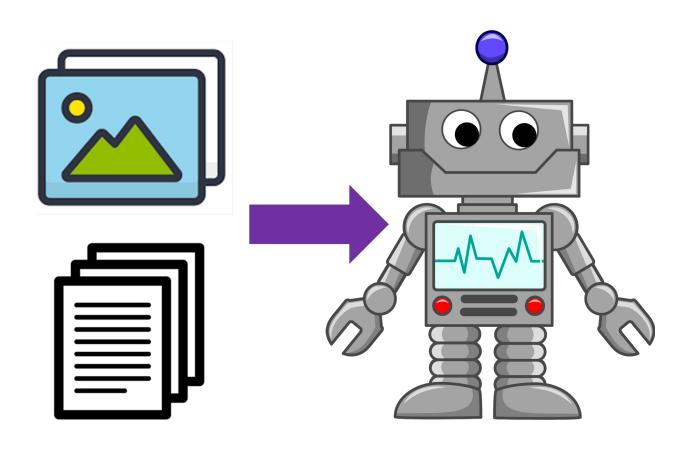
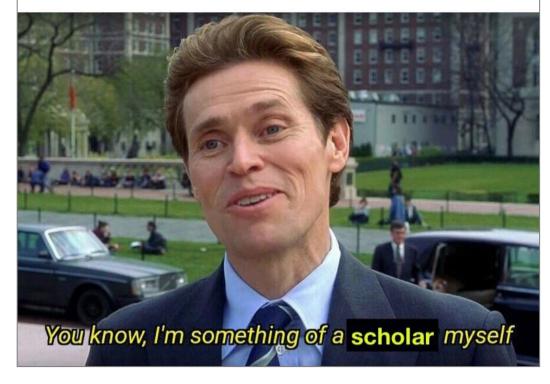
# Lecture 20a: CMSC 472 / 672 Computer Vision "Multimodal" Computer Vision

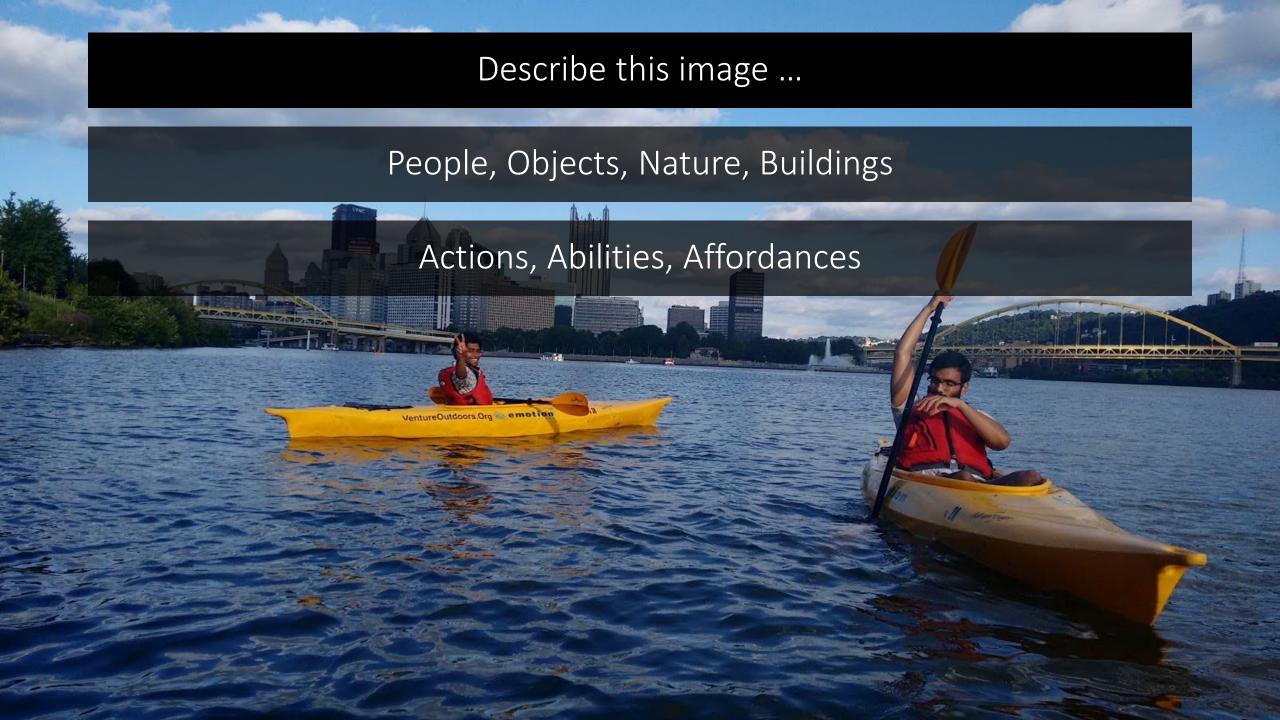


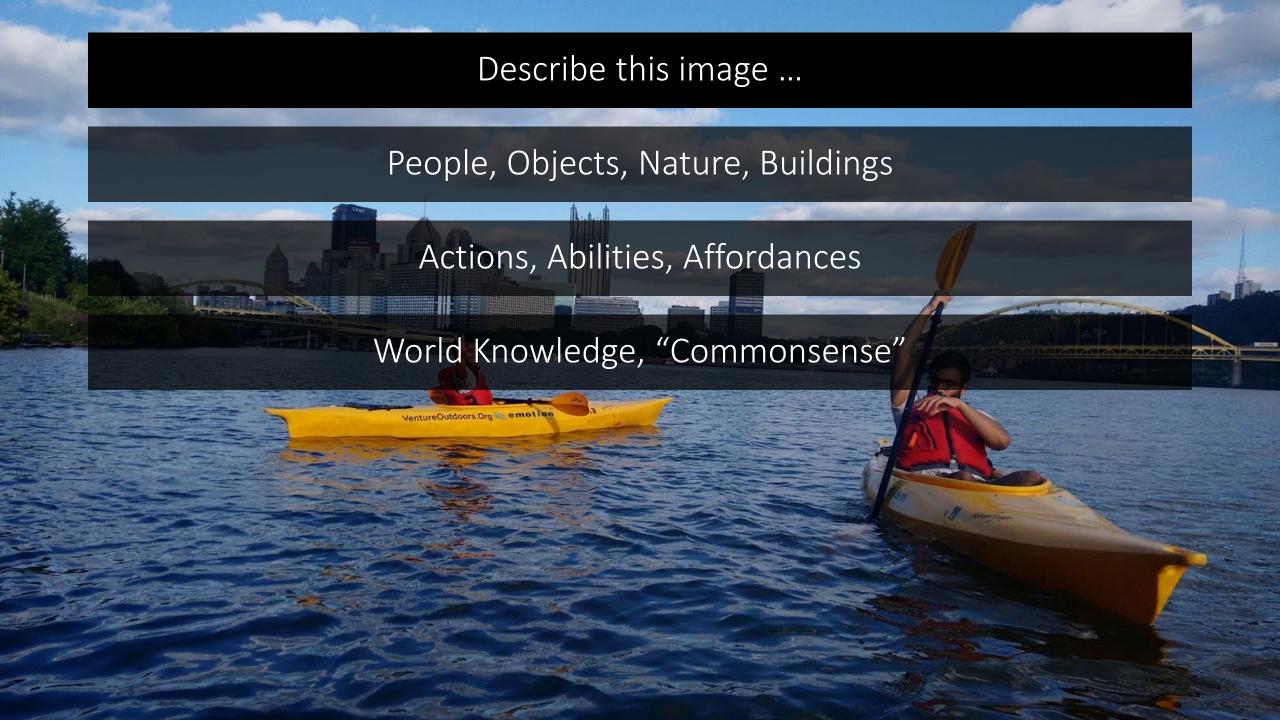
When you realize that memes are multimodal texts, making them a form of literature

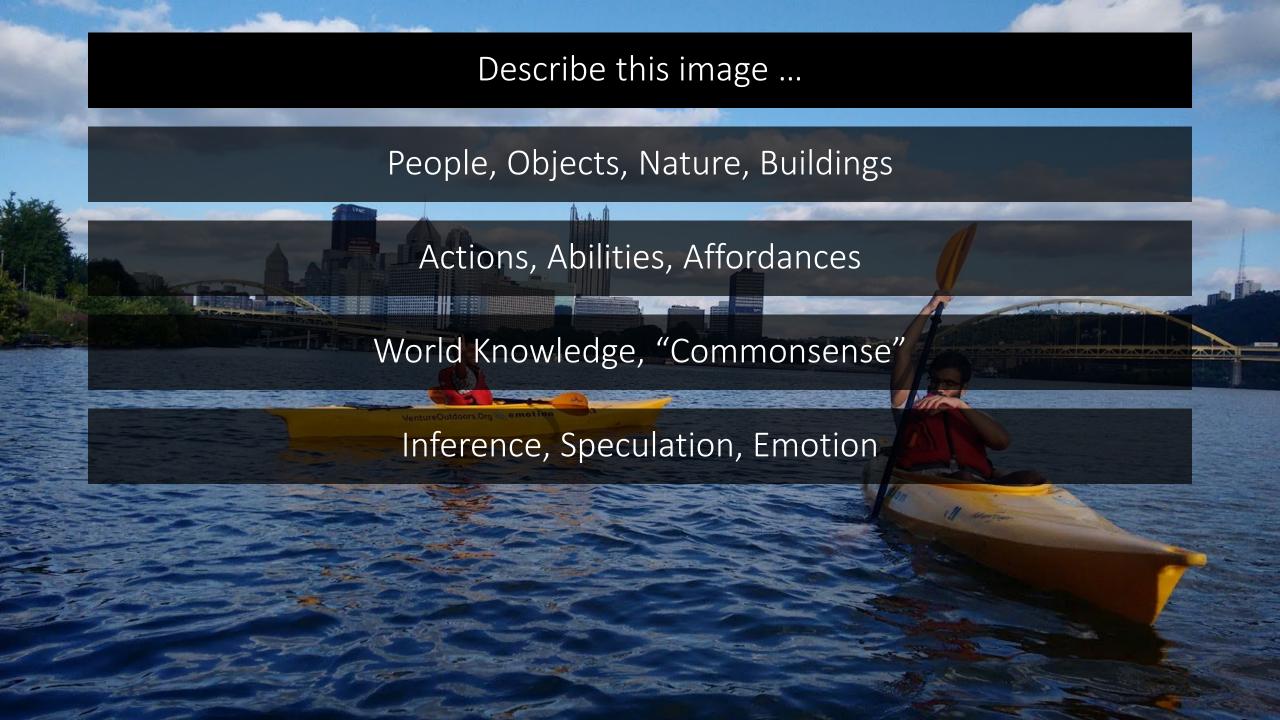










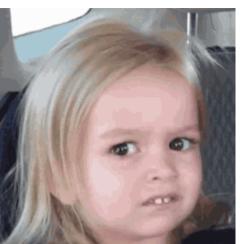


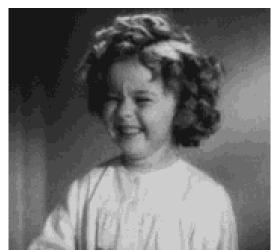






# Images convey emotions









# Vision + Language

A brand new era for computer vision!



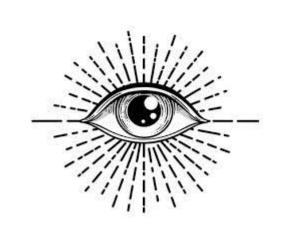
# DOG

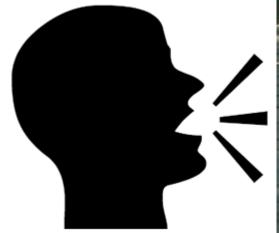


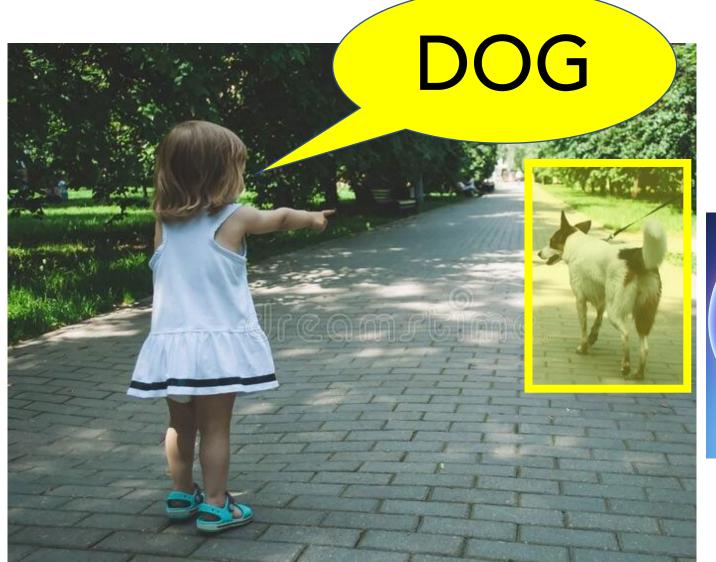
# Vision + Language

A brand new era for computer vision!

# Perception+Reasoning needs Vision+Language





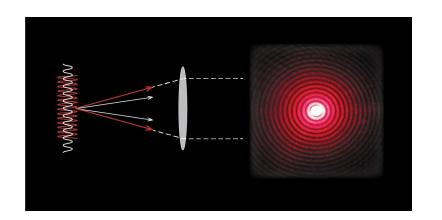




# Computer Vision: A Pyramid

#### PHYSICS-BASED

- Optics
- Computational Imaging

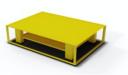


#### **GEOMETRIC**

- 3D Reconstruction
- Shape, Depth, ...



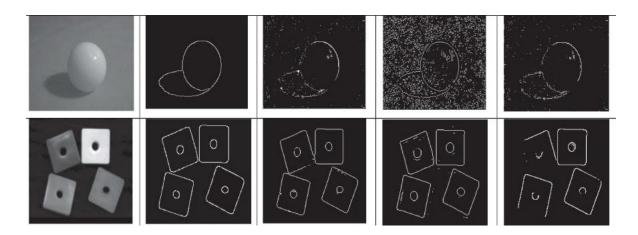


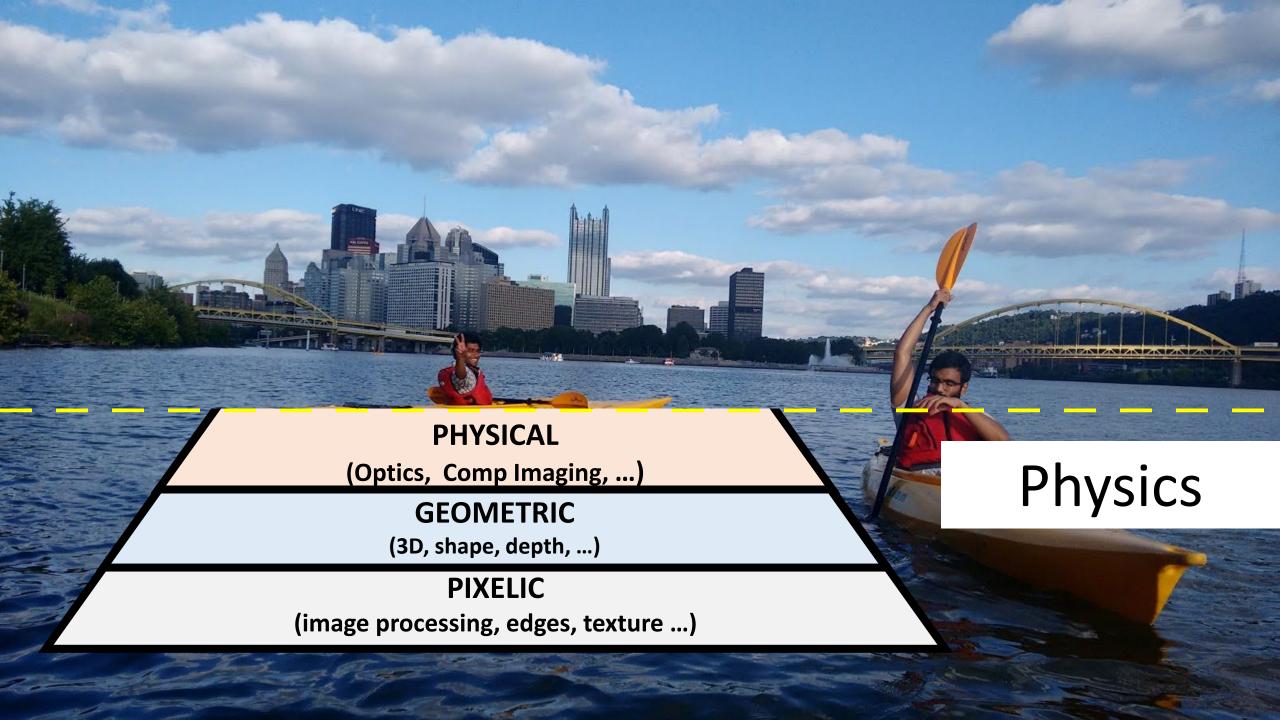


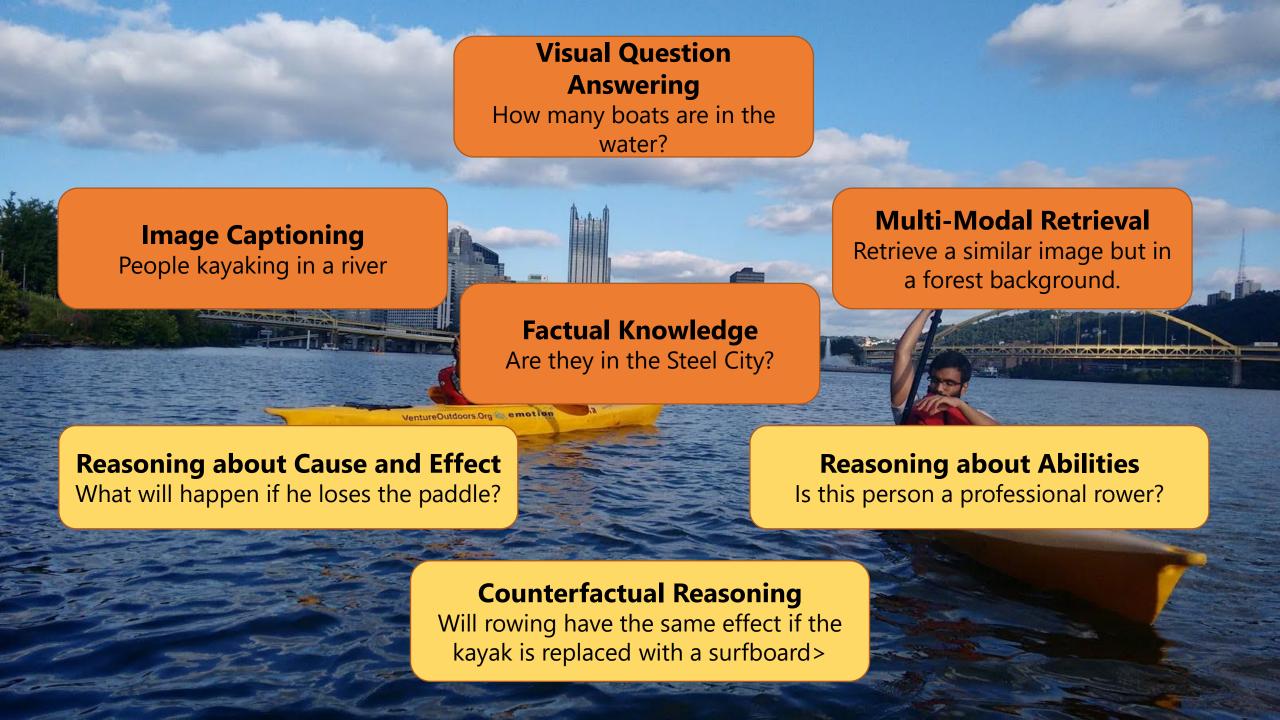


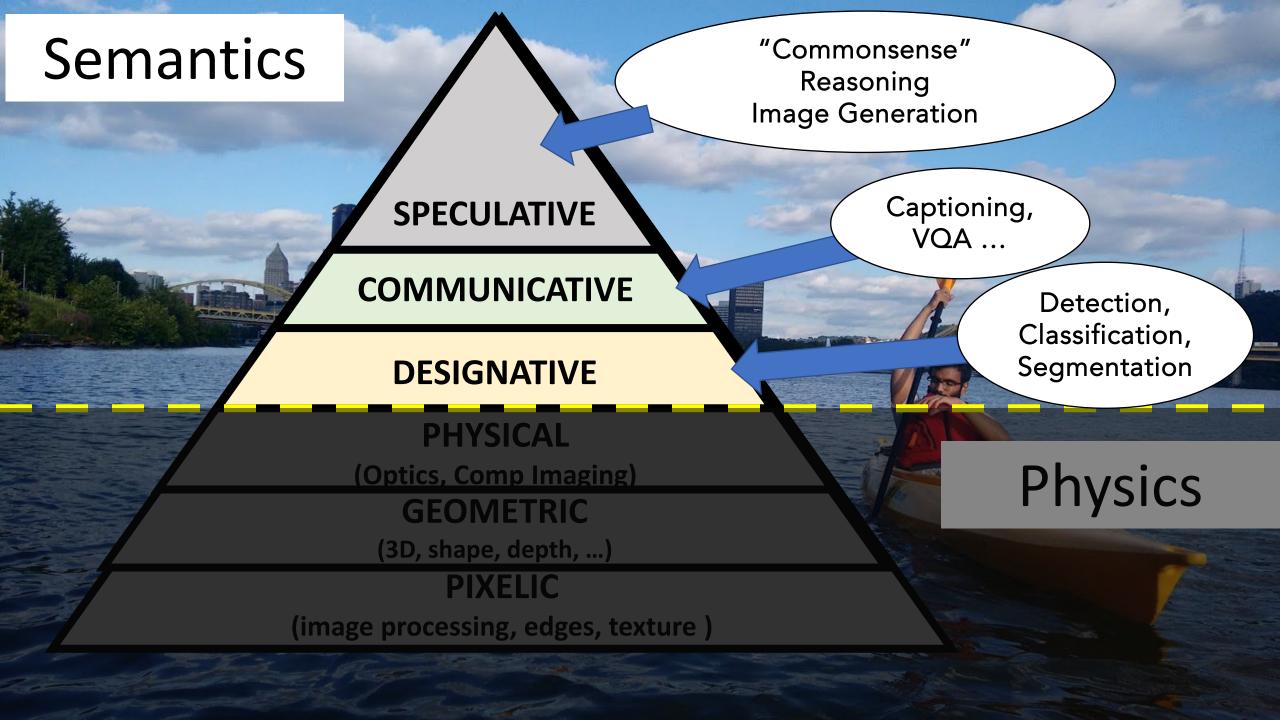
#### **PIXELIC**

- Image Processing
- Edge Detection











A New Paradigm

**Automatically captioned** 



An astronaut Teddy bears A bowl of soup

riding a horse lounging in a tropical resort in space playing basketball with cats in space

in a photorealistic style in the style of Andy Warhol as a pencil drawing

**COMMUNICATIVE** 

**DESIGNATIVE** 

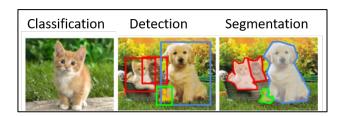


What is the mustache made of?

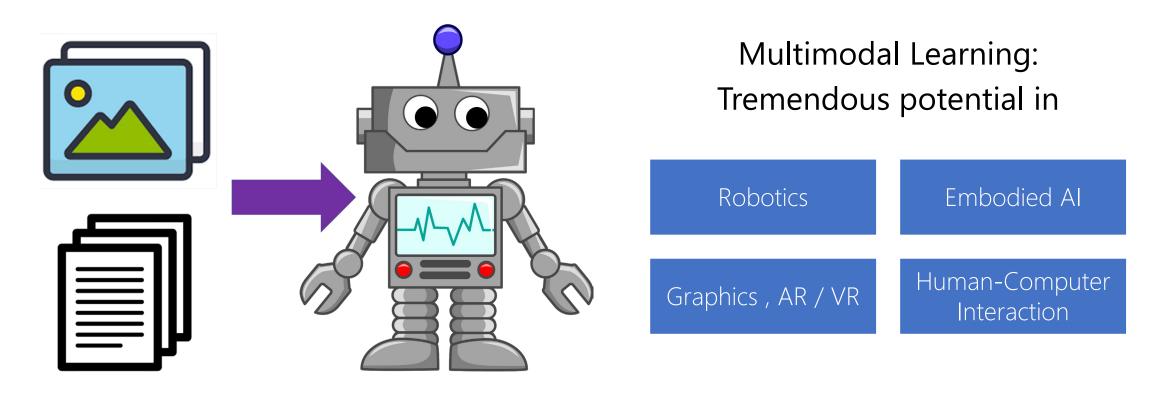
Al System

DALL-E 2

bananas



#### Multi-Modal (Vision + Language) Learning



Learning jointly from images and text has caused a paradigm shift in Al



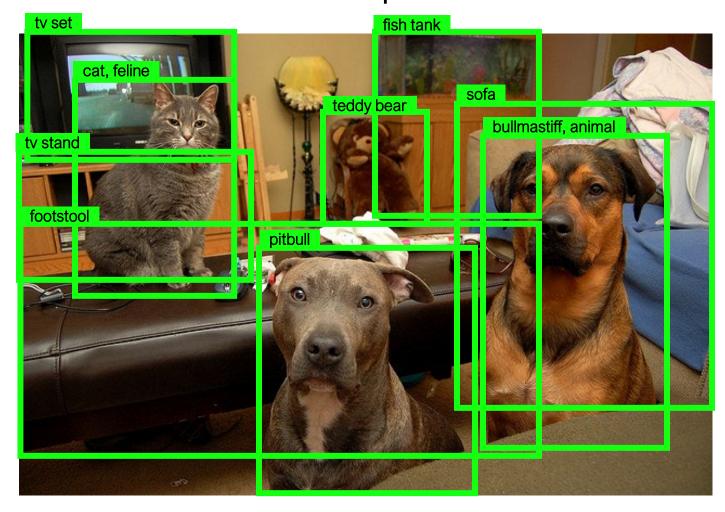
#### Old Computer Vision



Image tagging / Image classification

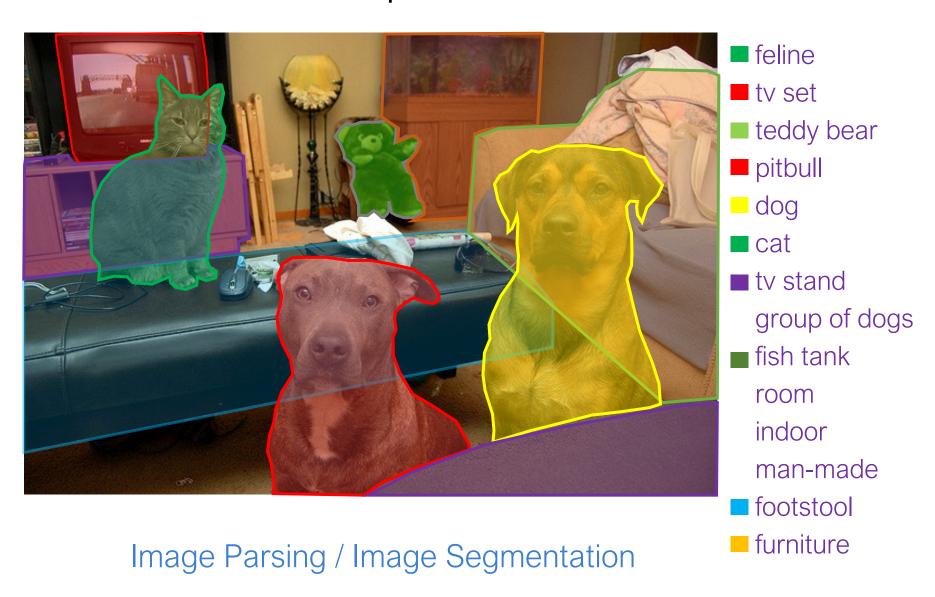
feline tv set teddy bear pitbull bullmastiff cat tv stand group of dogs fish tank room indoor man-made footstool furniture

#### Old Computer Vision



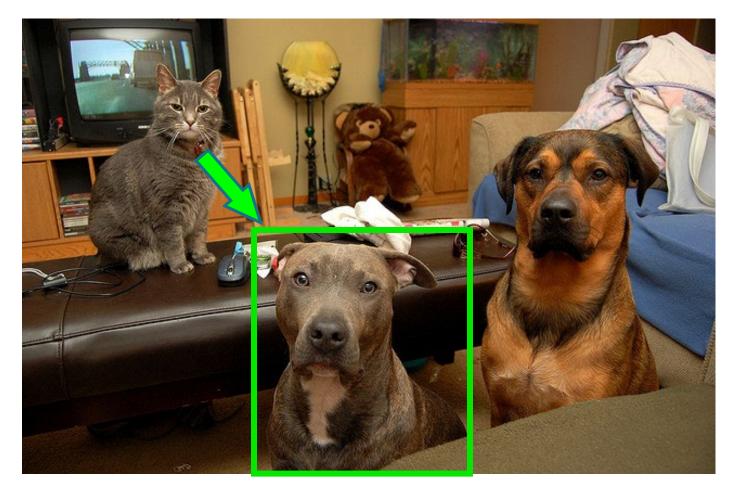
feline tv set teddy bear pitbull bullmastiff cat tv stand group of dogs fish tank room indoor man-made footstool furniture

### Old Computer Vision



New tasks with Vision + Language

### Referring to objects

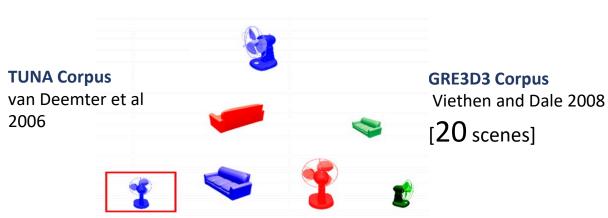


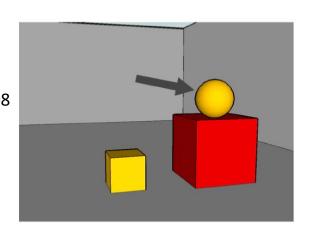
The dog in the middle

The gray dog in the middle

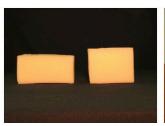
The gray dog

# Work on Referring Expression





Size Corpus
Mitchell et al 2011
[96 scenes]





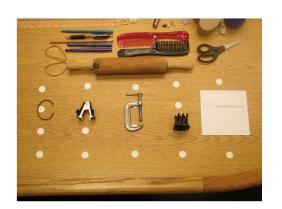




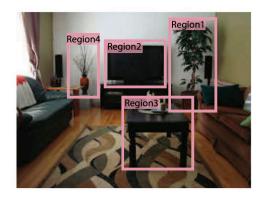
**GenX Corpus**FitzGerald et al 2013
[269 scenes]



Typicality Corpus
Mitchell et al 2013
[35 scenes]

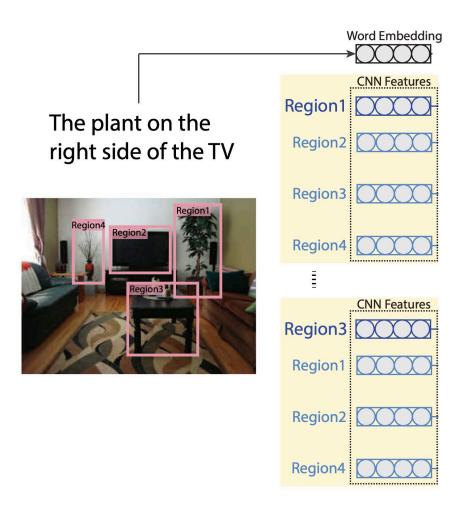


The plant on the right side of the TV

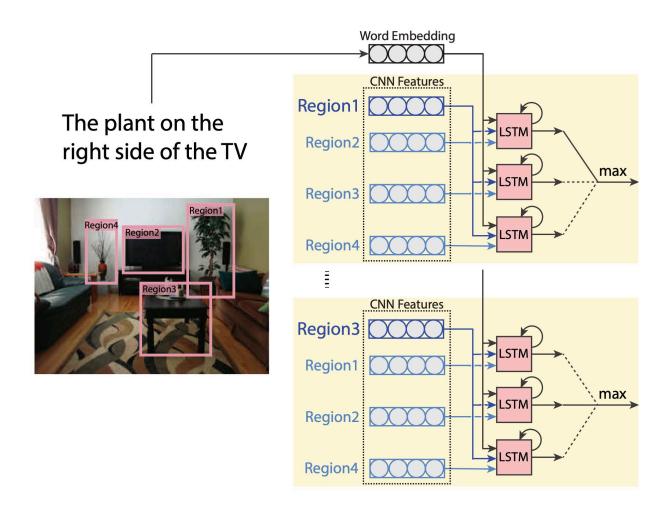


Modeling Context Between Objects for Referring Expression Understanding

Varun K. Nagaraja Vlad I. Morariu Larry S. Davis

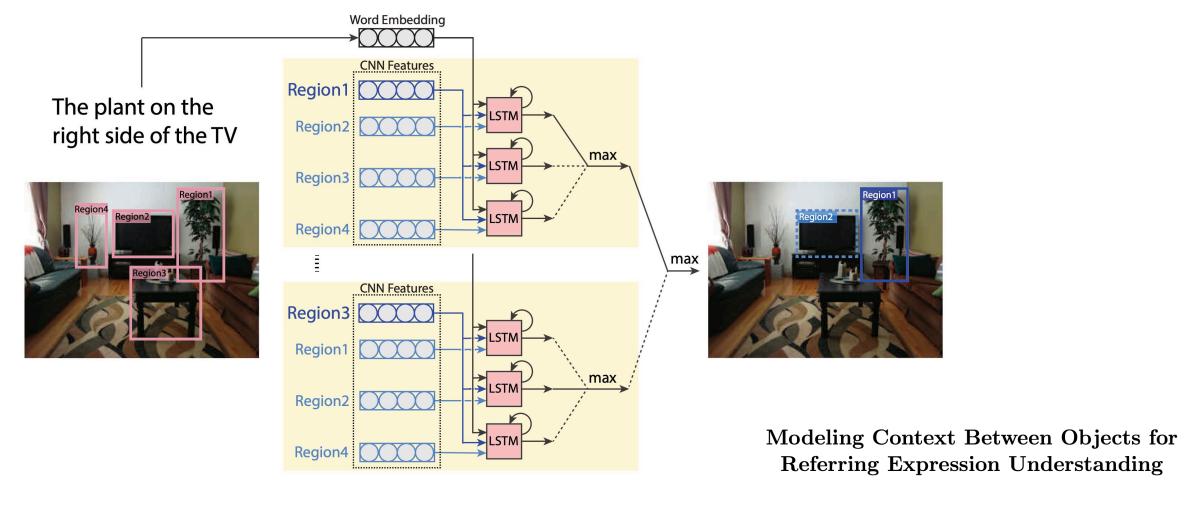


Modeling Context Between Objects for Referring Expression Understanding



#### Modeling Context Between Objects for Referring Expression Understanding

Varun K. Nagaraja Vlad I. Morariu Larry S. Davis



## Visual Question Answering

Given an image and a question about it, produce an answer to that question.



yes 100.000% no 0.000%

Is the food made of eggs?

# VQA: Visual Question Answering www.visualqa.org

Aishwarya Agrawal\*, Jiasen Lu\*, Stanislaw Antol\*, Margaret Mitchell, C. Lawrence Zitnick, Dhruv Batra, Devi Parikh

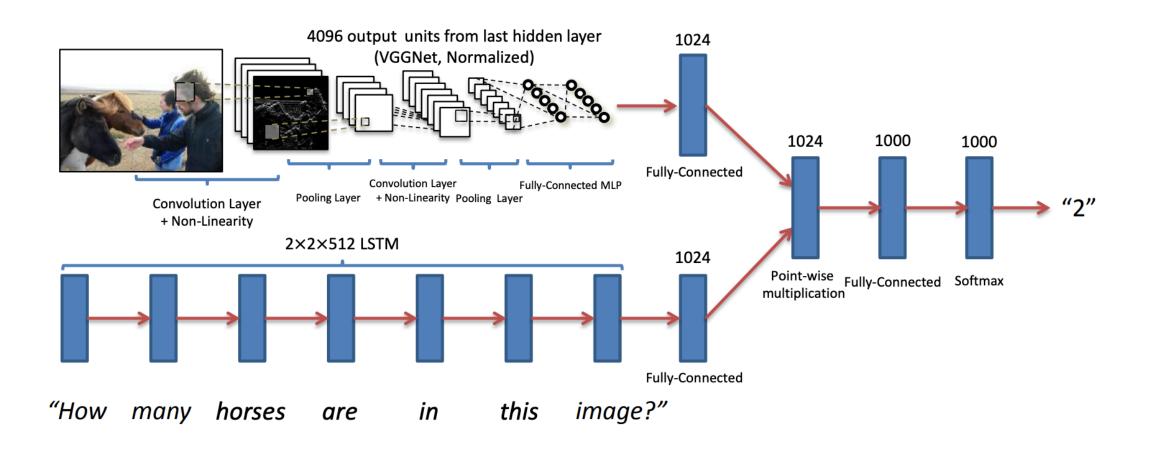


Is this person trying to hit a ball?	yes yes yes	yes yes
What is the person hitting the ball with?	frisbie racket round paddle	bat bat racket

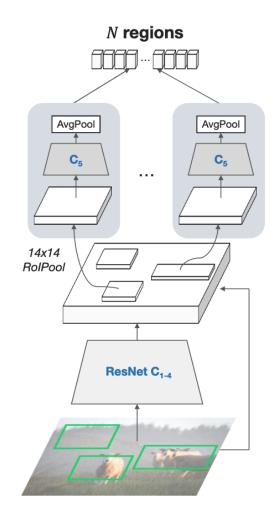


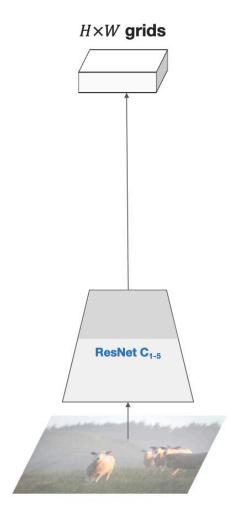
What is the guy	phone	reading
doing as he sits	taking picture	reading
on the bench?	taking picture with phone	smokes
What color are his shoes?	blue blue blue	black black brown

## Visual Question Answering: Naïve Approach



### What Features to use as input visual features?





#### CVPR 2017

## **Bottom-Up and Top-Down Attention for Image Captioning and Visual Question Answering**

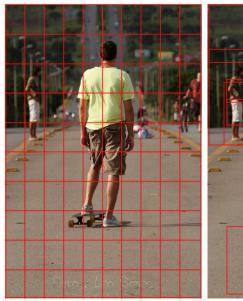
Peter Anderson<sup>1\*</sup> Xiaodong He<sup>2</sup> Chris Buehler<sup>3</sup> Damien Teney<sup>4</sup>

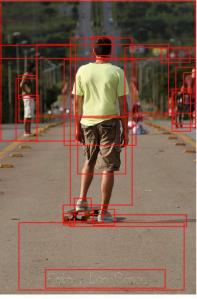
Mark Johnson<sup>5</sup> Stephen Gould<sup>1</sup> Lei Zhang<sup>3</sup>

<sup>1</sup>Australian National University <sup>2</sup>JD AI Research <sup>3</sup>Microsoft Research <sup>4</sup>University of Adelaide <sup>5</sup>Macquarie University

<sup>1</sup>firstname.lastname@anu.edu.au, <sup>2</sup>xiaodong.he@jd.com, <sup>3</sup>{chris.buehler,leizhang}@microsoft.com

<sup>4</sup>damien.teney@adelaide.edu.au, <sup>5</sup>mark.johnson@mg.edu.au



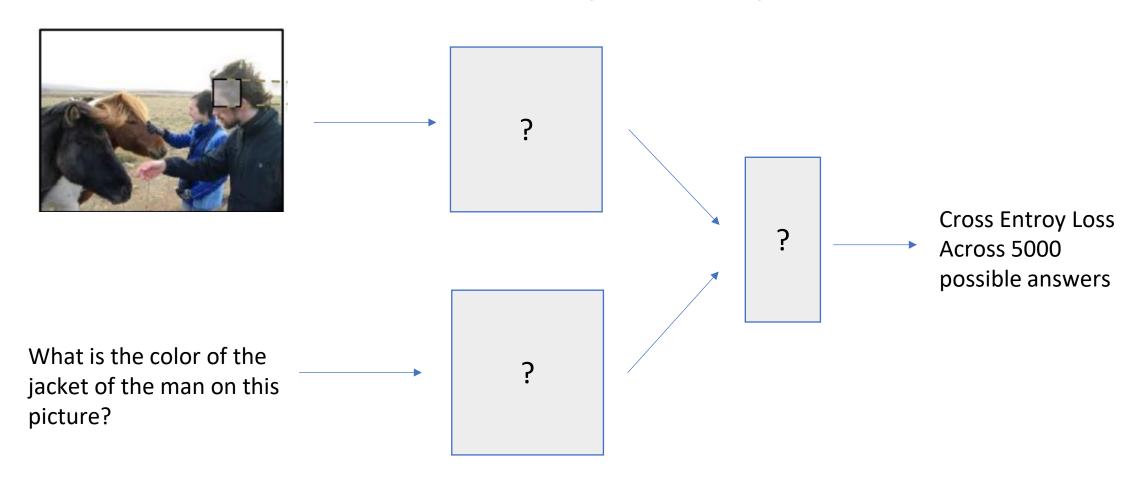




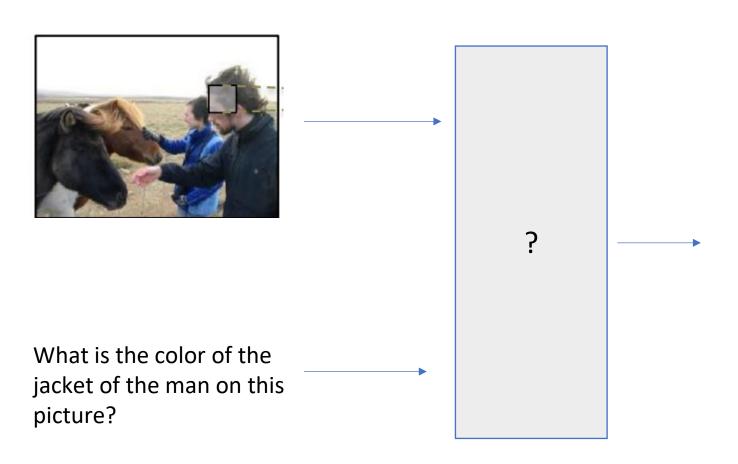


Question: What room are they in? Answer: kitchen

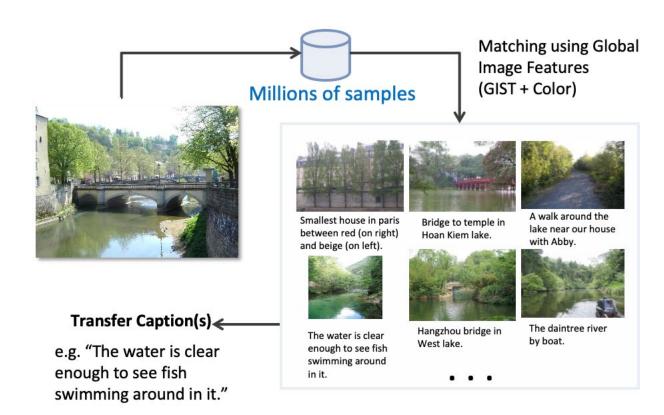
## VQA Solution 5 years ago: Learn V and L features separately, and fuse.



# VQA Solution from 2019 ... Multimodal Pretraining (typically using masked language modeling)



# Describing images with language (Image Captioning)





Im2Text: Describing Images Using 1 Million Captioned Photographs
Vicente Ordonez, Girish Kulkarni, Tamara L. Berg.
Advances in Neural Information Processing Systems. NIPS 2011. Granada, Spain.

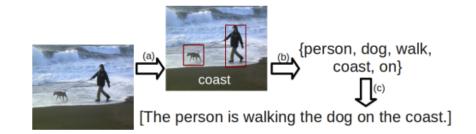


Figure 3: Overview of our approach. (a) Detect objects and scenes from input image. (b) Estimate optimal sentence structure quadruplet  $\mathcal{T}^*$ . (c) Generating a sentence from  $\mathcal{T}^*$ .

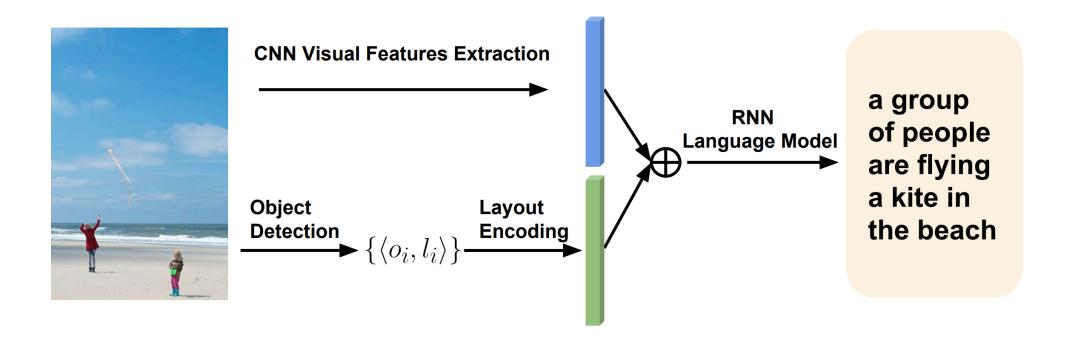
**Corpus-Guided Sentence Generation of Natural Images** 



#### **EMNLP 2011**

Yezhou Yang † and Ching Lik Teo † and Hal Daumé III and Yiannis Aloimonos
University of Maryland Institute for Advanced Computer Studies
College Park, Maryland 20742, USA
{yzyang, cteo, hal, yiannis}@umiacs.umd.edu

# One method for image captioning ...

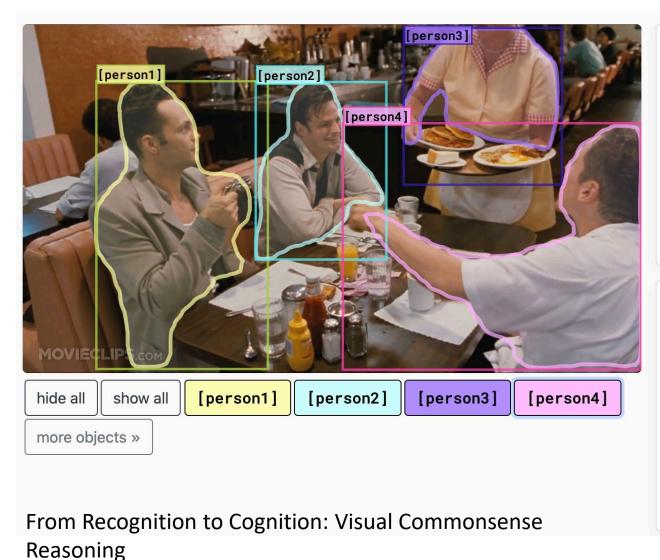


Obj2Text: Generating Visually Descriptive Language from Object Layouts
Xuwang Yin, Vicente Ordonez. Empirical Methods in Natural Language Processing.

EMNLP 2017. Copenhagen, Denmark. September 2017. [pdf] [arxiv] [code] [bibtex]

(~Oral presentation)

# Visual Common Sense Reasoning



Why is [person4 ] pointing at [person1 ]? a) He is telling [person3 that [person1 ordered] ordered the pancakes. b) He just told a joke. c) He is feeling accusatory towards [person1 ]. d) He is giving [person1 directions. Rationale: I think so because... a) [person1 has the pancakes in front of him. b) [person4 is taking everyone's order and asked for clarification.

https://visualcommonsense.com/

c) [person3[3]] is looking at the pancakes both she and

d) [person3 ] is delivering food to the table, and she

[person2 ] are smiling slightly.

might not know whose order is whose.

### **Enriching Video Captioning with Commonsense Descriptions**



### **Standard Caption**

A band is playing at a concert

### **Generated Commonsense Descriptions**

Intention

to entertain the audience

Effect

will get standing ovation

### **Video2Commonsense Dataset**

- Videos of agents doing actions
- Annotations for intentions of agents, effect of actions

### **Benchmarking Video Captioning**

- Existing models found lacking
- Guidance from commonsense knowledge bases required













# Video2Commonsense Enriching Video Captioning with Commonsense Descriptions



**Conventional Caption** 

Commonsense-Enriched Caption

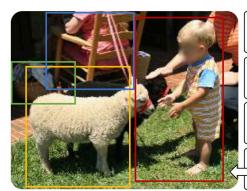
Commonsense Question Answering Group of runners get prepared to run a race.

In order to win a medal, a group of runners get prepared to run a race. As a result they are congratulated at the finish line. They are athletic.

What happens next to the runners?

Are congratulated at the finish line become tired

# Multi-task Learning / More General Models



Visual Question Answering

What color is the child's outfit? Orange

Referring Expressions

sheep basket people sitting on chair

Multi-modal Verification

The child is petting a dog. false

Caption-based Image Retrieval

A child in orange clothes plays with sheep.

12-in-1: Multi-task Vision and Language Representation Learning https://arxiv.org/abs/1912.02315 **CVPR 2020** 

#### **Ouestion**

What is a major importance of Southern California in relation to California and the US?

What is the translation from English to German?

What is the summary?

Hypothesis: Product and geography Premise: Conceptually cream are what make cream skimming work. Entailment, neutral, or contradiction?

positive or negative?

#### Context

economic center for the state of California and the US....

Most of the planet is ocean water.

Radcliffe gains access to a reported £320 million fortune...

Is this sentence

...Southern California is a major

Harry Potter star Daniel

skimming has two basic dimensions - product and geography.

A stirring, funny and finally transporting re-imagining of Beauty and the Beast and 1930s horror film.

#### Answer

major economic center

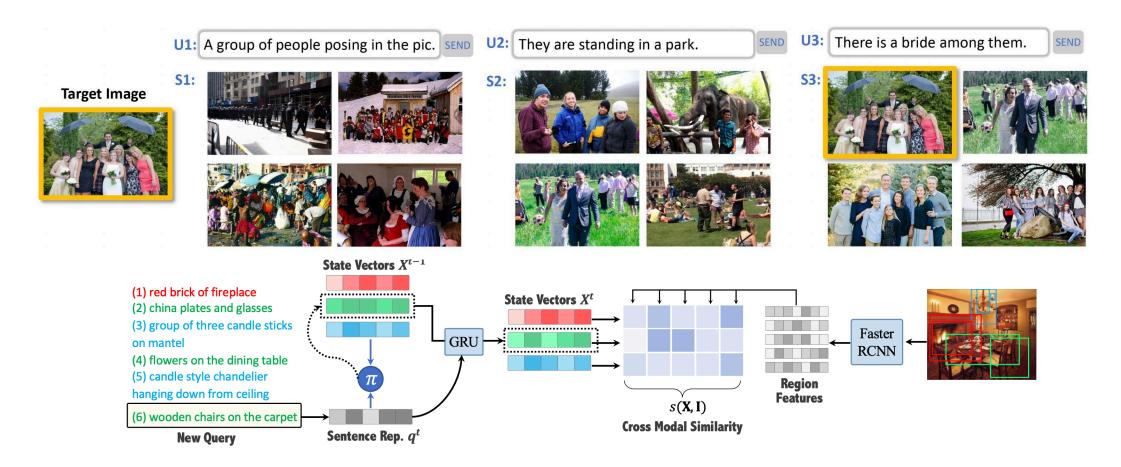
Der Großteil der Erde ist Meerwasser

Harry Potter star **Daniel Radcliffe gets** £320M fortune...

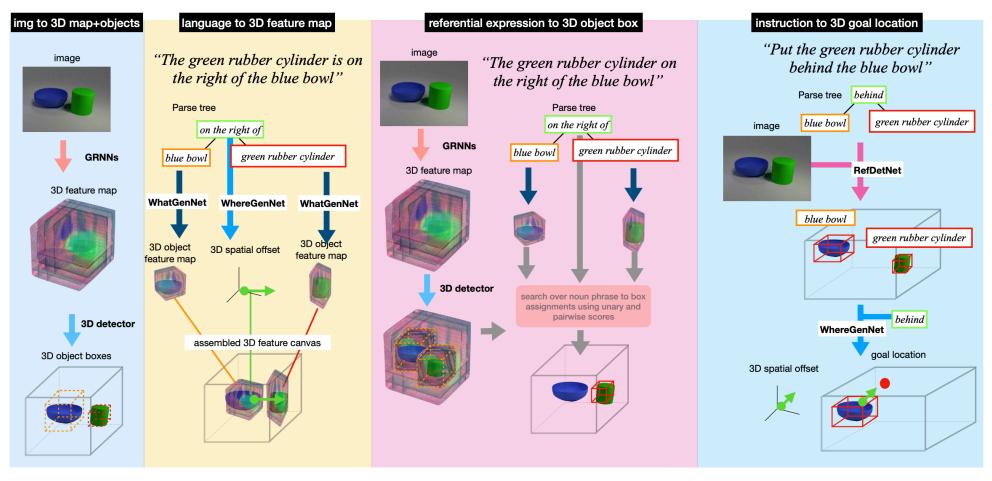
Entailment

positive

## Interactivity + Language and Vision

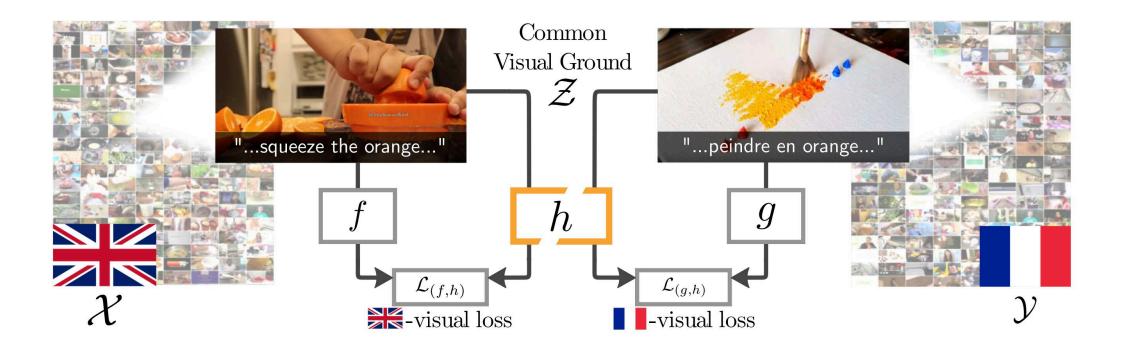


# Vision + Language + 3D



https://arxiv.org/abs/1910.01210

# Multiple Languages and Vision



# Video + Language Tasks













00:00:03,576 --> 00:00:05,697 Gavin Mitchell's office. Rachel Green's office. 00:00:05,870 --> 00:00:07,409 Give me that phone.

00:00:08,873 --> 00:00:12,293 Hello, this is Rachel Green. How can I help you? 00:00:12,460 --> 00:00:17,629 Uh-huh. Okay, then. I'll pass you back to your son. 00:00:18,800 --> 00:00:21,639 Hey, Mom. No, that's just my secretary.

(positive) The woman becomes upset when the man answers the phone because he pretends it is his own office.

(negative) The woman becomes upset when the man answers the phone because she is expecting a phone call from her mom.

Inferring reasons

(positive) The woman realizes it is the man's mother who is calling and she passes the phone back to the man.

(negative) The man realizes it is the woman's mother who is calling and he passes the phone back to the woman.

Identifying characters

(positive) The phone rings, a man picks it up, and a woman slams her hand on the desk and demands the man give her the phone.

(negative) The two people that the man in the glasses is talking to need to be briefed on something.

Global video understanding

# Counterfactuals in Vision and Language

Question Image	Counterfactual Questions	Counterfactual Images
Is this in Australia?	1. Is the grass green? 2. Is there grass on the ground? 3. Are they standing on a green grass field? 4. Is the stop light green?	
What color is the person's helmet?	1. What color jacket is the girl wearing? 2. What color jacket is the person wearing? 3. What color is the jacket? 4. What color is the woman's jacket?	
Where did the shadow on the car come from?	1. What kind of dog is this? 2. What type of dog is this? 3. What kind of dog is shown? 4. What is the breed of dog?	

https://openaccess.thecvf.com/content\_CVPR\_2020/papers/Abbasnejad\_Counterfactual\_Vision\_and\_Language\_Learning\_CVPR\_2020\_paper.pdf

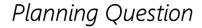
### Asking Counterfactual Questions to Reason about Physical Properties



**Input Video** 

Counterfactual Question
What will happen if the yellow cube is removed?

(A) Purple Cube will collide with brown cube



How can the collision between yellow and purple cube be stopped?

(A) **Add** teal sphere to the right of purple sphere









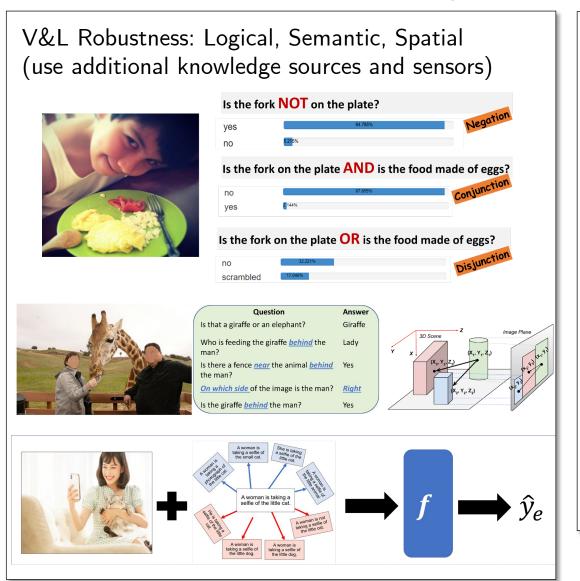


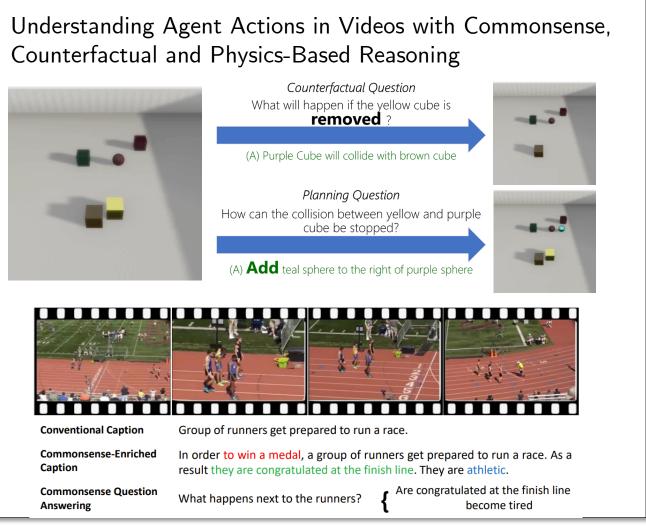


**Effect of Action** 

## My work: Perception & Reasoning with Robustness

### Robust Visual Reasoning (Visual QA, Video Captioning, V&L Inference)





Gokhale ECCV '20; Gokhale EMNLP'20; Gokhale ACL'21; Fang EMNLP'20; Banerjee ICCV'21; Patel EMNLP'22

All of this research by hundreds of people in hundreds of places ...

culminated in "Vision-Language Models (VLMs)" that you have today (in ChatGPT, Gemini, etc.)

These ideas started as research projects led by **PhD students** and with enough academic evidence, led to companies and products.





### Isaac Newton remarked in a 1675 letter to his rival Robert Hooke:

"What Des-Cartes did was a good step.

You have added much several ways, & especially in taking the colours of thin plates into philosophical consideration.

If I have seen further, it is by standing on the shoulders of Giants."