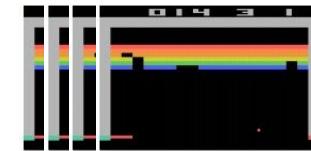
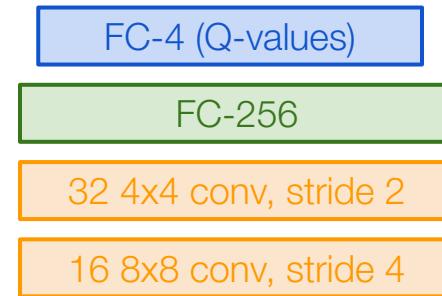
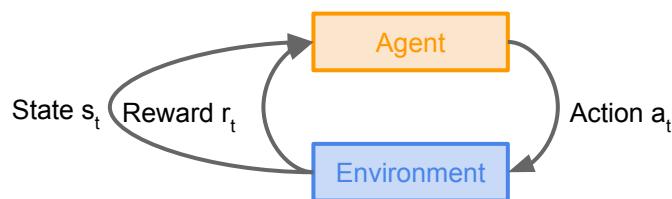


Lecture 18:

Scene Graphs and Graph Convolutions

Last time: Reinforcement learning



Reinforce
algorithm

$$\nabla_{\theta} J(\theta) \approx \sum_{t \geq 0} r(\tau) \nabla_{\theta} \log \pi_{\theta}(a_t | s_t)$$

Policy gradients
Q-learning
Actor-Critic
Soft action critic
Proximal policy gradients

Today's agenda

- Beyond objects
- Scene Graphs
- Scene Graph Generation
- Graph Convolution Networks

Computer vision was focused on disconnected objects

Image Classification



Object Detection



Instance Segmentation



Shilane et al, 2004; Fei-Fei et al, 2004; Griffin et al, 2006; Russell et al, IJCV 2007; Torralba et al, TPAMI 2008; Chen et al, SIGGRAPH 2009; Quattoni and Torralba, CVPR 2009; Deng et al, CVPR 2009; Xiao et al, CVPR 2010; Everingham, IJCV 2010; Silberman et al, ECCV 2012; Xiao et al, ICCV 2013; Lim et al, ICCV 2013; Lin et al, ECCV 2014; Zhou et al, NIPS 2014; Russakovsky et al, IJCV 2015; Chen et al, arXiv 2015; Chang et al, 2015

image #1

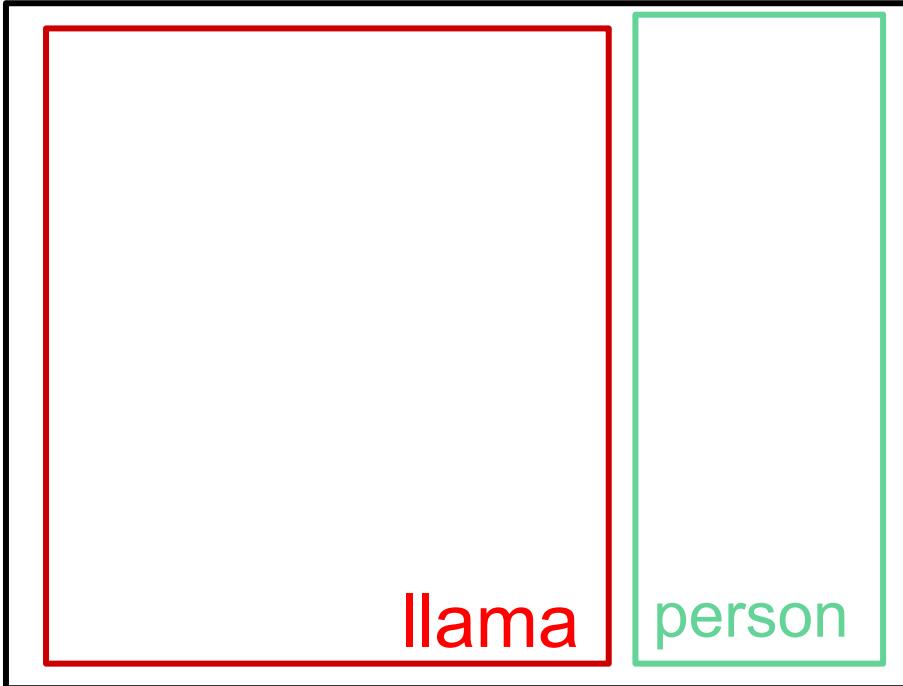
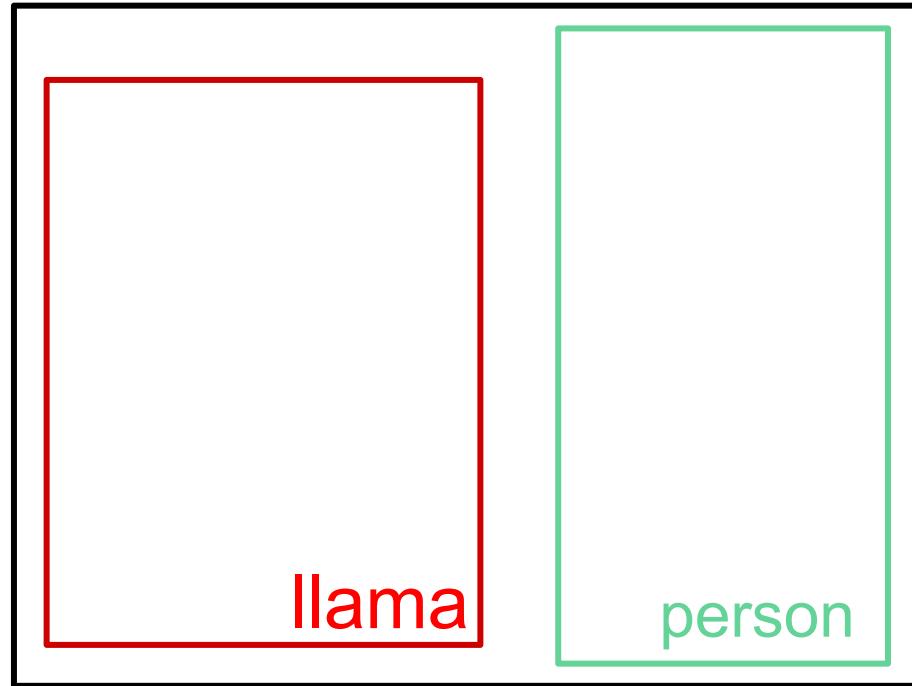


image #2

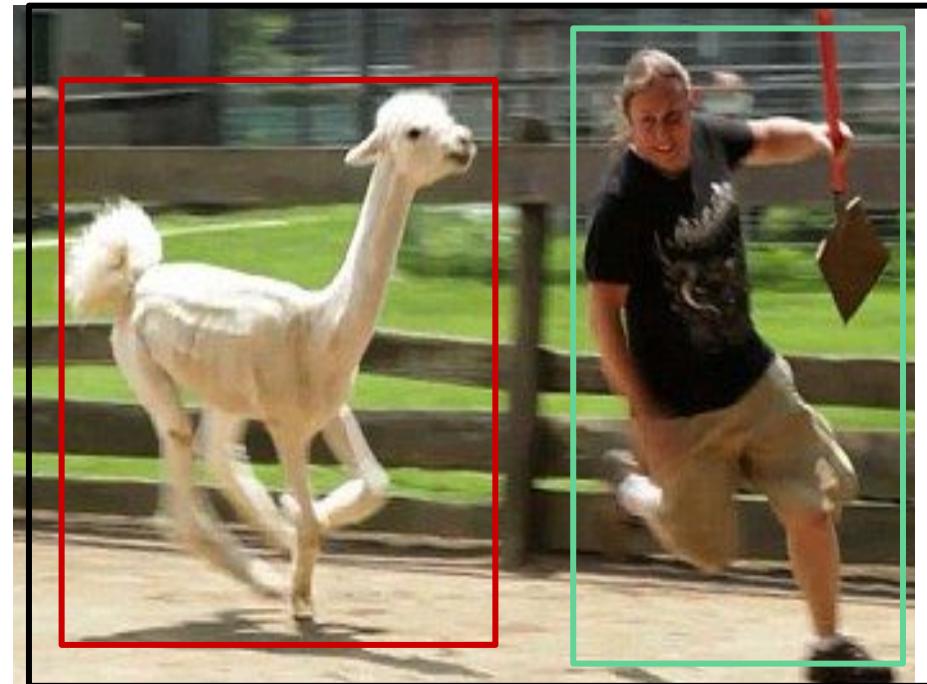


Fei-Fei et al, 2004; Griffin et al, 2006; Torralba et al, TPAMI 2008; Quattoni and Torralba, CVPR 2009; Deng et al, CVPR 2009; Xiao et al, CVPR 2010; Zhou et al, NIPS 2014; Russakovsky et al, IJCV 2015



next to

Fei-Fei et al, 2004; Griffin et al, 2006; Torralba et al, TPAMI 2008; Quattoni and Torralba, CVPR 2009; Deng et al, CVPR 2009; Xiao et al, CVPR 2010; Zhou et al, NIPS 2014; Russakovsky et al, IJCV 2015



chasing

Can image captioning models capture this information?



A man walking a dog

- Wrong! Not a dog
- Wrong! Not walking
- Missed ribbon held by person
- Missed any descriptions of the llama (the model could have said that they are next to one another or that they are in front of the wall).

Lin et al, ECCV 2014
Chen et al, arXiv 2015

What information would people convey if asked to caption?



A llama standing next to a person

White llama in front of a blue wall

A huacaya alpaca held by a person
who is holding a big ribbon

What information would people convey if asked to caption?



Objects

A **llama** standing next to a **person**

White **llama** in front of a blue wall

A huacaya **alpaca** held by a **person**
who is holding a big **ribbon**

What information would people convey if asked to caption?



Objects Attributes

A llama standing next to a person

White llama in front of a blue wall

A huacaya alpaca held by a person who is holding a big ribbon

What information would people convey if asked to caption?



Objects Attributes Relationships

A llama **standing next to** a person

White llama **in front of** a blue wall

A huacaya alpaca **held by** a person
who is **holding** a big ribbon

Many Vision tasks share a similar underlying structure

Action classification



action: drinking from a cup
action: take notebook from somewhere

Grounding objects

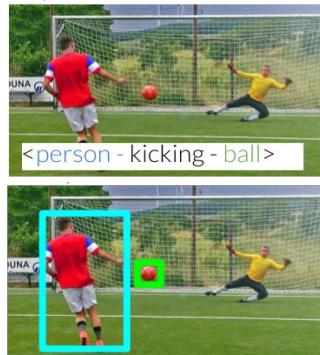
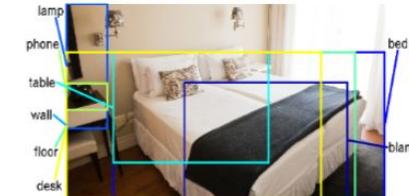


Image retrieval

Black phone is on top of white, wooden desk. The desk is next to a clean white bed that has a black blanket and is next to a white table. The lamp is on a tan wall. The table is by the bed, which is next to the phone. The floor is under the bed, table, lamp and blanket.



Question answering



how many types of vegetables are there? is the food in the foreground prickly?



how many people are in the photo? is this a busy street?
how many skateboards are there? is the man wearing a hat?

Agrawal, et al. Vqa: Visual question answering, ICCV 2015
Swets et al. Using discriminant eigenfeatures for image retrieval, TPAMI 1996
Yu et al. Modeling context in referring expressions, ECCV 2016
Simonyan et al. Two-stream convolutional networks for action recognition in videos ,NeurIPS 2014

Ji, Krishna et al. Action Genome: Actions as Compositions of Spatio-Temporal Scene Graphs, CVPR 2020
Krishna et al. Information Maximizing Visual Question Answering, CVPR 2019
Krishna et al. Referring Relationships, CVPR 2018
Johnson, Krishna et al. Image Retrieval with Scene Graphs, CVPR 2015

Many Vision tasks share a similar underlying structure

Action classification



Grounding objects

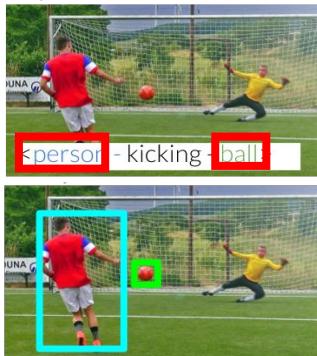
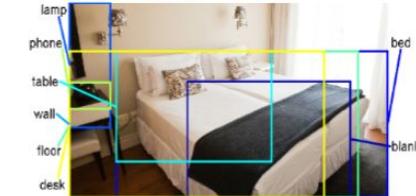


Image retrieval

Black phone is on top of white, wooden desk. The desk is next to a clean white bed that has a black blanket and is next to a white table. The table is by the bed, which is next to the phone. The floor is under the bed, table, lamp and blanket.



Question answering



Agrawal, et al. Vqa: Visual question answering, ICCV 2015

Swets et al. Using discriminant eigenfeatures for image retrieval, TPAMI 1996

Yu et al. Modeling context in referring expressions, ECCV 2016

Simonyan et al. Two-stream convolutional networks for action recognition in videos ,NeurIPS 2014

Ji, Krishna et al. Action Genome: Actions as Compositions of Spatio-Temporal Scene Graphs, CVPR 2020

Krishna et al. Information Maximizing Visual Question Answering, CVPR 2019

Krishna et al. Referring Relationships, CVPR 2018

Johnson, Krishna et al. Image Retrieval with Scene Graphs, CVPR 2015

Many Vision tasks share a similar underlying structure

objects
attributes

Action classification



Grounding objects

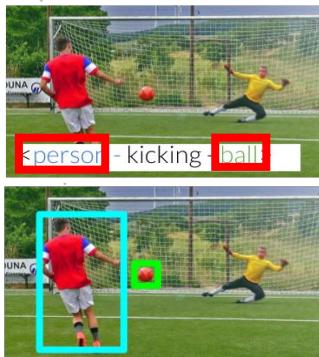


Image retrieval



Question answering



Agrawal, et al. Vqa: Visual question answering, ICCV 2015

Swets et al. Using discriminant eigenfeatures for image retrieval, TPAMI 1996

Yu et al. Modeling context in referring expressions, ECCV 2016

Simonyan et al. Two-stream convolutional networks for action recognition in videos ,NeurIPS 2014

Ji, Krishna et al. Action Genome: Actions as Compositions of Spatio-Temporal Scene Graphs, CVPR 2020

Krishna et al. Information Maximizing Visual Question Answering, CVPR 2019

Krishna et al. Referring Relationships, CVPR 2018

Johnson, Krishna et al. Image Retrieval with Scene Graphs, CVPR 2015

Many Vision tasks share a similar underlying structure

objects
attributes
relationships

Action classification



Grounding objects

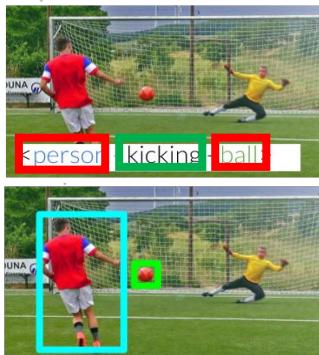
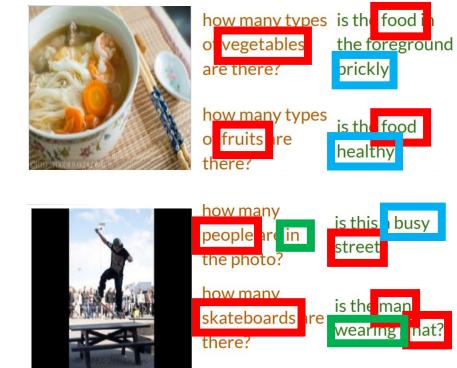


Image retrieval



Question answering



Agrawal, et al. Vqa: Visual question answering, ICCV 2015

Swets et al. Using discriminant eigenfeatures for image retrieval, TPAMI 1996

Yu et al. Modeling context in referring expressions, ECCV 2016

Simonyan et al. Two-stream convolutional networks for action recognition in videos ,NeurIPS 2014

Ji, Krishna et al. Action Genome: Actions as Compositions of Spatio-Temporal Scene Graphs, CVPR 2020

Krishna et al. Information Maximizing Visual Question Answering, CVPR 2019

Krishna et al. Referring Relationships, CVPR 2018

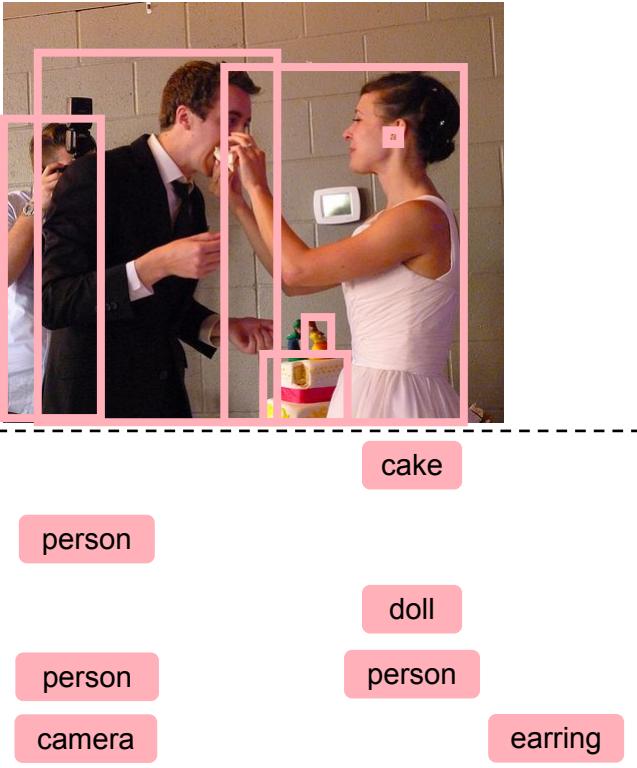
Johnson, Krishna et al. Image Retrieval with Scene Graphs, CVPR 2015

The scene graph representation



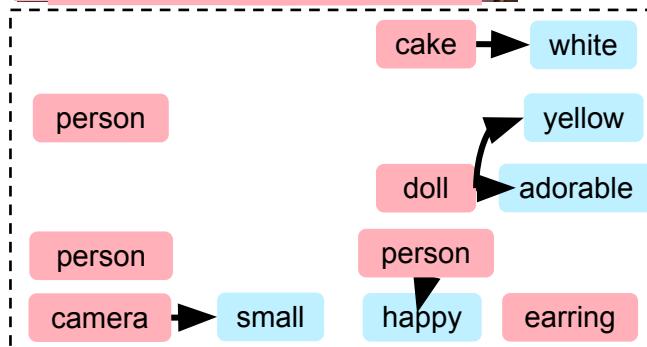
Krishna et al., Visual Genome: Connecting Vision and Language using
Crowdsourced Image Annotations, IJCV 2017

The scene graph representation



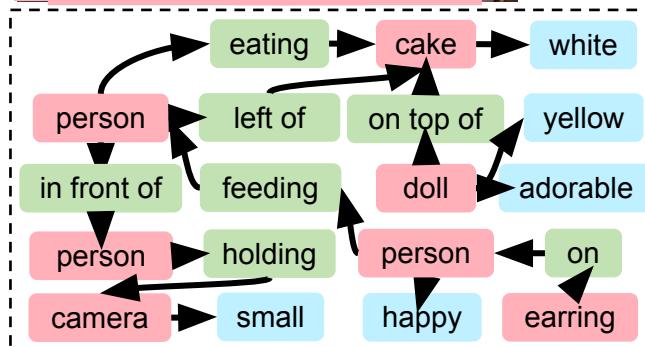
Krishna et al., Visual Genome: Connecting Vision and Language using
Crowdsourced Image Annotations, IJCV 2017

The scene graph representation



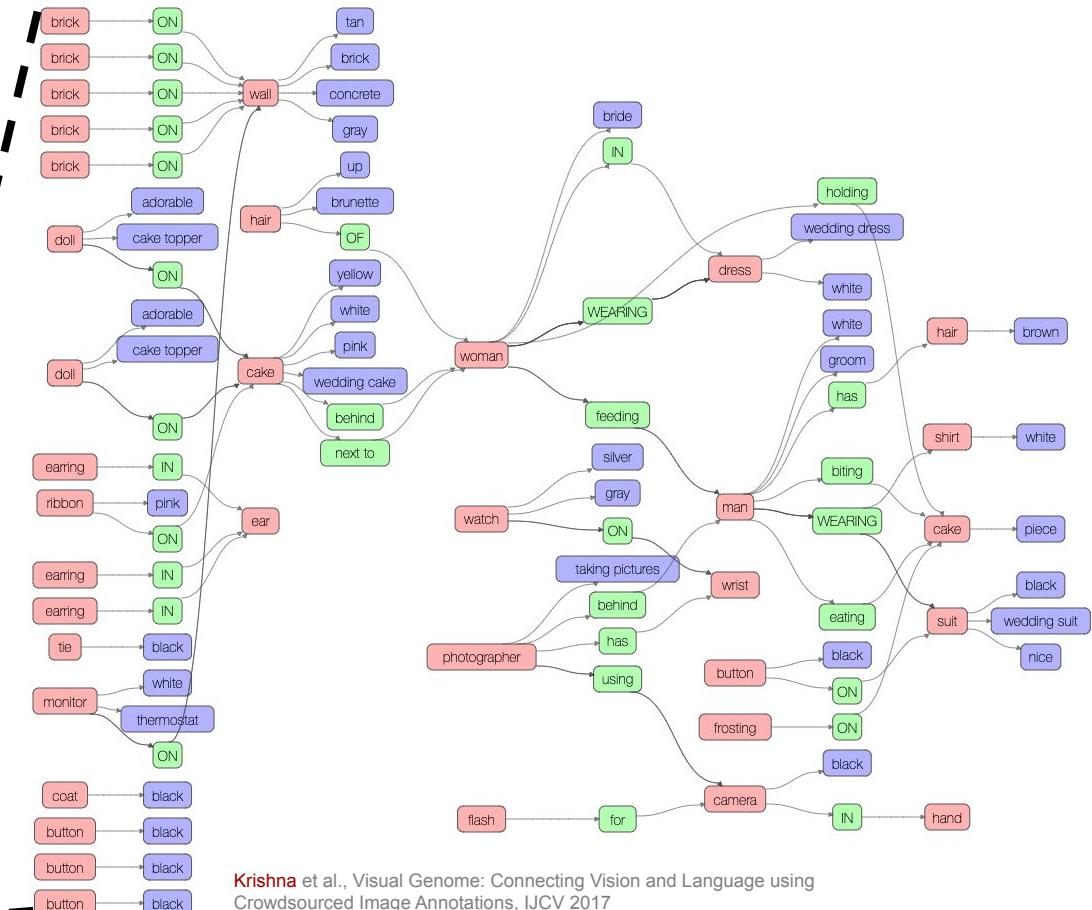
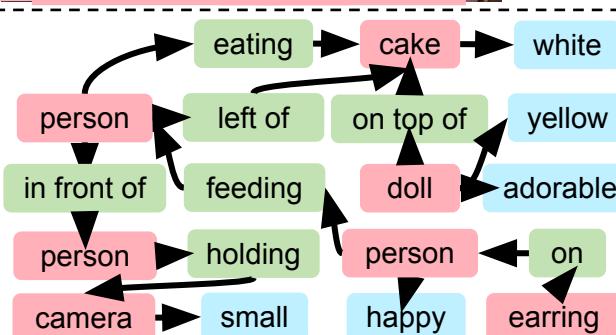
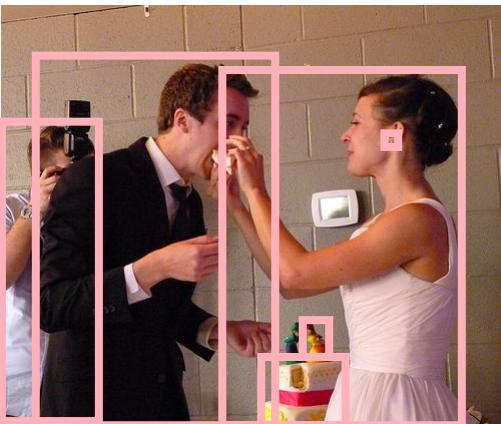
Krishna et al., Visual Genome: Connecting Vision and Language using
Crowdsourced Image Annotations, IJCV 2017

The scene graph representation



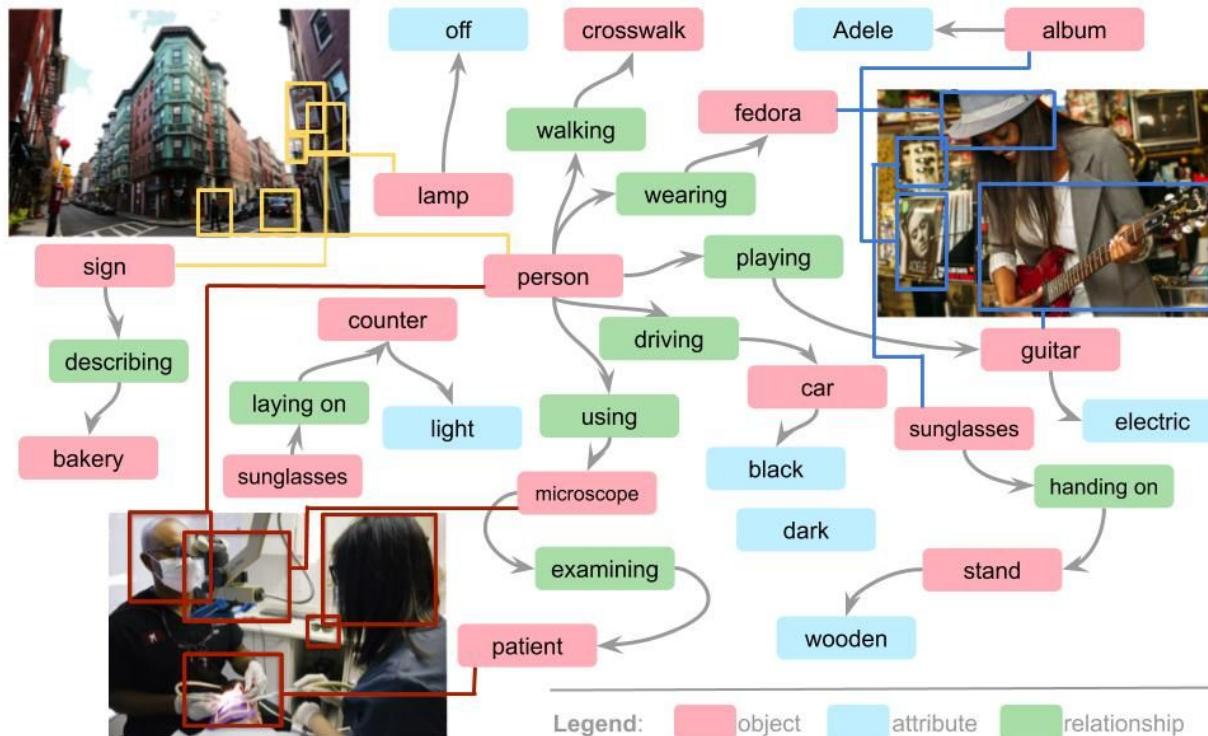
Krishna et al., Visual Genome: Connecting Vision and Language using
Crowdsourced Image Annotations, IJCV 2017

The scene graph representation



Krishna et al., Visual Genome: Connecting Vision and Language using
Crowdsourced Image Annotations, IJCV 2017

Visual Genome – connects images together with scene graphs



108K images
3.8 Million Objects
2.8 Million Attributes
2.3 Million Relationships
1.7 Million question answers
5.4 Millions descriptions

Everything Mapped to Wordnet Synsets

Code and dataset available:
<http://visualgenome.org>
Visualization code:
<https://github.com/ranjaykrishna/graphviz>

Krishna et al., Visual Genome: Connecting Vision and Language using Crowdsourced Image Annotations, IJCV 2017

But why is scene graph the right representation?

Try and remember all these images



All images are CC0 1.0 public domain. sources: [1](#), [2](#), [3](#), [4](#), [5](#), [6](#), [7](#), [8](#), [9](#), [10](#), [11](#), [12](#), [13](#), [14](#), [15](#), [16](#), [17](#), [18](#), [19](#), [20](#), [21](#), [22](#), [23](#), [24](#)

Do you remember seeing this image?

a



b



c



d



The difficulty with the appealing idea that we remember the gist of a scene is that there is no consensus about the contents of a 'gist'. Intuition suggests that an inventory of some of the objects in the scene should be at least a part of the gist.

Wolfe, Visual Memory: What do you know about what you saw?
Biology, 1998

R304 Current Biology, Vol 8 No 9

to one object than to spatially equivalent properties spread over two or more objects [17–19].

needed to show that changes in the gist were necessary and sufficient for efficient change detection. Given the

Where does this leave us in the situation? Evidence from visual search suggests that objects can be identified per second [20]. It is possible that a any sort of stable memory runs at the 'rapid serial visual presentation' example). In either case, a relatively scene would allow several objects passed to memory. Is that listing the names of N objects? Or is it the whole scene, but a series of thought experiments must be coded into the gist being poured from a carton into a glass picture of milk being poured from next to a glass, even if all of the. Moreover, even if all the properties between objects remain the same, we

Dispatch R303

Visual memory: What do you know about what you saw?

Jeremy M. Wolfe

Recent studies of visual perception are bringing us closer to an understanding of what we remember – and what we forget – when we recall a scene.

Address: Center for Ophthalmic Research, Brigham and Women's Hospital, 221 Longwood Avenue, Boston, Massachusetts 02115, USA.

Current Biology 1998, 8:R303–R304
<http://biomednet.com/electr/09609822008R0303>

© Current Biology Ltd ISSN 0960-9822

d is your memory? One line of research starting years ago shows that your memory for visually material is quite remarkably good [1,2]. In a picture recognition study, subjects are shown a of scenes—such as images cut from a glossy magazine—each of which is presented for a second the test phase, subjects are shown a second set half of them from the first set and the other half for the first time. The task is to identify of the second set as old or new. Subjects very well on such a task, even when thousands of

they were not—but because books are part of the schema for what should be in an office. People routinely remember seeing more of a scene than was presented [13,14]. On a more sinister note, memory for scenes can be colored by the biases of the observer [15].

The difficulty with the appealing idea that we remember the gist of a scene is that there is no consensus about the contents of a 'gist'. Intuition suggests that an inventory of some of the objects in the scene should be at least a part of the gist. If you have seen a scene with a house, you would be surprised if the description named no objects but relied only on a description of features, such as color or size. A recent experiment by Luck and Vogel [16] seems to show this coding into memory for objects, rather than simple features. They performed a variation of a change-detection experiment. Two arrays of items were presented to subjects; on half the trials, the second array contained one item that was changed. If one-to-three colored squares were presented, subjects could perfectly detect the change; performance fell off with larger sets of items. These results suggest that subjects keep track of four objects. Now, suppose that each item on the screen could vary in color, orientation, size and the presence or absence of a gap. Would subjects be able to keep track of just four individual features, or would they be able to keep track of up to four objects with all of their associated features? The answer, in a variety of versions of this experiment, is that subjects kept track of objects. They could detect any single feature change in any of up to four objects, even though that meant keeping track of more than a dozen individual features.

There is a bottleneck between vision and memory. If you close your eyes, you will immediately lose access to many of the details that were obvious a moment ago. The results of Luck and Vogel [16] show that it is objects and not raw features, that move through that bottleneck. The selection of objects is governed by attention. There is copious evidence that it is easier to attend to properties belonging

In an earlier paper, Wolfe found that of subjects could identify the memory, and you might just have it represented in this way, it will not a brain can remember thousands of objects in a scene. Another program

We encode more than objects



Visual memory: What do you know about what you saw?

Jeremy M. Wolfe

Recent studies of
closer to an under-
what we forget – w

Address: Center for Opt
Hospital, 221 Longwood

Current Biology 1998,
<http://biomednet.com/lni>

© Current Biology I

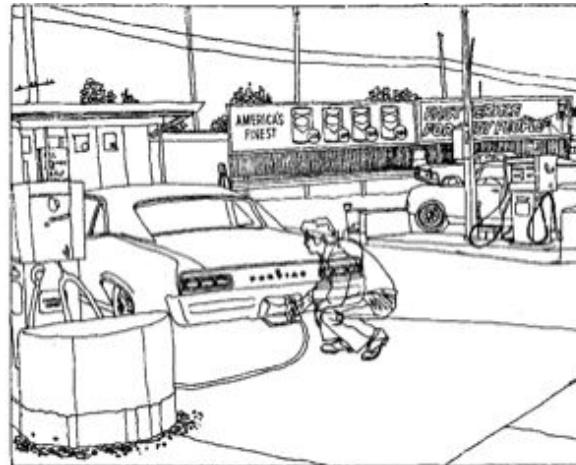
How good is your
short-term visual and

Attributes and Relationships are processed independent of Objects

Attribute and relationship violations are noticed within 150ms.

Relationship violations slow down object identification.

Biederman, *Visual Memory: What do you know about what you saw?* Cognitive Psychology, 1982



Scene Graph Generation - Problem formulation

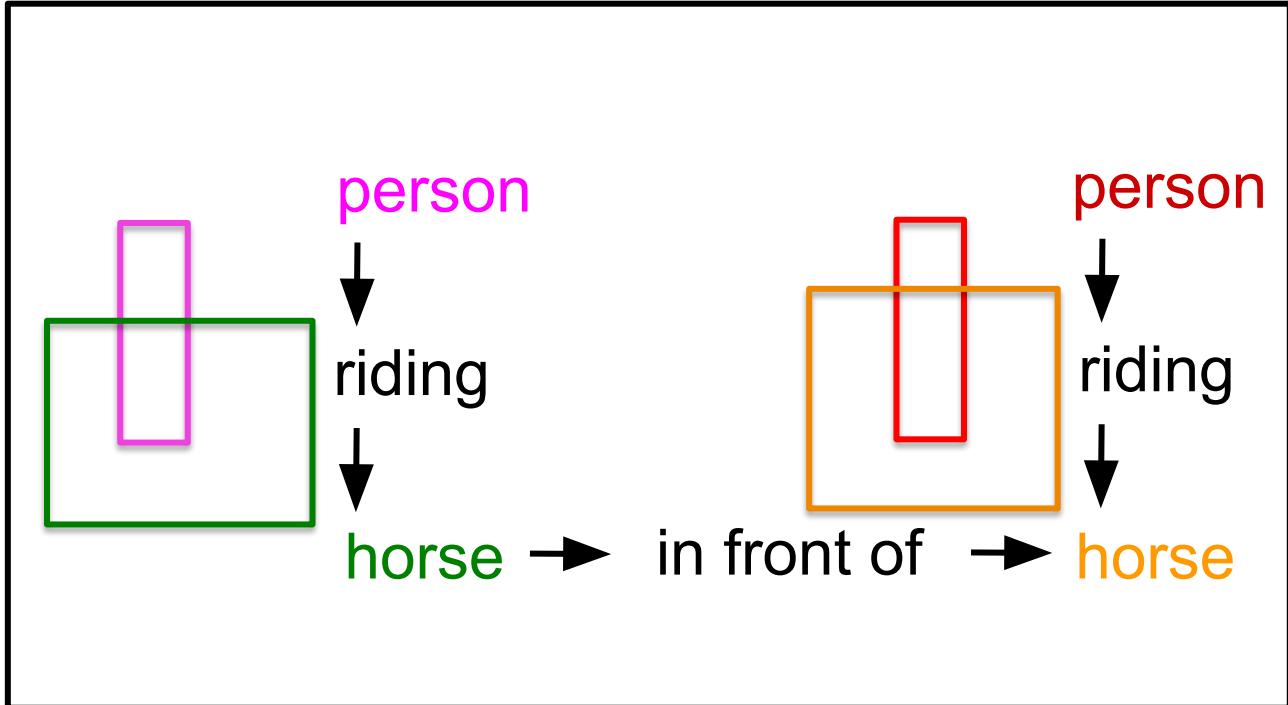


Input
(image only)

Scene Graph Generation - Problem formulation



Input
(image only)



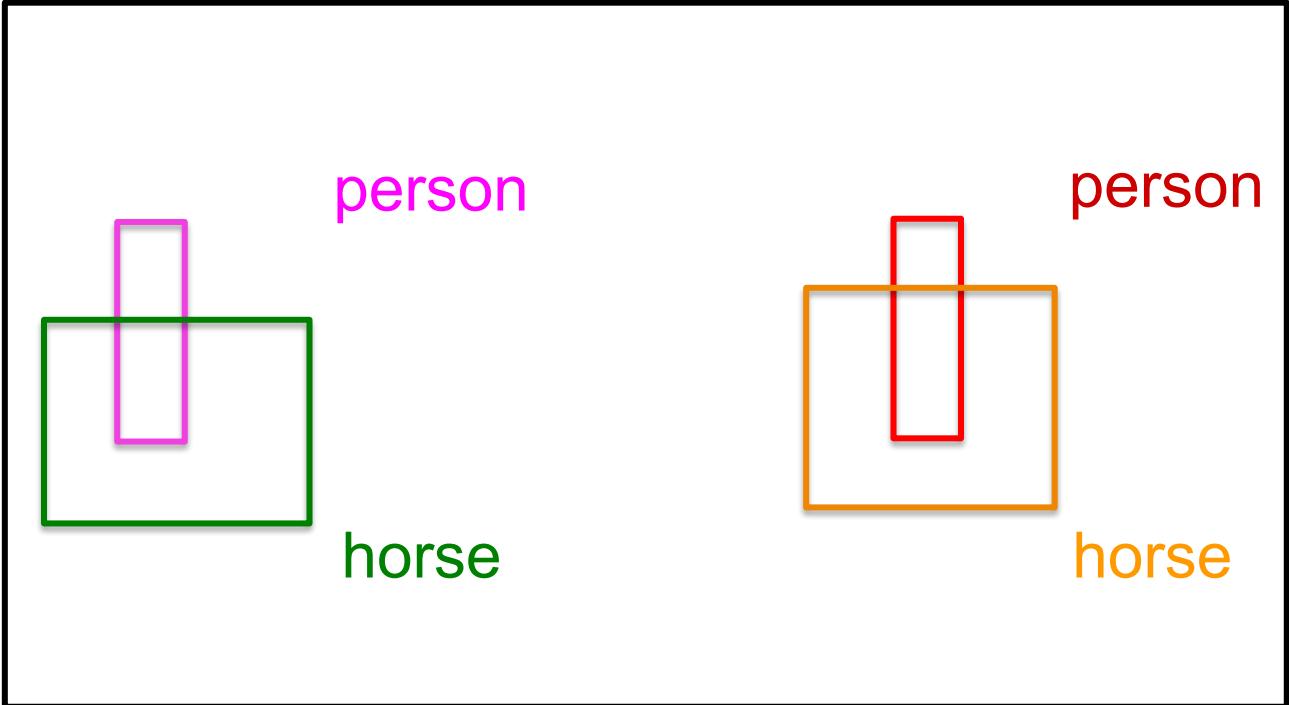
Output

Lu, Krishna et al., Visual Relationship Detection with Language Priors, ECCV 2016

Scene Graph Generation - Problem formulation



Input
(image only)

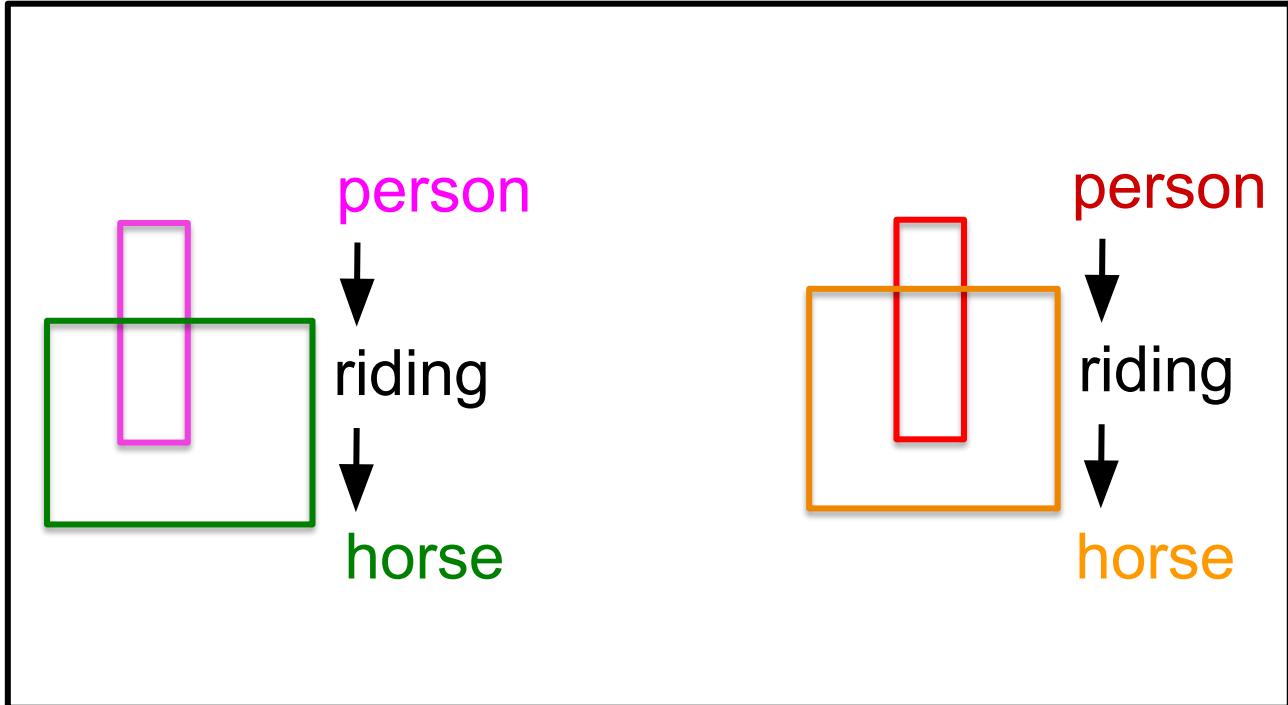


Lu, Krishna et al., Visual Relationship Detection with Language Priors, ECCV 2016

Scene Graph Generation - Problem formulation



Input
(image only)



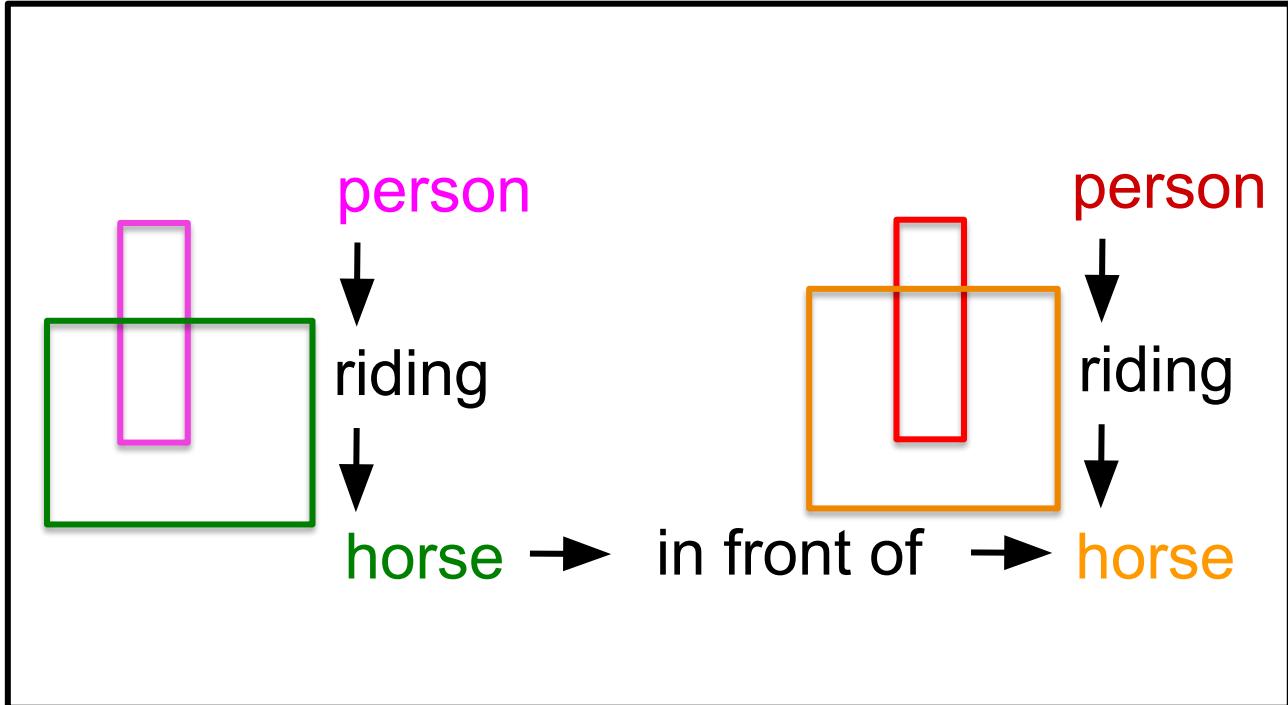
Output

Lu, Krishna et al., Visual Relationship Detection with Language Priors, ECCV 2016

Scene Graph Generation - Problem formulation



Input
(image only)



Output

Challenge 1:

Quadratic explosion of

- N objects,
- K relationships

leading to N^2K detectors



ride



falling off



next to



carry



lying



resting on



drag



throw

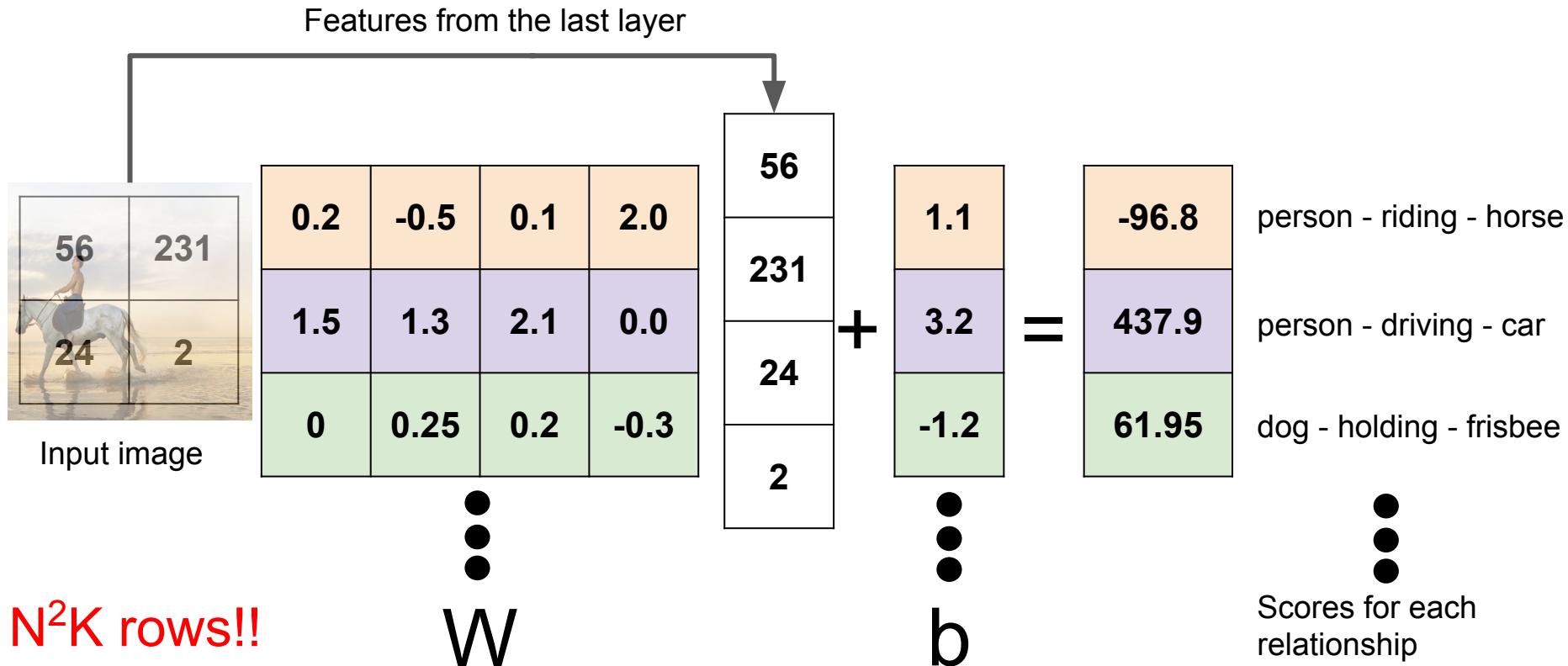
Visual Genome dataset

$N = 33K$

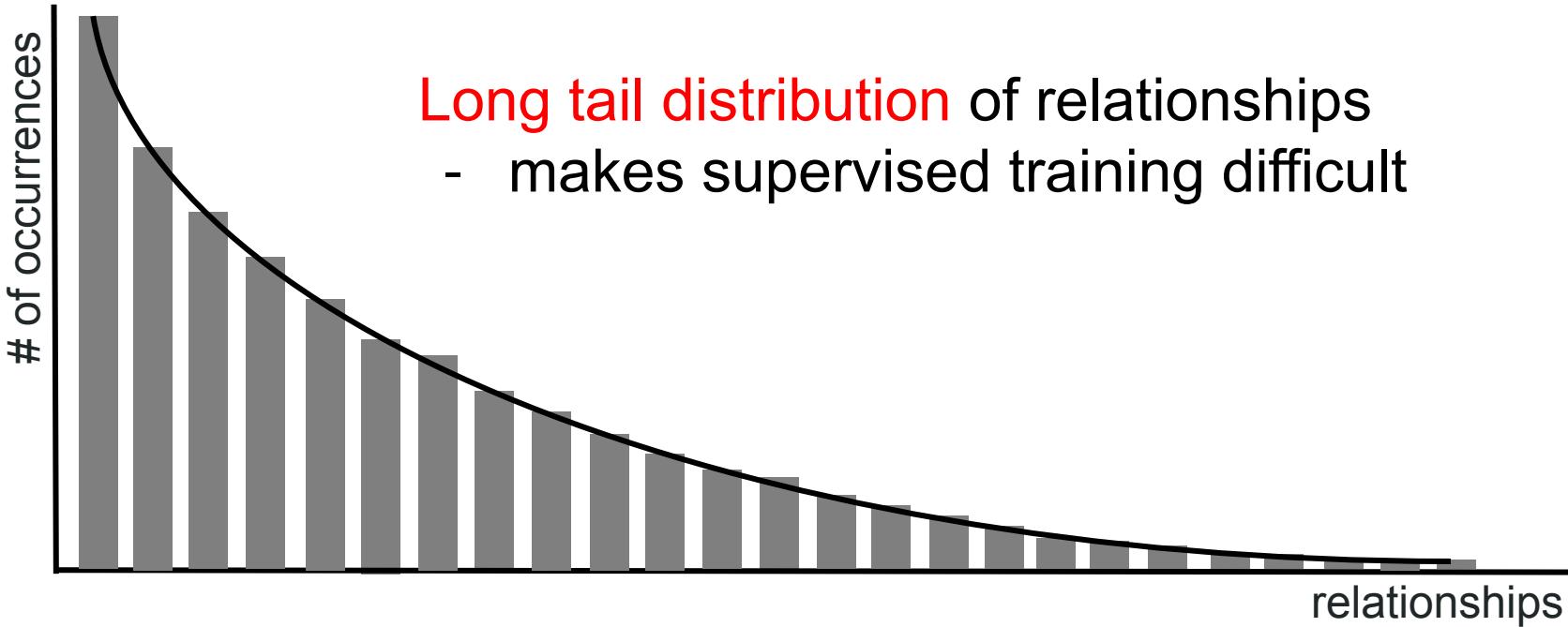
$K = 42K$

Lu, Krishna et al., Visual Relationship Detection with Language Priors, ECCV 2016

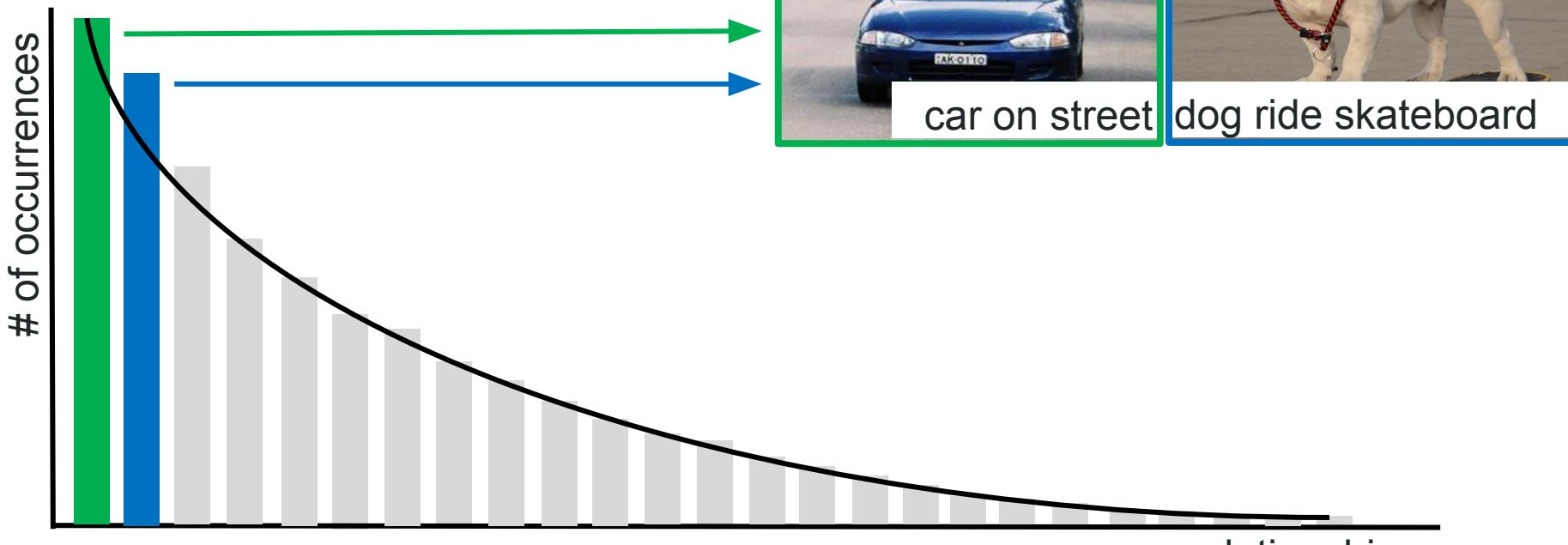
Recall the algebraic interpretation of linear models:



Challenge #2



Challenge #2

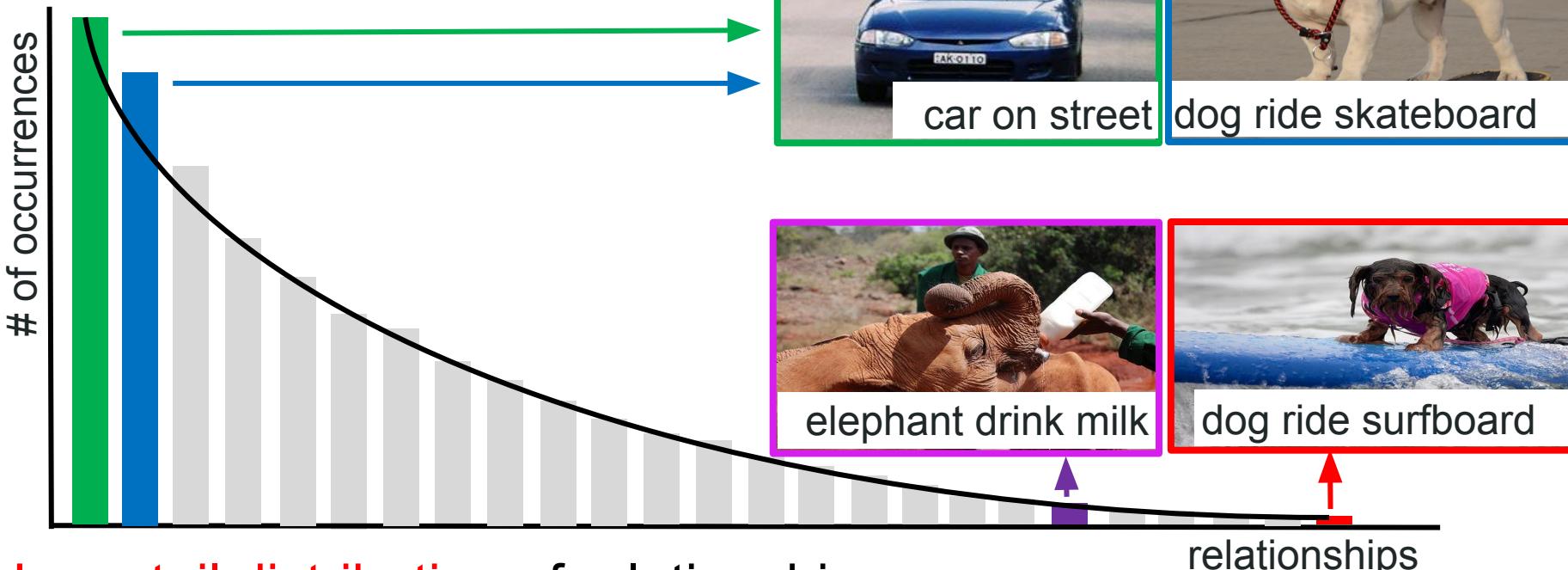


Long tail distribution of relationships

- makes supervised training difficult

Lu, Krishna et al., Visual Relationship Detection with Language Priors, ECCV 2016

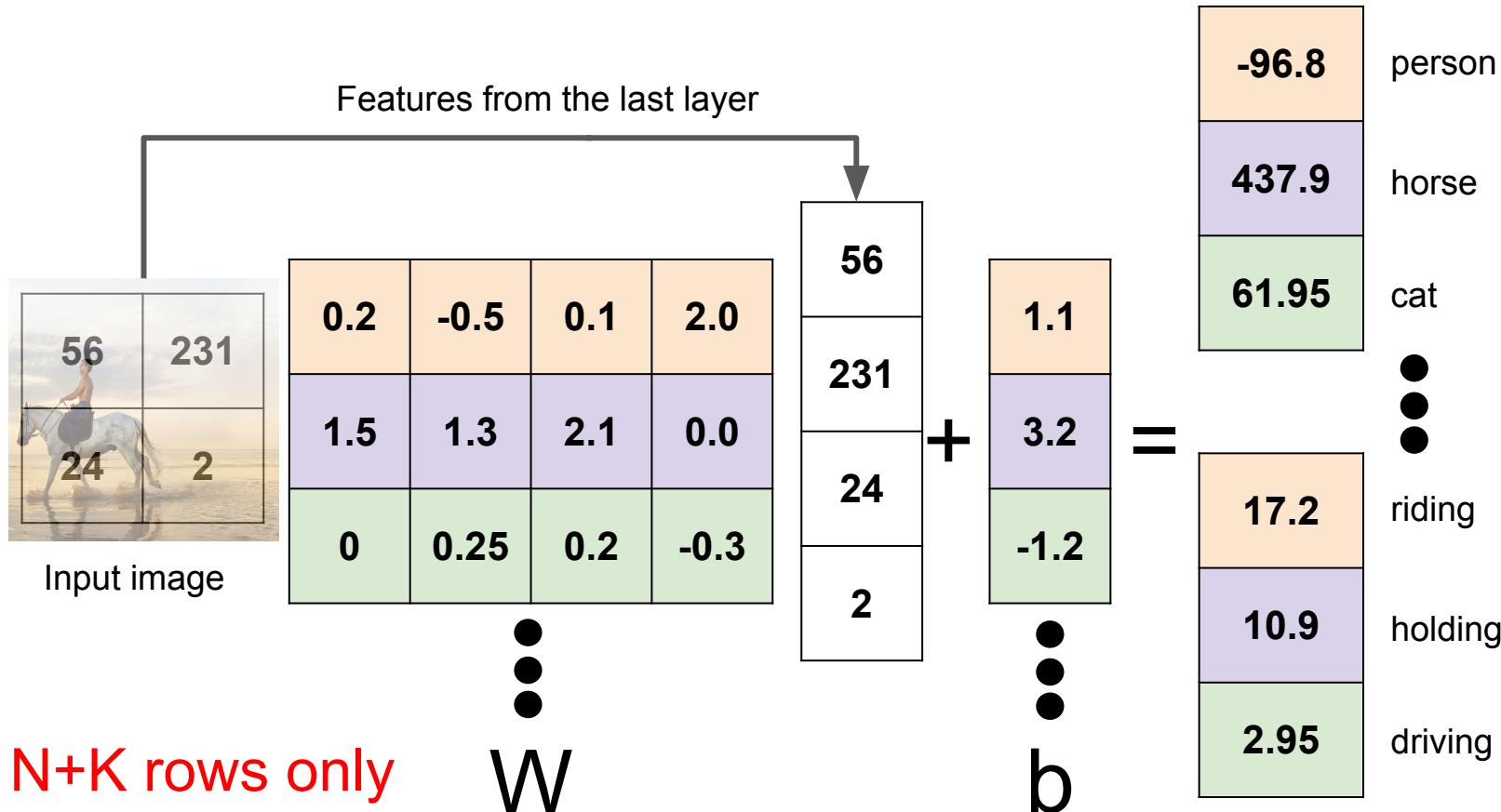
Challenge #2



Long tail distribution of relationships
- makes supervised training difficult

Lu, Krishna et al., Visual Relationship Detection with Language Priors, ECCV 2016

Intuition: Compose visual relationships from objects and predicates



N+K rows only

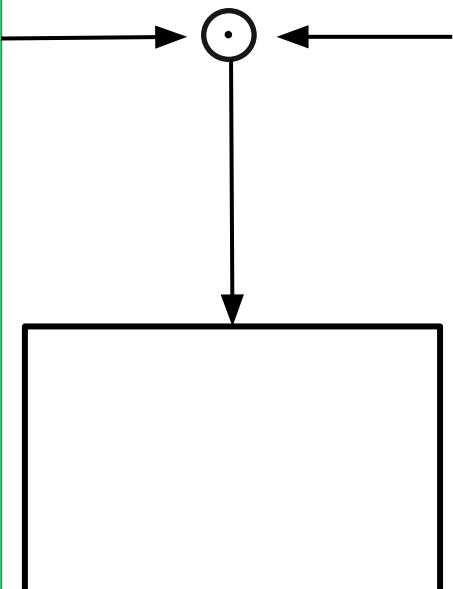
W

b

Visual module

Tackles:

Quadratic explosion of N^2K detectors



Language module

Tackles:

Long tail distribution of relationships

Visual module



Input

Definitions:

Visual module

Proposals:



Input

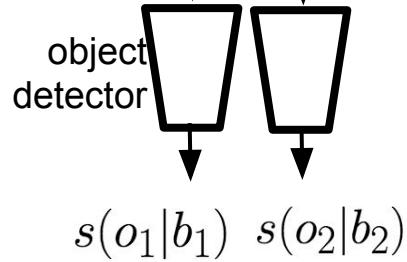
Definitions:
 b_1, b_2 are object proposals

Visual module

Proposals:



Sample: b_1 b_2

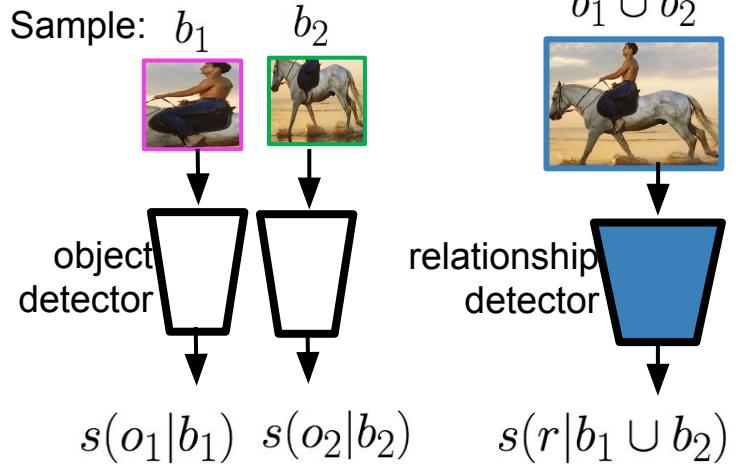


Input

Definitions:
 b_1, b_2 are object proposals
 $o_1, o_2 \in [\text{person, horse, ...}]$

Visual module

Proposals:

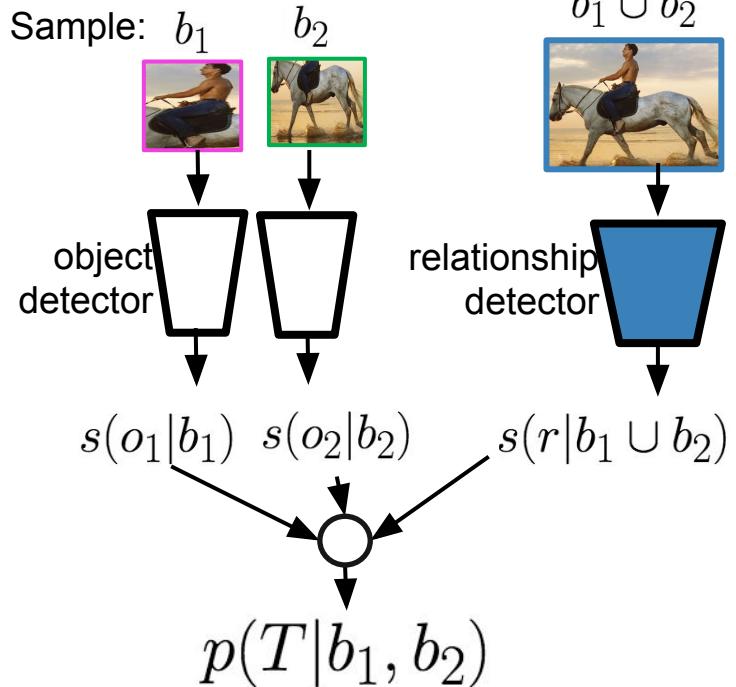


Input

Definitions:
 b_1, b_2 are object proposals
 $o_1, o_2 \in [\text{person, horse, ...}]$
 $r \in [\text{on, in, ride, front of, ...}]$

Visual module

Proposals:

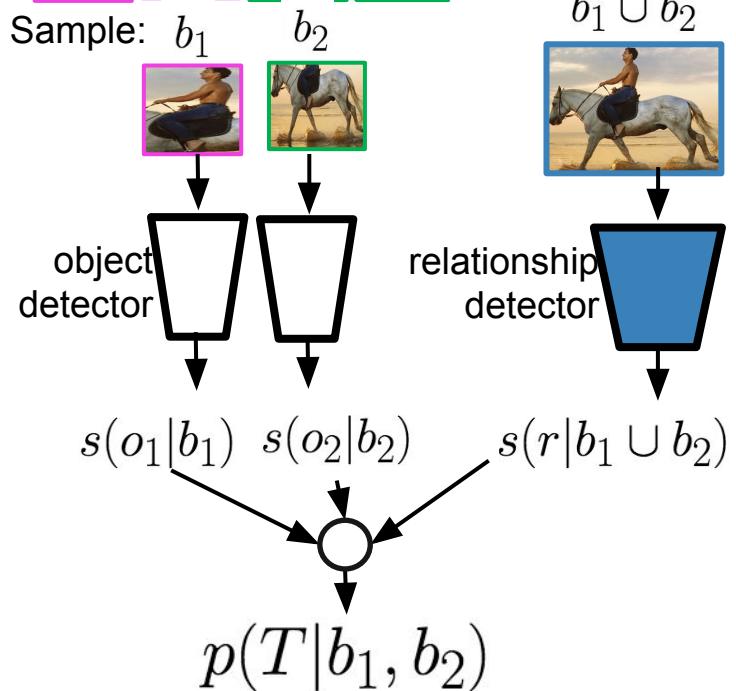


Definitions:

b_1, b_2 are object proposals
 $o_1, o_2 \in [\text{person, horse, ...}]$
 $r \in [\text{on, in, ride, front of, ...}]$
 T is a $\langle o_1, r, o_2 \rangle$ triple

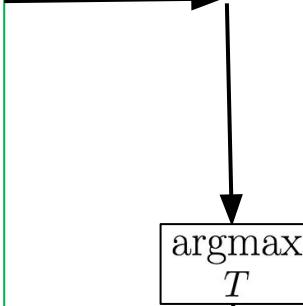
Visual module

Proposals:



$$p(T|b_1, b_2)$$

Input



Definitions:

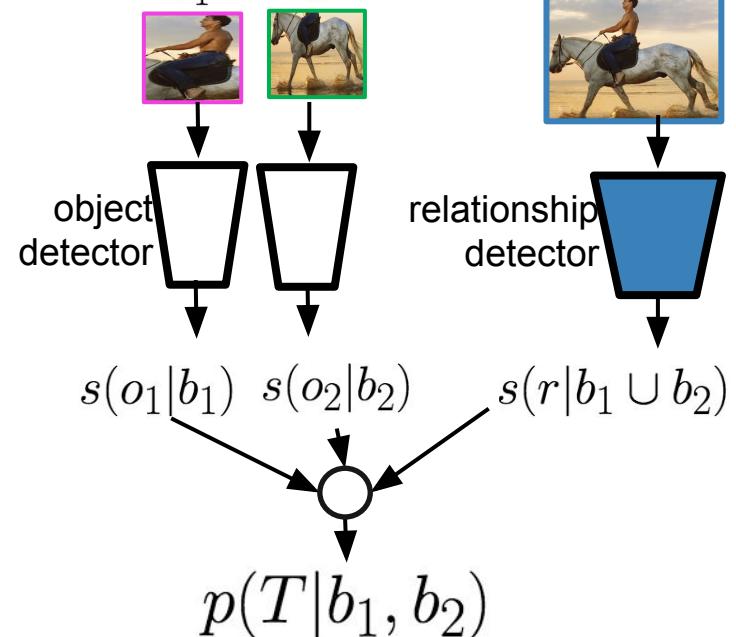
b_1, b_2 are object proposals
 $o_1, o_2 \in [\text{person}, \text{horse}, \dots]$
 $r \in [\text{on}, \text{in}, \text{ride}, \text{front of}, \dots]$
 T is a $\langle o_1, r, o_2 \rangle$ triple

Visual module

Proposals:



Sample: b_1 b_2



$$p(T|b_1, b_2)$$

Input

Language module

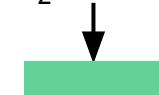
o_1 : person



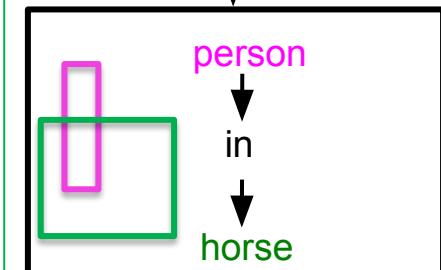
r : ride



o_2 : horse



$$\text{argmax}_T$$



Definitions:

b_1, b_2 are object proposals

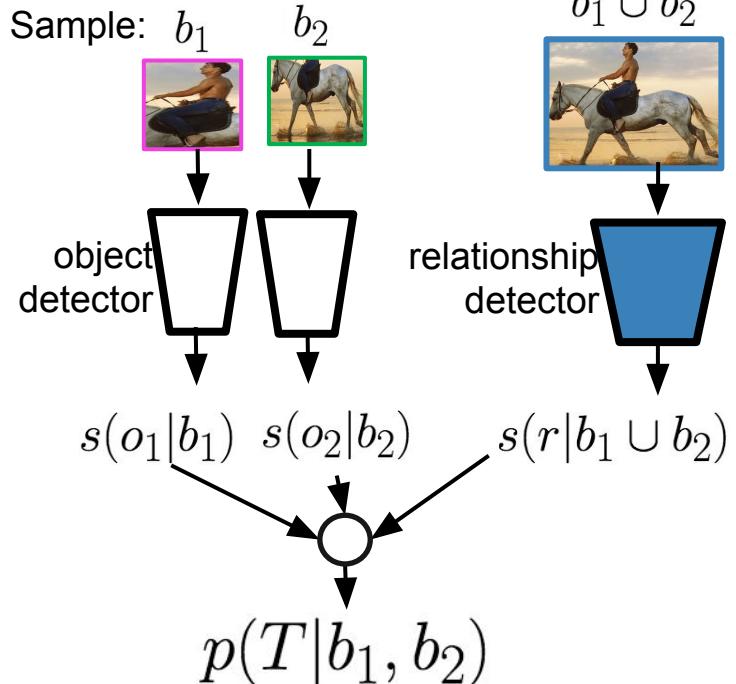
$o_1, o_2 \in [\text{person}, \text{horse}, \dots]$

$r \in [\text{on}, \text{in}, \text{ride}, \text{front of}, \dots]$

T is a $\langle o_1, r, o_2 \rangle$ triple

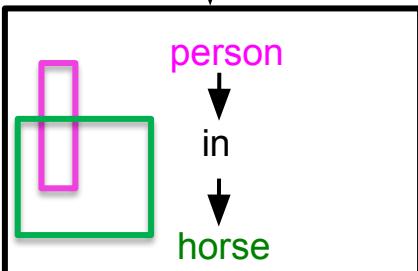
Visual module

Proposals:

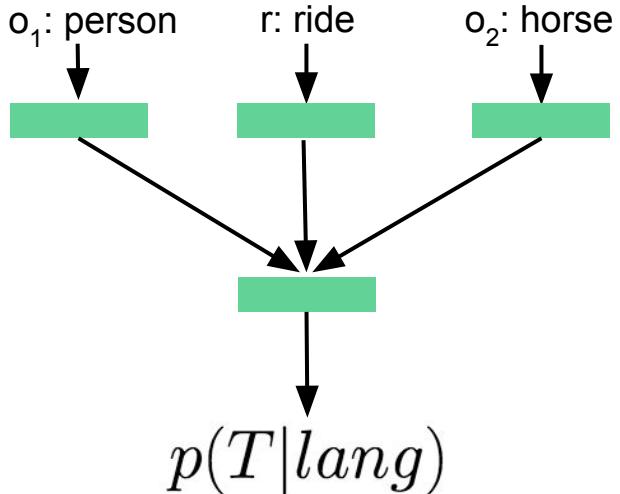


$$p(T|b_1, b_2)$$

argmax
 T



Language module



Definitions:

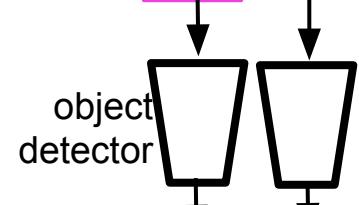
b_1, b_2 are object proposals
 $o_1, o_2 \in [\text{person}, \text{horse}, \dots]$
 $r \in [\text{on}, \text{in}, \text{ride}, \text{front of}, \dots]$
 T is a $\langle o_1, r, o_2 \rangle$ triple

Visual module

Proposals:



Sample: b_1 b_2



$$b_1 \cup b_2$$



$$s(o_1|b_1) \quad s(o_2|b_2)$$

relationship
detector

$$s(r|b_1 \cup b_2)$$

$$p(T|b_1, b_2)$$

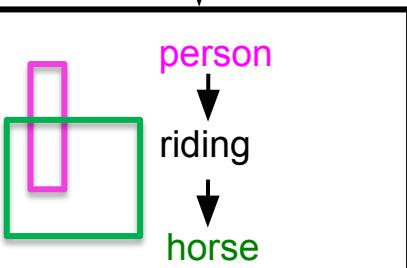


$$p(T|b_1, b_2)$$

$$\odot$$

$$p(T|b_1, b_2, lang)$$

$$\underset{T}{\operatorname{argmax}}$$

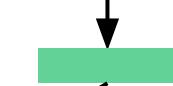
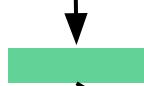


Language module

o_1 : person

r : ride

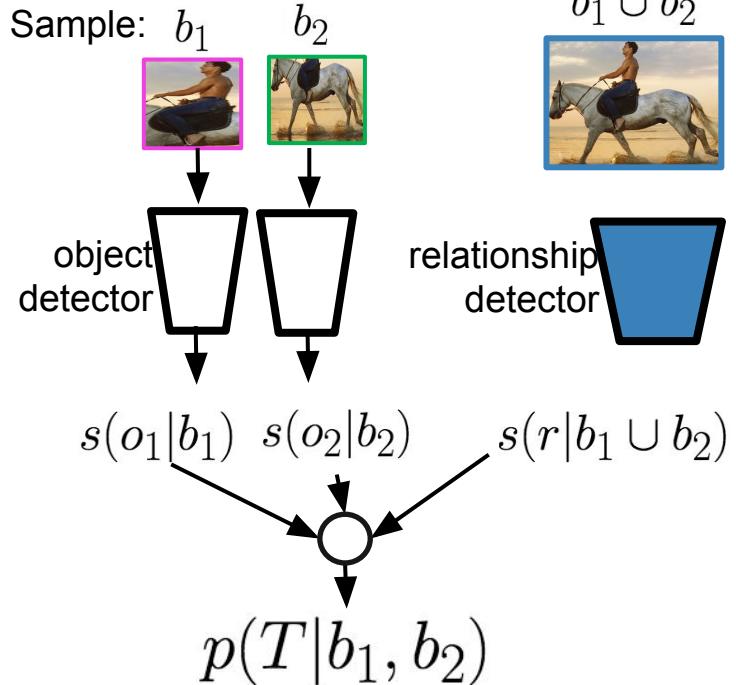
o_2 : horse



$$p(T|lang)$$

Visual module

Proposals:



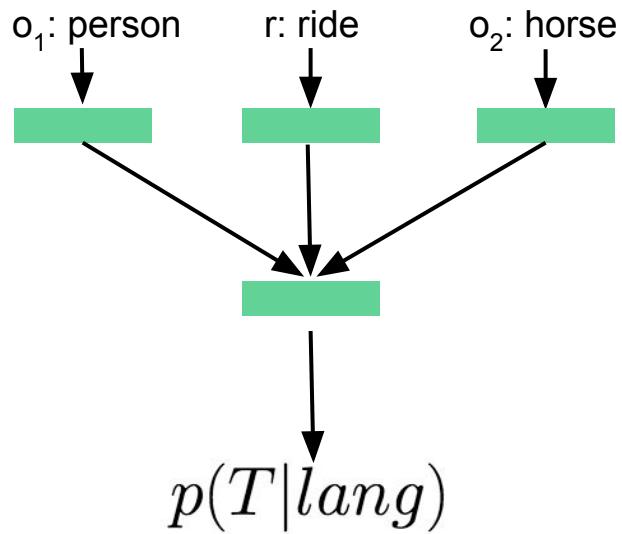
Tackles:

**Quadratic explosion
only requires $N+K$
detectors**

Tackles:

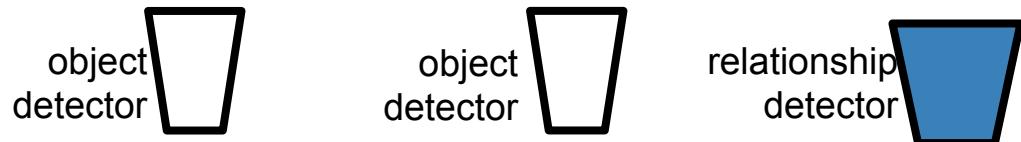
Long tail distribution
can predict rare
relationships

Language module



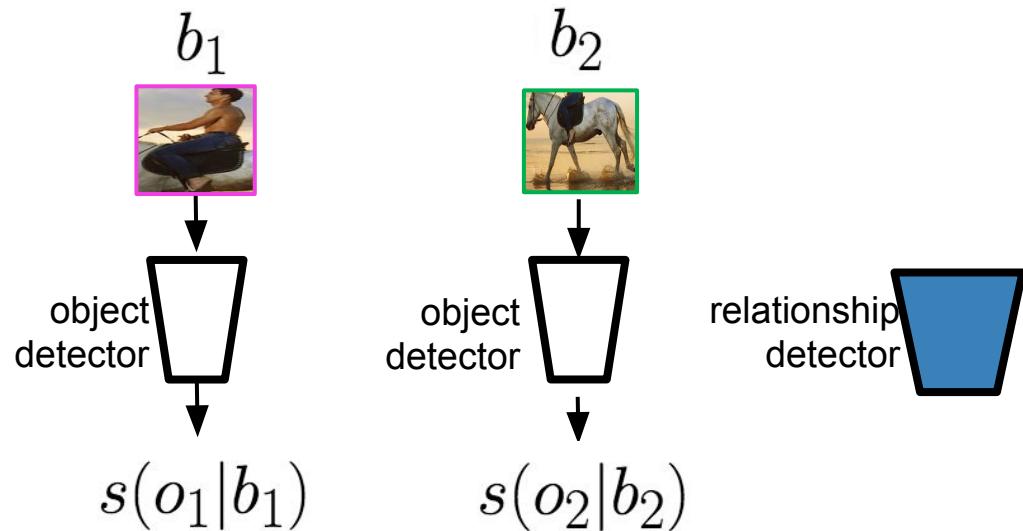
Training the visual module

1. Pre-train using ImageNet



Definitions:

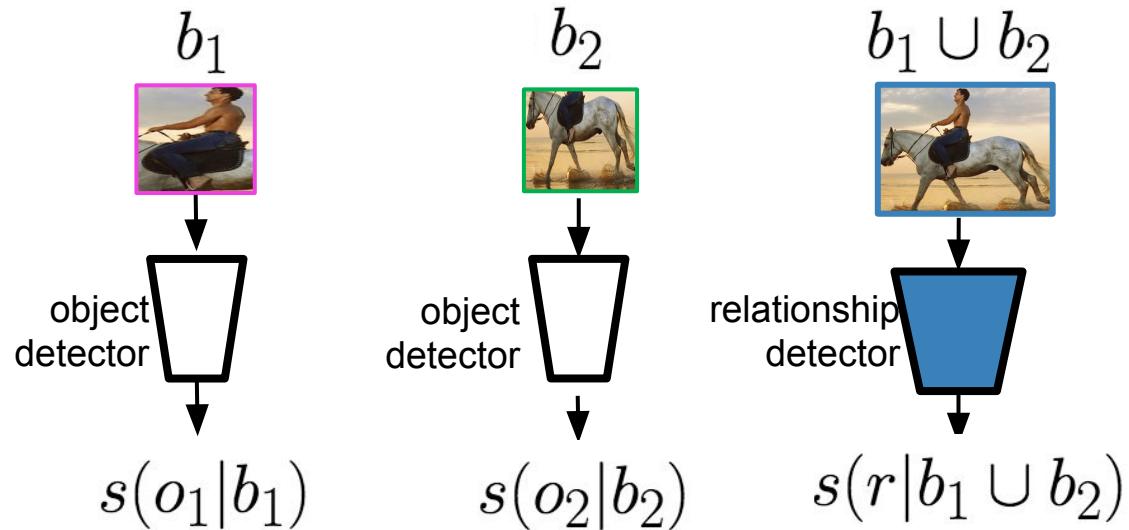
Training the visual module



1. Pre-train using ImageNet
2. Train object detector

Definitions:
 b_1, b_2 are object proposals
 $o_1, o_2 \in [\text{person}, \text{horse}, \dots]$

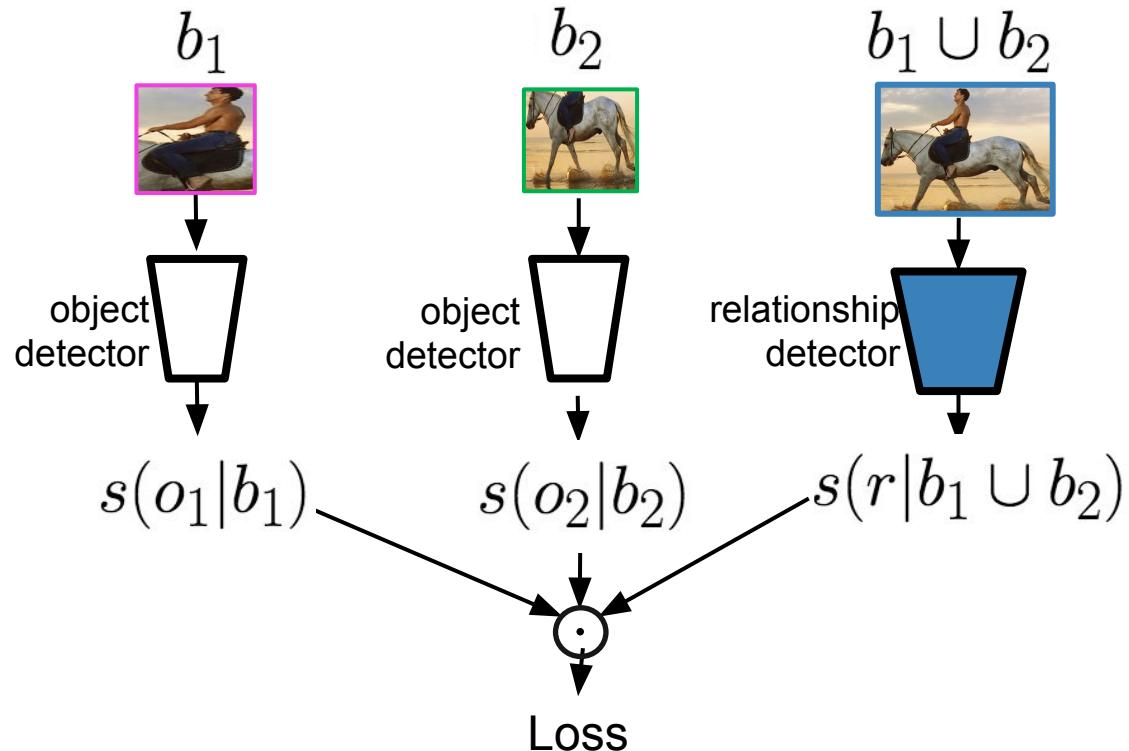
Training the visual module



1. Pre-train using ImageNet
2. Train object detector
3. Train relationship detector

Definitions:
 b_1, b_2 are object proposals
 $o_1, o_2 \in [\text{person, horse, ...}]$
 $r \in [\text{on, in, ride, front of, ...}]$

Training the visual module

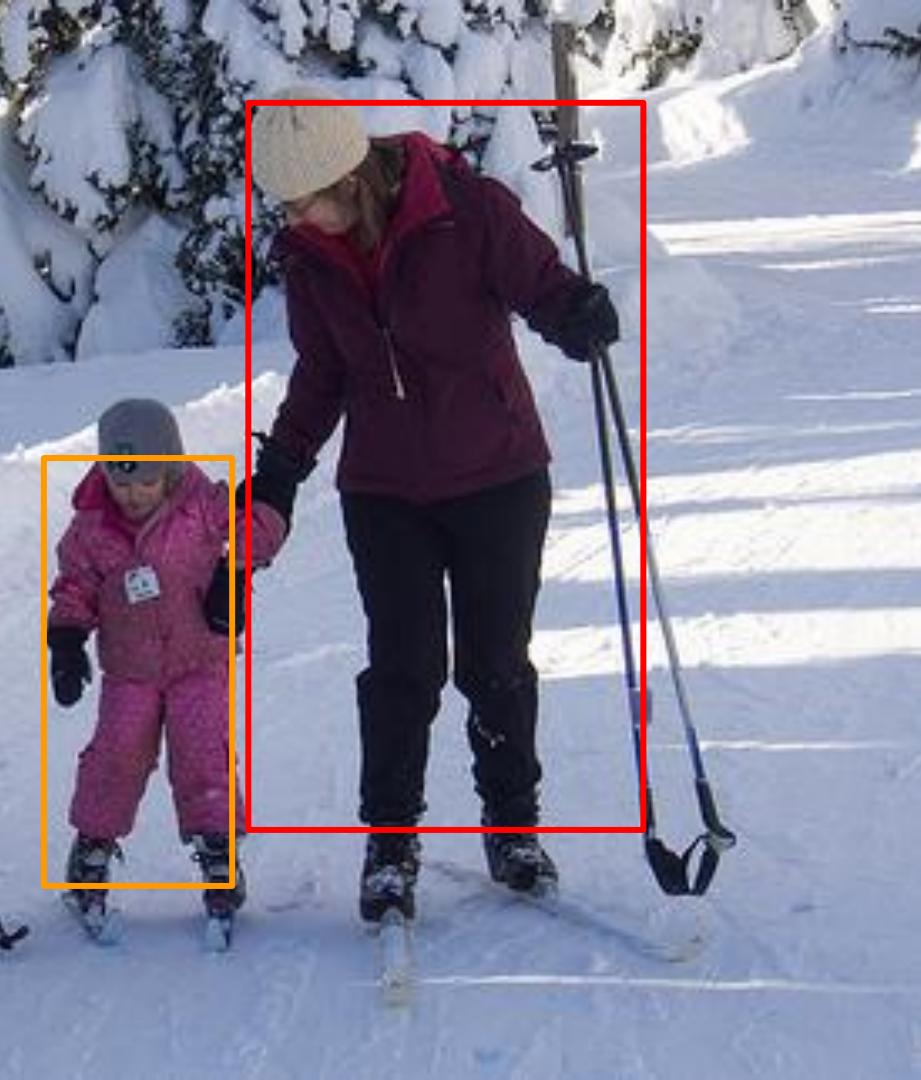


1. Pre-train using ImageNet
2. Train object detector
3. Train relationship detector
4. Fine-tune both jointly

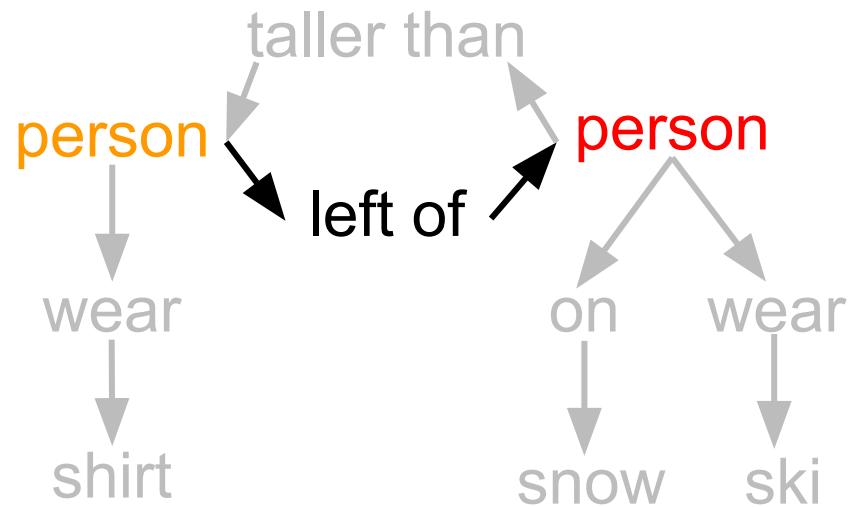
Definitions:
 b_1, b_2 are object proposals
 $o_1, o_2 \in [\text{person, horse, ...}]$
 $r \in [\text{on, in, ride, front of, ...}]$



Our results:



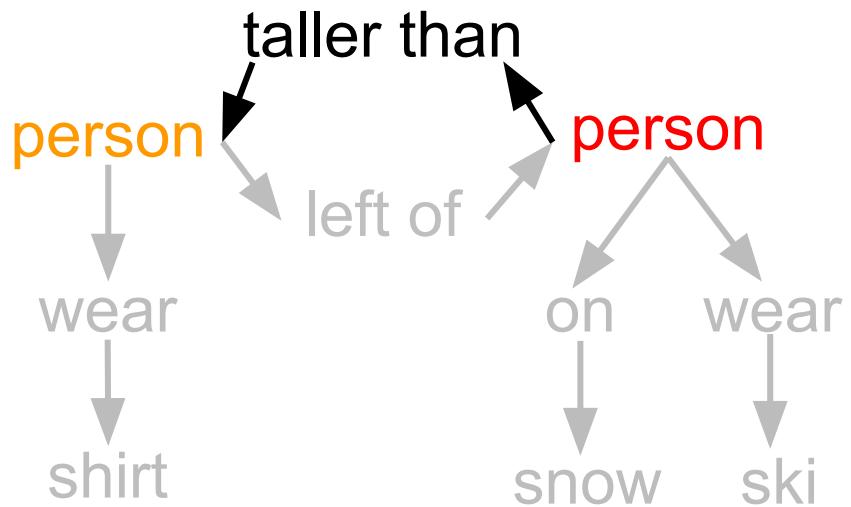
Our results:
spatial, comparative, asymmetrical,
verb, prepositional

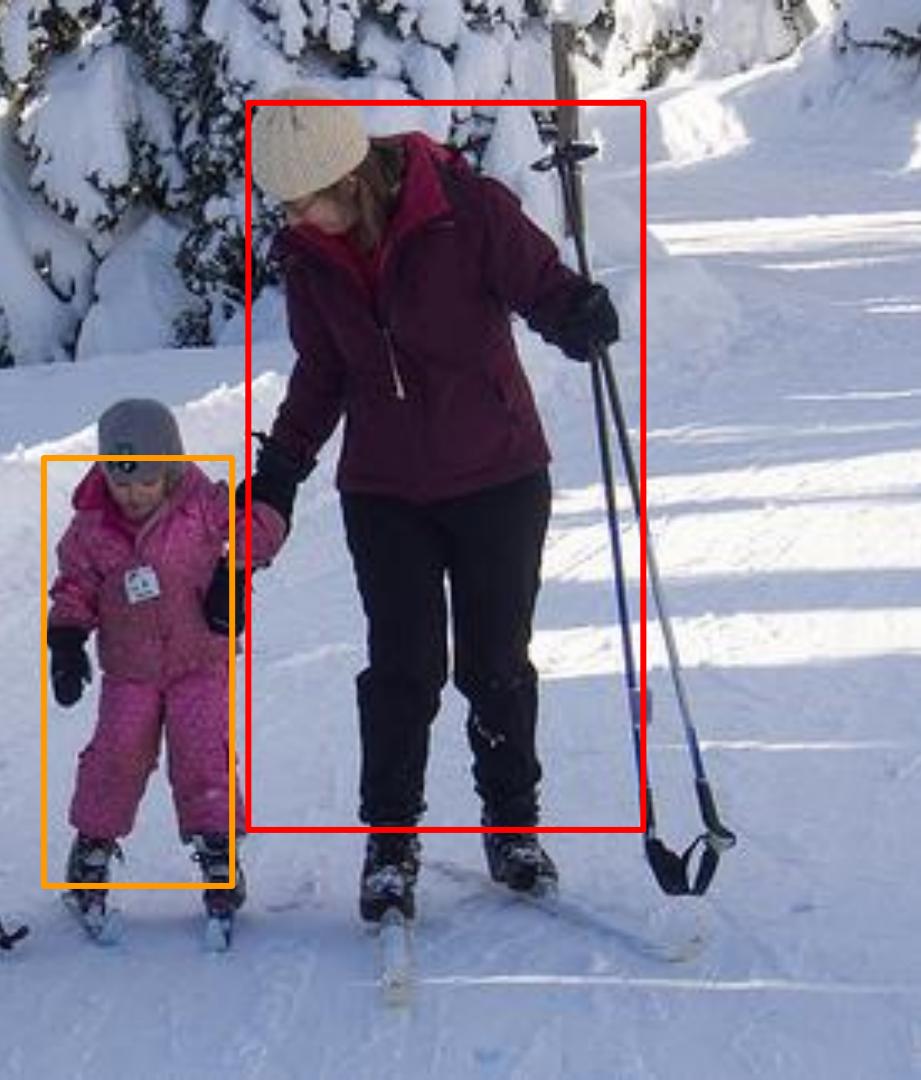




Our results:

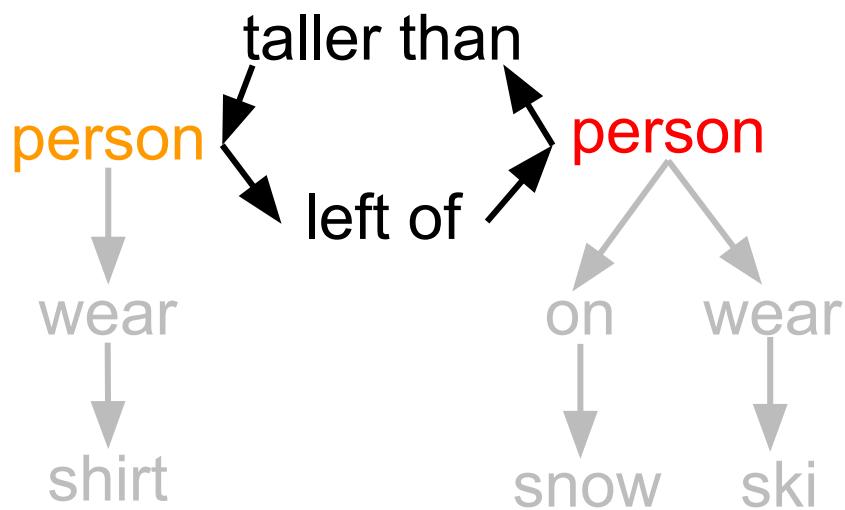
spatial, comparative, asymmetrical,
verb, prepositional

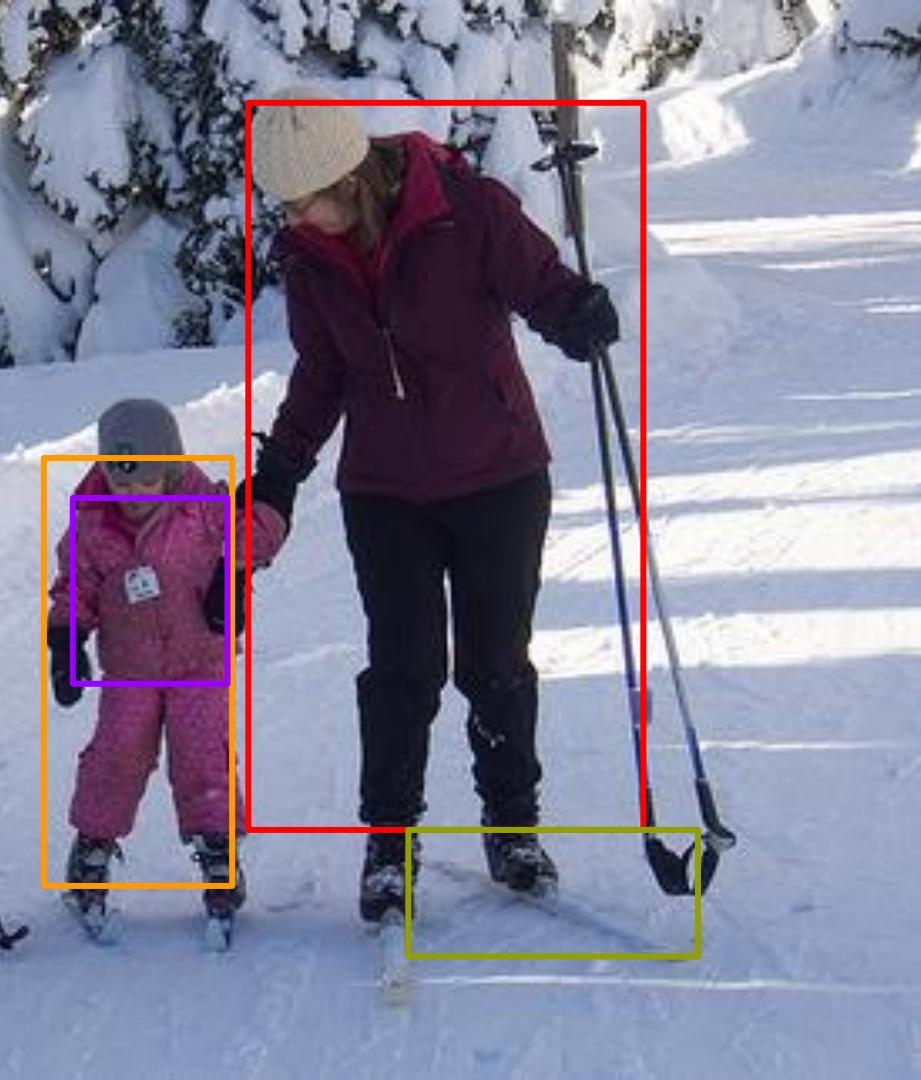




Our results:

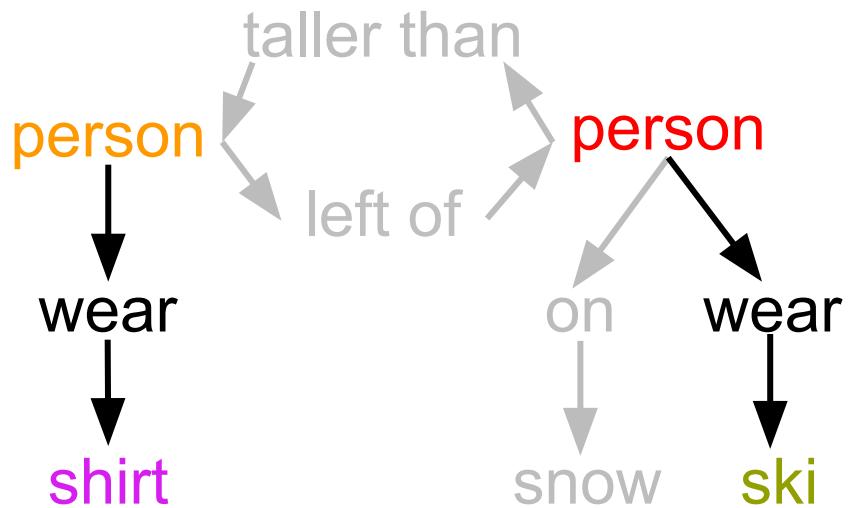
spatial, comparative, asymmetrical,
verb, prepositional





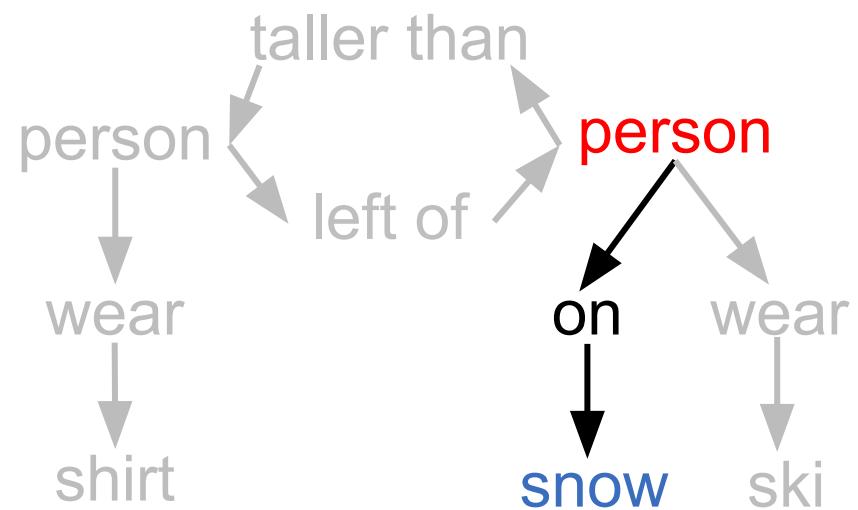
Our results:

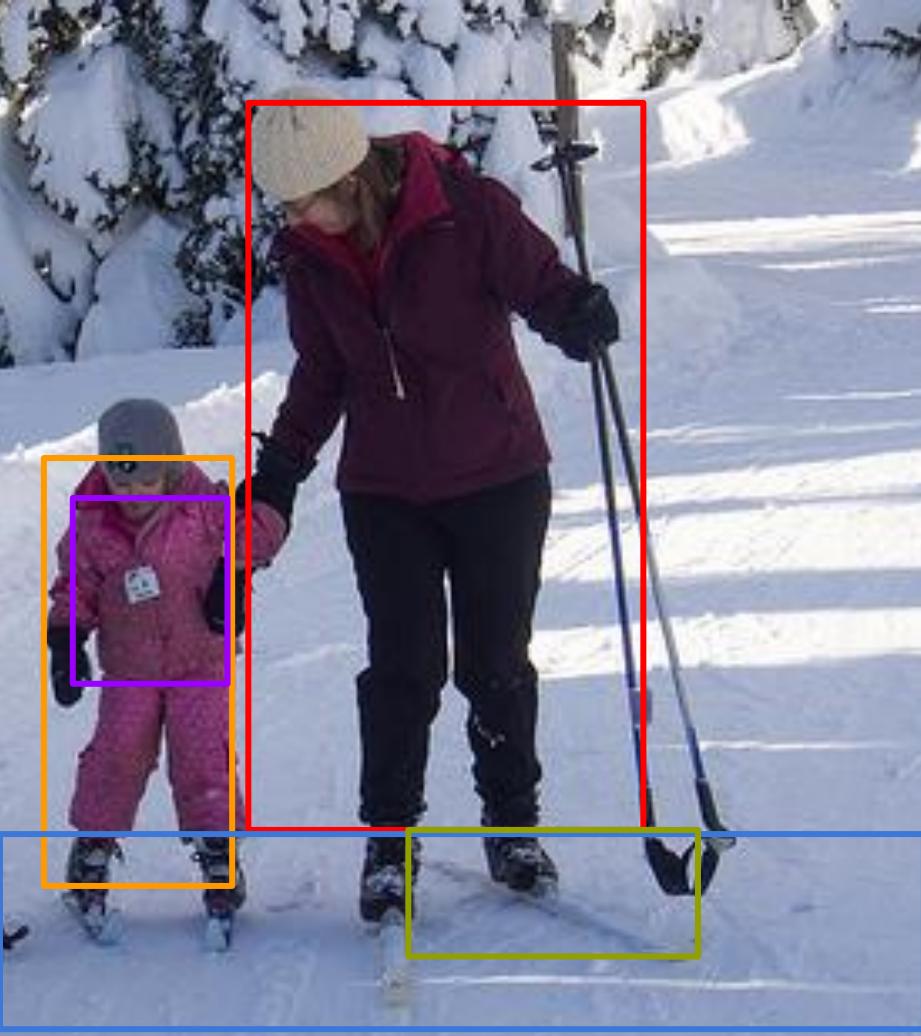
spatial, comparative, asymmetrical,
verb, prepositional



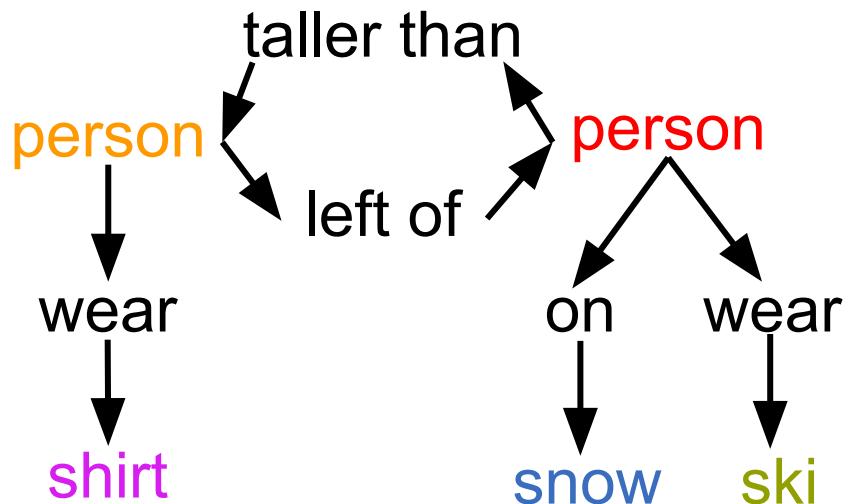


Relationship types:
spatial, comparative, asymmetrical,
verb, prepositional





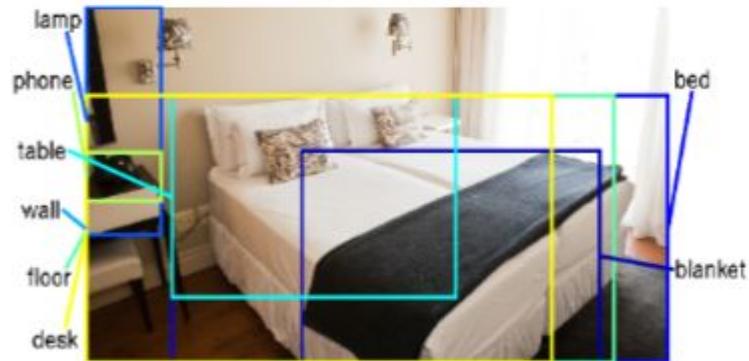
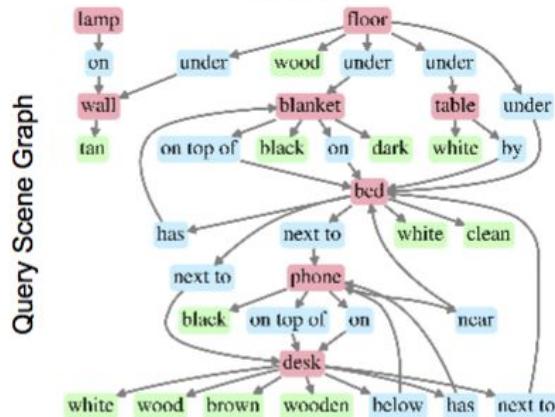
Our results:
spatial, comparative, asymmetrical,
verb, prepositional



Scene graphs can improve image retrieval

Sentence Description

Black phone is on top of white, wooden desk. The desk is next to a clean white bed that has a black blanket and is next to a white table. The lamp is on a tan wall. The table is by the bed, which is next to the phone. The floor is under the bed, table, lamp and blanket.



Johnson, Krishna et al., Image Retrieval using Scene Graphs CVPR, 2015

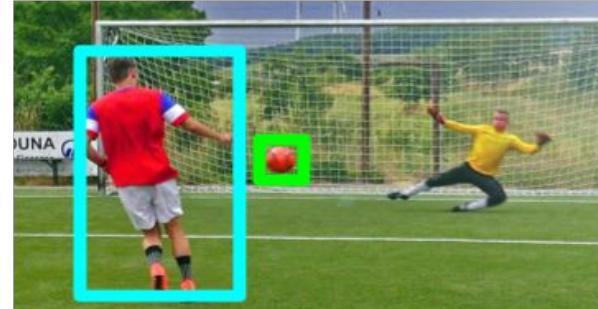
Schuster, Krishna, et al., Generating Semantically Precise Scene Graphs from Textual Descriptions for Improved Image Retrieval, EMNLP 2015 workshop

Modeling relationships can improve existing vision tasks like **object localization**

Input



Output



Krishna et al., Referring Relationships CVPR, 2018

Zero shot detection



person sit chair
948 training examples

hydrant on ground
29 training examples

Zero shot detection



person sit chair
948 training examples

hydrant on ground
29 training examples



person sit hydrant
0 training examples

Zero shot detection



person ride horse
578 training examples

person wear hat
1023 training examples

Zero shot detection



person ride horse
578 training examples

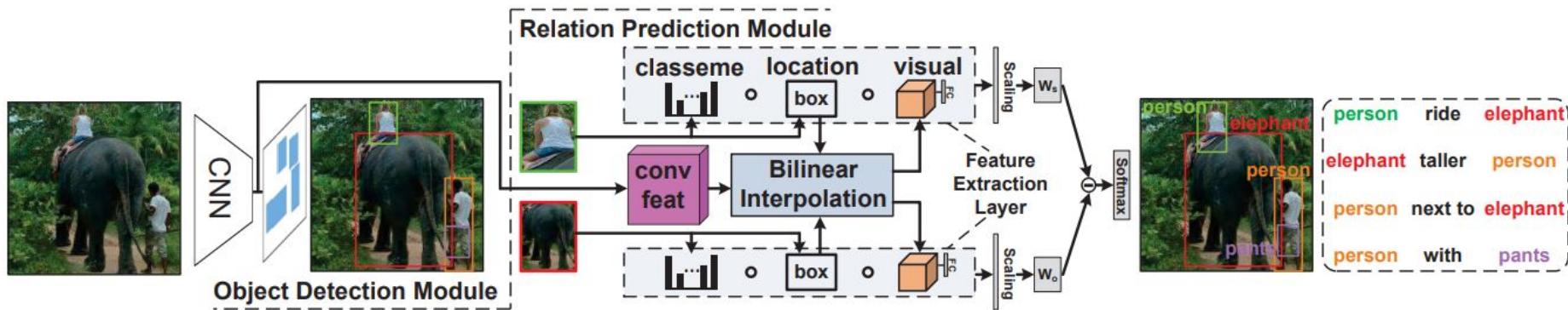


person wear hat
1023 training examples



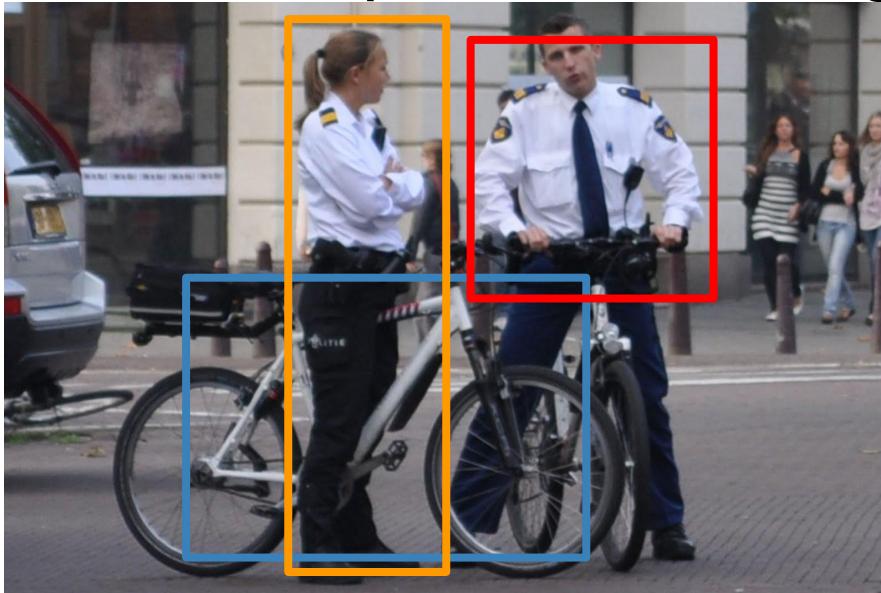
horse wear hat
0 training examples

Incorporating spatial features and classemes



Zhang, Hanwang, et al. "Visual translation embedding network for visual relation detection CVPR 2017
Copyright Zellers. Reproduced with permission.

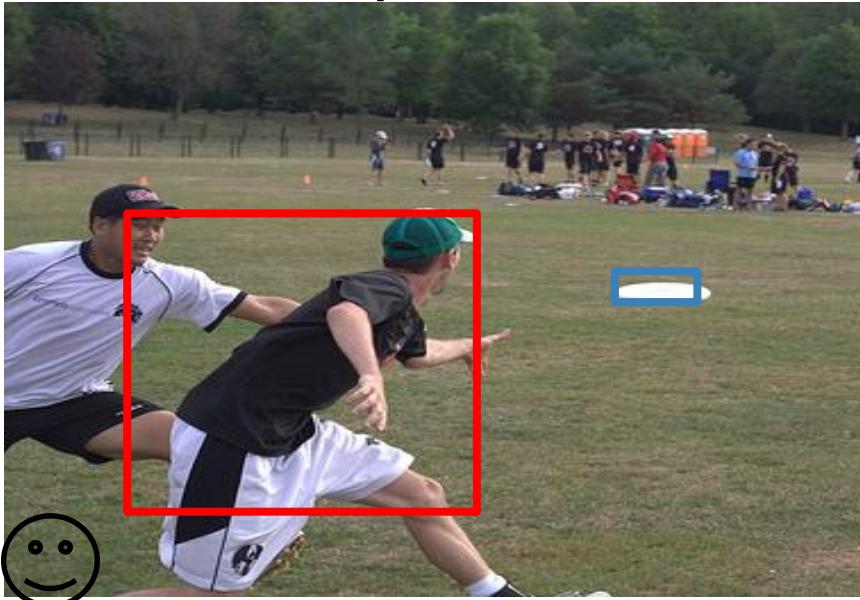
Problem with current method: Doesn't consider other relationships when making predictions



person ride bicycle
person ride bicycle



Problem with current method: Doesn't consider other relationships when making predictions



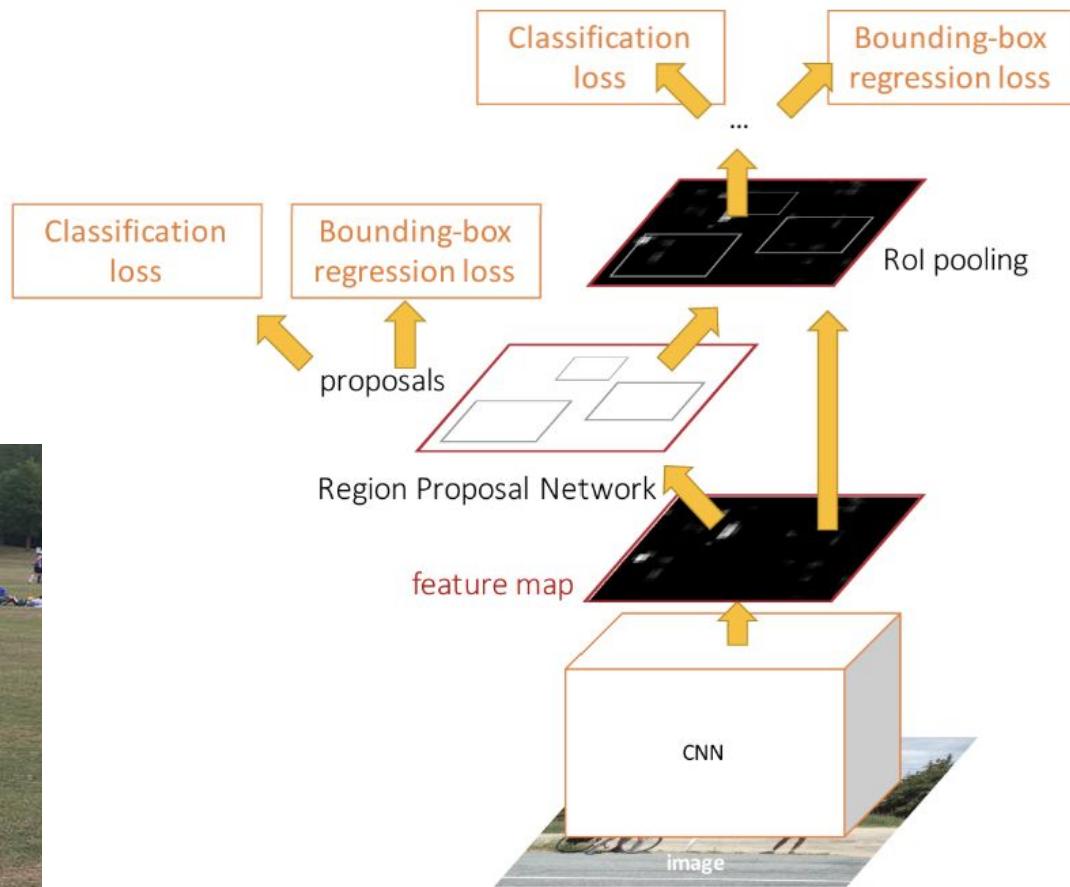
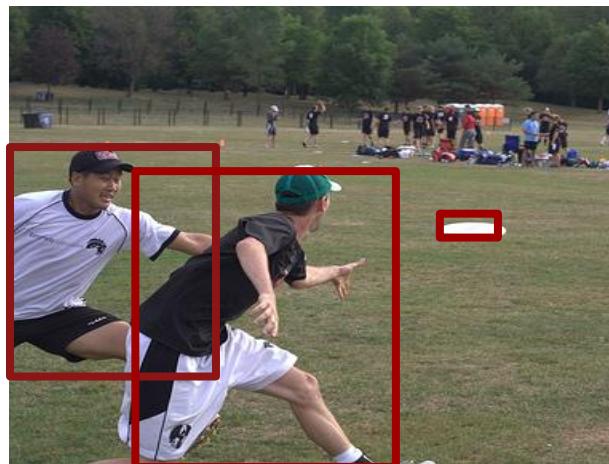
person throw frisbee



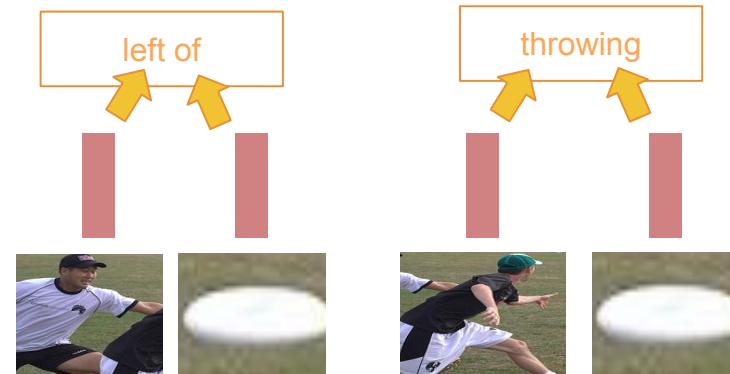
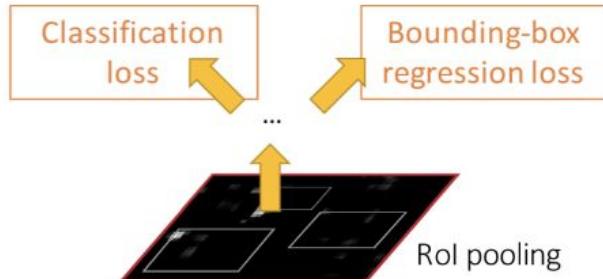
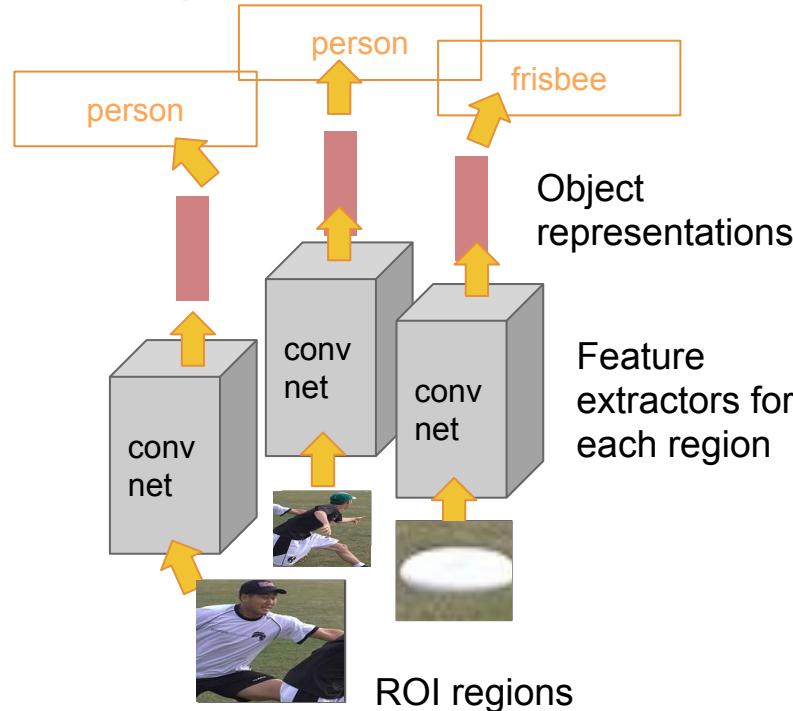
person throw frisbee

How do we model the other relationships in the image when making a prediction for a given relationship?

Recall Faster RCNN

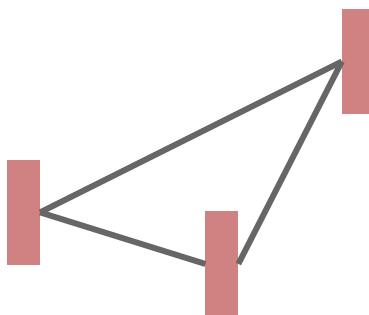


Each prediction in isolation



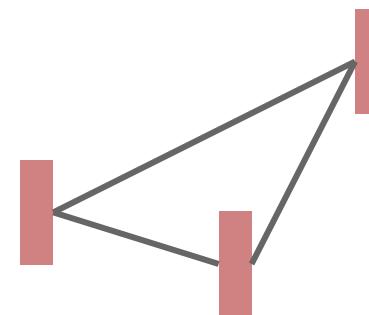
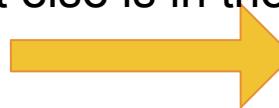
Representing objects as a graph with pairwise connections

But this graph doesn't encode
the different kinds of
relationships



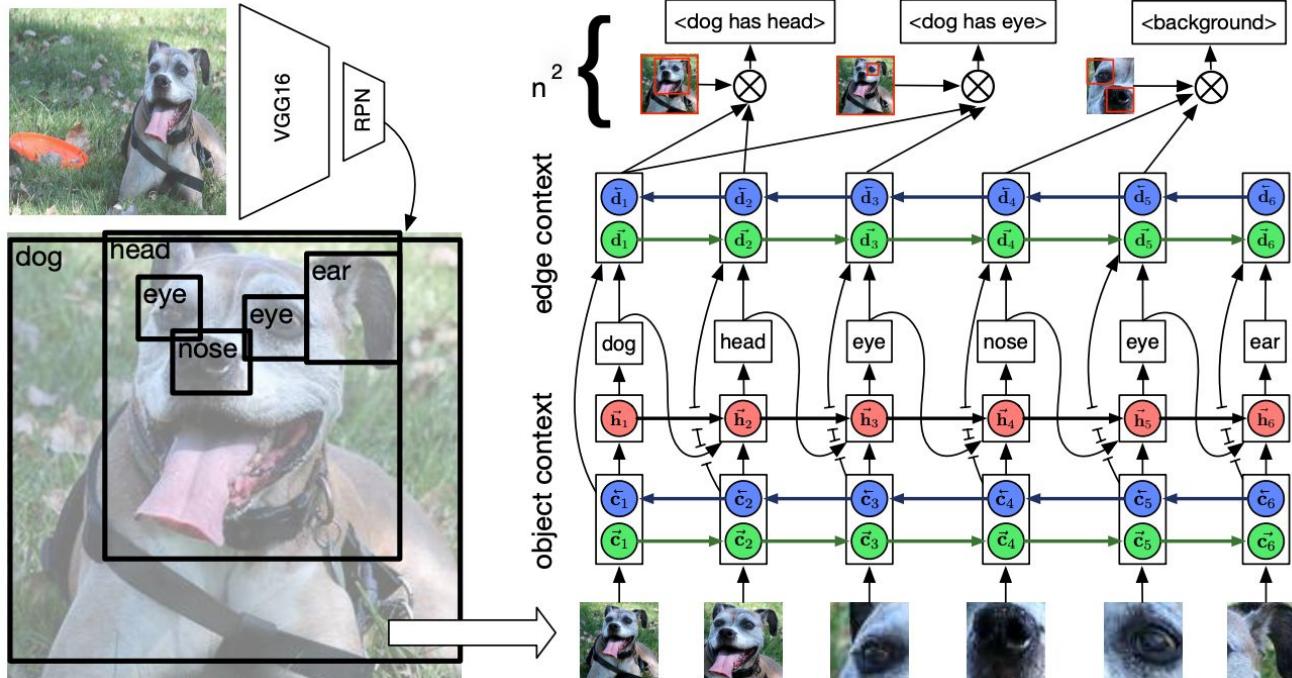
Each node contains features
from individual regions

Perform some operation that
allows each node to encode
what else is in the image.



Each node now also contains
features from all other regions

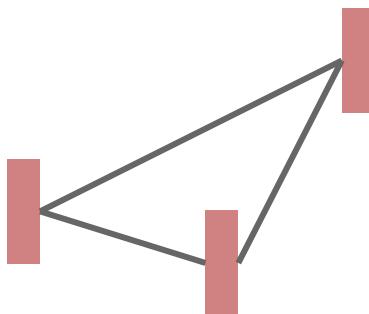
Use an RNN to collect information? But order of objects impacts predictions



Zellers et al. "Neural motifs: Scene graph parsing with global context." CVPR 2018
Copyright Zellers. Reproduced with permission.

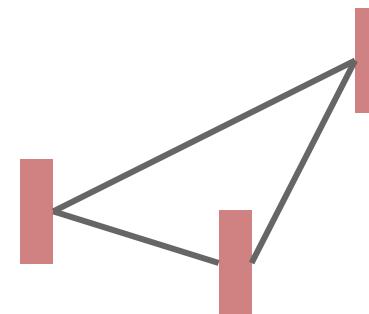
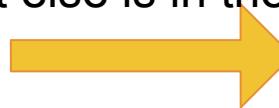
Representing objects as a graph with pairwise connections

But this graph doesn't encode
the different kinds of
relationships



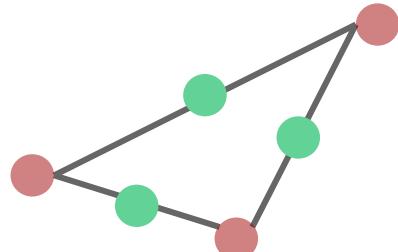
Each node contains features
from individual regions

Perform some operation that
allows each node to encode
what else is in the image.



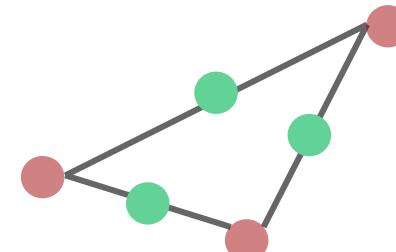
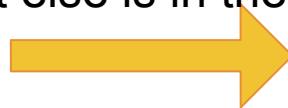
Each node now also contains
features from all other regions

Graph representation with relationships included as nodes



Each node contains features
from individual regions

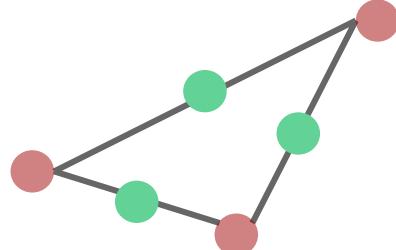
Perform some operation that
allows each node to encode
what else is in the image.



Each node now also contains
features from all other regions

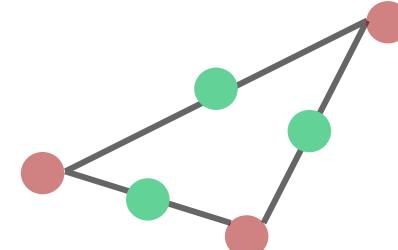
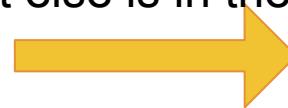
Graph representation with edges included as nodes

What operation have we already seen that updates features in a graph?



Each node contains features
from individual regions

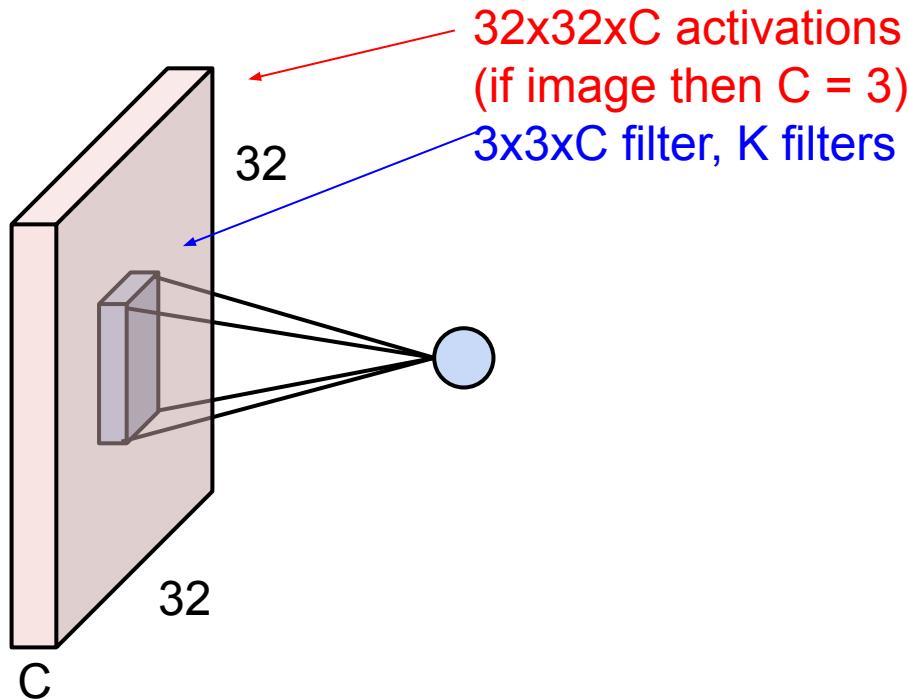
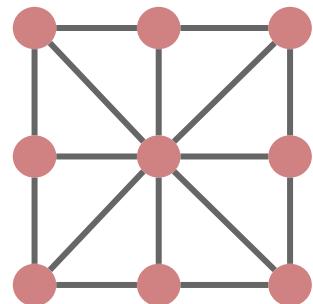
Perform some operation that
allows each node to encode
what else is in the image.



Each node now also contains
features from all other regions

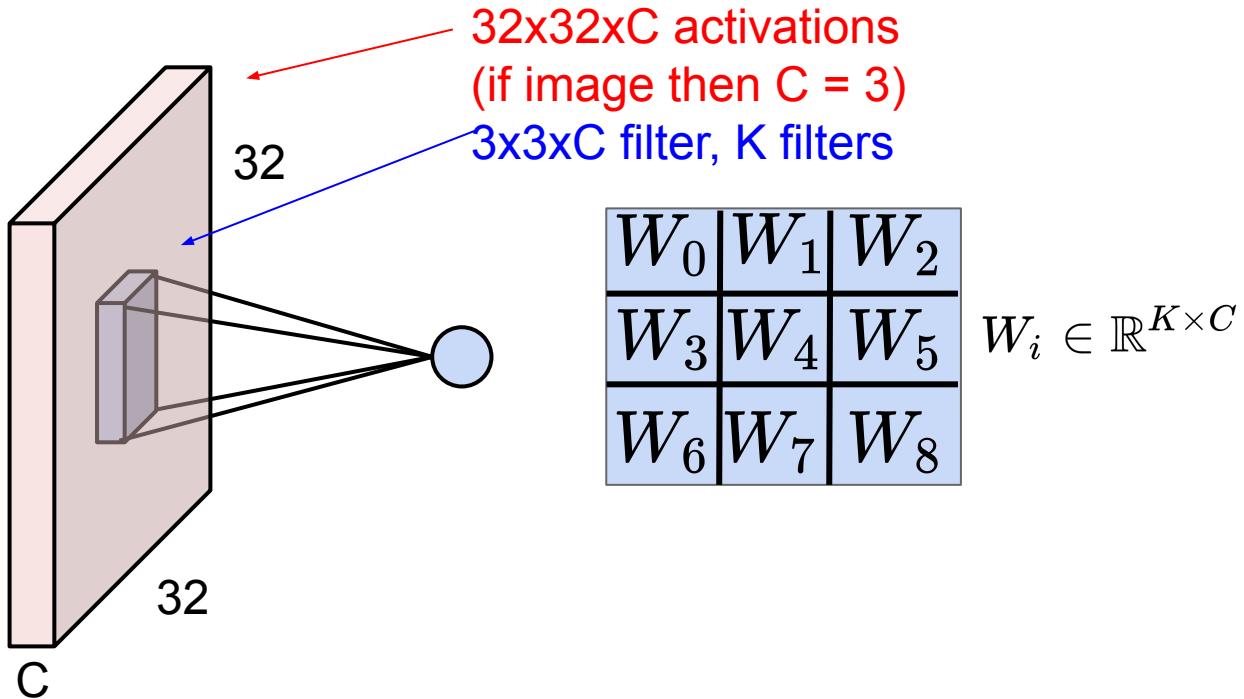
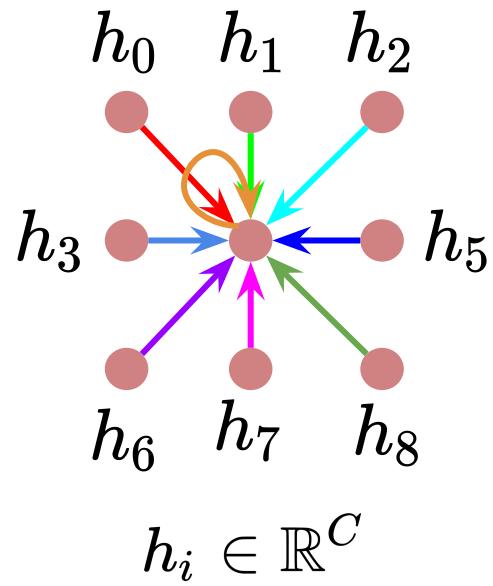
Recall Convolutions

Images are a **structured graph of pixels!**



Recall Convolutions

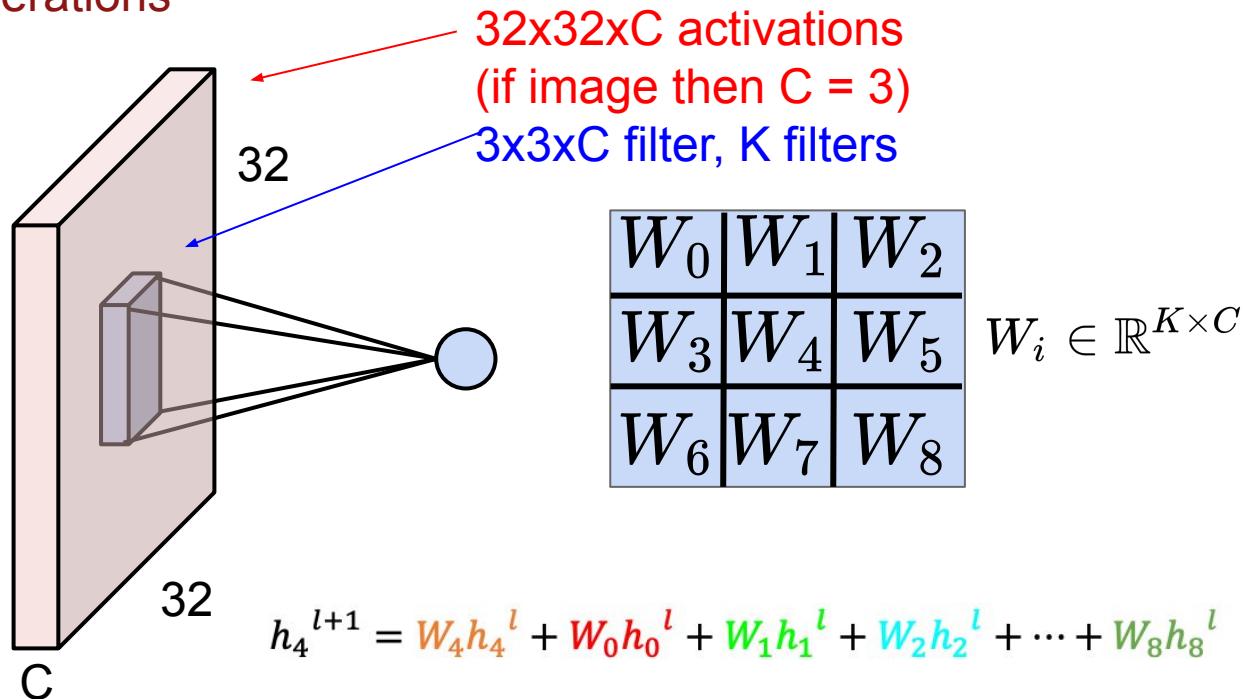
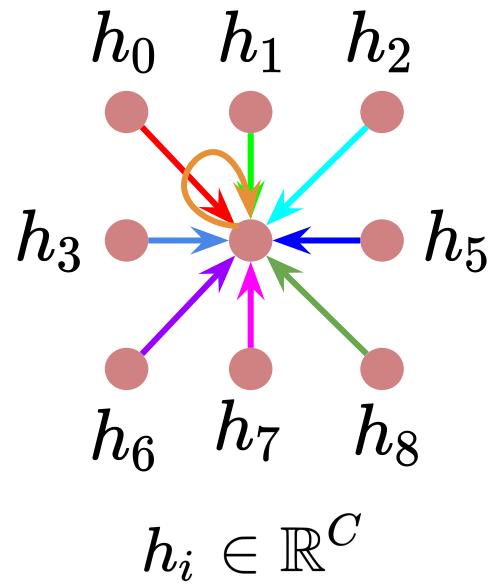
Images are a **structured graph of pixels!**



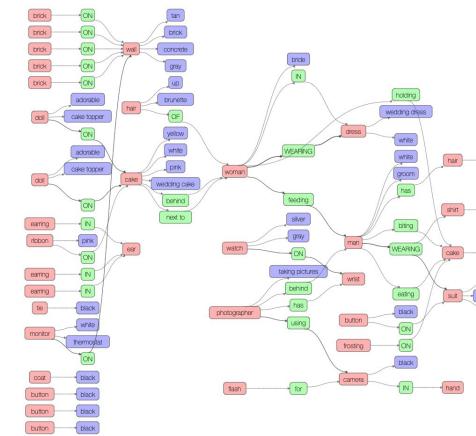
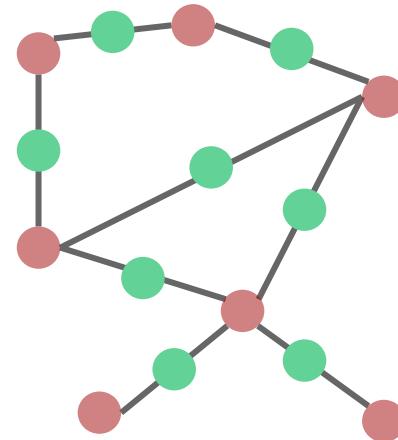
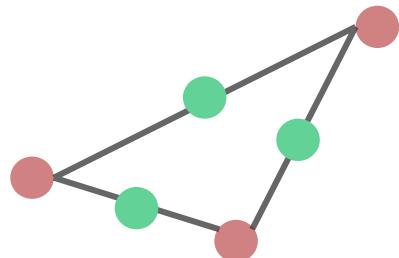
Recall Convolutions

Images are a **structured graph of pixels!**

Convolutions are **local operations**

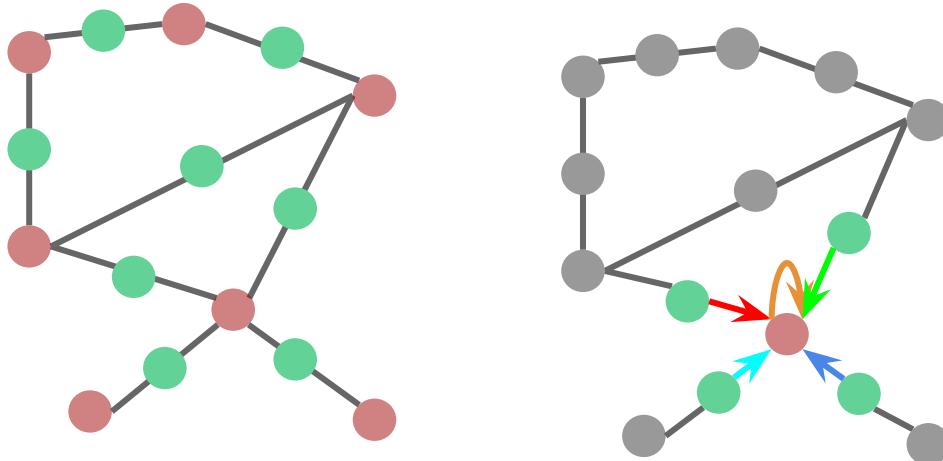


In comparison, scene graphs are not uniformly structured



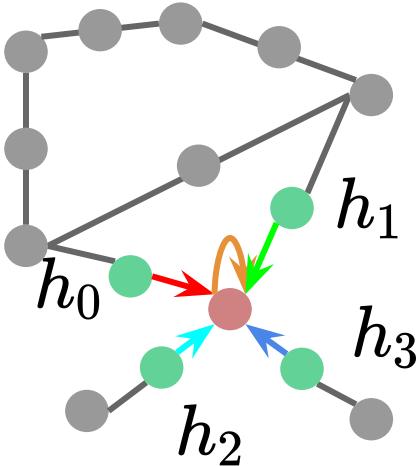
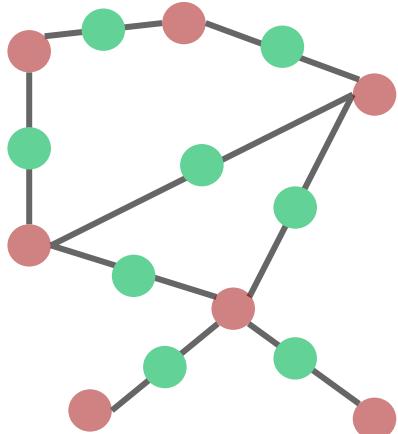
Objects have varying number of relationships

Generalizing 2D convolutions to Graph Convolutions



- Graph convolutions involve similar **local operations** on nodes.
- The **ordering of neighbors** should not matter.
- The **number of neighbors** should not matter.

Generalizing 2D convolutions to Graph Convolutions

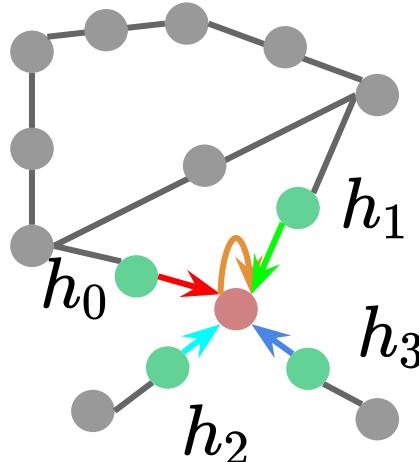
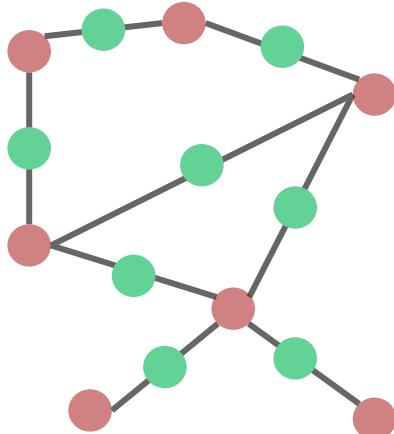


- Graph convolutions involve similar **local operations** on nodes.
- The **ordering of neighbors** should not matter.
- The **number of neighbors** should not matter.

$$h_4^{l+1} = W_4 h_4^l + W_0 h_0^l + W_1 h_1^l + W_2 h_2^l + W_3 h_3^l$$

But in this formulation the ordering matters

Generalizing 2D convolutions to Graph Convolutions

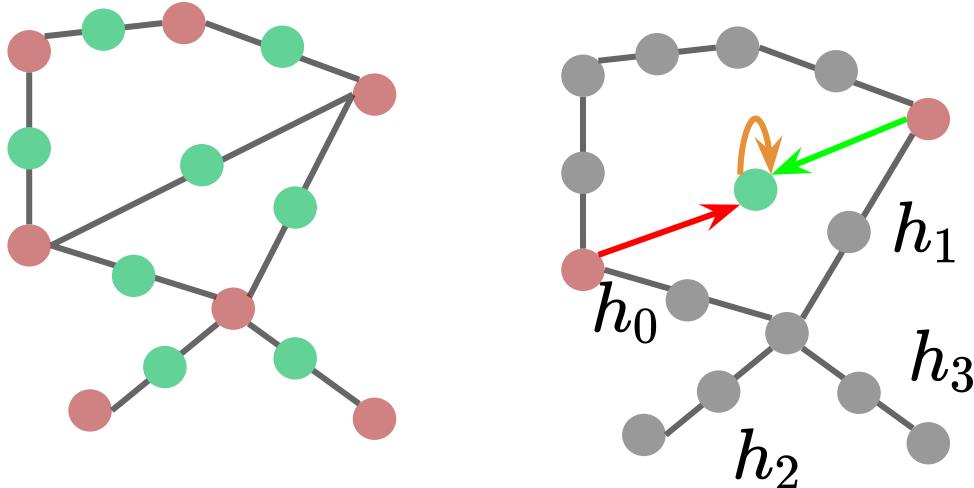


- Graph convolutions involve similar **local operations** on nodes.
- Nodes are now object representations and not activations
- The **ordering of neighbors** should not matter.
- The **number of neighbors** should not matter.
- $N(i)$ are the neighbors of node i
- c_{ij} is a normalization constant

$$h_4^{l+1} = W_4 h_4^l + W_0 h_0^l + W_1 h_1^l + W_2 h_2^l + W_3 h_3^l$$

$$h_i^{l+1} = W_i h_i^l + \sum_{j \in N(i)} \frac{1}{c_{ij}} W_j h_j^l$$

Generalizing 2D convolutions to Graph Convolutions

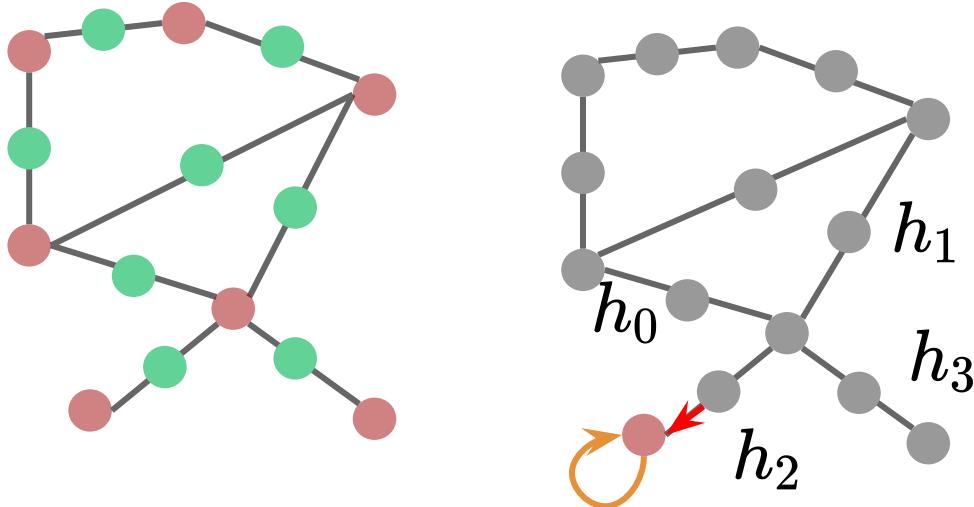


- Graph convolutions involve similar **local operations** on nodes.
- Nodes are now object representations and not activations
- The **ordering of neighbors** should not matter.
- The **number of neighbors** should not matter.
- $N(i)$ are the neighbors of node i

$$h_i^{l+1} = \textcolor{brown}{W_i h_i^l} + \sum_{j \in \mathcal{N}(i)} \frac{1}{c_{ij}} \textcolor{red}{W_j h_j^l}$$

Kipf & Welling (ICLR 2017)

Generalizing 2D convolutions to Graph Convolutions

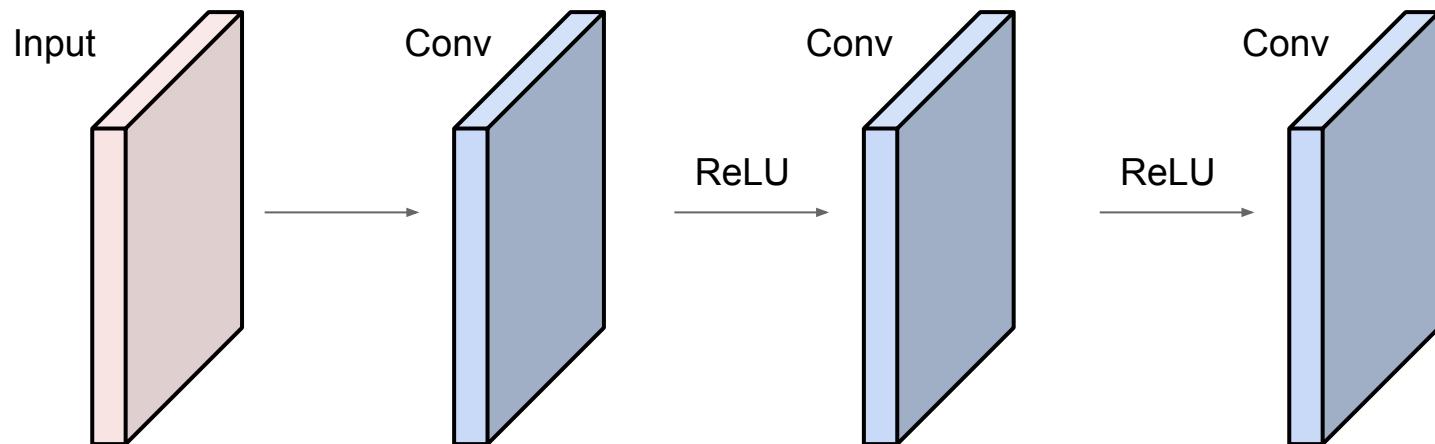


- Graph convolutions involve similar local operations on nodes.
- Nodes are now object representations and not activations
- The ordering of neighbors should not matter.
- The number of neighbors should not matter.
- $N(i)$ are the neighbors of node i

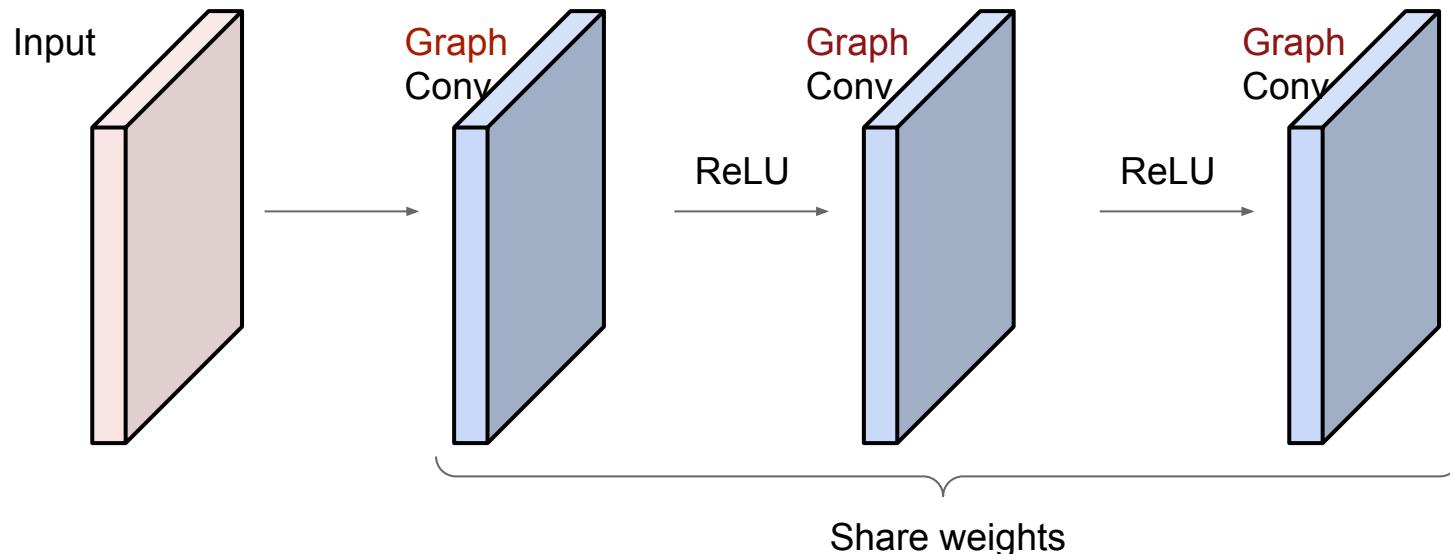
$$h_i^{l+1} = \mathbf{W}_i h_i^l + \sum_{j \in \mathcal{N}(i)} \frac{1}{c_{ij}} \mathbf{W}_j h_j^l$$

Kipf & Welling (ICLR 2017)

To increase receptive field of CNNs: increase depth



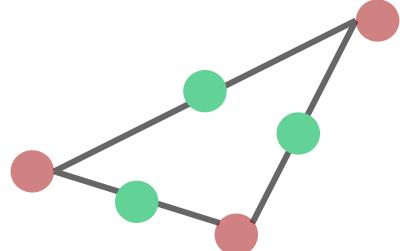
To increase receptive field of GCNs: increase depth



GCNs: Graph Convolutional Networks

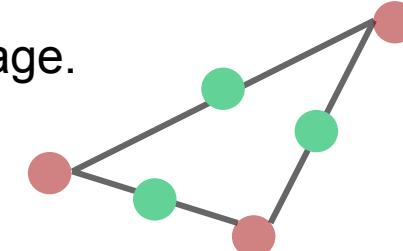
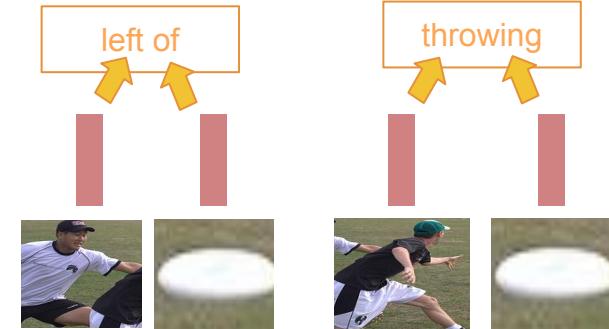
Kipf & Welling (ICLR 2017)

Graph representation with edges included as nodes



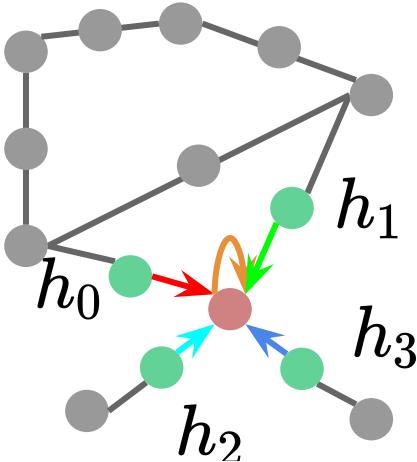
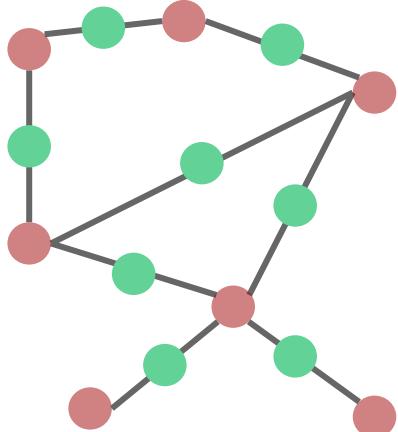
Each node contains features from individual regions

Perform **Graph Convolutions**,
which allows each node to
encode what else is in the image.



Each node now also contains features from all other regions

Graph Convolutions with Attention



- Updates from some neighbors can be more important than others.
- Attention over neighbors allows graph convolutions to focus on specific neighbors
- σ is a non-linearity, usually ReLU or LeakyReLU.

$$\text{Without attention: } h_i^{l+1} = W_i h_i^l + \sum_{j \in \mathcal{N}(i)} \frac{1}{c_{ij}} W_j h_j^l$$

$$\text{With attention: } h_i^{l+1} = W_i h_i^l + \sum_{j \in \mathcal{N}(i)} \frac{1}{c_{ij}} \alpha_{ij} W_j h_j^l$$

$$\text{where } \alpha_{ij} = \frac{e^{\sigma(a^T [W h_i || W h_j])}}{\sum_{k \in \mathcal{N}(i)} e^{\sigma(a^T [W h_i || W h_k])}}$$

How is it actually implemented?

For loops iterating over all the neighbors is expensive

$$h_i^{l+1} = W_i h_i^l + \sum_{j \in \mathcal{N}(i)} \frac{1}{c_{ij}} W_j h_j^l$$

Formalizing a graph representation

For loops iterating over all the neighbors is expensive

$$h_i^{l+1} = W_i h_i^l + \sum_{j \in \mathcal{N}(i)} \frac{1}{c_{ij}} W_j h_j^l$$

Let's define a graph with nodes and edges: $\mathcal{G} = (\mathcal{V}, \mathcal{E})$ with N nodes

Formalizing a graph representation

For loops iterating over all the neighbors is expensive

$$h_i^{l+1} = W_i h_i^l + \sum_{j \in \mathcal{N}(i)} \frac{1}{c_{ij}} W_j h_j^l$$

Let's define a graph with nodes and edges: $\mathcal{G} = (\mathcal{V}, \mathcal{E})$ with N nodes

Let's define the adjacency matrix of a graph as: $A \in \mathbb{R}^{N \times N}$ $A_{ij} = \begin{cases} 1 & \text{if } e_{ij} \in \mathcal{E} \\ 0 & \text{otherwise} \end{cases}$

Formalizing a graph representation

For loops iterating over all the neighbors is expensive

$$h_i^{l+1} = W_i h_i^l + \sum_{j \in \mathcal{N}(i)} \frac{1}{c_{ij}} W_j h_j^l$$

Let's define a graph with nodes and edges: $\mathcal{G} = (\mathcal{V}, \mathcal{E})$ with N nodes

Let's define the adjacency matrix of a graph as: $A \in \mathbb{R}^{N \times N}$ $A_{ij} = \begin{cases} 1 & \text{if } e_{ij} \in \mathcal{E} \\ 0 & \text{otherwise} \end{cases}$

Finally, let's define the degree matrix: $D \in \mathbb{R}^{N \times N}$ $D_{ij} = \begin{cases} \mathcal{N}(i) & \text{if } i = j \\ 0 & \text{otherwise} \end{cases}$

Vectorized graph convolution

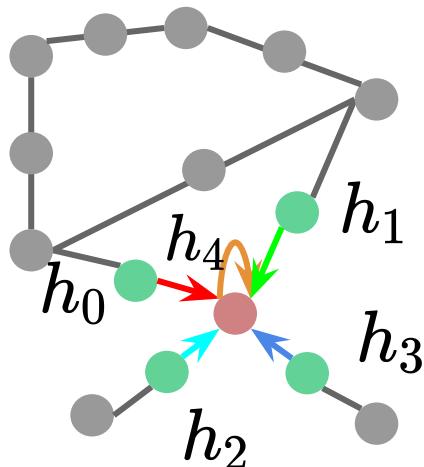
Examples:

$$D_{00} = 2$$

$$D_{44} = 4$$

$$A_{04} = A_{40} = 1$$

$$A_{01} = A_{10} = 0$$



For loops iterating over all the neighbors is expensive

$$h_i^{l+1} = W_i h_i^l + \sum_{j \in \mathcal{N}(i)} \frac{1}{c_{ij}} W_j h_j^l$$

$$A \in \mathbb{R}^{N \times N} \quad A_{ij} = \begin{cases} 1 & \text{if } e_{ij} \in \mathcal{E} \\ 0 & \text{otherwise} \end{cases}$$

$$D \in \mathbb{R}^{N \times N} \quad D_{ij} = \begin{cases} \mathcal{N}(i) & \text{if } i = j \\ 0 & \text{otherwise} \end{cases}$$

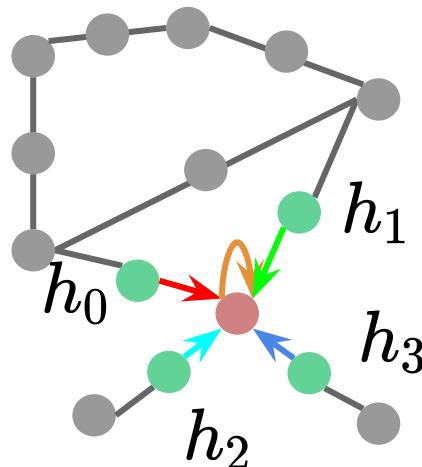
Vectorized graph convolution

First, let's stack all the node representations in a matrix H :

$$H^l \in \mathbb{R}^{N \times C}$$

Such that every row is a node:

$$h_i \in \mathbb{R}^C$$



For loops iterating over all the neighbors is expensive

$$h_i^{l+1} = W_i h_i^l + \sum_{j \in \mathcal{N}(i)} \frac{1}{c_{ij}} W_j h_j^l$$

$$A \in \mathbb{R}^{N \times N} \quad A_{ij} = \begin{cases} 1 & \text{if } e_{ij} \in \mathcal{E} \\ 0 & \text{otherwise} \end{cases}$$

$$D \in \mathbb{R}^{N \times N} \quad D_{ij} = \begin{cases} \mathcal{N}(i) & \text{if } i = j \\ 0 & \text{otherwise} \end{cases}$$

Vectorized graph convolution

First, let's stack all the node representations in a matrix H :

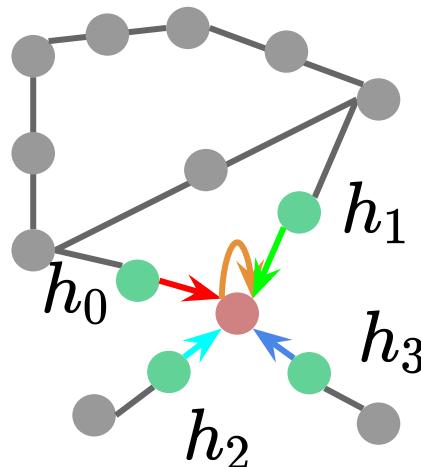
$$H^l \in \mathbb{R}^{N \times C}$$

Such that every row is a node:

$$h_i \in \mathbb{R}^C$$

The vectorized computation of graph convolution is:

$$H^{l+1} = D^{-1/2} \hat{A} D^{-1/2} H^l W$$



For loops iterating over all the neighbors is expensive

$$h_i^{l+1} = W_i h_i^l + \sum_{j \in \mathcal{N}(i)} \frac{1}{c_{ij}} W_j h_j^l$$

$$A \in \mathbb{R}^{N \times N} \quad A_{ij} = \begin{cases} 1 & \text{if } e_{ij} \in \mathcal{E} \\ 0 & \text{otherwise} \end{cases}$$

$$D \in \mathbb{R}^{N \times N} \quad D_{ij} = \begin{cases} \mathcal{N}(i) & \text{if } i = j \\ 0 & \text{otherwise} \end{cases}$$

$$\hat{A} = A + I$$

Vectorized graph convolution

First, let's stack all the node representations in a matrix H :

$$H^l \in \mathbb{R}^{N \times C}$$

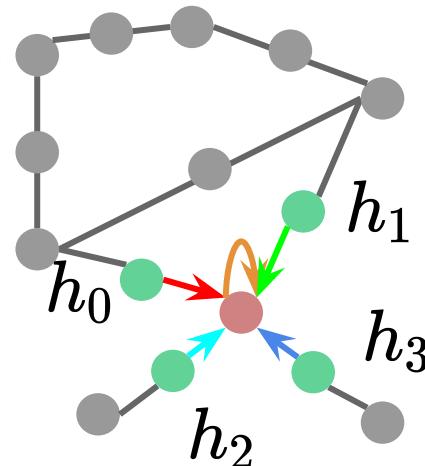
Such that every row is a node:

$$h_i \in \mathbb{R}^C$$

Can be pre-calculated once per graph:

$$H^{l+1} = D^{-1/2} \hat{A} D^{-1/2} H^l W$$

Linear layer weights



For loops iterating over all the neighbors is expensive

$$h_i^{l+1} = W_i h_i^l + \sum_{j \in \mathcal{N}(i)} \frac{1}{c_{ij}} W_j h_j^l$$

$$A \in \mathbb{R}^{N \times N} \quad A_{ij} = \begin{cases} 1 & \text{if } e_{ij} \in \mathcal{E} \\ 0 & \text{otherwise} \end{cases}$$

$$D \in \mathbb{R}^{N \times N} \quad D_{ij} = \begin{cases} \mathcal{N}(i) & \text{if } i = j \\ 0 & \text{otherwise} \end{cases}$$

$$\hat{A} = A + I$$

Aside: Grounding to spectral convolutions with graph laplacian

Convolutions in the spectral domain:

$$W * h = U \text{diag}(W) U^T h$$

Where U is the eigenvectors of the graph laplacian:

$$L = I + D^{-1/2}AD^{-1/2} = U \Lambda U^T$$

You can approximate spectral graph convolutions as 1st order Chebyshev polynomials to get:

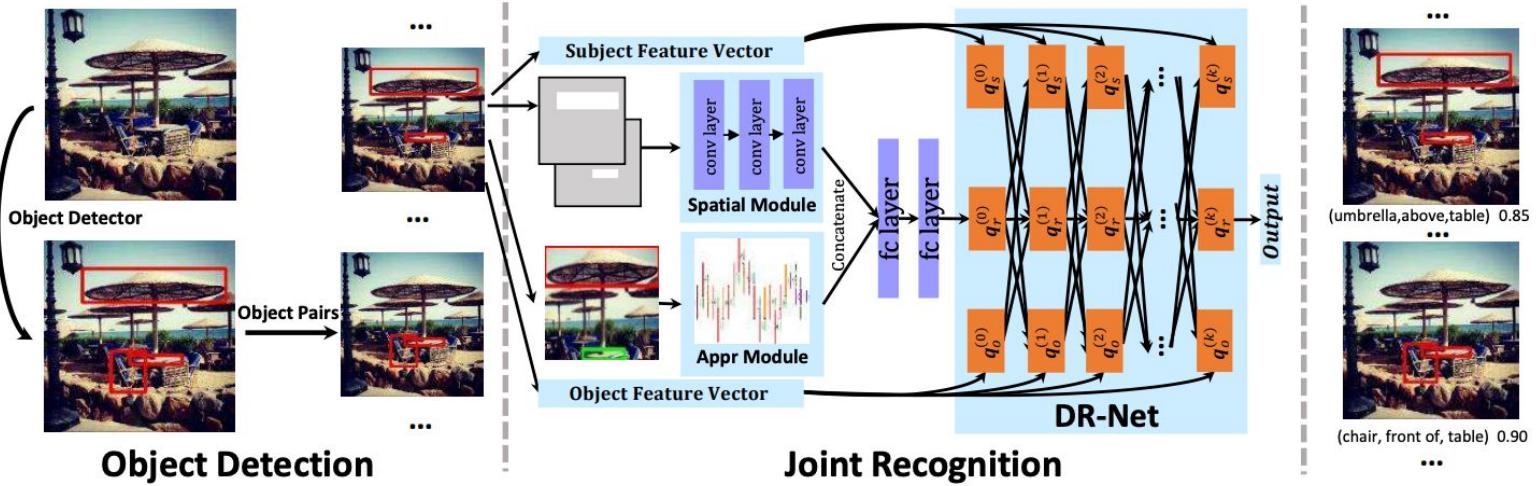
$$W * h = W(I + D^{-\frac{1}{2}}AD^{-\frac{1}{2}})h$$

Renormalize the weights to get our spatial
graph convolutions: $I + D^{-1/2}AD^{-1/2} \rightarrow D^{-1/2}\hat{A}D^{-1/2}$

$$H^{l+1} = D^{-1/2}\hat{A}D^{-1/2}H^lW$$

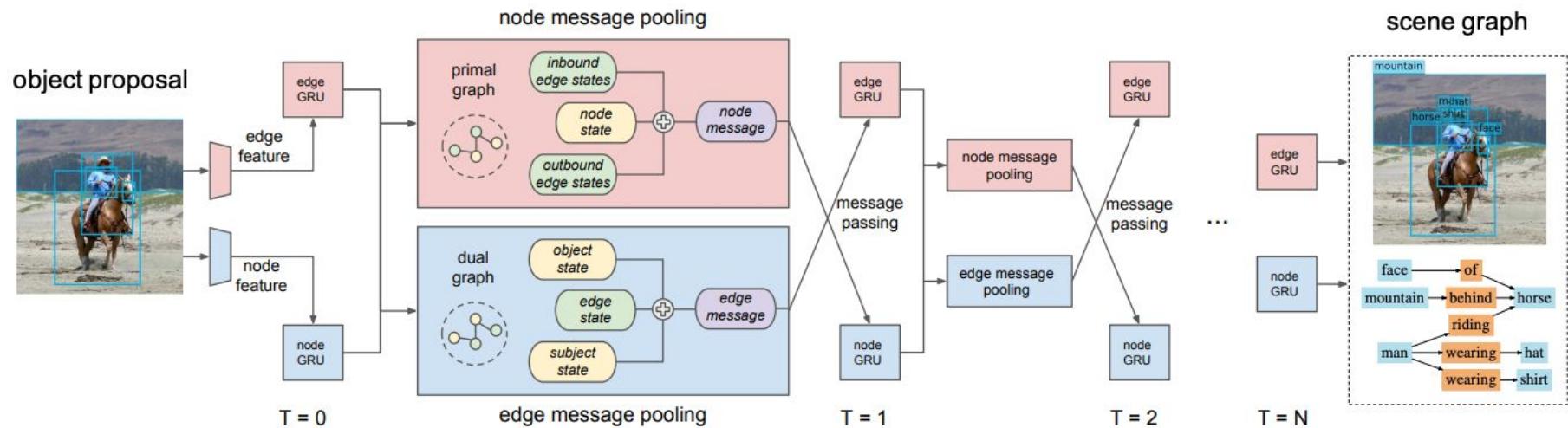
$$\hat{A} = A + I$$

Scene Graph Generation with Graph Convolution methods



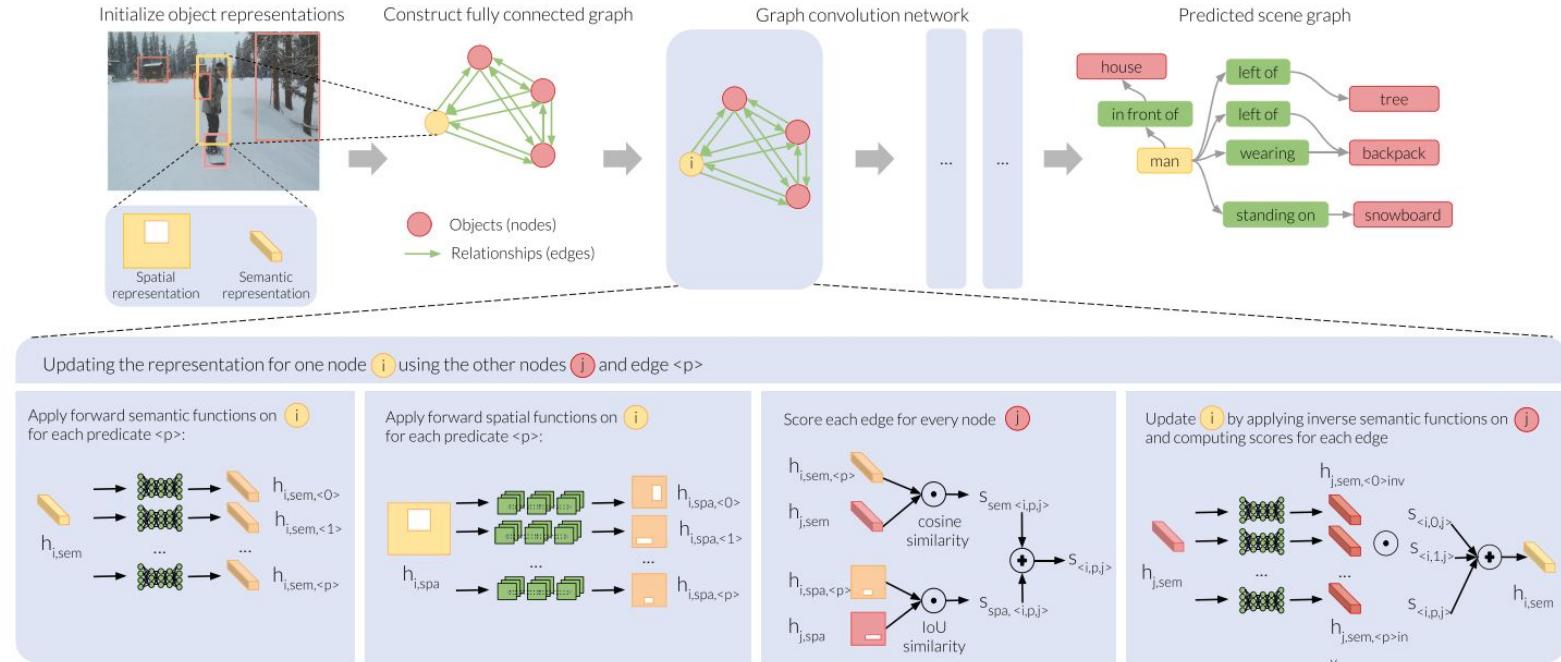
Liang et al. Deep variation-structured reinforcement learning for visual relationship and attribute detection, CVPR 2017

Scene Graph Generation with node and edge Graph Convolution methods



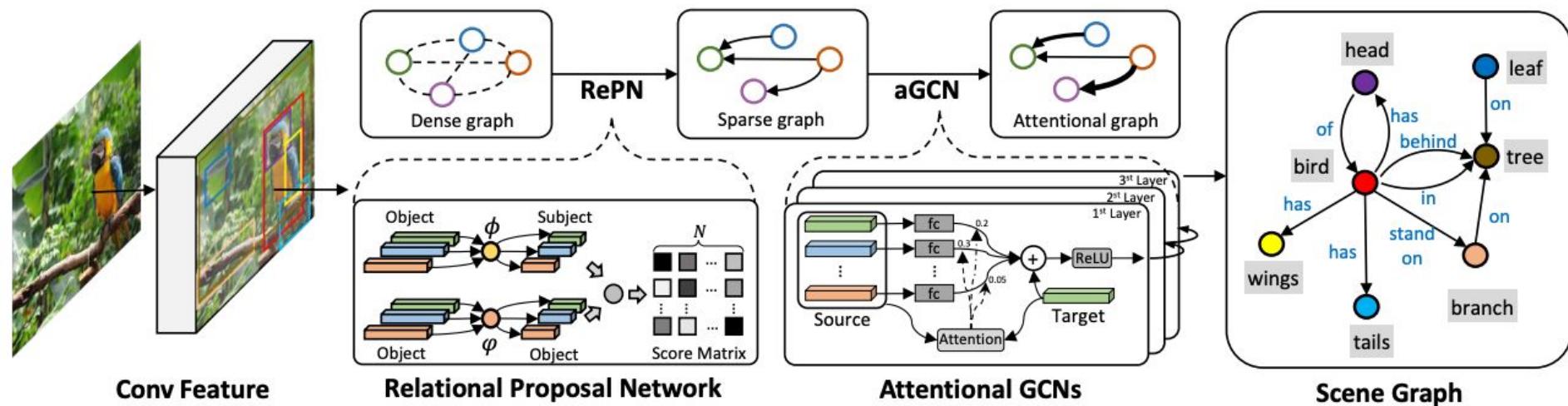
Xu et al. "Scene graph generation by iterative message passing, CVPR 2017"

Few shot scene graph generation with graph convolution methods



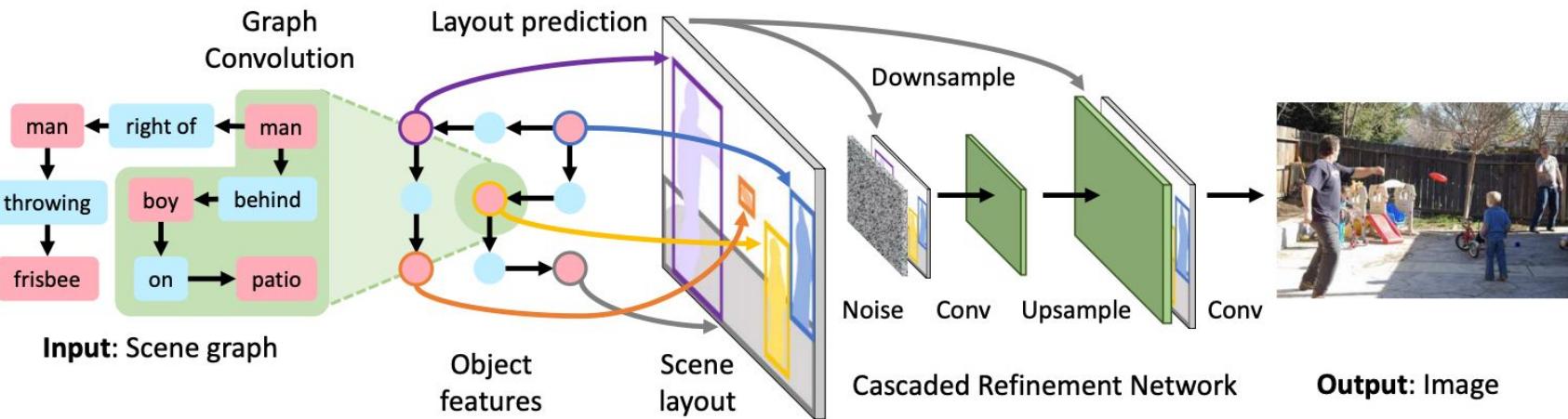
Dornadula, Narcomey, Krishna, et al. "Visual Relationships as Functions: Enabling Few-Shot Scene Graph Prediction." Proceedings of the IEEE International Conference on Computer Vision Workshops. 2019

Scene graph generation with graph attention convolutions



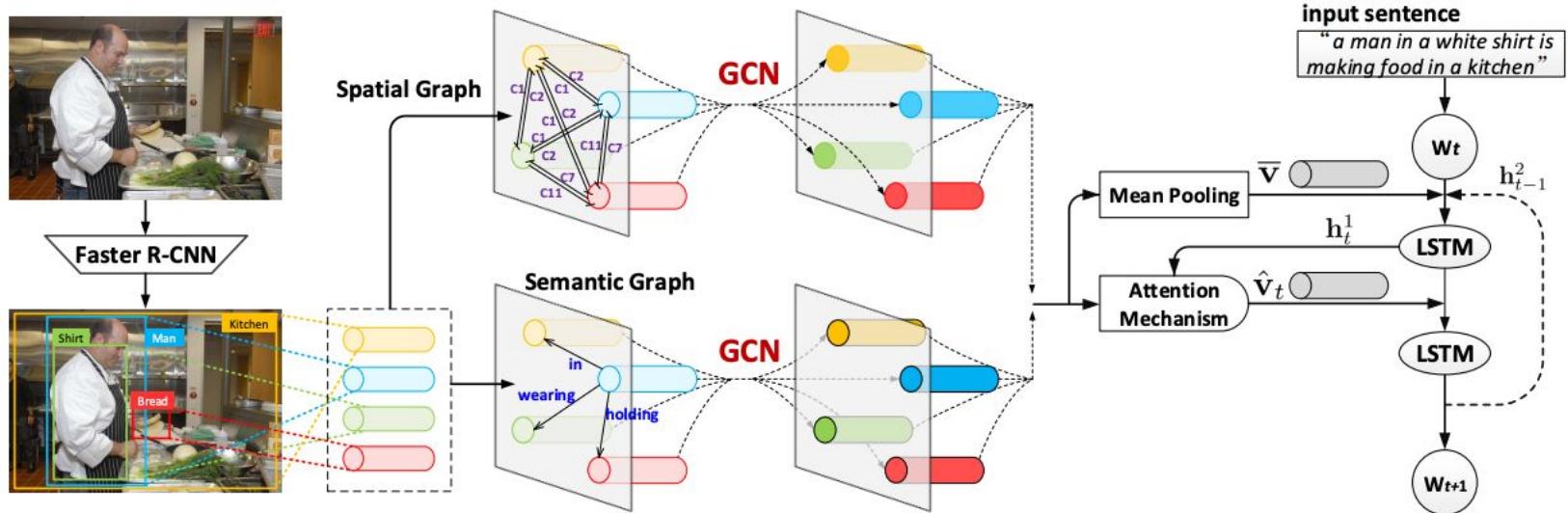
Yang, et al. "Graph r-cnn for scene graph generation." Proceedings of the European conference on computer vision ECCV 2018

Image generation from scene graphs



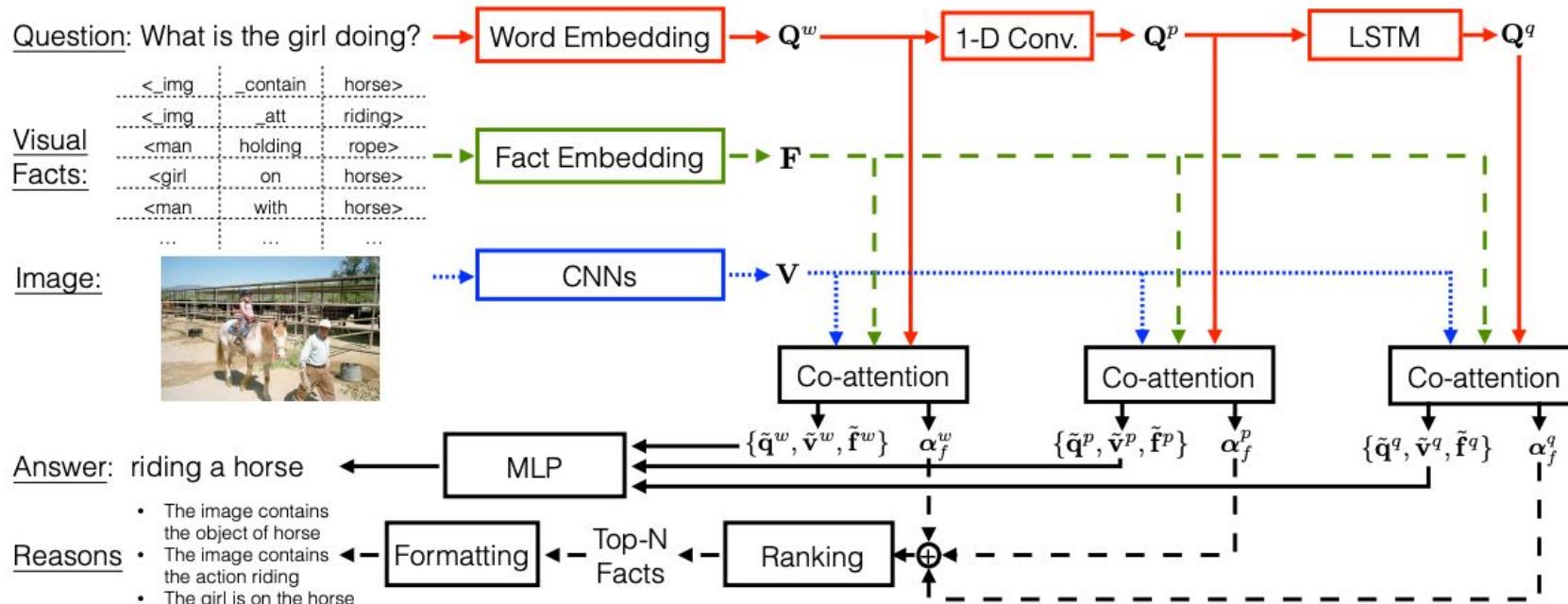
Johnson et al. Image generation from scene graphs, CVPR 2019

Scene graphs as intermediate representation for image captioning



Yao et al. Exploring Visual Relationship for Image Captioning, ECCV 2018

Scene graphs as intermediate representation for visual question answering



Wang et al. The vqa-machine: Learning how to use existing vision algorithms to answer new questions CVPR 2017

So what's next for scene graphs?

Action Genome: Understanding Actions with Spatio-Temporal Scene Graphs

action: take a bag from somewhere

action: drinking from a cup

action: take notebook from somewhere



Krishna et al. Dense Captioning Events in Videos, CVPR 2017

Ji, Krishna et al. Action Genome: Actions as Compositions of Spatio-Temporal Scene Graphs, CVPR 2020

Action Genome: Understanding Action with Spatio-Temporal Scene Graphs

action: take a bag from somewhere

action: drinking from a cup

action: take notebook from somewhere



Krishna et al. Dense Captioning Events in Videos, CVPR 2017

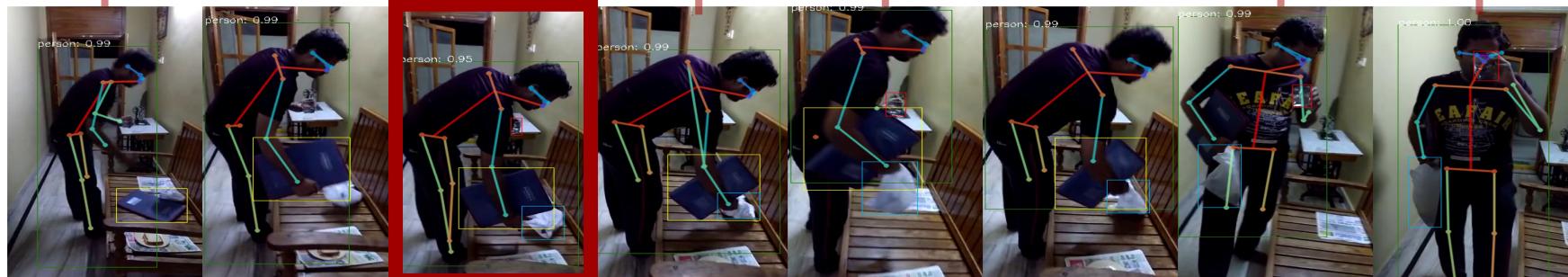
Ji, Krishna et al. Action Genome: Actions as Compositions of Spatio-Temporal Scene Graphs, CVPR 2020

Action Genome: Understanding Action with Spatio-Temporal Scene Graphs

action: take a bag from somewhere

action: drinking from a cup

action: take notebook from somewhere



Krishna et al. Dense Captioning Events in Videos, CVPR 2017

Ji, Krishna et al. Action Genome: Actions as Compositions of Spatio-Temporal Scene Graphs, CVPR 2020

Action Genome: Understanding Action with Spatio-Temporal Scene Graphs

action: take a bag from somewhere

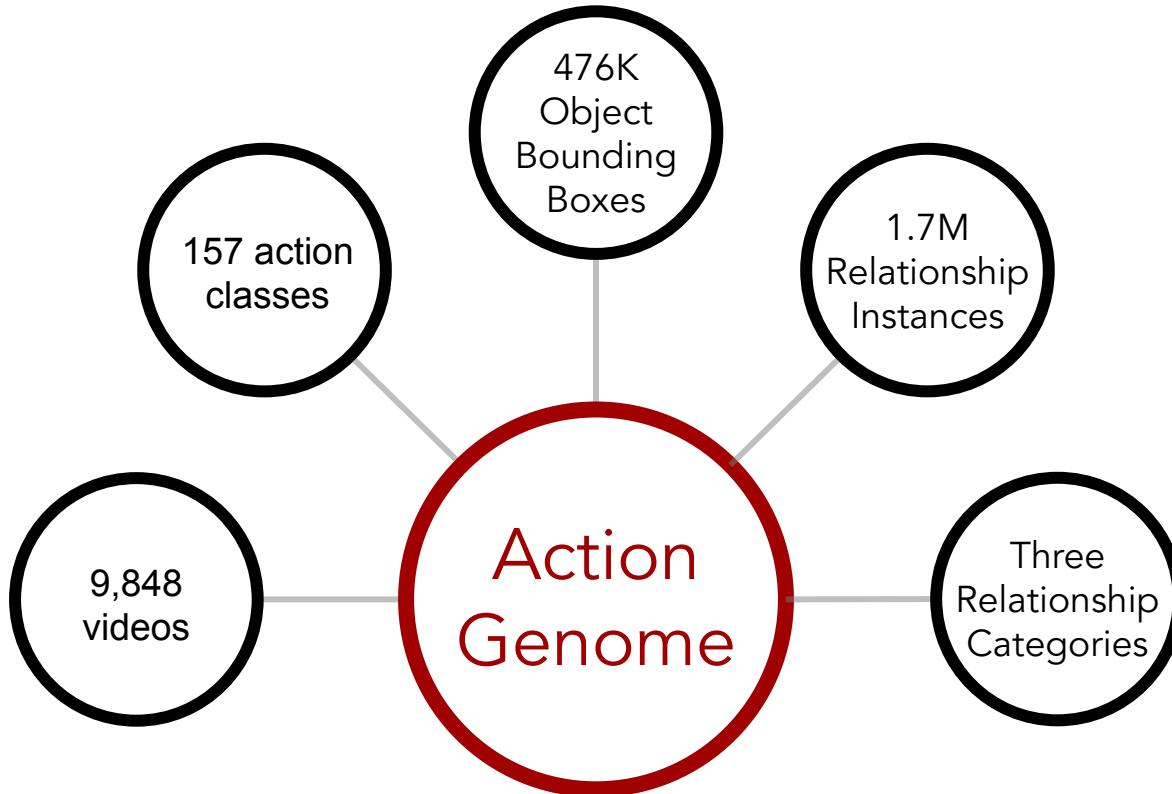
action: drinking from a cup

action: take notebook from somewhere



Krishna et al. Dense Captioning Events in Videos, CVPR 2017

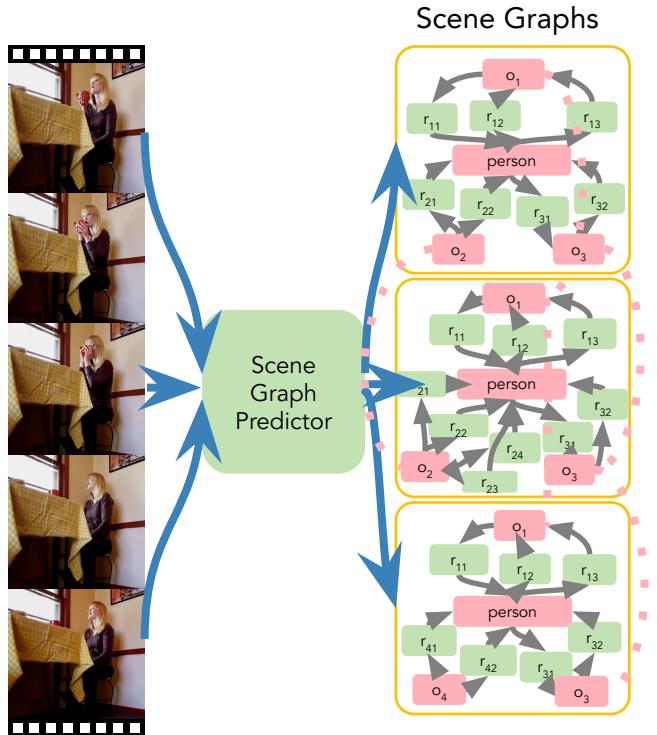
Ji, Krishna et al. Action Genome: Actions as Compositions of Spatio-Temporal Scene Graphs, CVPR 2020



Code and dataset available: <http://actiongenome.org>

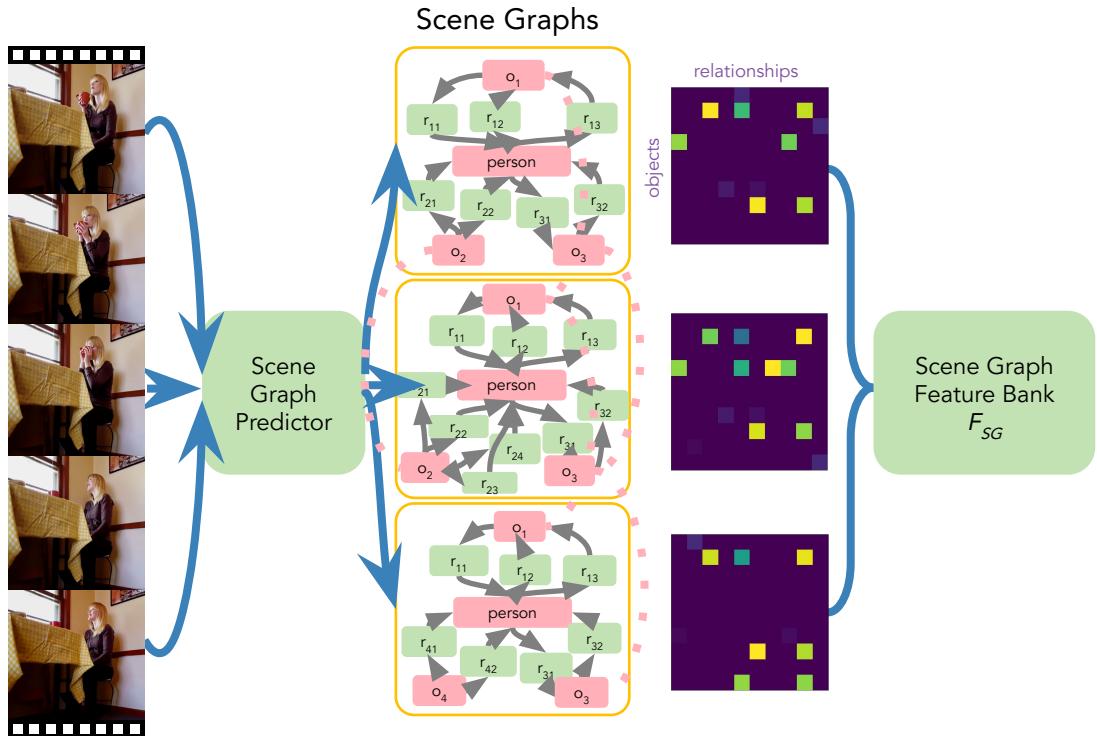
Ji, Krishna et al. Action Genome: Actions as Compositions of Spatio-Temporal Scene Graphs, CVPR 2020

Spatio-temporal Scene Graph Feature Banks (SGFB) for Action Recognition



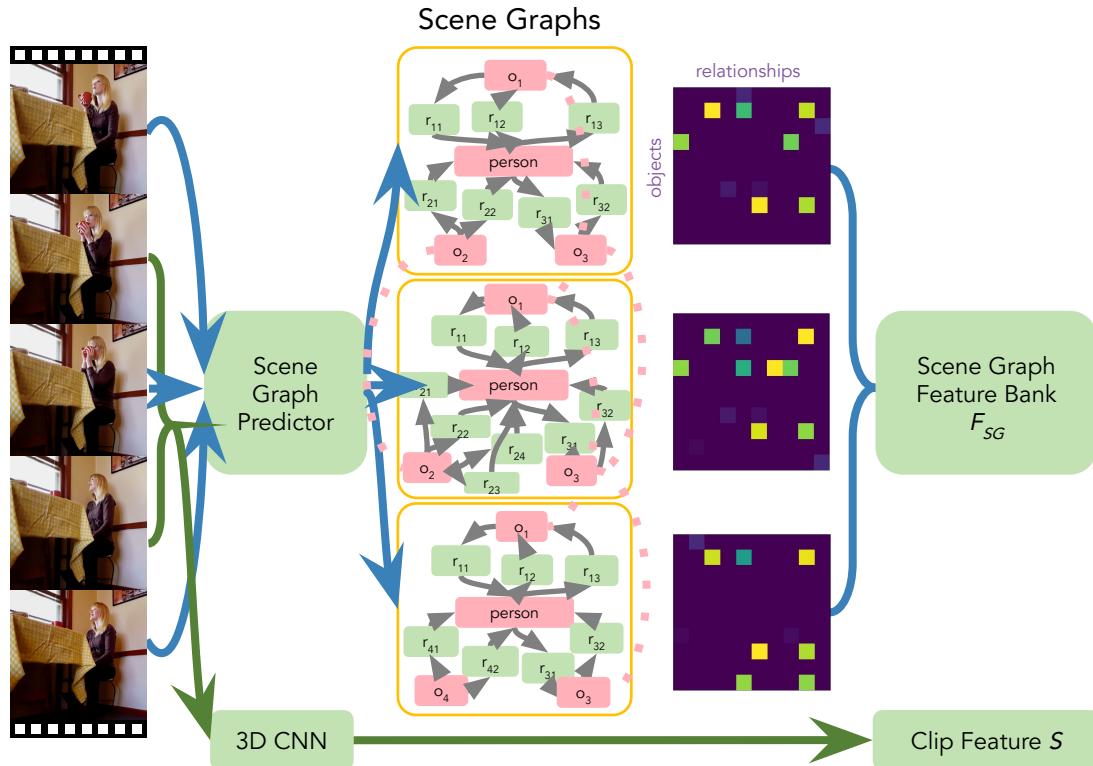
Ji, Krishna et al. Action Genome: Actions as Compositions of Spatio-Temporal Scene Graphs, CVPR 2020

Spatio-temporal Scene Graph Feature Banks (SGFB) for Action Recognition



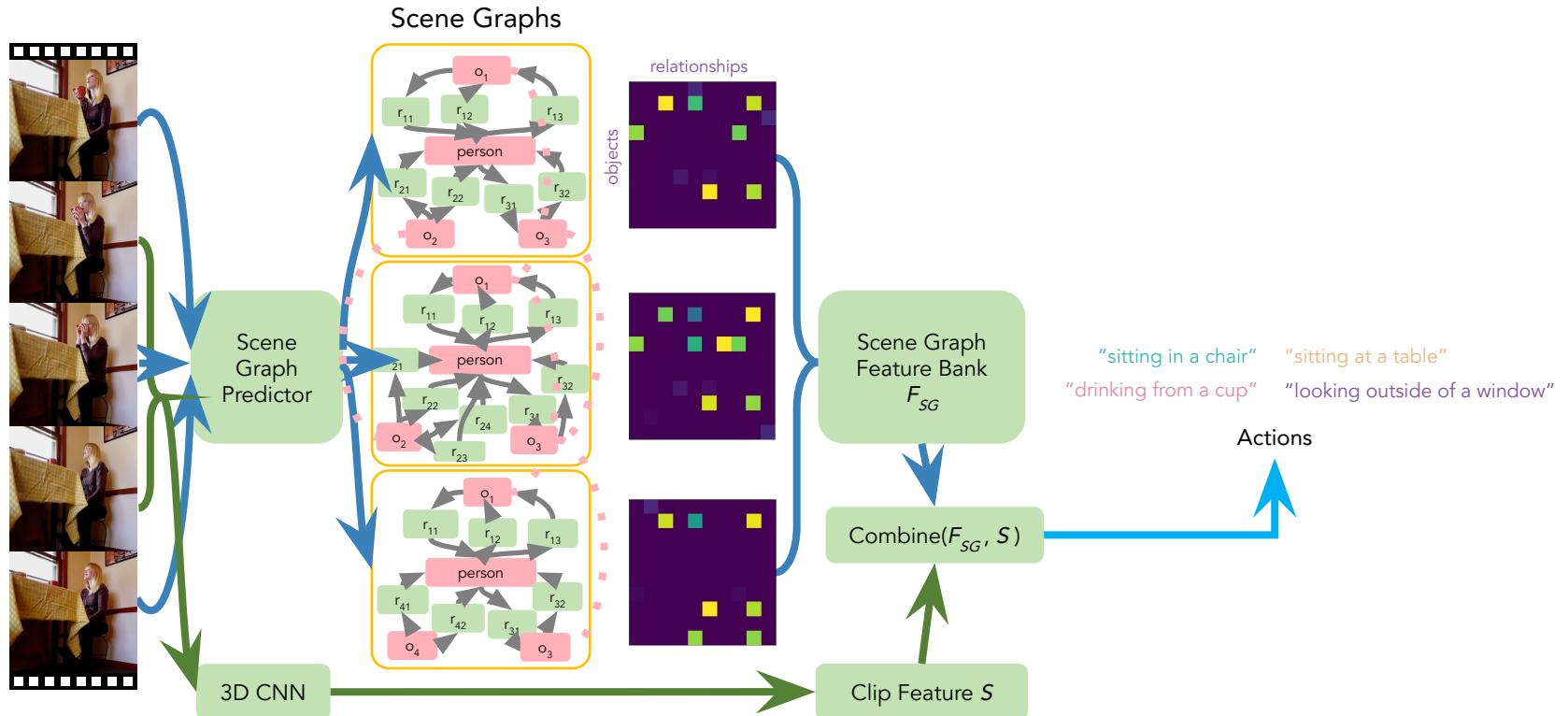
Ji, Krishna et al. Action Genome: Actions as Compositions of Spatio-Temporal Scene Graphs, CVPR 2020

Spatio-temporal Scene Graph Feature Banks (SGFB) for Action Recognition



Ji, Krishna et al. Action Genome: Actions as Compositions of Spatio-Temporal Scene Graphs, CVPR 2020

Spatio-temporal Scene Graph Feature Banks (SGFB) for Action Recognition



Ji, Krishna et al. Action Genome: Actions as Compositions of Spatio-Temporal Scene Graphs, CVPR 2020

From Scene Graphs to Action Recognition



Ground truth action labels:
Lying on a bed,
Awakening in bed,
Holding a pillow

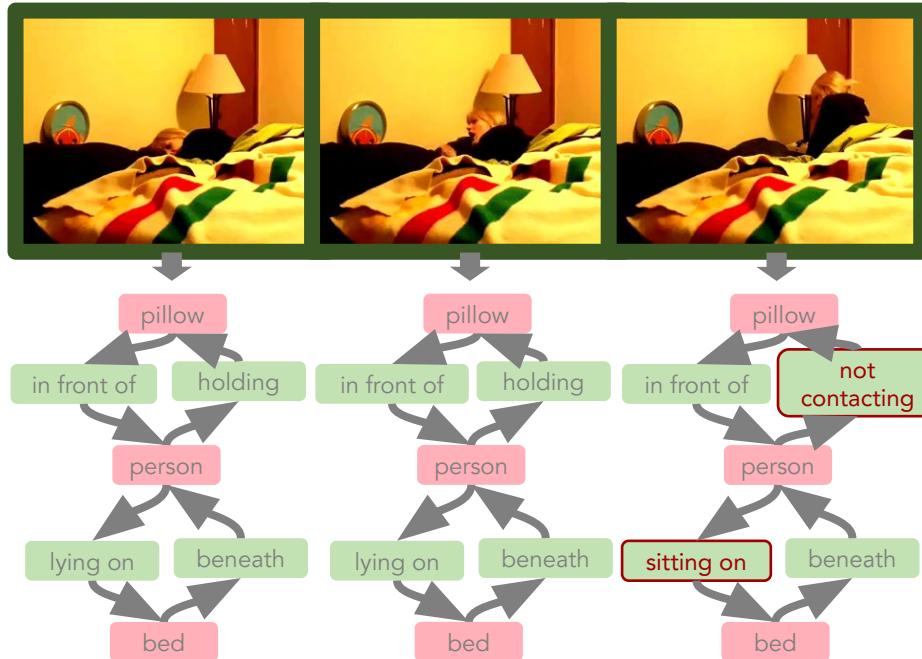
Baselines rely heavily on training set priors



Ground truth:
Lying on a bed,
Awakening in bed,
Holding a pillow

Baseline (LFB)
predictions:
Lying on a bed,
Watching television,
Holding a pillow

Modeling temporal changes in relationships lead to improved inference



Ground truth:
Lying on a bed,
Awakening in bed,
Holding a pillow

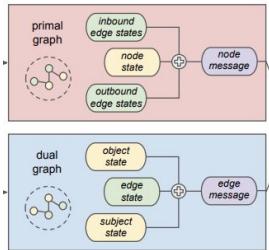
Baseline (LFB) predictions:
Lying on a bed,
Watching television,
Holding a pillow

Our top-3 predictions:
Lying on a bed,
Awakening in bed,
Holding a pillow

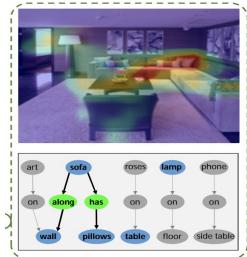
Wu et al. Long-term feature banks for detailed video understanding, CVPR 2019
Ji, Krishna et al. Action Genome: Actions as Compositions of Spatio-Temporal Scene Graphs, CVPR 2020

The community has published hundreds of scene graph papers

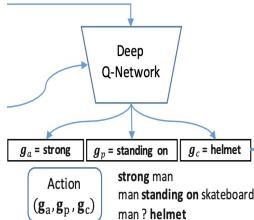
Message passing



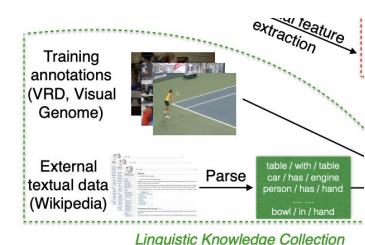
Attention



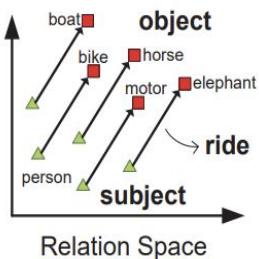
Reinforce



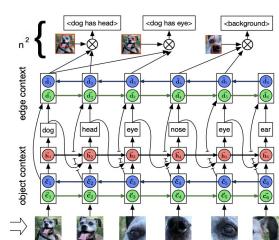
External knowledge



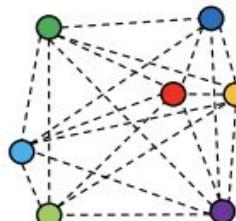
Transformations



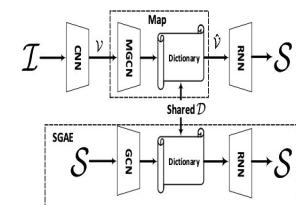
Recurrent networks



Graph Convolutions



Auto-encoders



Zellers et al. Neural motifs: Scene graph parsing with global context CVPR. 2018

Yang et al. Graph r-cnn for scene graph generation ECCV 2018

Yang et al. Shuffle-then-assemble: Learning object-agnostic visual relationship features ECCV 2018

Zhang et al. Visual translation embedding network for visual relation detection CVPR 2017

Liang et al. Deep variation-structured reinforcement learning for visual relationship and attribute detection CVPR 2017

Dornadula et al. Visual Relationships as Functions: Enabling Few Shot Scene Graph Generation ICCV SGRL 2019

Xu et al. Scene graph generation by iterative message passing CVPR 2017

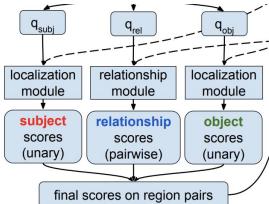
Yu et al. Visual relationship detection with internal and external linguistic knowledge distillation ICCV 2017

Scene graphs have achieved state of the art in many tasks

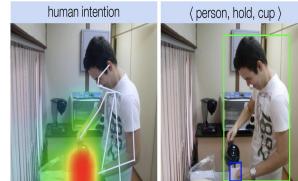
3D scene graphs



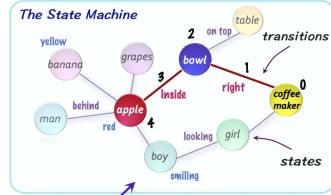
Explainable AI



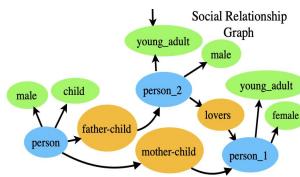
Human intentions



VQA



Social relationships



Fashion



Image generation

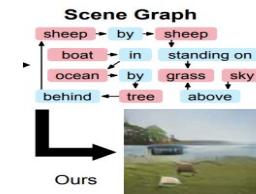
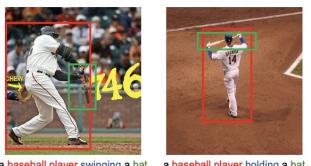
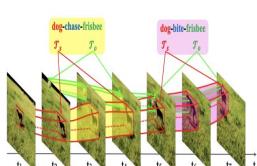


Image captioning



Video understanding



Armeni et al. 3D Scene Graph: A Structure for Unified Semantics, 3D Space, and Camera ICCV 2019

Hu et al. Modeling relationships in referential expressions with compositional modular networks, CVPR 2017

Xu et al. Interact as you intend: Intention-driven human-object interaction detection, Transactions on Multimedia 2019

Hudson et al. Neural State Machine, NeurIPS 2019

Hu, Ronghang, et al. Learning to reason: End-to-end module networks for visual question answering /ICCV 2017

Johnson et al. Image generation from scene graphs CVPR 2018

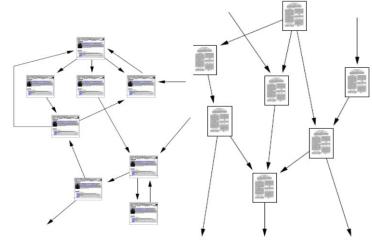
Yu et al. Layout-graph reasoning for fashion landmark detection CVPR 2019

Goel et al. An End-to-End Network for Generating Social Relationship Graphs CVPR 2019

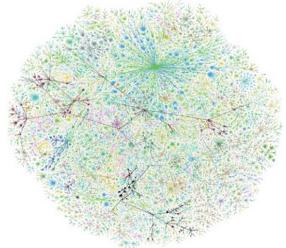
Kim et al. Dense relational captioning: Triple-stream networks for relationship-based captioning CVPR 2019

Tsai et al. Video relationship reasoning using gated spatio-temporal energy graph CVPR 2019

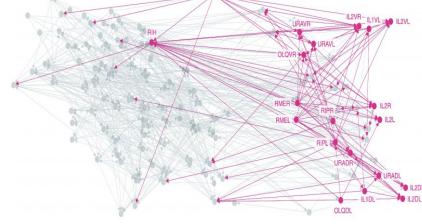
Graphs are everywhere – in numerous fields



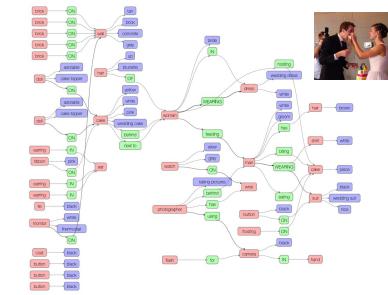
Information networks: Web & citations



Internet



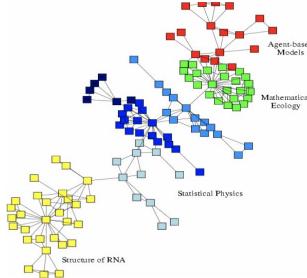
Networks of neurons



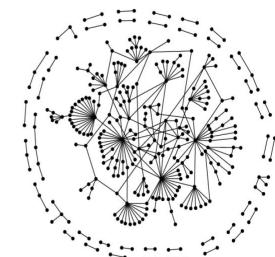
Scene Graphs



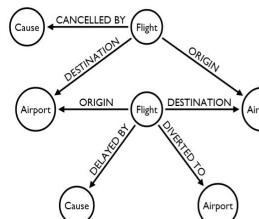
Social networks



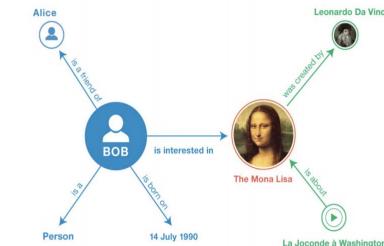
Economic networks



Communication networks



Event Graphs



Knowledge Graphs

Summary

- Scene graphs are a symbolic, compositional, knowledge representation inspired by Cognitive Science and is a common underlying structure in many Computer Vision tasks.
- The task of Scene Graph Generation requires more complex structured prediction models
- GCNs are a generalization of the CNNs you have already learned about.
 - Use them when you work with graph-related data
- This is a relatively new sub-field and there is a lot of work left to do and a lot of promise for future research.

What have we learned this quarter?

Neural Network Fundamentals

Data-driven learning
Linear classification & kNN
Loss functions
Optimization
Backpropagation
Multi-layer perceptrons
Neural Networks

Convolutional Neural Networks

Convolutions
Pytorch 1.4 / Tensorflow 2.0
Activation functions
Batch normalization
Transfer learning
Data augmentation
Momentum / RMSProp / Adam
Architecture design

Computer Vision Applications

RNNs / LSTMs
Image captioning
Interpreting neural networks
Style transfer
Adversarial examples
Fairness & ethics
Human-centered AI
3D vision
Deep reinforcement learning
Scene graphs
Graph Convolutions

Instructors



Fei-Fei Li

Teaching Assistants



William Shen
(Head TA)



Jonathan Braatz



Daniel Cai



JunYoung Gwak



De-An Huang



Ranjay Krishna



Andrew Kondrich



Fang-Yu Lin



Damian Mrowca



Boxiao Pan



Chris Waites



Danfei Xu



Rui Wang



Yi Wen



Karen Yang



Brent Yi



Christina Yuan

Course Coordinator



Amelie Byun



Kevin Zakka



Yiheng Zhang