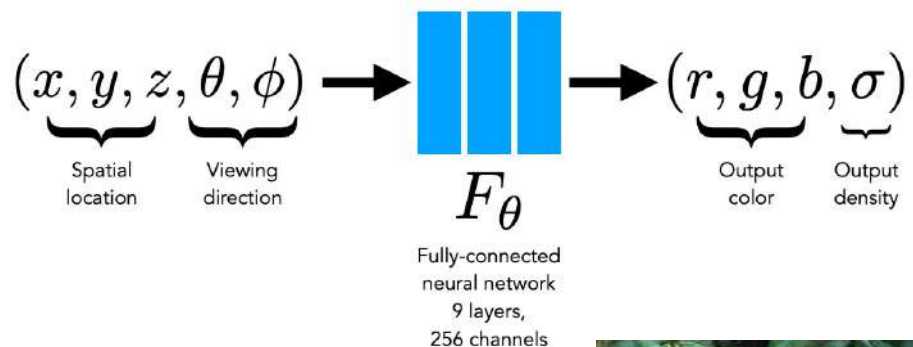# Lecture 17:

## Scene Graphs and Graph Convolutions

# Administrative

- A3 grades will be released next week

- Final project reports due June 3rd
    - Final project video due June 4th
    - No late days for final project

- 2 guest lectures next week:
    - First on multimodal learning combining vision and sound
    - Second on combining vision with action

# Last time: NeRF



$(x, y, z, \theta, \phi)$ → $F_\theta$ → $(r, g, b, \sigma)$

Spatial location | Viewing direction

Fully-connected neural network
9 layers,
256 channels

Output color | Output density

# Today's agenda

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

# Computer vision was focused on disconnected objects

Image Classification



red panda

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
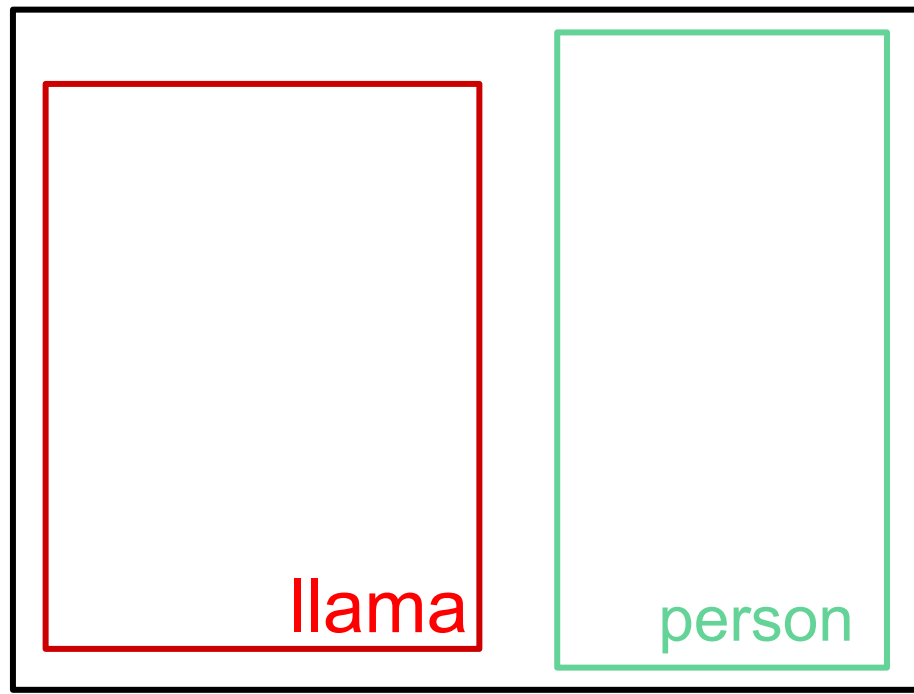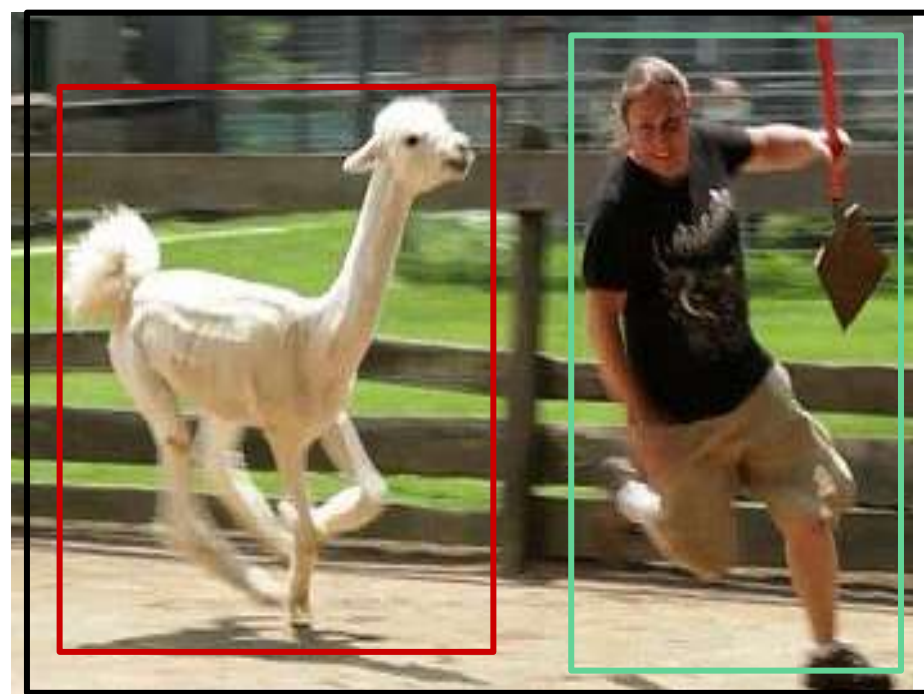
# image #1



llama

person

# image #2



llama

person

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

chasing

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

# 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



<person - kicking - ball>

### 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 types of fruits are there?

is the food healthy?

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

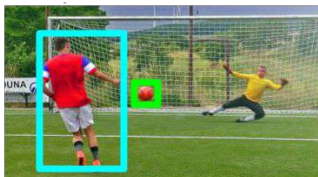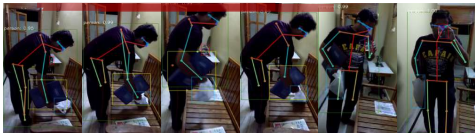# Many Vision tasks share a similar underlying structure
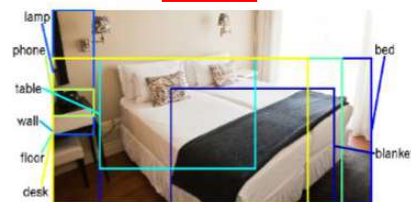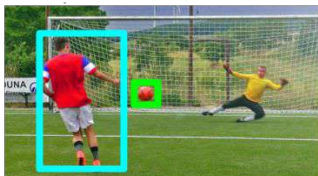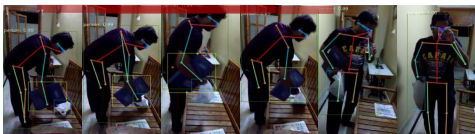
objects

## Action classification

action: drinking from a cup

action: take notebook from somewhere



## Grounding objects



person - kicking - ball

## 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 table 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 types of fruits are there?

is the food healthy?

how many people are in the photo?

is this a busy street

how many skateboards are there?

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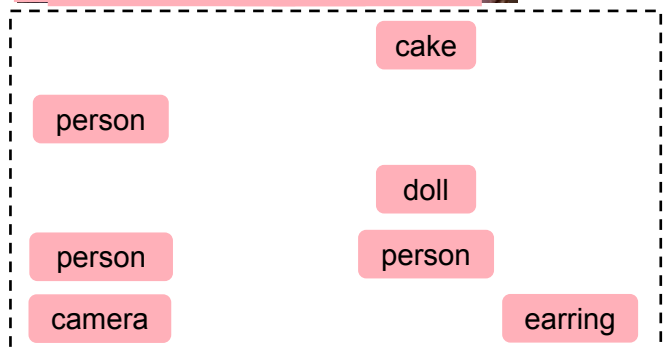
Ji, Krishna et al. Action Genome: Actions as Compositions of Spatio-Temporal Scene Graphs, CVPR 2020
Krishna et al. Information Maximizing Visual Question Answering, CVPR 2019
Krishna et al. Referring Relationships, CVPR 2018
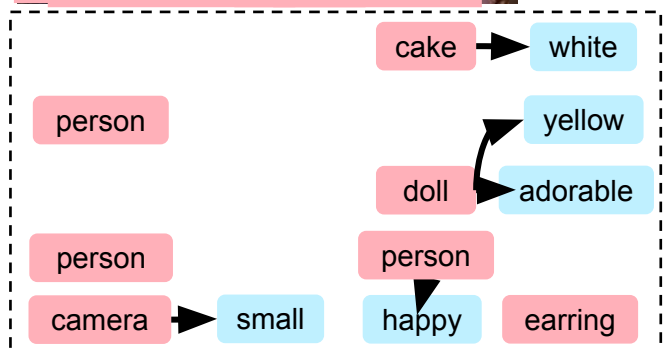Johnson, Krishna et al. Image Retrieval with Scene Graphs, CVPR 2015

# Many Vision tasks share a similar underlying structure



**objects** (red)
**attributes** (blue)

**Action classification**

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

**Grounding objects**

<person> - kicking - ball

**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.

lamp
phone
table
wall
floor
desk
bed
blanket

**Question answering**

how many types of vegetables are there?
is the food in the foreground prickly
how many types of fruits are there?
is the food healthy
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   Question answering



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

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

Fei-Fei Li, Ranjay Krishna, Danfei Xu          Lecture 17 - 17          May 27, 2021

# The scene graph representation



person
person
camera
cake
doll
person
earring

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

# The scene graph representation

# The scene graph representation

# 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

Legend: object attribute relationship

But why is scene graph the right representation?

# Try and remember all these images

# 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

# We encode more than objects

Some relationships between objects must be coded into the gist. A picture of milk being poured from a carton into a glass is not the same as a picture of milk being poured from a carton into the space next to a glass, even if all of the objects are the same.

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

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)

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

# Scene Graph Generation - Problem formulation



Input
(image only)

person
↓
riding
↓
horse → in front of → horse
↑
riding
↑
person

Output

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

# Scene Graph Generation - Problem formulation



Input
(image only)

person

horse

person

horse

Output

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

# Scene Graph Generation - Problem formulation



Input
(image only)

person

riding

horse

person

riding

horse

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

Visual Genome dataset
N = 100
K = 50

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


ride


falling off


next to


carry


lying


resting on


drag


throw

# Recall the algebraic interpretation of linear models:

Features from the last layer



Input image

| 0.2 | -0.5 | 0.1 | 2.0 |
| 1.5 | 1.3 | 2.1 | 0.0 |
| 0 | 0.25 | 0.2 | -0.3 |

| 56 |
| 231 |
| 24 |
| 2 |

+

| 1.1 |
| 3.2 |
| -1.2 |

=

| -96.8 |
| 437.9 |
| 61.95 |

person - riding - horse

person - driving - car

dog - holding - frisbee

W

b

Scores for each relationship

$N^2K$ rows!!

# 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



car on street | dog ride skateboard

**Long tail distribution** of relationships
- makes supervised training difficult

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

# Challenge #2



*# of occurrences* (y-axis)

*relationships* (x-axis)

car on street

dog ride skateboard

elephant drink milk

dog ride surfboard

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

Features from the last layer



Input image

N+K rows only

W

b

person -96.8

horse 437.9

cat 61.95

riding 17.2

holding 10.9

driving 2.95

Visual module

Language module

Input

Tackles:
Quadratic explosion of $N^2K$ detectors

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: $b_1$     $b_2$



object
detector

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

Input

Definitions:
$b_1$, $b_2$ are object proposals
$o_1$, $o_2 \in$ [person, horse, …]

# Visual module

Proposals:

Sample: $b_1$    $b_2$      $b_1 \cup b_2$

Input

object detector    relationship detector

$s(o_1|b_1)$   $s(o_2|b_2)$     $s(r|b_1 \cup b_2)$

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

# Visual module

Proposals:



Sample: $b_1$   $b_2$

$b_1 \cup b_2$

object detector

relationship detector

$s(o_1|b_1)$   $s(o_2|b_2)$   $s(r|b_1 \cup b_2)$

$p(T|b_1, b_2)$

Input

Definitions:
$b_1$, $b_2$ are object proposals
$o_1$, $o_2 \in$ [person, horse, …]
$r \in$ [on, in, ride, front of, …]
T is a <$o_1$, r, $o_2$> triple

# Visual module

Proposals:

Sample: $b_1$ $b_2$

$b_1 \cup b_2$

Input

$p(T|b_1, b_2)$

object detector

relationship detector

$s(o_1|b_1)$  $s(o_2|b_2)$  $s(r|b_1 \cup b_2)$

$p(T|b_1, b_2)$

$\underset{T}{\mathrm{argmax}}$

person

in

horse

Definitions:
$b_1$, $b_2$ are object proposals
$o_1$, $o_2 \in$ [person, horse, …]
$r \in$ [on, in, ride, front of, …]
T is a $<o_1, r, o_2>$ triple

# Visual module

Proposals:



Sample: $b_1$    $b_2$



object detector

relationship detector

$b_1 \cup b_2$



$p(T|b_1, b_2)$

$s(o_1|b_1)$    $s(o_2|b_2)$    $s(r|b_1 \cup b_2)$

$p(T|b_1, b_2)$

Input

$\underset{T}{\mathrm{argmax}}$

person

in

horse

# Language module

$o_1$: person        r: ride        $o_2$: horse

Definitions:

$b_1$, $b_2$ are object proposals

$o_1$, $o_2$ $\in$ [person, horse, …]

r $\in$ [on, in, ride, front of, …]

T is a <$o_1$, r, $o_2$> triple

Visual module

Proposals:

Sample: $b_1$  $b_2$

object detector

relationship detector

$b_1 \cup b_2$

$s(o_1|b_1)$  $s(o_2|b_2)$  $s(r|b_1 \cup b_2)$

$p(T|b_1, b_2)$

Input

$p(T|b_1, b_2)$

$\underset{T}{\text{argmax}}$

person

in

horse

Language module

$o_1$: person  r: ride  $o_2$: horse

$p(T|lang)$

Definitions:
$b_1$, $b_2$ are object proposals
$o_1$, $o_2 \in$ [person, horse, ...]
$r \in$ [on, in, ride, front of, ...]
T is a <$o_1$, r, $o_2$> triple

Fei-Fei Li, Ranjay Krishna, Danfei Xu          Lecture 17 - 49          May 27, 2021

Visual module

Proposals:

Sample: $b_1$ $b_2$

$b_1 \cup b_2$

object detector

relationship detector

$s(o_1|b_1)$ $s(o_2|b_2)$ $s(r|b_1 \cup b_2)$

$p(T|b_1, b_2)$

Input

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

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

$\underset{T}{\mathrm{argmax}}$

person

riding

horse

Language module

$o_1$: person     r: ride     $o_2$: horse

$p(T|lang)$

Definitions:
$b_1$, $b_2$ are object proposals
$o_1$, $o_2 \in$ [person, horse, …]
$r \in$ [on, in, ride, front of, …]
T is a <$o_1$, r, $o_2$> triple

# Visual module

Proposals:

Sample: $b_1$ $b_2$

$b_1 \cup b_2$

object detector

relationship detector

$s(o_1|b_1)$ $s(o_2|b_2)$ $s(r|b_1 \cup b_2)$

$p(T|b_1, b_2)$

Tackles:

Quadratic explosion

only requires N+K detectors

Tackles:

Long tail distribution

can predict rare
relationships



Language module

$o_1$: person    r: ride    $o_2$: horse

$p(T|lang)$

# Training the visual module

1. Pre-train using ImageNet

object detector

object detector

relationship detector

Definitions:

# Training the visual module



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

Definitions:
$b_1$, $b_2$ are object proposals
$o_1$, $o_2 \in$ [person, horse, …]

# 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
$o_1$, $o_2 \in$ [person, horse, …]
$r \in$ [on, in, ride, front of, …]

# Training the visual module



$$b_1 \qquad b_2 \qquad b_1 \cup b_2$$

object detector    object detector    relationship detector

$$s(o_1|b_1) \qquad s(o_2|b_2) \qquad s(r|b_1 \cup b_2)$$

Loss

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
$o_1$, $o_2 \in$ [person, horse, …]
$r \in$ [on, in, ride, front of, …]

Our results:

Our results:

spatial, comparative, asymmetrical, verb, prepositional

person — taller than — person

person — left of — person

person → wear → shirt

person → on → snow

person → wear → ski

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

taller than

person              person

left of

wear            on     wear

shirt            snow     ski

Our results:

spatial, comparative, **asymmetrical**, verb, prepositional

taller than

person → person

left of

person → wear → shirt

person → on → snow

person → wear → ski

Our results:

spatial, comparative, asymmetrical, **verb**, prepositional

taller than

person          left of        person

wear             on       wear

shirt            snow      ski

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

taller than

person    person

left of

on    wear

wear

shirt    snow    ski

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

# Scene graphs can improve image retrieval



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

`<person - kicking - ball>`

`<person, guarding - goal>`

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

# 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

person

person

frisbee

Object representations

conv net

conv net

conv net

Feature extractors for each region

ROI regions

Classification loss

Bounding-box regression loss

...

RoI pooling

left of

throwing

# Representing objects as a graph with pairwise connections

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

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

Each node contains features from individual regions

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

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

Each node contains features from individual regions

Each node now also contains features from all other regions

# Graph representation with relationships included as nodes



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

Each node contains features from individual regions

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?

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



Each node contains features from individual regions

Each node now also contains features from all other regions

# Recall Convolutions

Images are a structured graph of pixels!

32x32xC activations
(if image then C = 3)
3x3xC filter, K filters

32

32

C

# Recall Convolutions

Images are a structured graph of pixels!



32x32xC activations
(if image then C = 3)
3x3xC filter, K filters

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

$W_i \in \mathbb{R}^{K \times C}$

# Recall Convolutions

Images are a structured graph of pixels!
Convolutions are local operations



32x32xC activations
(if image then C = 3)
3x3xC filter, K filters

$h_0 \quad h_1 \quad h_2$

$h_3 \qquad\qquad h_5$

$h_6 \quad h_7 \quad h_8$

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

32

32

C

| $W_0$ | $W_1$ | $W_2$ |
|-------|-------|-------|
| $W_3$ | $W_4$ | $W_5$ |
| $W_6$ | $W_7$ | $W_8$ |

$W_i \in \mathbb{R}^{K \times C}$

$h_4{}^{l+1} = W_4 h_4{}^l + W_0 h_0{}^l + W_1 h_1{}^l + W_2 h_2{}^l + \cdots + W_8 h_8{}^l$

# In comparison, scene graphs are not uniformly structured



Objects have 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

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

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

# 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
- $c_{ij}$ is a normalization constant

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

$$h_i^{l+1} = W_{self}\, h_i^l + \sum_{j \in N(i)} \frac{1}{c_{ij}} W_{other} h_j^l$$

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

# 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

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

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

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

Kipf & Welling (ICLR 2017)

To increase receptive field of CNNs: increase depth



Input       Conv          Conv          Conv

ReLU         ReLU

# To increase receptive field of GCNs: increase depth



GCNs: Graph Convolutional Networks

Kipf & Welling (ICLR 2017)

# Graph representation with edges included as nodes

left of

throwing

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

Each node contains features from individual regions

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
- $\sigma$ is a non-linearity, usually ReLU or LeakyReLU.

Without attention: $\quad h_i^{l+1} = W \; h_i^l + \sum_{j \in N(i)} \frac{1}{c_{ij}} W \; h_j^l$

With attention: $\quad h_i^{l+1} = W \; h_i^l + \sum_{j \in N(i)} \frac{1}{c_{ij}} \boxed{\alpha_{ij}} W \; h_j^l$

where $\quad \alpha_{ij} = \dfrac{e^{\sigma(a^T[Wh_i \| Wh_j])}}{\sum_{k \in \mathcal{N}(i)} e^{\sigma(a^T[Wh_i \| Wh_k])}}$

# How is it actually implemented?

For loops iterating over all the neighbors is expensive

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

# Formalizing a graph representation

For loops iterating over all the neighbors is expensive

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

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

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

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

# Formalizing a graph representation

For loops iterating over all the neighbors is expensive

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

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



For loops iterating over all the neighbors is expensive

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

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

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

Examples:

$$D_{00} = 2$$
$$D_{44} = 4$$

$$A_{04} = A_{40} = 1$$
$$A_{01} = A_{10} = 0$$

# 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

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

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

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

# 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

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

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

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

The vectorized computation of graph convolution is:

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

# 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:

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



For loops iterating over all the neighbors is expensive

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

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

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

Can be pre-calculated once per graph:

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

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

# Aside: Grounding to spectral convolutions with graph laplacian

Convolutions in the spectral domain:

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

Where U is the eigenvectors of the graph laplacian:

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

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

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

Renormalize the weights to get our spatial graph convolutions:

$$I + D^{-1/2} A D^{-1/2} \rightarrow D^{-1/2} \hat{A} D^{-1/2}$$

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

# 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 a. 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 a. 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 a. 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 a. 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

Wu et al. Long-term feature banks for detailed video understanding, CVPR 2019

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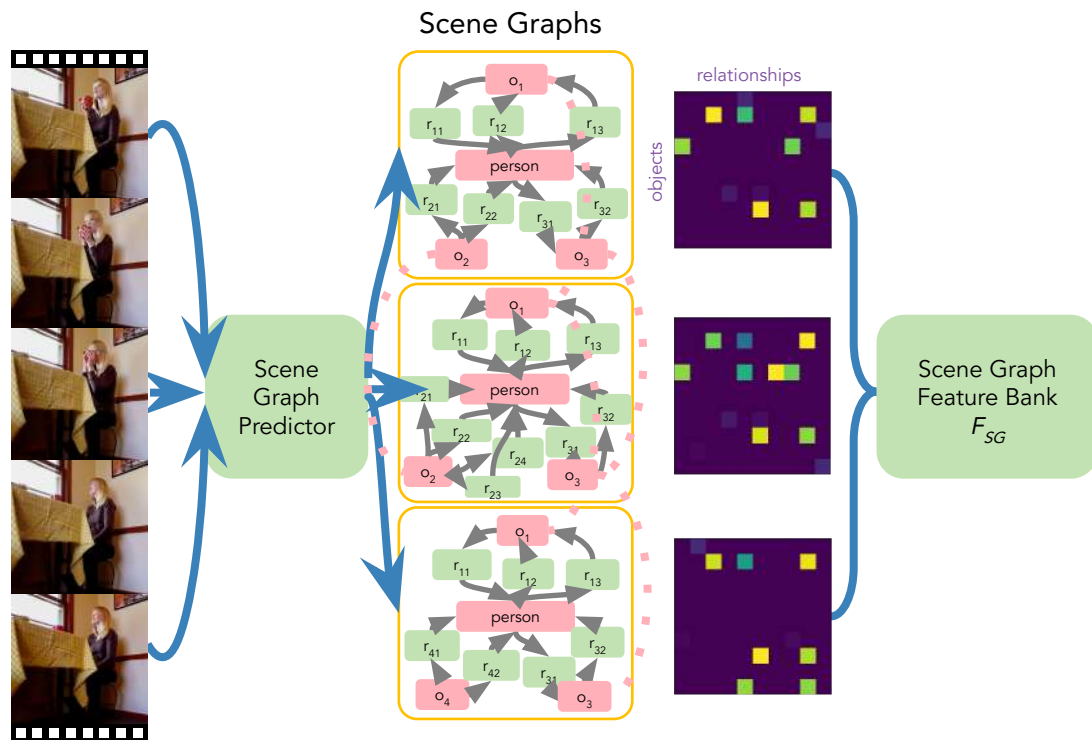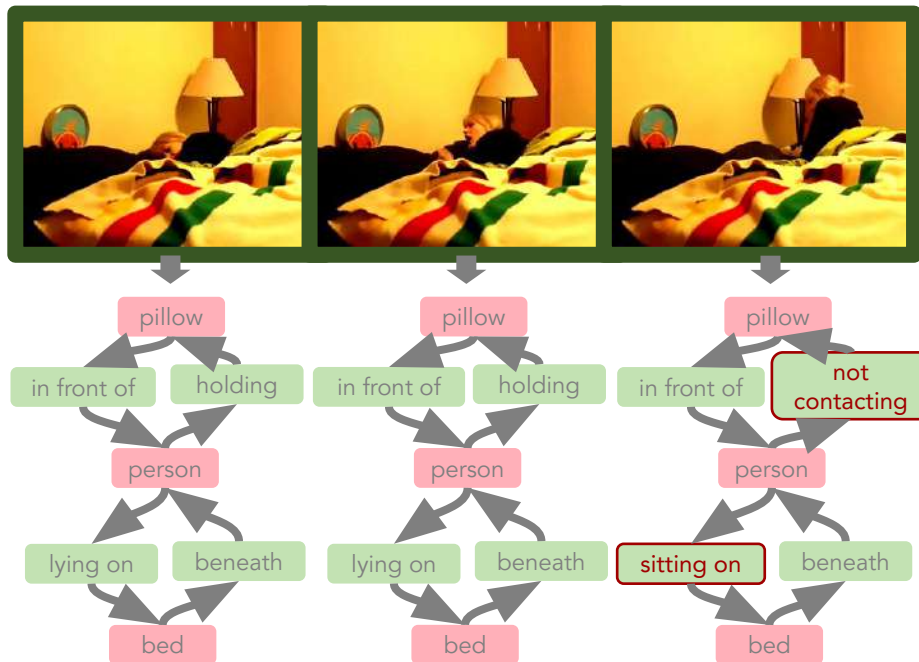Ji, Krishna et al. Action Genome: Actions as Compositions of Spatio-Temporal Scene Graphs, CVPR 2020

# 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



### Program synthesis



### Image captioning



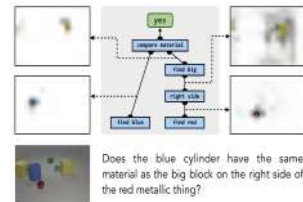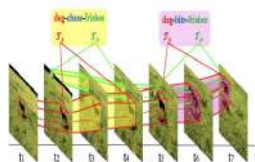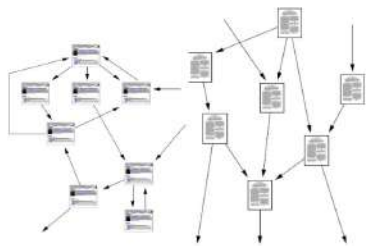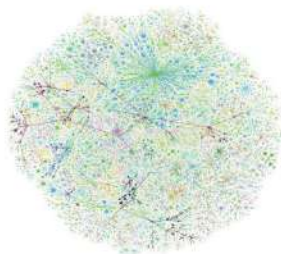### Video understanding

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
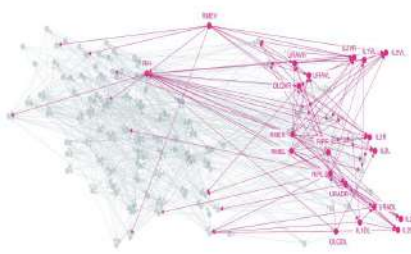
125

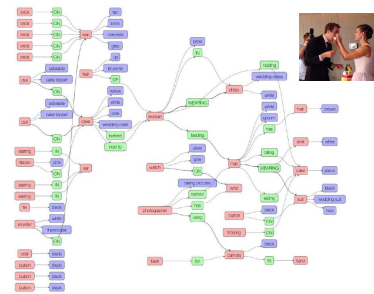# Graphs are everywhere – in numerous fields



Information networks:
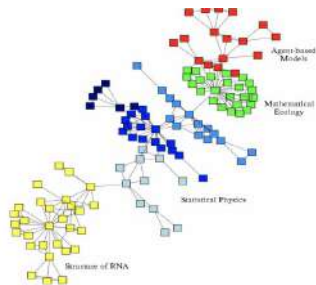Web & citations

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**   **Convolutional Neural Networks**   **Computer Vision Applications**

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

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

RNNs / LSTMs
Attention & Transformers
Image captioning
Interpreting neural networks
Style transfer
Adversarial examples
NeRF
Scene graphs
Graph Convolutions
Self-supervised learning
Multimodal learning
Perception & Action

Instructors

Teaching Assistants

Fei-Fei Li    Kevin Zakka (Head TA)    Haofeng Chen    Rachel Gardner    Samuel Kwong    Yichen Li

Ranjay Krishna    Sean Liu    Mandy Lu    Nishant Rai    Geet Sethi    Lin Shao

Danfei Xu    Guanzhi Wang    Chris Waites    Jiequan Zhang    Russel Xie

Course Coordinator

Yosefa Gilon

Instructors

Fei-Fei Li

Ranjay Krishna

Danfei Xu

Teaching Assistants

Kevin Zakka (Head TA)

Haofeng Chen

Rachel Gardner

Samuel Kwong

Yichen Li

Sean Liu

Mandy Lu

Nishant Rai

Geet Sethi

Lin Shao

Guanzhi Wang

Chris Waites

Jiequan Zhang

Russel Xie

Course Coordinator

Yosefa Gilon