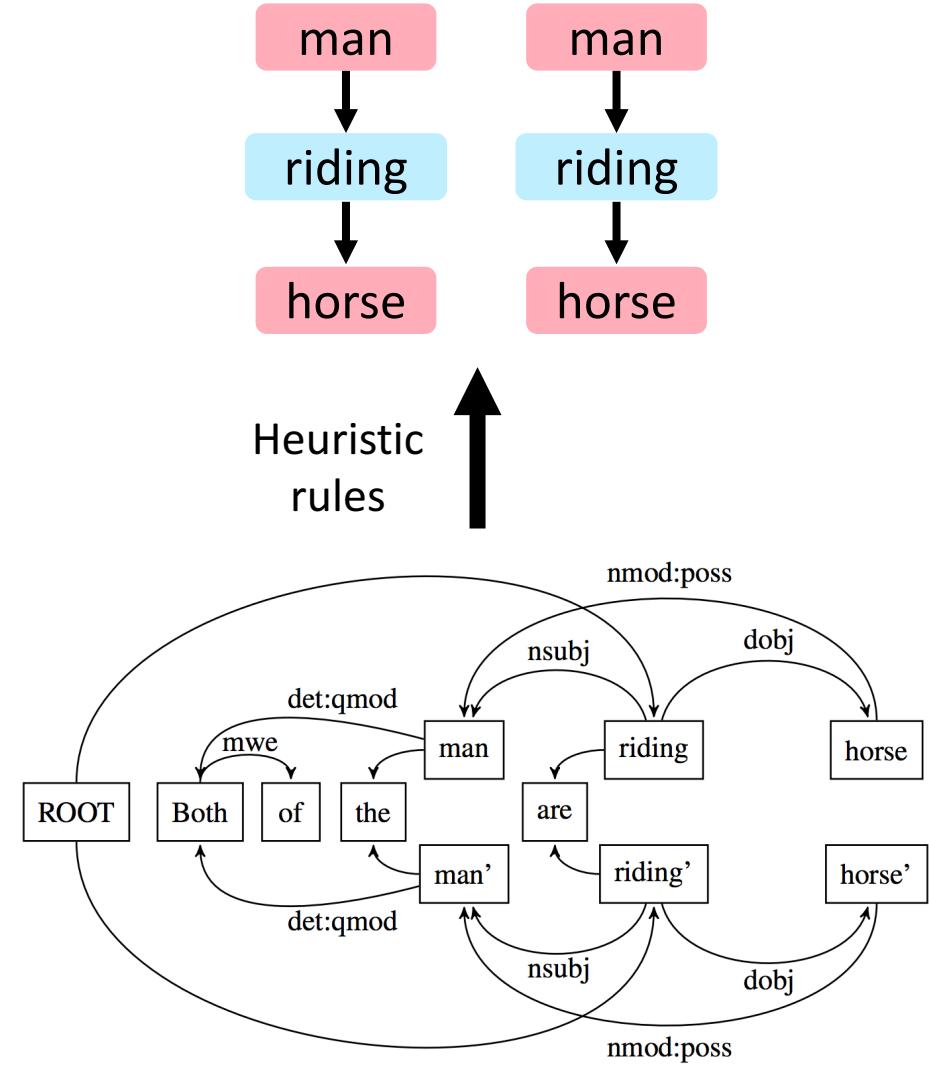
Both of the men are riding their horses

Dependency parser



Tree  
transforms

Heuristic  
rules



Image

# Scene Graph

Sentence

Anderson et al, ECCV 2016  
Liu et al, ICCV 2017

Lu et al, ECCV 2016  
Xu et al, CVPR 2017  
Newell et al, NIPS 2017  
Zellers et al, CVPR 2018

Schuster et al, EMNLP-VLW 2015

# Image Captioning (with Scene Graphs)

## SPICE: Semantic Propositional Image Caption Evaluation



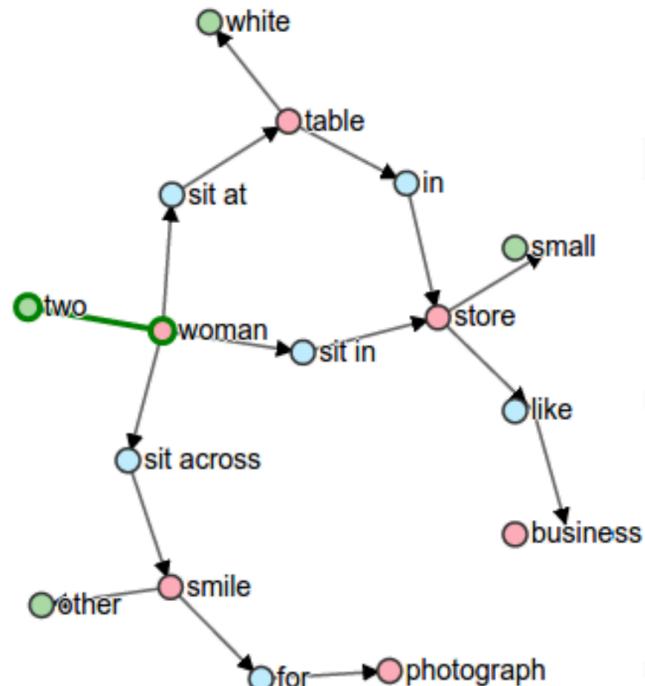
"two women are sitting at a white table"

"two women sit at a table in a small store"

"two women sit across each other at a table smile for the photograph"

"two women sitting in a small store like business"

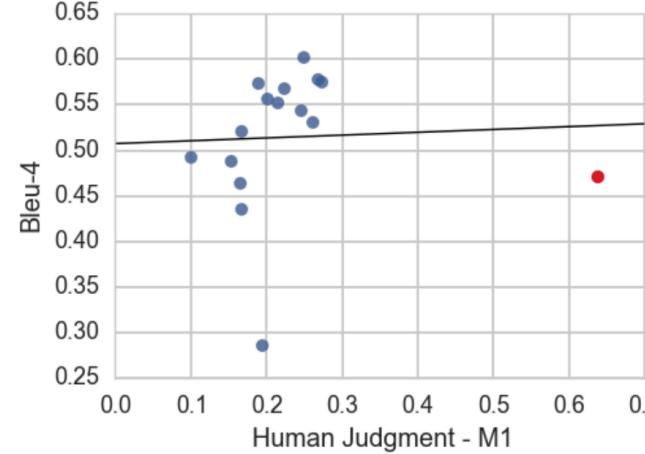
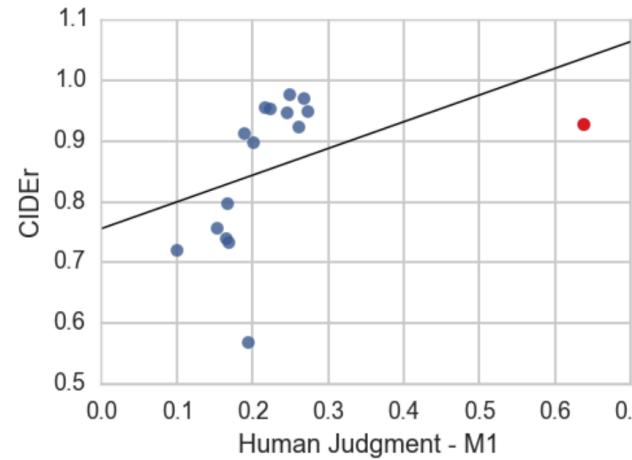
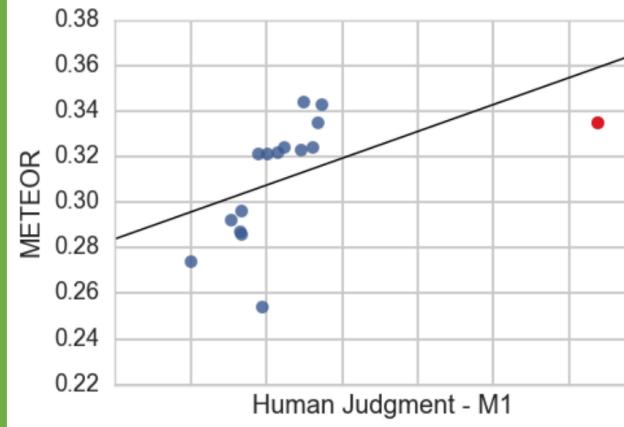
"two woman are sitting at a table"



Evaluate captioning  
by comparing **scene graphs** instead of  
comparing **sentences**

# Image Captioning (with Scene Graphs)

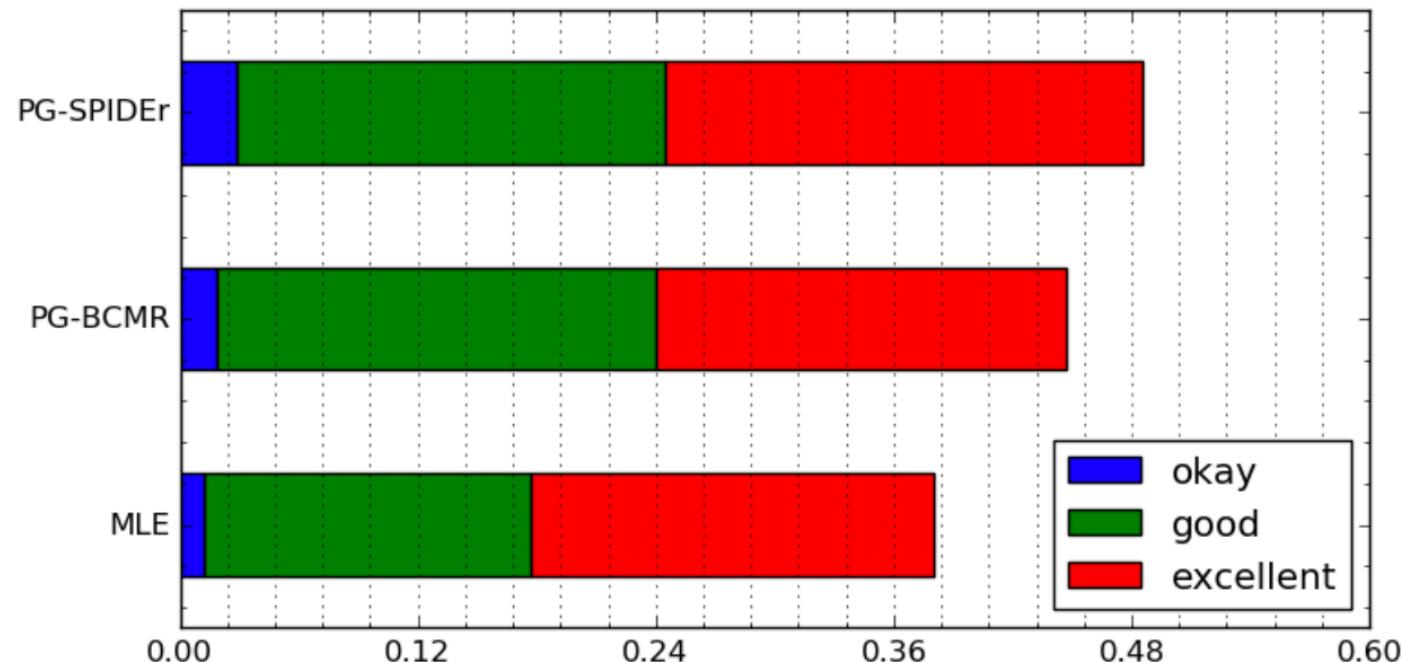
## SPICE: Semantic Propositional Image Caption Evaluation



Better correlation  
with human  
judgement

# Image Captioning (with Scene Graphs)

**Idea:** Train captioning model to maximize  
 $\text{SPIDEr} = \text{SPICE} + \text{CIDEr}$  using policy gradient



**Result:** People like the generated captions better!

Image

# Scene Graph

Sentence

Lu et al, ECCV 2016  
Xu et al, CVPR 2017  
Newell et al, NIPS 2017  
Zellers et al, CVPR 2018

Anderson et al, ECCV 2016  
Liu et al, ICCV 2017

Schuster et al, EMNLP-VLW 2015

**Image**

# Scene Graph

**Sentence**

Anderson et al, ECCV 2016  
Liu et al, ICCV 2017

Lu et al, ECCV 2016  
Xu et al, CVPR 2017  
Newell et al, NIPS 2017  
Zellers et al, CVPR 2018

Schuster et al, EMNLP-VLW 2015

Mansimov et al, ICLR 2016; Reed et al, ICML 2016;  
Reed et al, NIPS 2016; Reed et al, NIPS 2017;  
Nguyen et al, CVPR 2017; Zhang et al, ICCV 2017

# Text to Image

A red bird with black wings



# Text to Image

A red bird with black wings

**Prior work:**

Short description, single object



A man in a purple shirt throwing a red frisbee over a boy in a blue shirt to a man wearing jeans in the backyard

**My goal:**

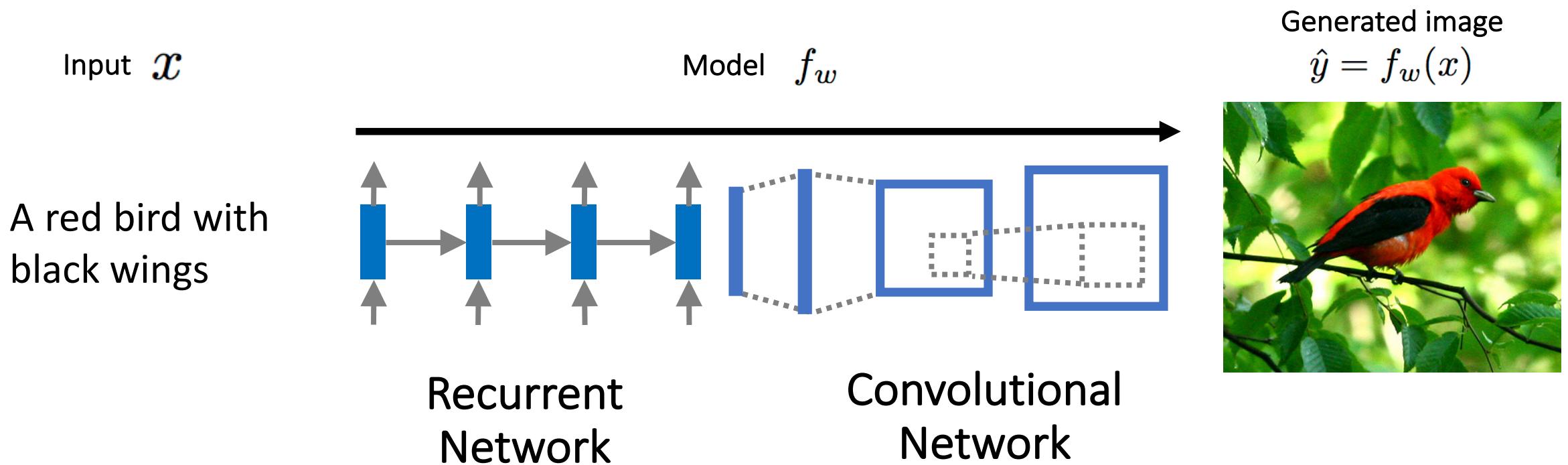
Complex description, many objects and relationships



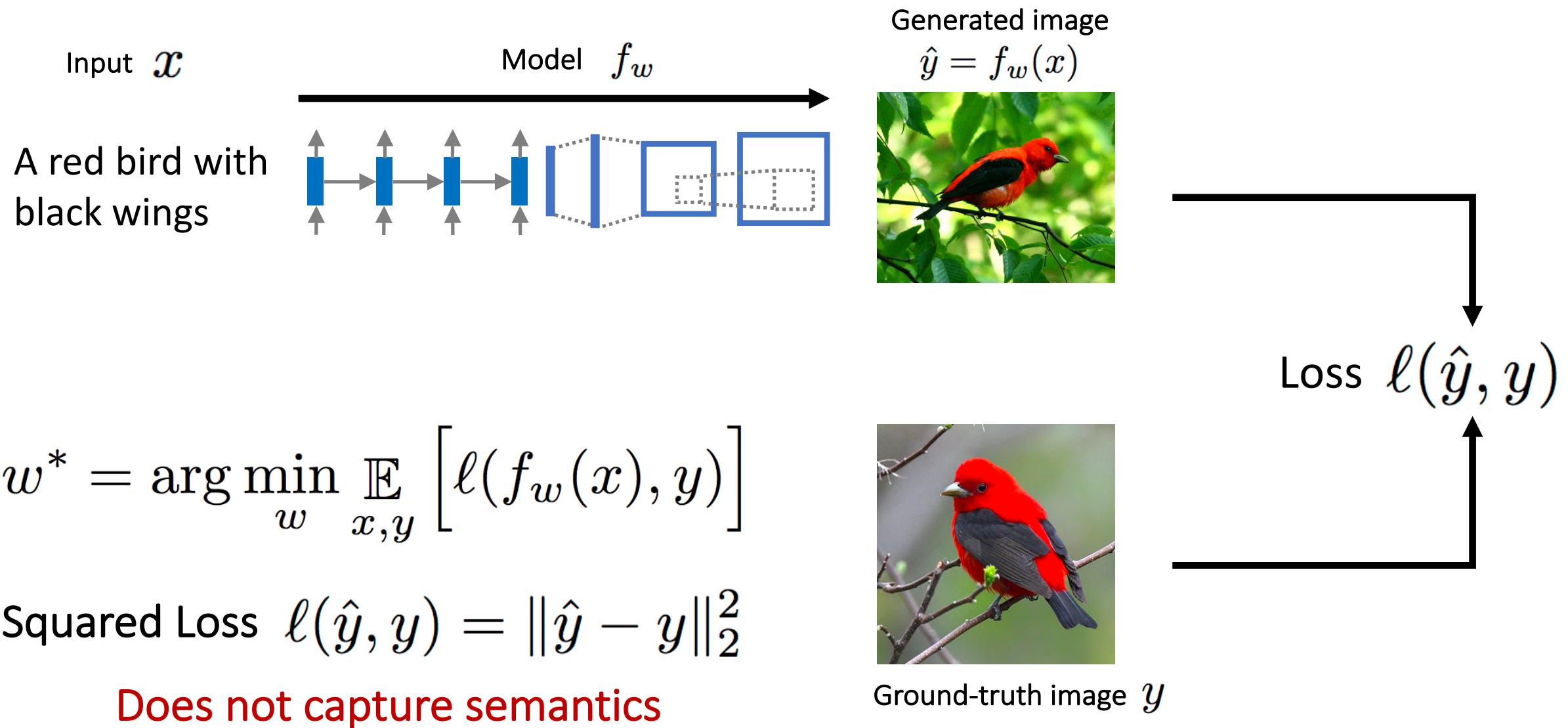
Johnson et al, CVPR 2018



# Text to Image

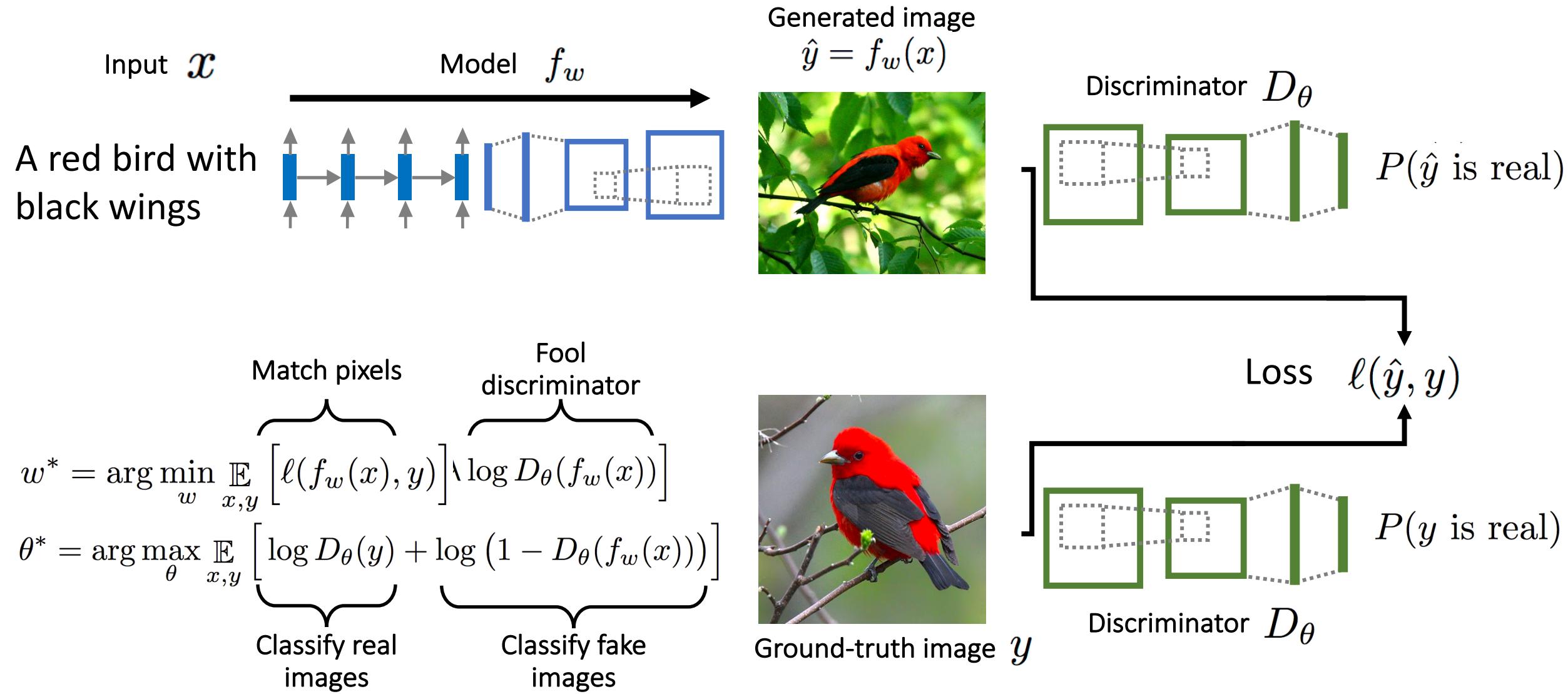


# Text to Image : Loss Function



# Text to Image: Adversarial Network

Goodfellow et al, NIPS 2014

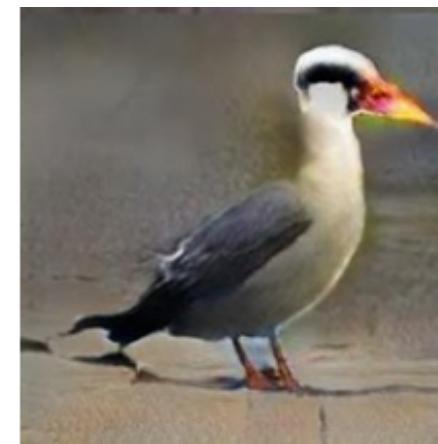


# Text to Image

This bird is blue with white  
and has a very short beak



A white bird with a black  
crown and yellow beak





# Text to Image: Complex Sentences

Schuster  
et al, 2015

## Natural Language

A sheep by another sheep standing on the grass with sky above and a boat in the ocean by a tree behind the sheep

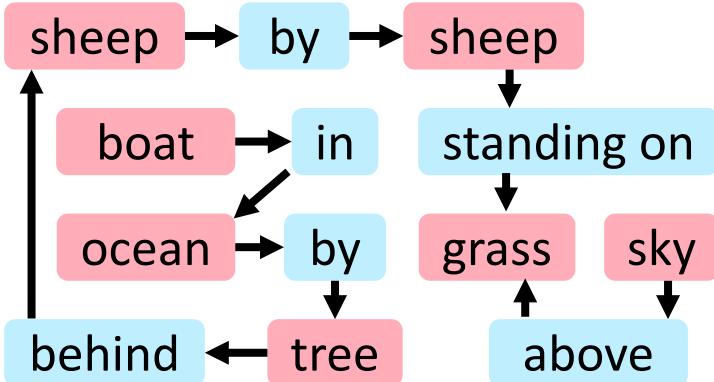
Better results: Zhang et al, Self-Attention Generative Adversarial Networks, arXiv 2018



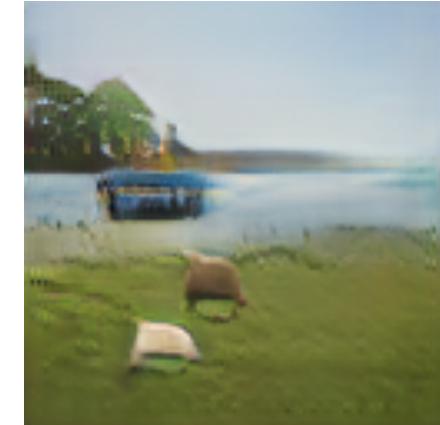
Zhang et al, ICCV 2017



## Scene Graph



Johnson, Gupta, and Fei-Fei, CVPR 2018



**Image**

# Scene Graph

**Sentence**

Anderson et al, ECCV 2016  
Liu et al, ICCV 2017

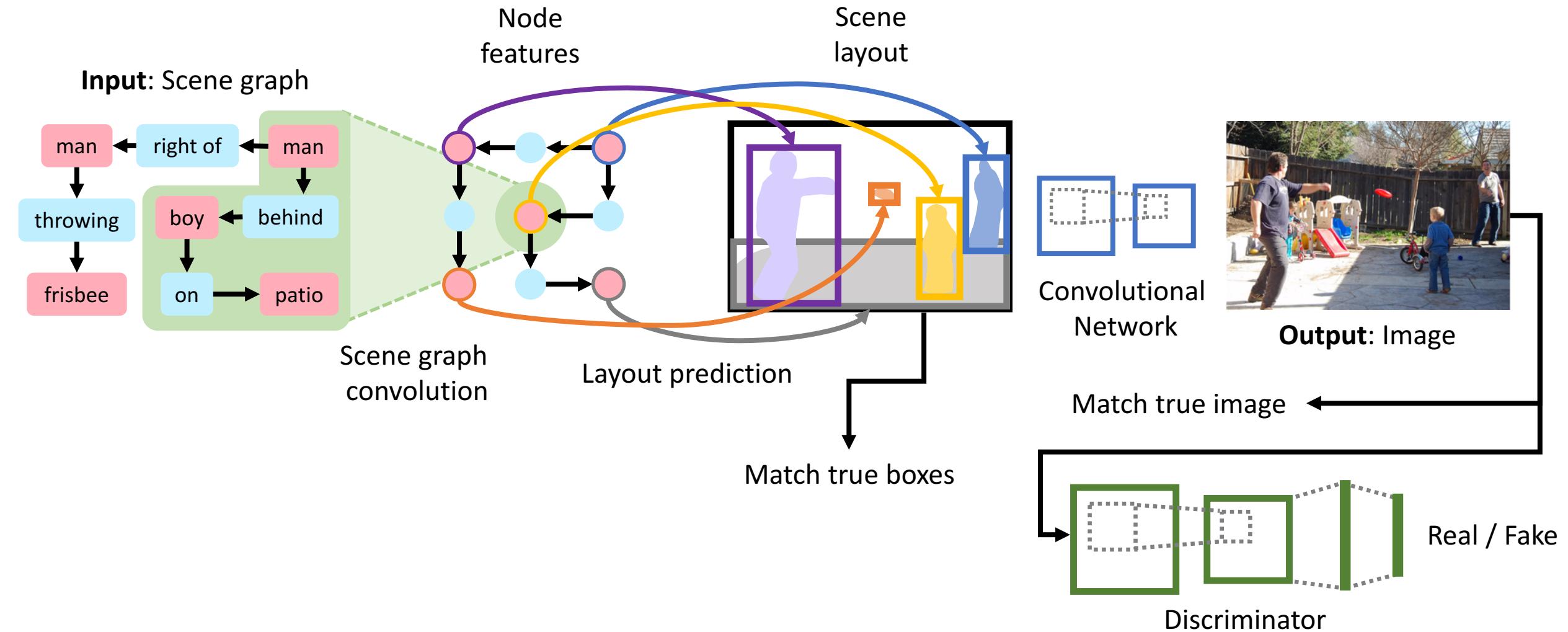
Lu et al, ECCV 2016  
Xu et al, CVPR 2017  
Newell et al, NIPS 2017  
Zellers et al, CVPR 2018

Johnson et al, CVPR 2018

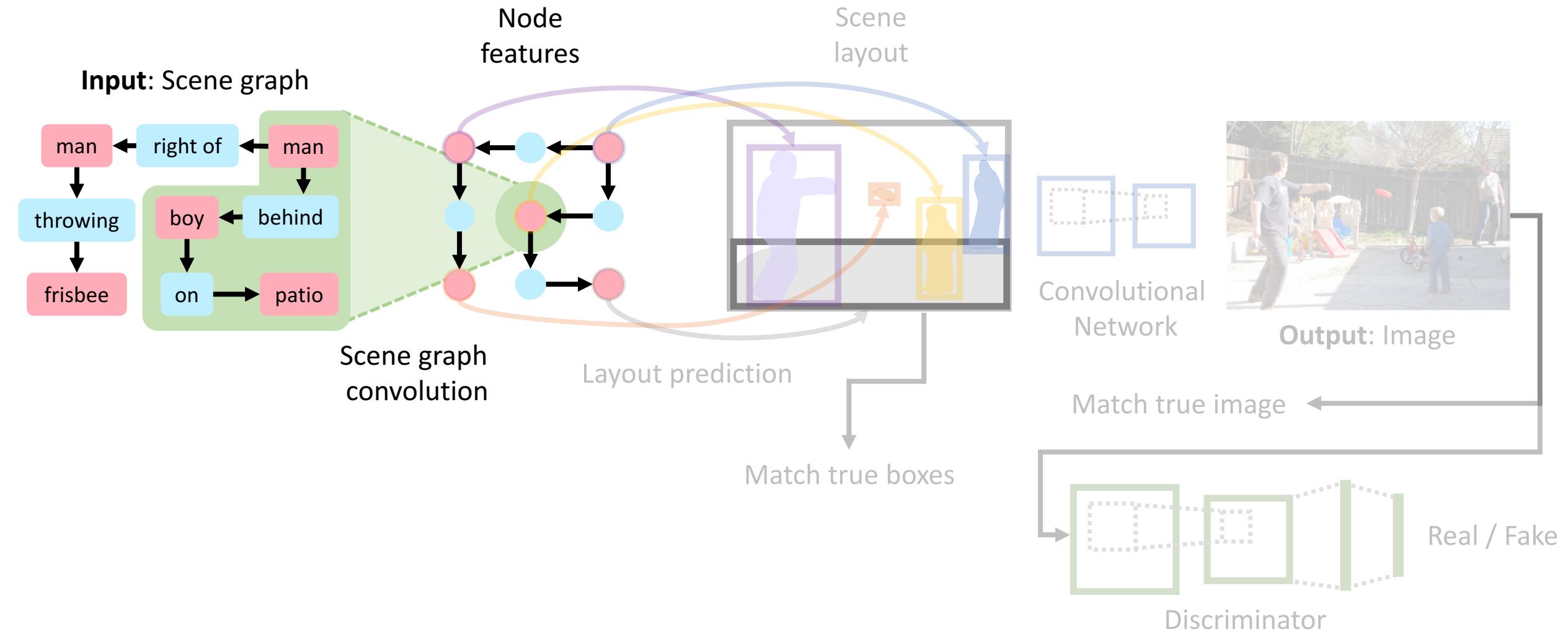
Schuster et al, EMNLP-VLW 2015

Mansimov et al, ICLR 2016; Reed et al, ICML 2016;  
Reed et al, NIPS 2016; Reed et al, NIPS 2017;  
Nguyen et al, CVPR 2017; Zhang et al, ICCV 2017

# Scene Graph to Image: Model

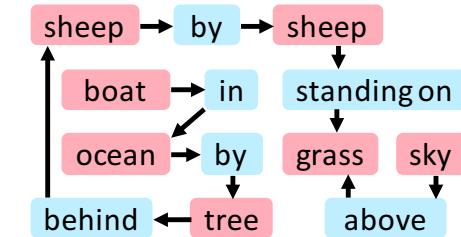


# Scene Graph to Image: Model

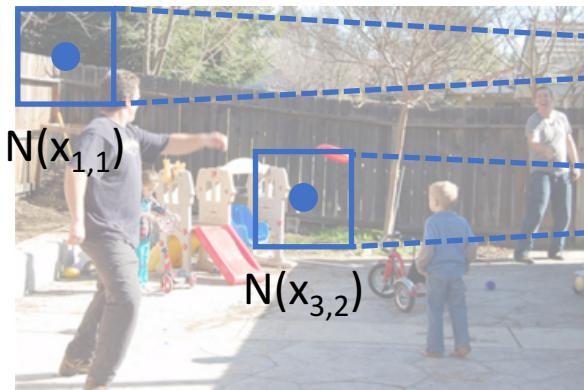


# Processing different input types

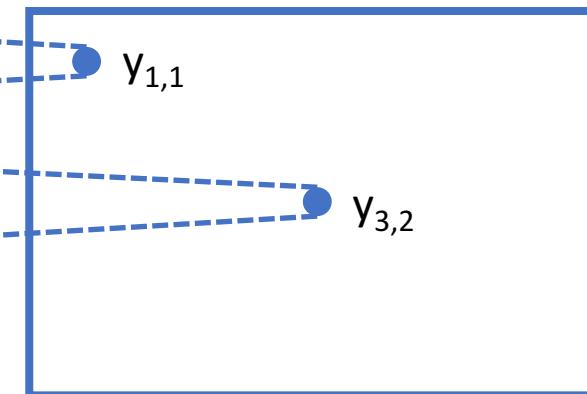
What about  
scene graphs?



## Images: Convolutional Networks



**Input:** Grid of vectors



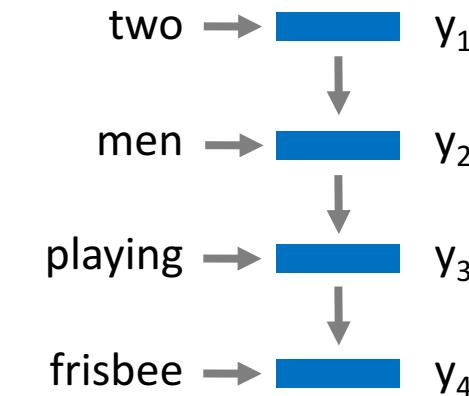
**Output:** Grid of vectors

Apply same function to all  
neighborhoods of input:

$$y_i = \boxed{f_W}(\mathcal{N}(x_i))$$

## Sequences: Recurrent Networks

**Input:**  
Sequence  
of vectors



**Output:**  
Sequence  
of vectors

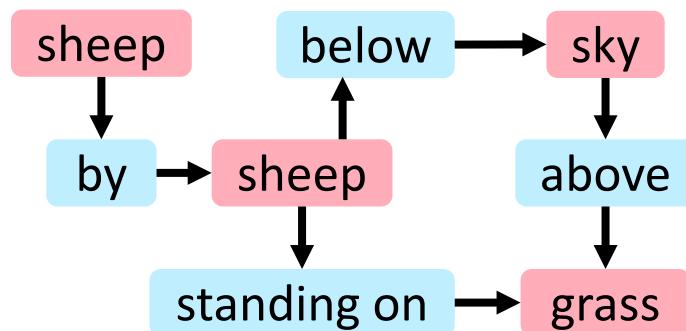
Apply same function to  
all timesteps of input:

$$y_t = \boxed{f_W}(x_t, y_{t-1})$$

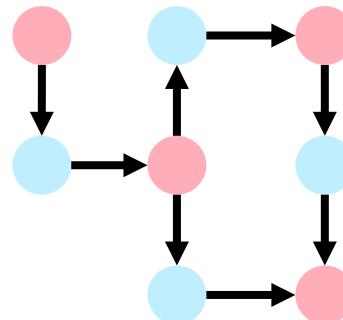
"Fully-connected" neural network: Matrix multiply and nonlinearity

# Graph Convolution Networks

**Input:** Graph of vectors



**Output:** Graph of vectors

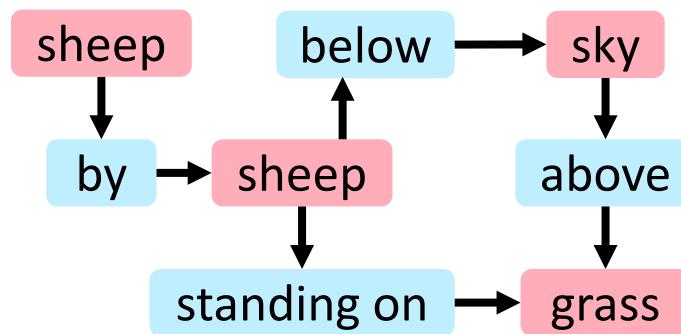


Apply same function to all neighborhoods of input:

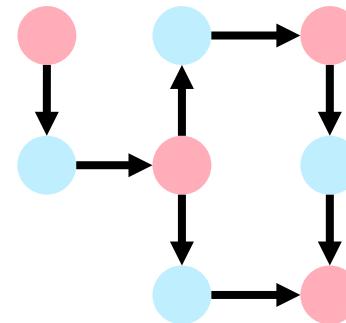
$$y_i = f_W(\mathcal{N}(x_i))$$

# Scene Graph Convolution Networks

**Input:** Graph of vectors



**Output:** Graph of vectors



**Relationship nodes**

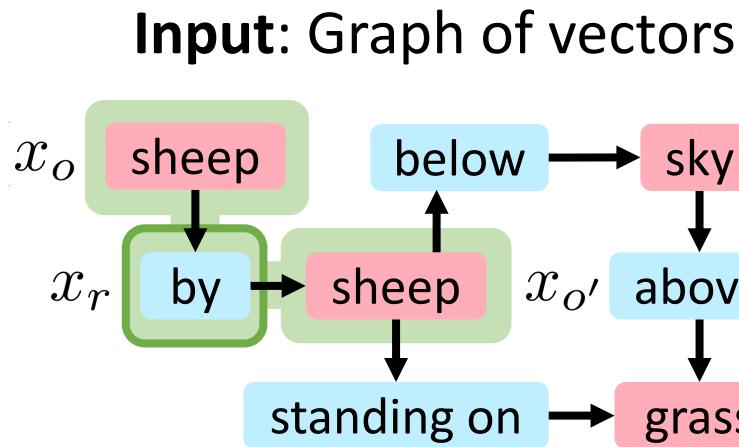
$$y_r = f_W^r(\mathcal{N}(x_r))$$

**Object nodes**

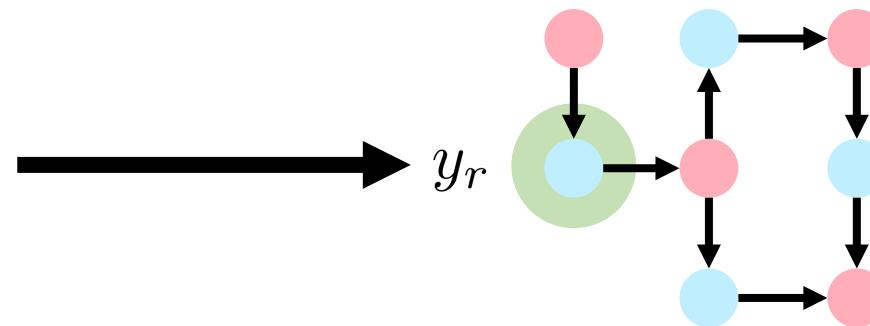
$$y_o = f_W^o(\mathcal{N}(x_o))$$

Two types of nodes: **relationships** and **objects**; use different functions for each

# Scene Graph Convolution Networks



**Output: Graph of vectors**



Two types of nodes: **relationships** and **objects**; use different functions for each

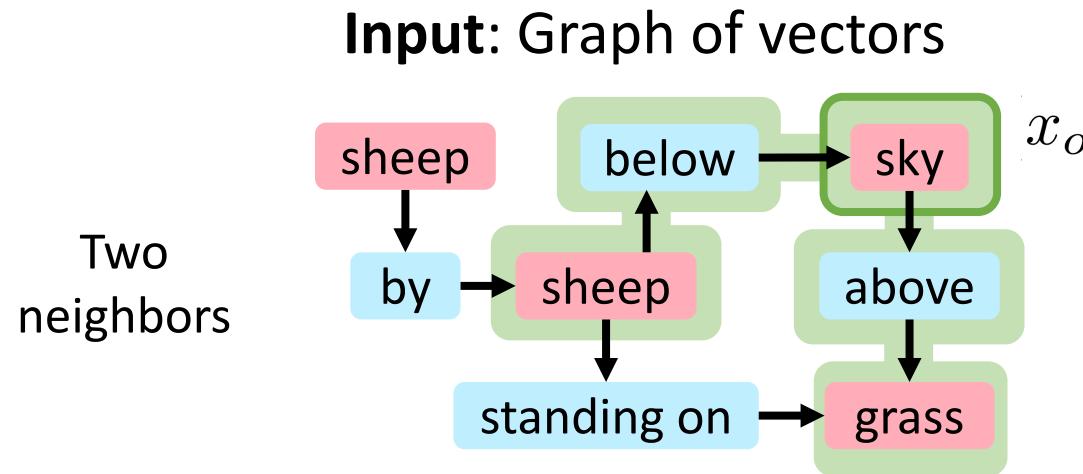
**Relationship** nodes have two neighbors

$$y_r = f_W^r(\mathcal{N}(x_r)) = f_W^r(x_o, x_r, x_{o'})$$

**Object** nodes

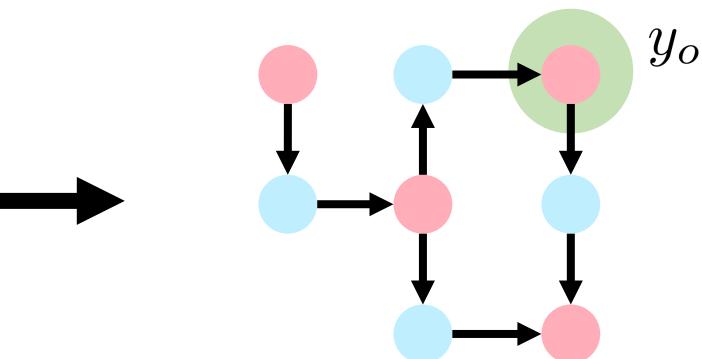
$$y_o = f_W^o(\mathcal{N}(x_o))$$

# Scene Graph Convolution Networks



Two types of nodes: **relationships** and **objects**; use different functions for each

**Output: Graph of vectors**



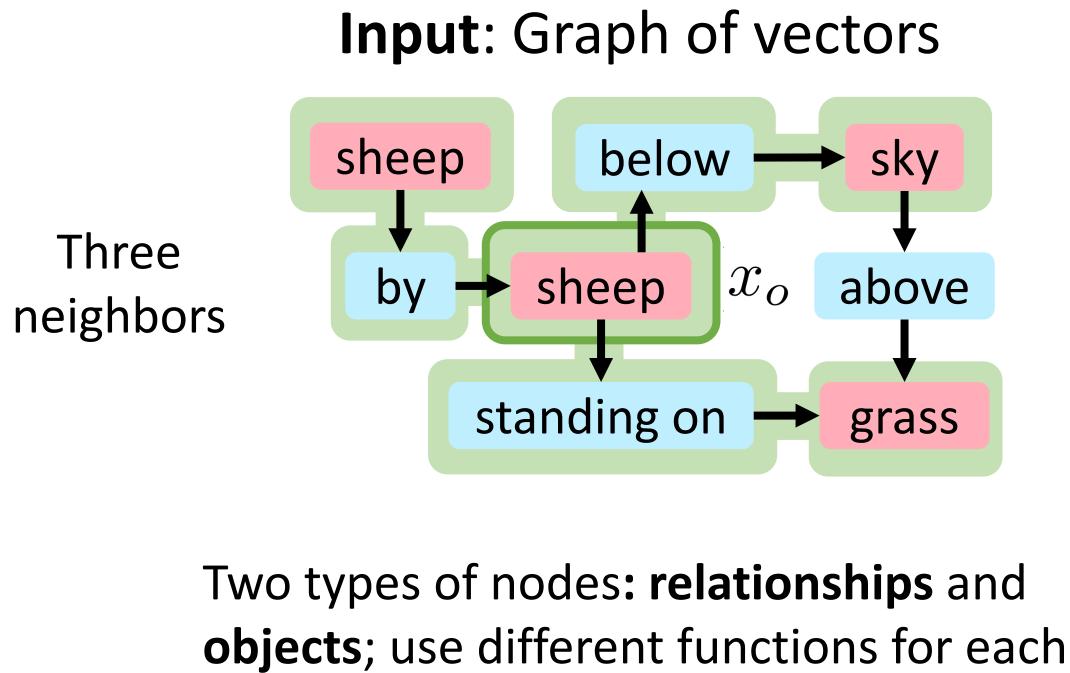
**Relationship** nodes have two neighbors

$$y_r = f_W^r(\mathcal{N}(x_r)) = f_W^r(x_o, x_r, x_{o'})$$

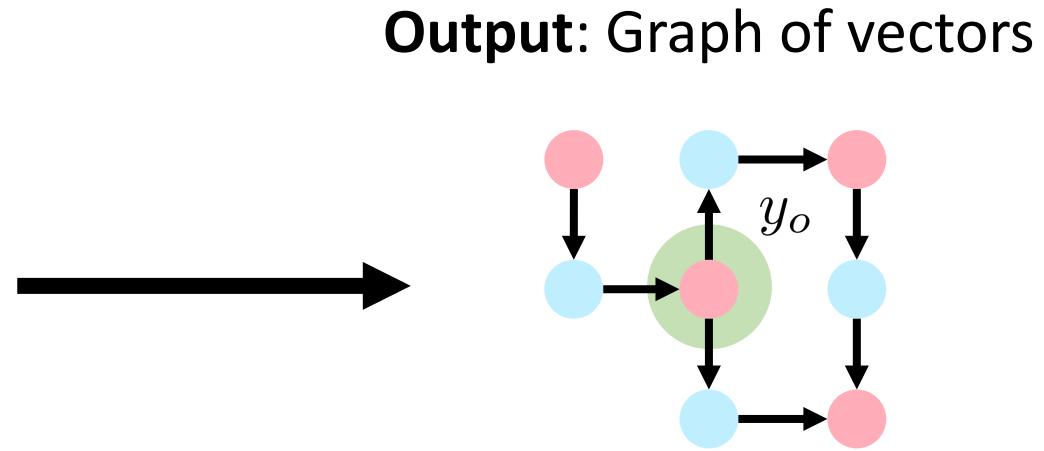
**Object** nodes have varying neighbors

$$y_o = f_W^o(\mathcal{N}(x_o))$$

# Scene Graph Convolution Networks



"Fully-connected" neural network:  
Matrix multiply and nonlinearity

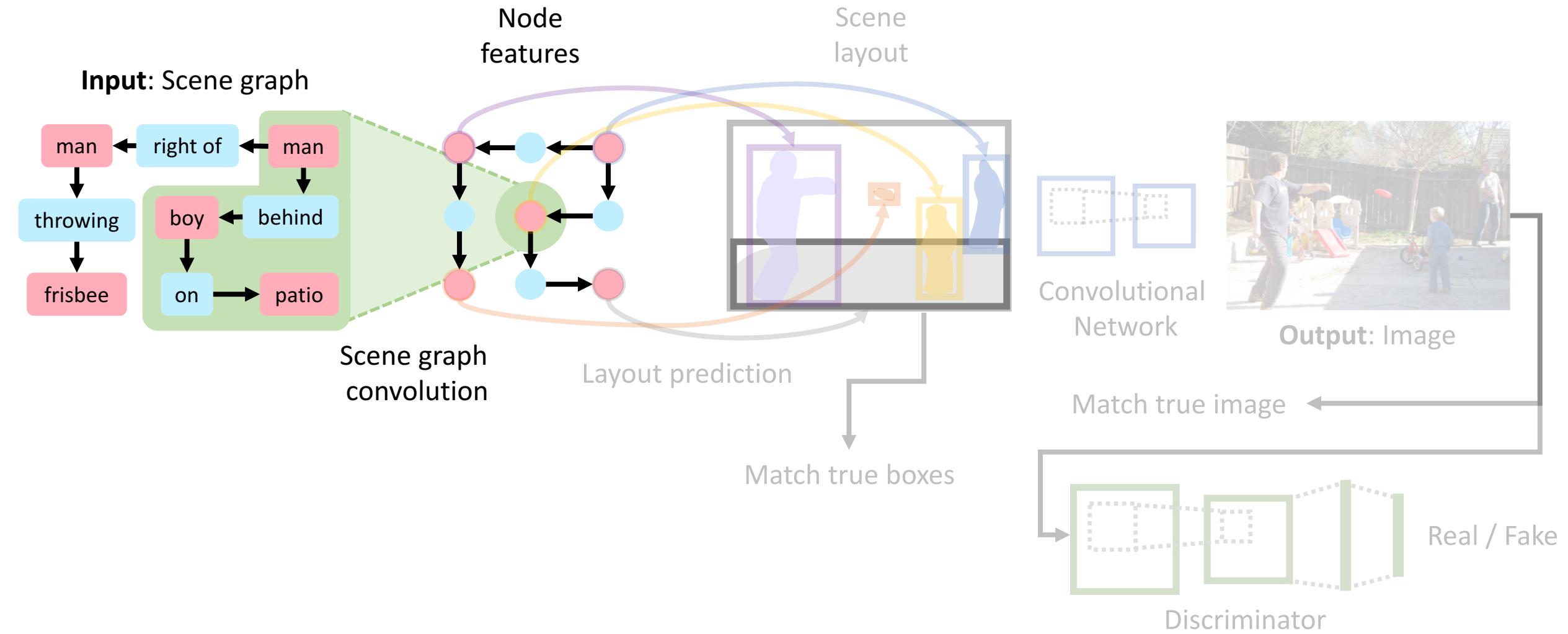


$$y_r = f_W^r(\mathcal{N}(x_r)) = f_W^r(x_o, x_r, x_{o'})$$

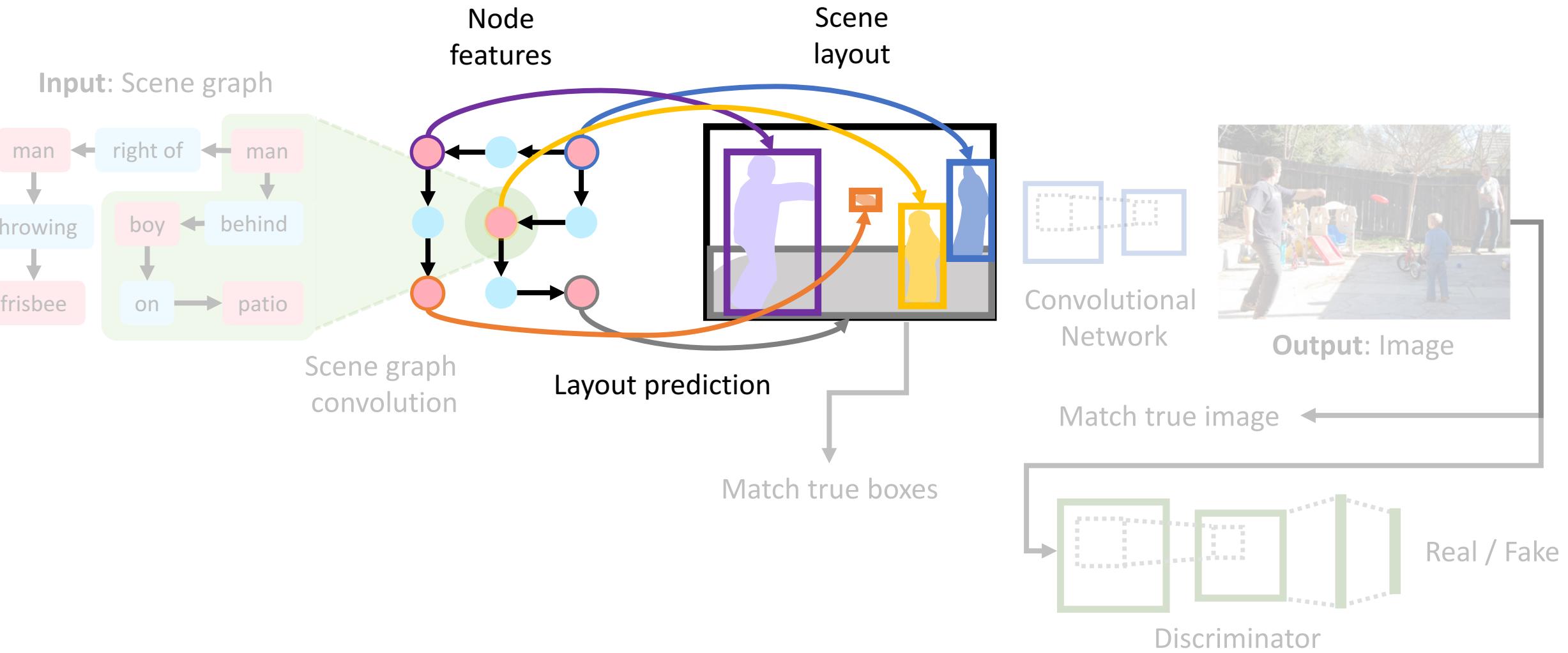
**Object** nodes have varying neighbors

$$y_o = f_W^o(\mathcal{N}(x_o)) \left[ \begin{array}{c} y_W \\ \diagdown \\ \cup \\ (o,e,o') \in G \end{array} \right] f_W^o(x_o, x_r, x_{o'}) \right)$$

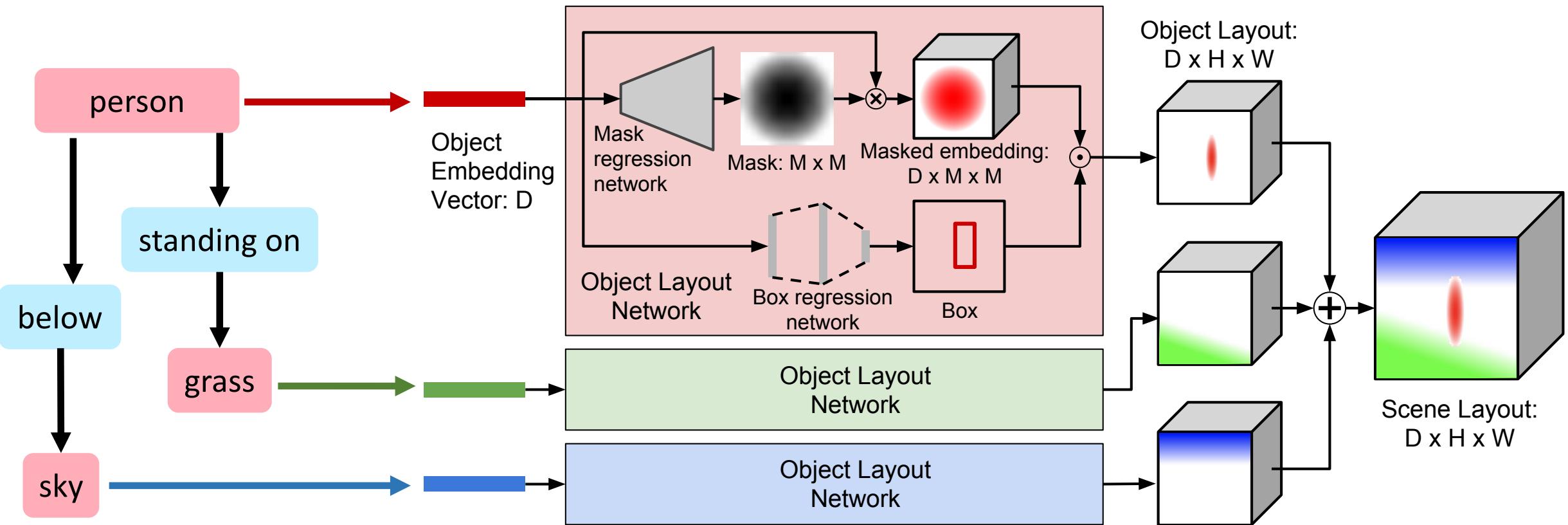
# Scene Graph to Image: Model



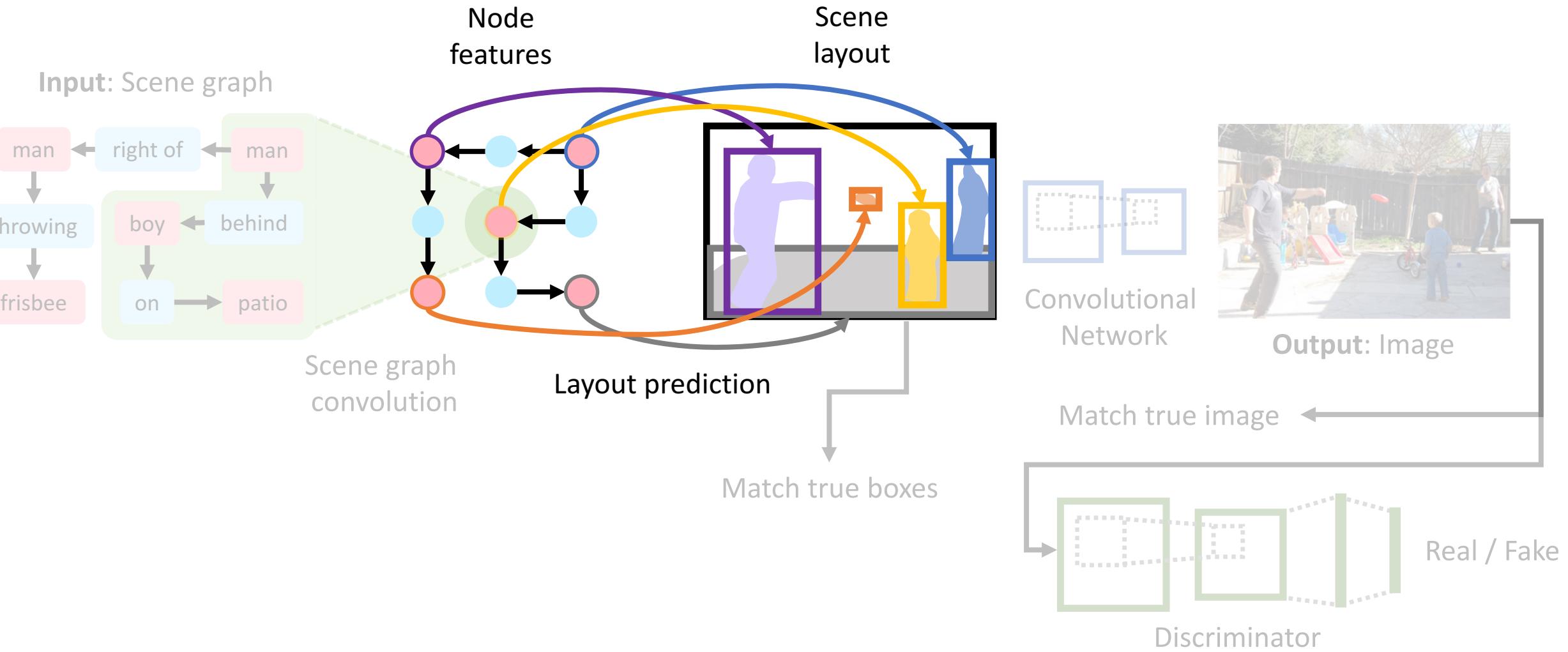
# Scene Graph to Image: Model



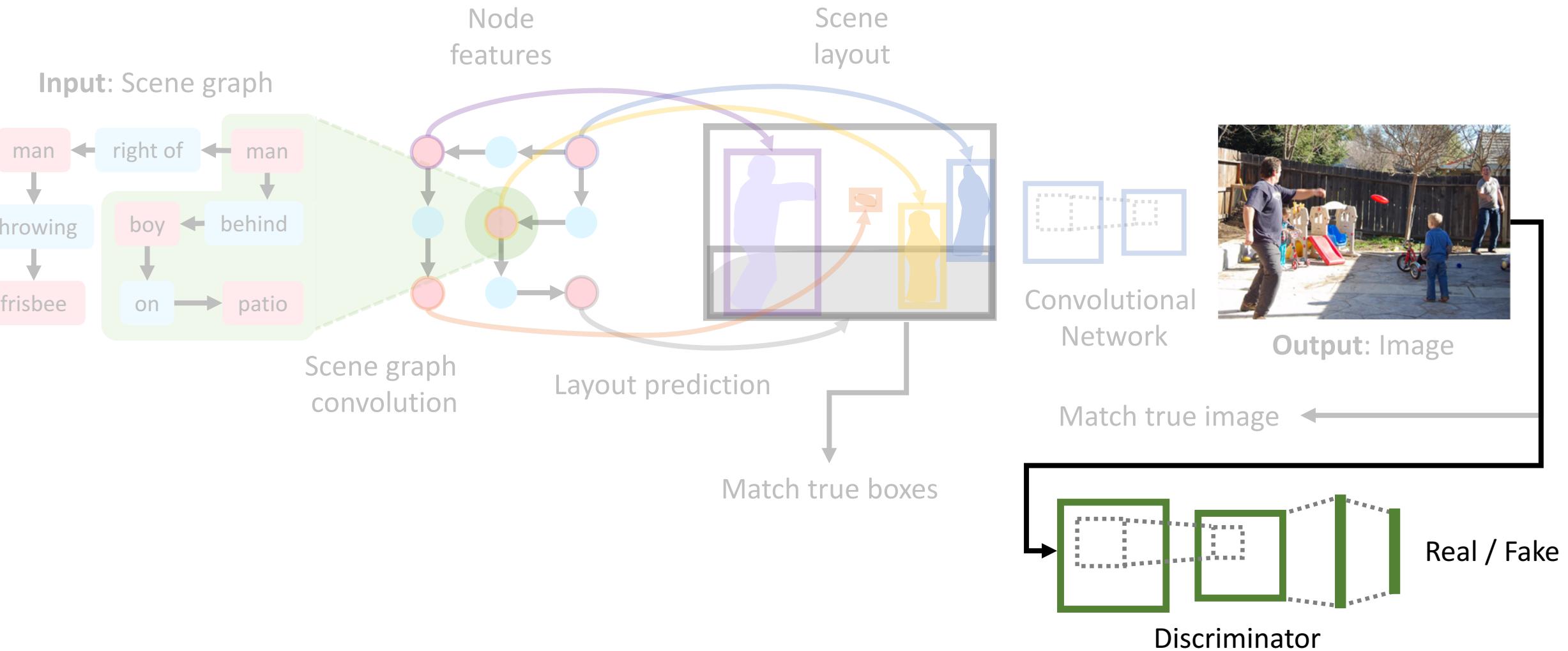
# Scene Layout



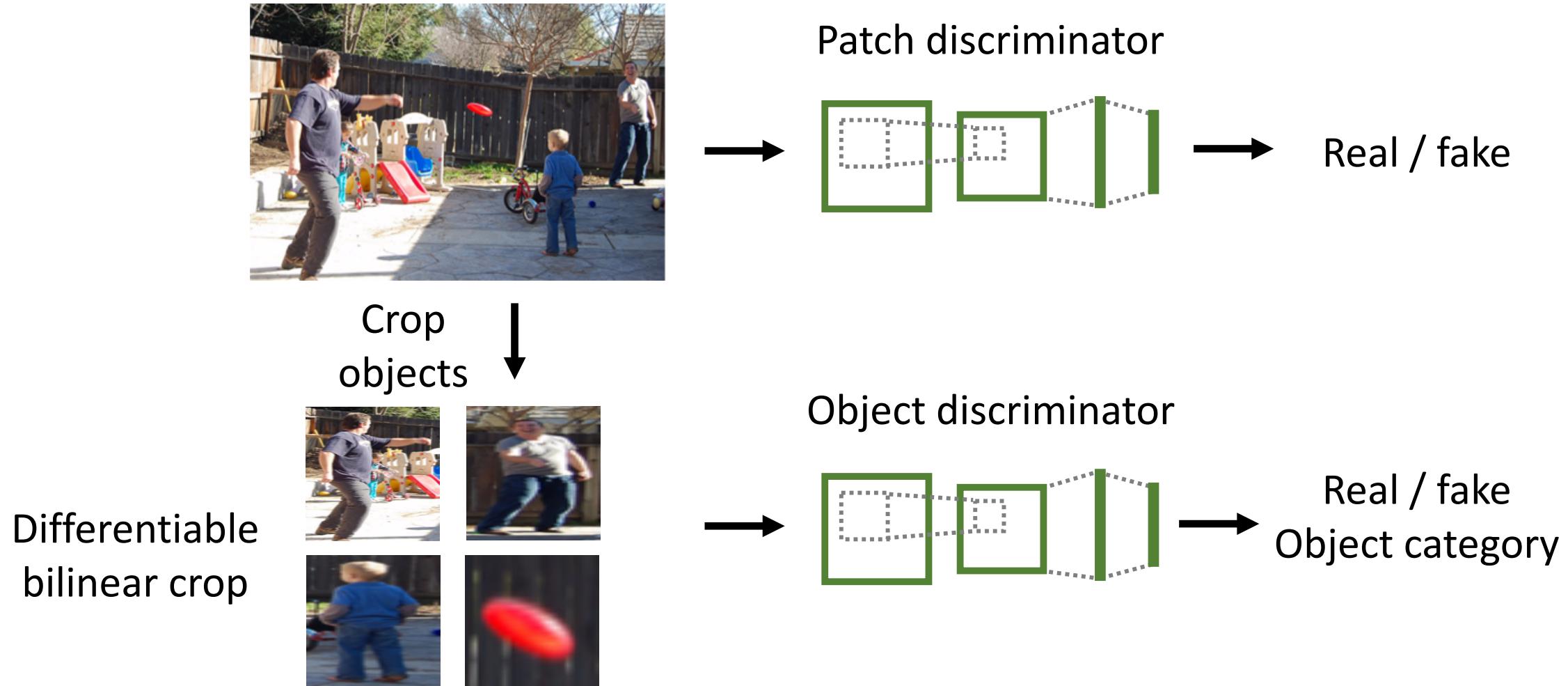
# Scene Graph to Image: Model



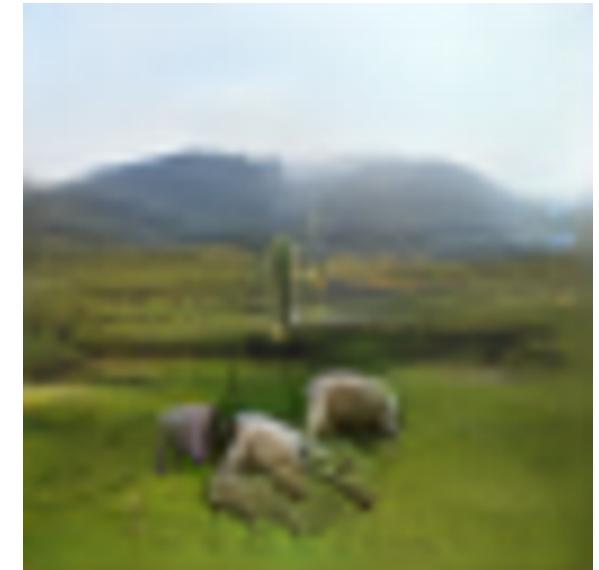
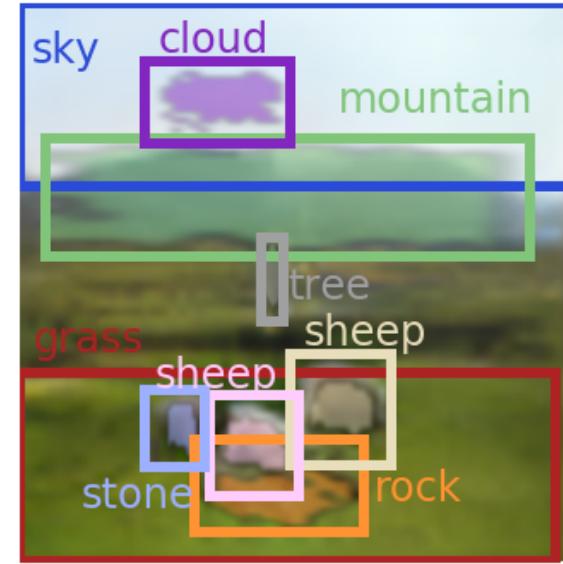
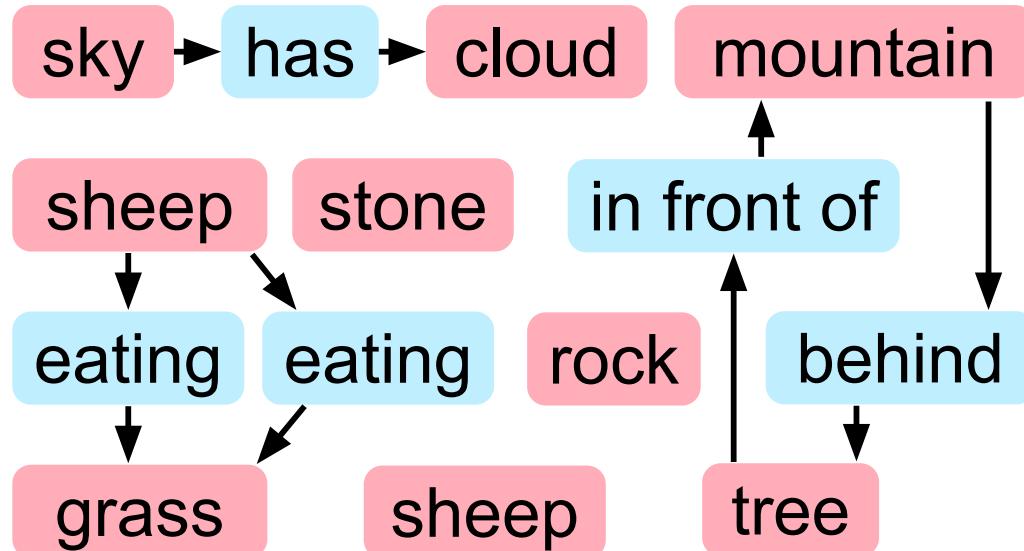
# Scene Graph to Image: Model



# Scene Graph to Image: Discriminator



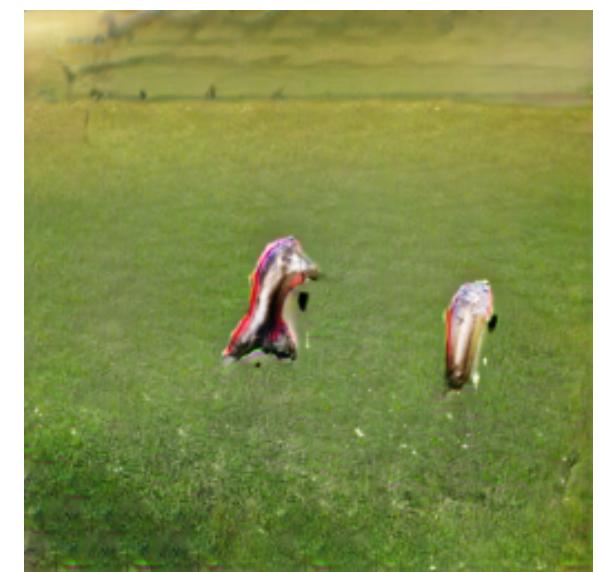
# Results on Visual Genome



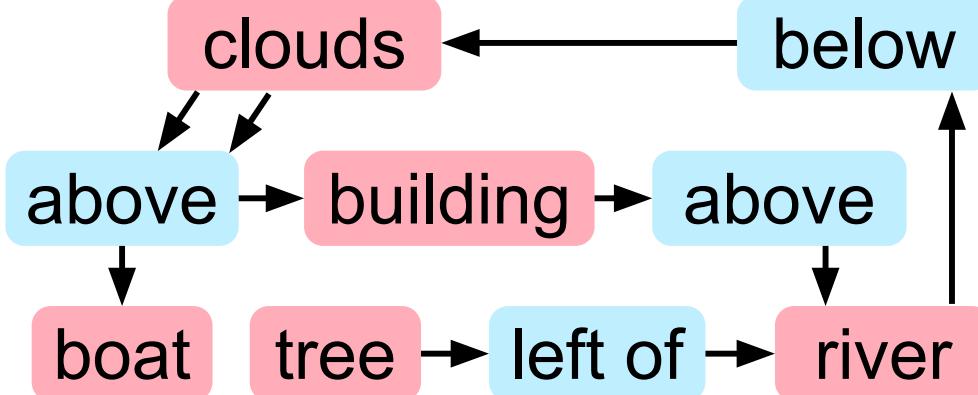
Two sheep, one eating grass  
with a tree in front of a  
mountain; the sky has a cloud



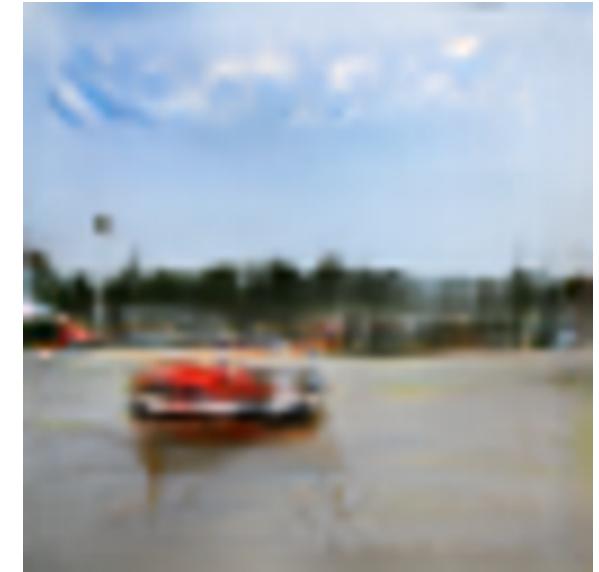
StackGAN  
Zhang et al, ICCV 2017



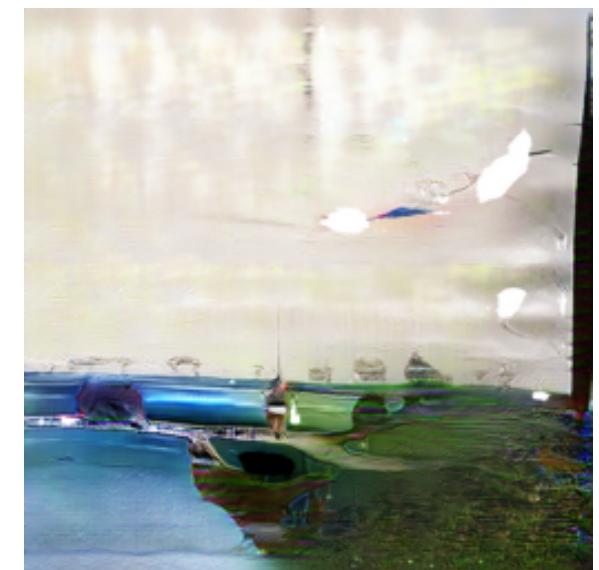
# Results on COCO-Stuff



Clouds above a boat and a building above a river, with trees left of the river



StackGAN  
Zhang et al, ICCV 2017



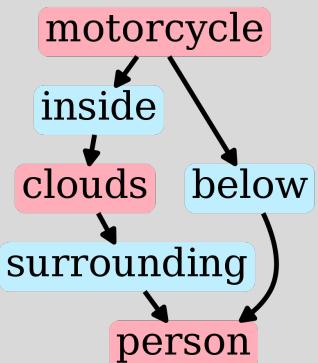
# Ablation Studies

Method	Inception	
	COCO	VG
Real Images ( $64 \times 64$ )	$16.3 \pm 0.4$	$13.9 \pm 0.5$
Ours (No gconv)	$4.6 \pm 0.1$	$4.2 \pm 0.1$
Ours (No relationships)	$3.7 \pm 0.1$	$4.9 \pm 0.1$
Ours (No discriminators)	$4.8 \pm 0.1$	$3.6 \pm 0.1$
Ours (No $D_{obj}$ )	$5.6 \pm 0.1$	$5.0 \pm 0.2$
Ours (No $D_{img}$ )	$5.6 \pm 0.1$	<b><math>5.7 \pm 0.3</math></b>
Ours (Full model)	<b><math>6.7 \pm 0.1</math></b>	$5.5 \pm 0.1$
Ours (GT Layout, no gconv)	$7.0 \pm 0.2$	$6.0 \pm 0.2$
Ours (GT Layout)	<b><math>7.3 \pm 0.1</math></b>	<b><math>6.3 \pm 0.2</math></b>
StackGAN [59] ( $64 \times 64$ )	<b><math>8.4 \pm 0.2</math></b>	-

# Results: User Study



## COCO Annotations



*A man flying  
through the  
air while  
riding a bike*

Ours



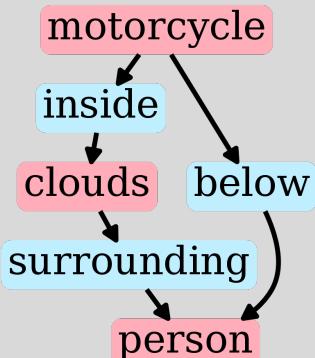
StackGAN



# Results: User Study



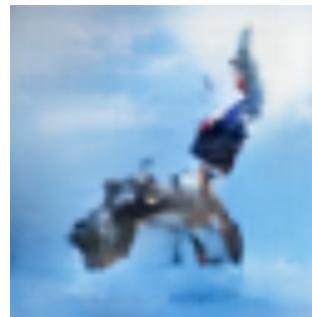
## COCO Annotations



*A man flying through the air while riding a bike*

Q: Which of the following objects are present?  
*motorcycle, clouds, person*

Ours

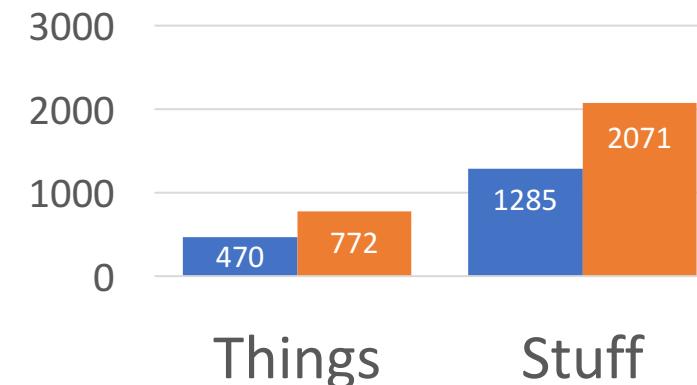


StackGAN



## Object Recall

■ StackGAN ■ Ours



Things

Stuff

# Results: User Study



## COCO Annotations

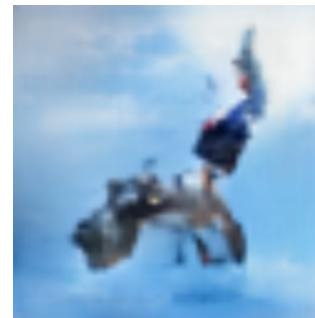
motorcycle  
inside  
clouds  
surrounding  
below  
person

A man flying through the air while riding a bike

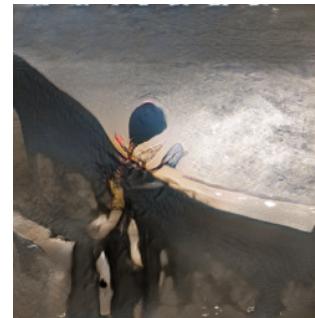
Q: Which image matches the caption better?



Ours

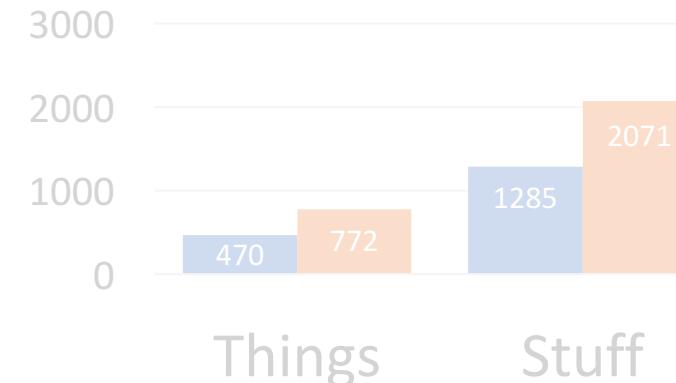


StackGAN

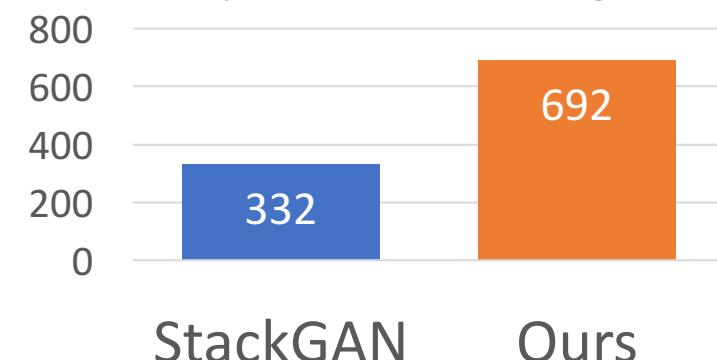


## Object Recall

■ StackGAN ■ Ours



## Caption Matching



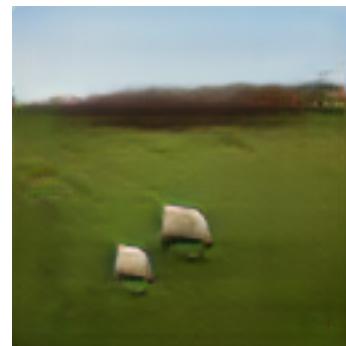
# Simple to Complex



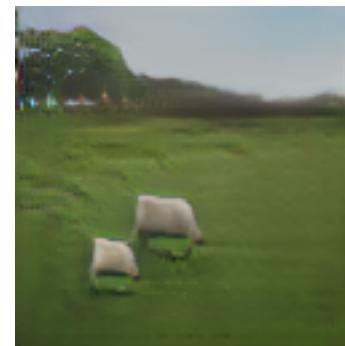
Sky above grass;  
zebra standing  
on grass



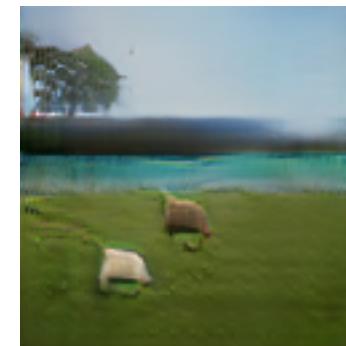
Sky above grass;  
**sheep** standing  
on grass



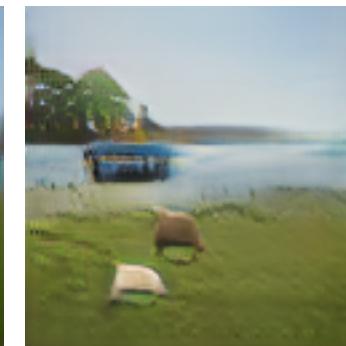
Sky above grass;  
sheep standing  
on grass; **sheep2**  
**by sheep**



Sky above grass;  
sheep standing  
on grass; **sheep2**  
**by sheep; tree**  
**behind sheep**



Sky above grass;  
sheep standing  
on grass; **sheep2**  
**by sheep; tree**  
**behind sheep;**  
**ocean by tree**



Sky above grass;  
sheep standing  
on grass; **sheep2**  
**by sheep; tree**  
**behind sheep;**  
**ocean by tree;**  
**boat in ocean**



Sky above grass;  
sheep standing  
on grass; **sheep2**  
**by sheep; tree**  
**behind sheep;**  
**ocean by tree;**  
**boat on grass**

Image

Scene  
Graph

Sentence

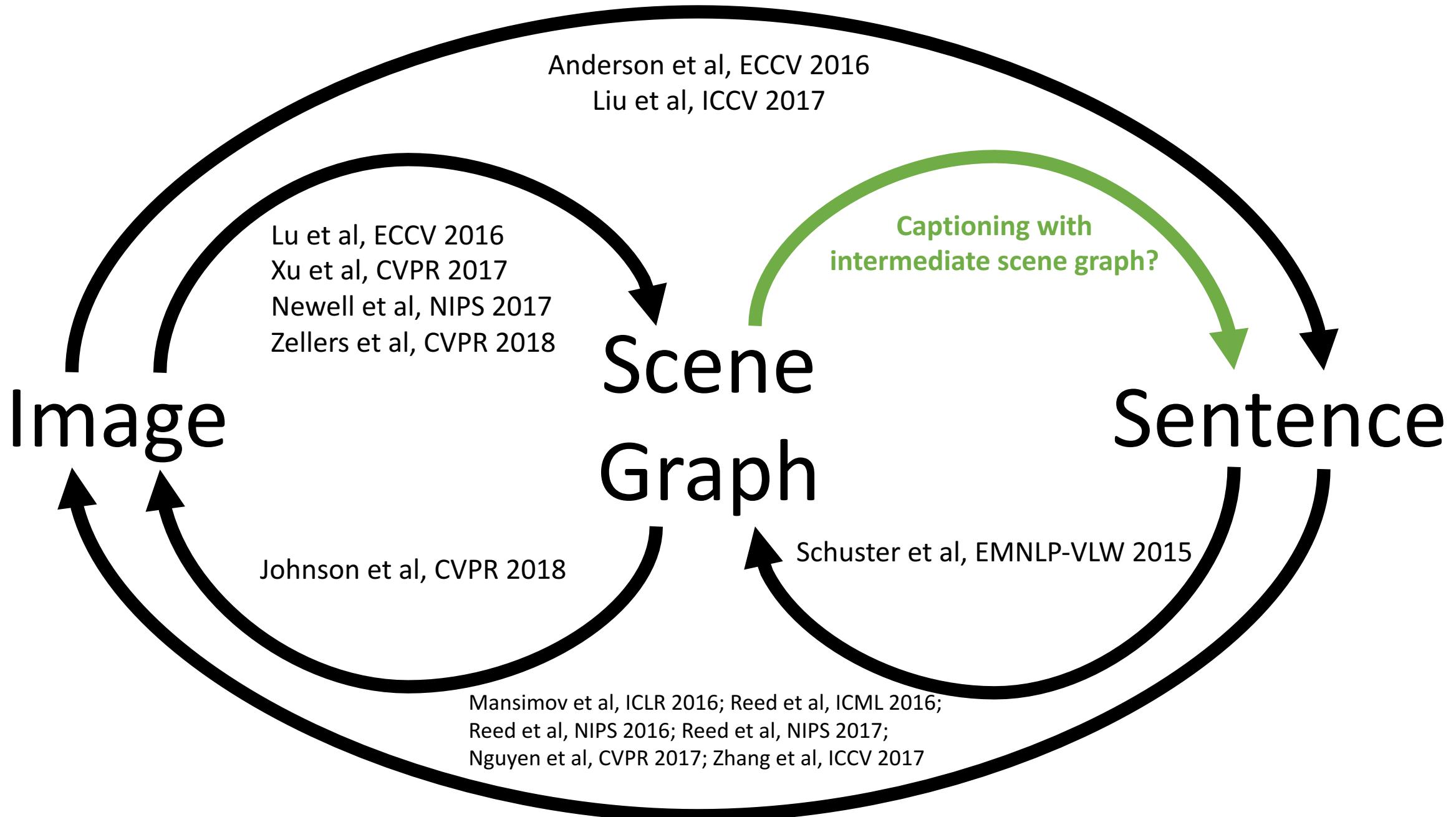
Anderson et al, ECCV 2016  
Liu et al, ICCV 2017

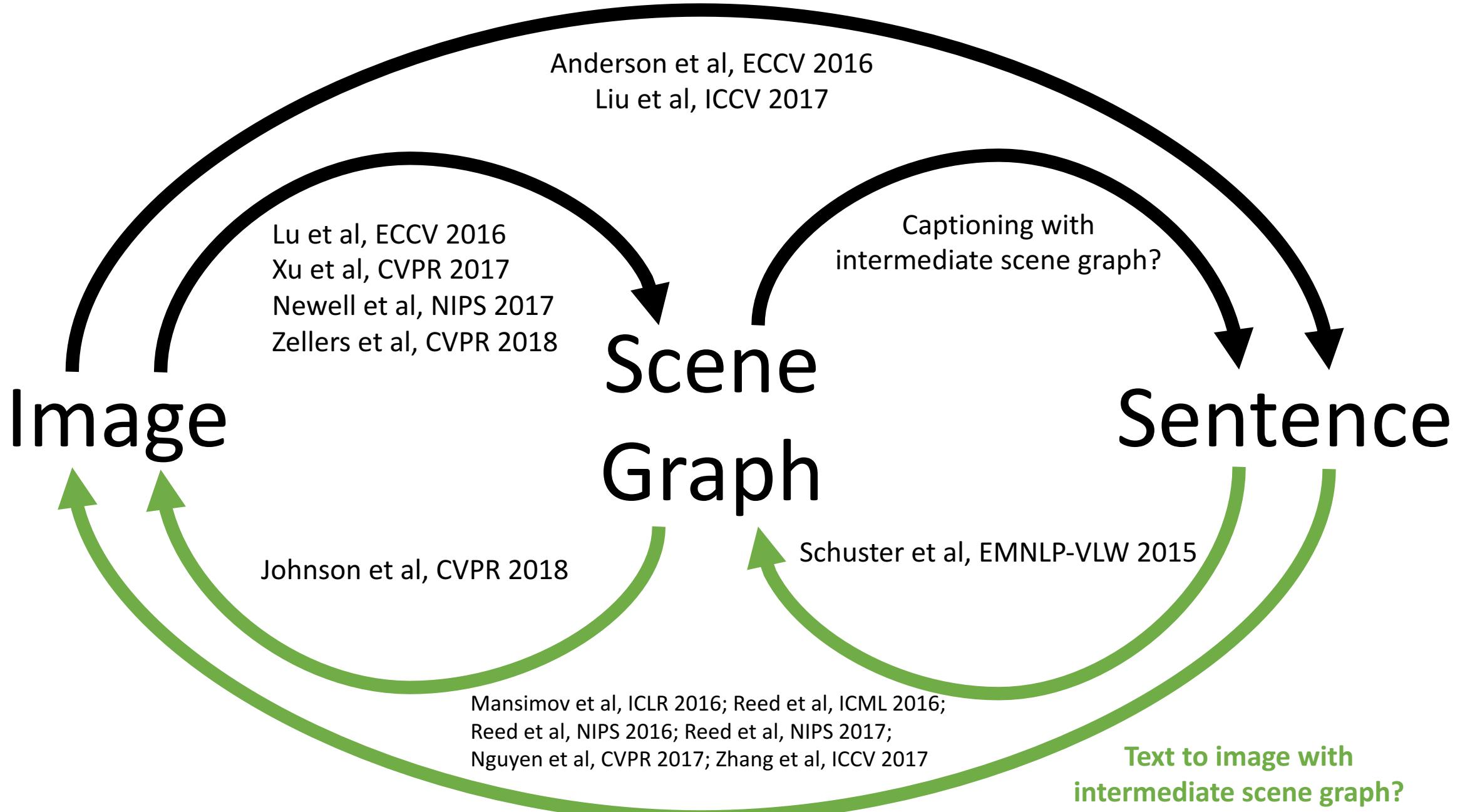
Lu et al, ECCV 2016  
Xu et al, CVPR 2017  
Newell et al, NIPS 2017  
Zellers et al, CVPR 2018

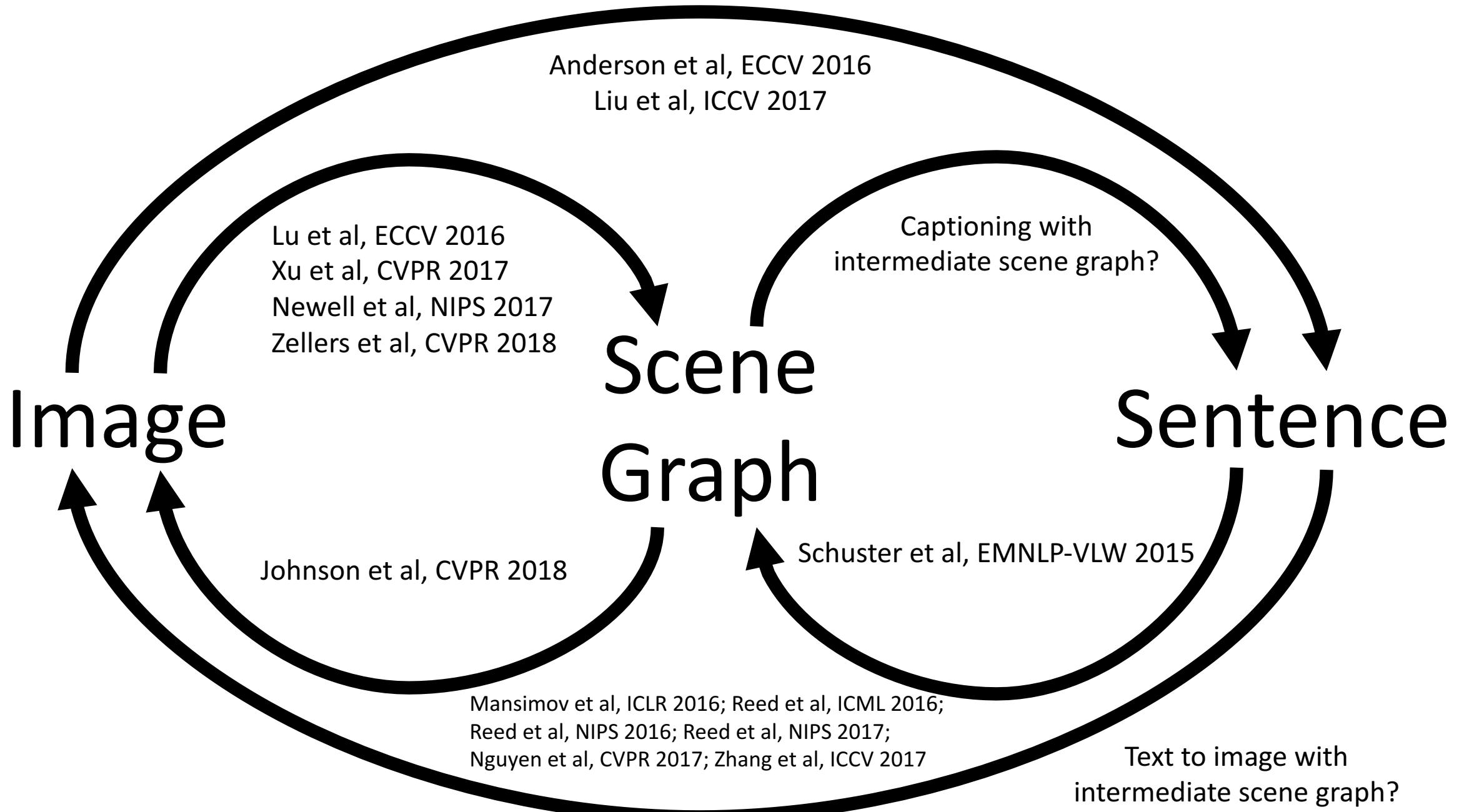
Johnson et al, CVPR 2018

Mansimov et al, ICLR 2016; Reed et al, ICML 2016;  
Reed et al, NIPS 2016; Reed et al, NIPS 2017;  
Nguyen et al, CVPR 2017; Zhang et al, ICCV 2017

Schuster et al, EMNLP-VLW 2015







Questions?