# Words, Pictures, and Common Sense

**Devi Parikh**

**Georgia Tech**

what i think

what i say

"Color College Avenue", Blacksburg, VA, May 2012

2

# People coloring a street in rural Virginia.

"Color College Avenue", Blacksburg, VA, May 2012

3

# It was a great event! It brought families out, and the whole community together.

Q. **What are they coloring the street with?**
A. **Chalk**

"Color College Avenue", Blacksburg, VA, May 2012

AI: What a nice picture! What event was this?

User: *"Color College Avenue". It was a lot of fun!*

AI: I am sure it was! Do they do this every year?

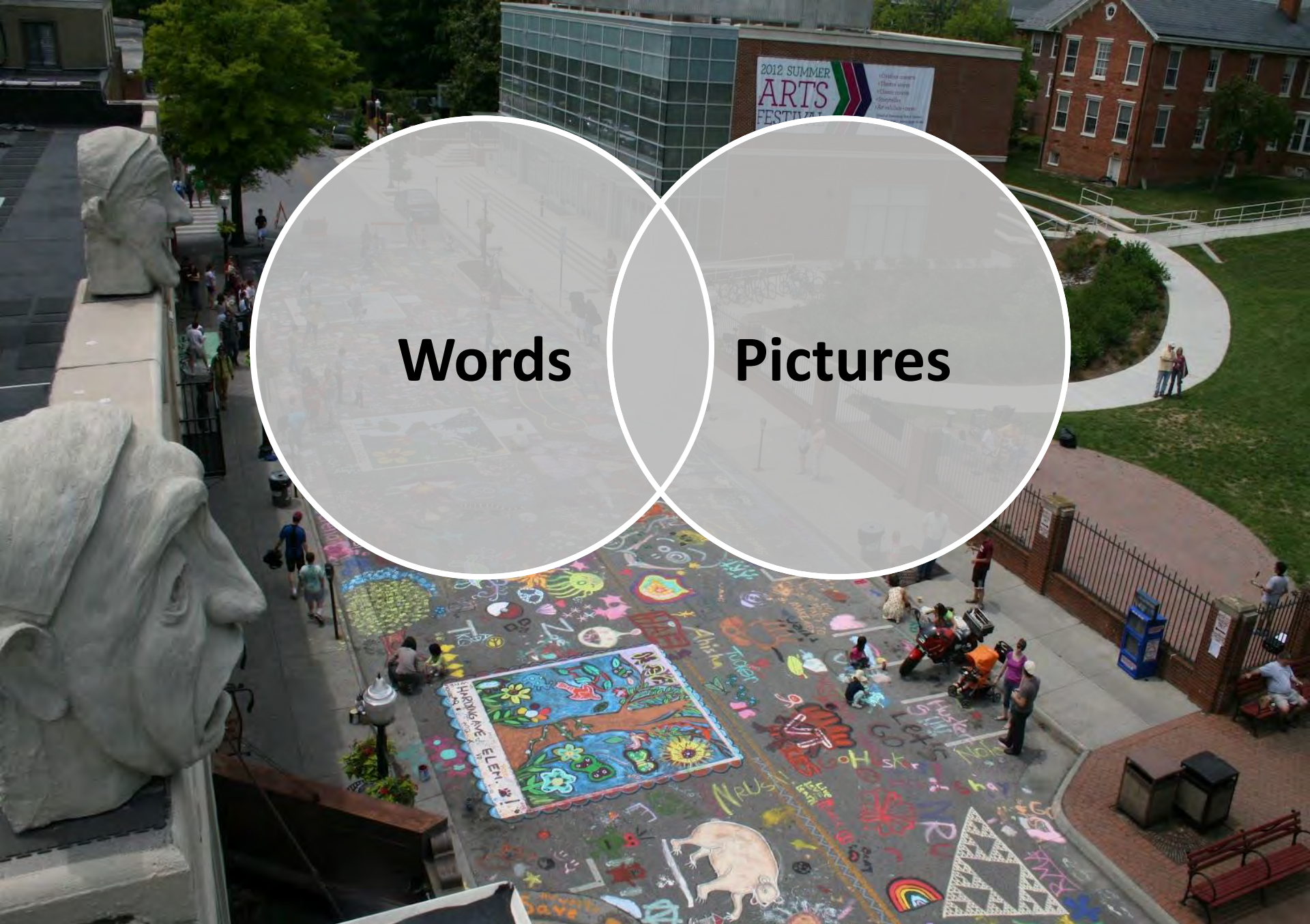User: *I wish they would. I don't think they've organized it again since 2012.*

...

"Color College Avenue", Blacksburg, VA, May 2012

# Pictures

"Color College Avenue", Blacksburg, VA, May 2012

**Words**     **Pictures**

Words

Pictures

Common Sense

"Color College Avenue", Blacksburg, VA, May 2012

# Applications

Pictures are everywhere

Words are how we communicate

# Applications

Interact with, organize, and navigate visual data

# Applications

## Leverage multi-modal information on the web

# Applications

## Aid visually-impaired users

# Applications

## Aid visually-impaired users



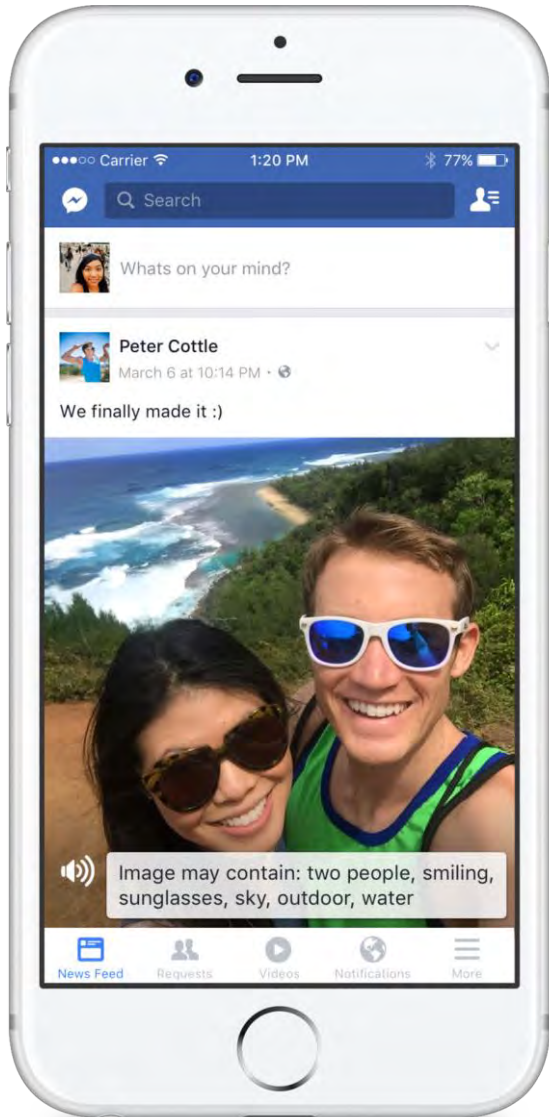FACEBOOK'S AI CAN CAPTION PHOTOS FOR THE BLIND ON ITS OWN

# Applications

# Applications

## Summarize visual data for analysts

Did anyone enter this room last week?

Yes, 127 instances logged on camera

Were any of them carrying a black bag?

...

# Applications

## Natural language instructions to an agent



Is there smoke in any room around you?
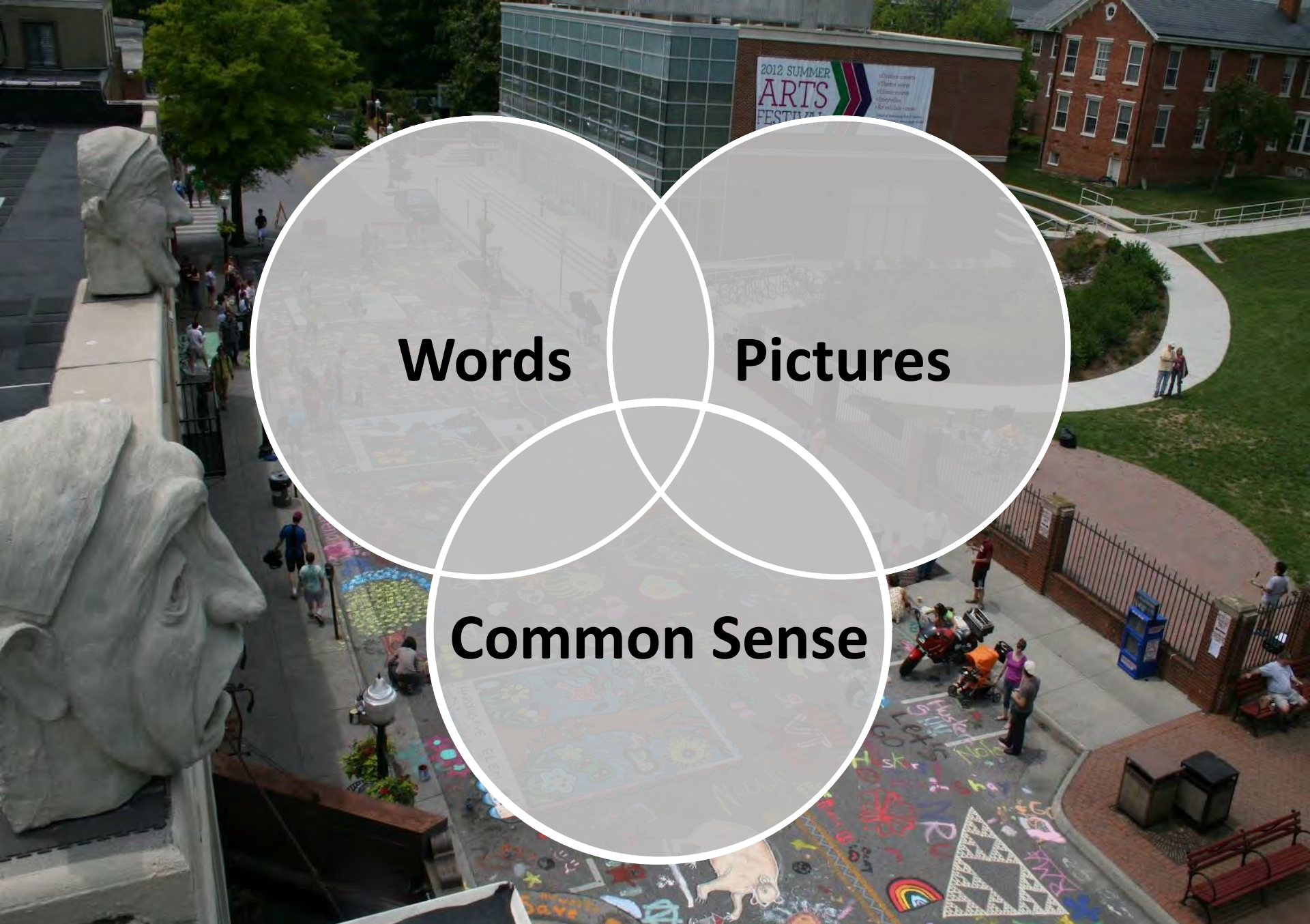
Yes, in one room

Go there and look for people

…

17

Words

Pictures

Common Sense

Words   Visually   Pictures

Grounded

Dialog

Common Sense

"Color College Avenue", Blacksburg, VA, May 2012

# Visual Question Answering (VQA)

# Visual Question Answering (VQA)



What is the mustache made of?

# Visual Question Answering (VQA)



What is the mustache made of?

AI System

# Visual Question Answering (VQA)



What is the mustache made of?
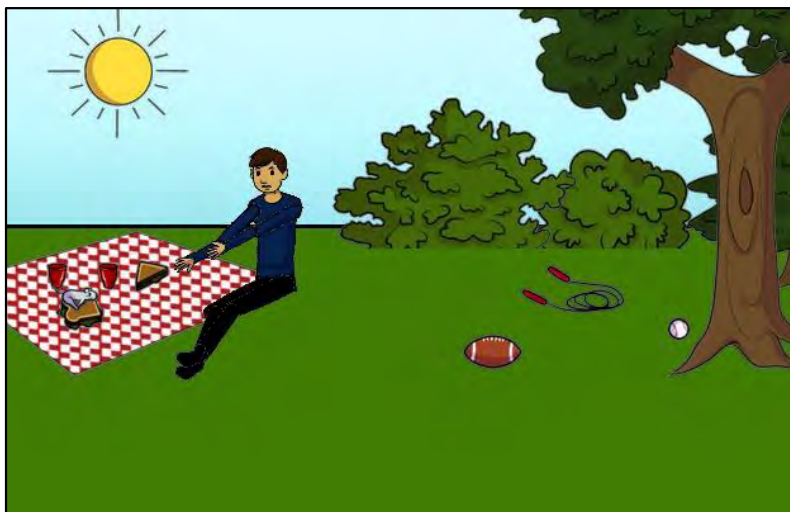
AI System

bananas

# Visual Question Answering (VQA)



What color are her eyes?
What is the mustache made of?



How many slices of pizza are there?
Is this a vegetarian pizza?



Is this person expecting company?
What is just under the tree?



Does it appear to be rainy?
Does this person have 20/20 vision?

>0.25 million images

# 254,721 images (MS COCO)

# 50,000 scenes

>0.25 million images

>0.76 million questions

# Questions



Stump a smart robot! Ask a question about this image that a human can answer, but a smart robot probably can't!

**Stump a smart robot!**
**Ask a question that a human can answer,**
**but a smart robot probably can't!**

We have built [...] kitchen, beach) [...] can recognize the scene (e.g, [...] smart robot!

Ask a question [...]
IMPORTANT: T[...] should not be able to answer
the question w[...]

[...]ns below:

- Do not repeat questions. Do not ask the same questions or the same questions with minor variations over and over again across images. Think of a new question each time specific to each image.

- Each question should be a single question. Do not ask questions that have multiple parts or multiple sub-questions in them.

- Do not ask generic questions that can be asked of many other images. Ask questions specific to each image.

Please ask a question about this image that a human can answer *if* looking at the image (and not otherwise), but would stump this smart robot:

Q1: Write your question here to stump this smart robot.

29

# >0.25 million images

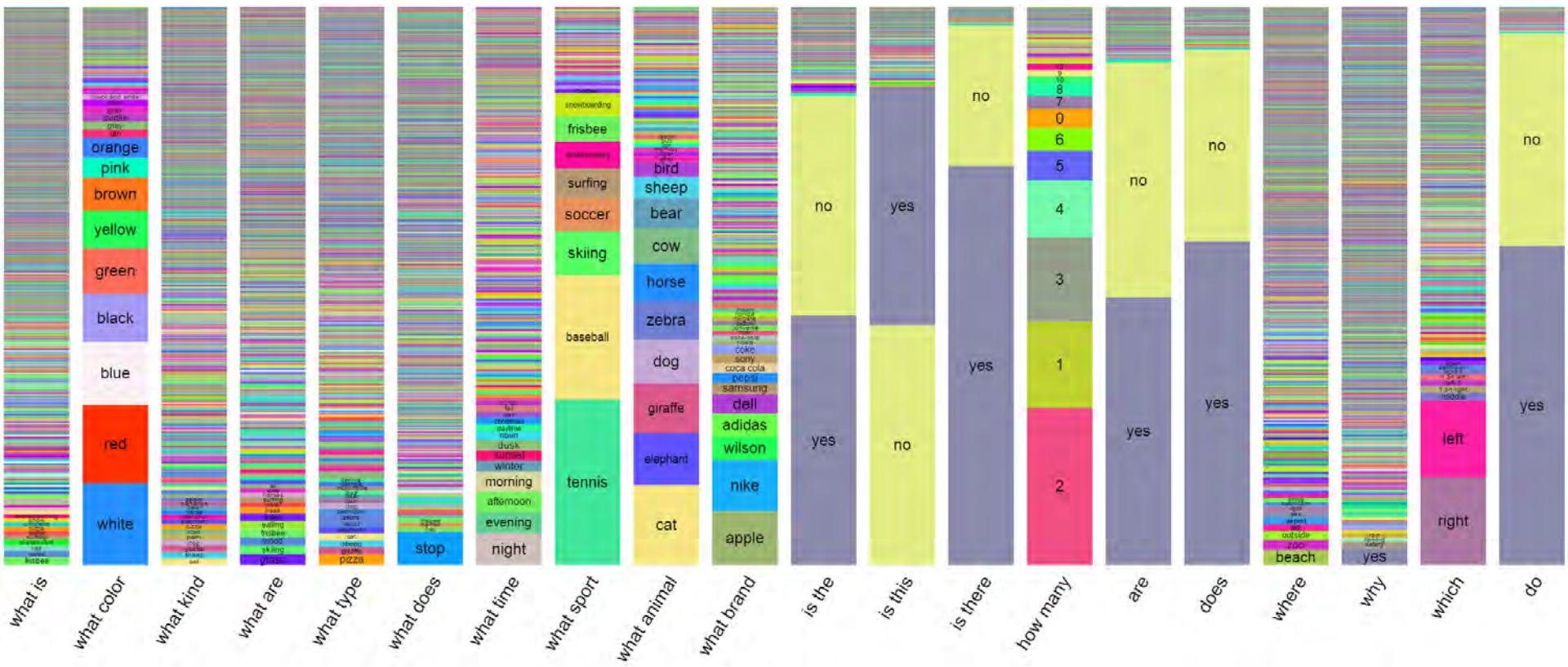# >0.76 million questions

# ~10 million answers

# >20 person-job-years

# Taxing the Turkers

- *Beware also the lasting effects of doing too many --for hours after the fact you will not be able to look at any photo without automatically generating a mundane question for it.*

- *If I were in possession of state secrets they could be immediately tortured out of me with the threat of being shown images of: skateboards, trains, Indian food and [long string of expletives] giraffes.*

- *(Please...I will tell you everything...just no more giraffes...)*

# Answers

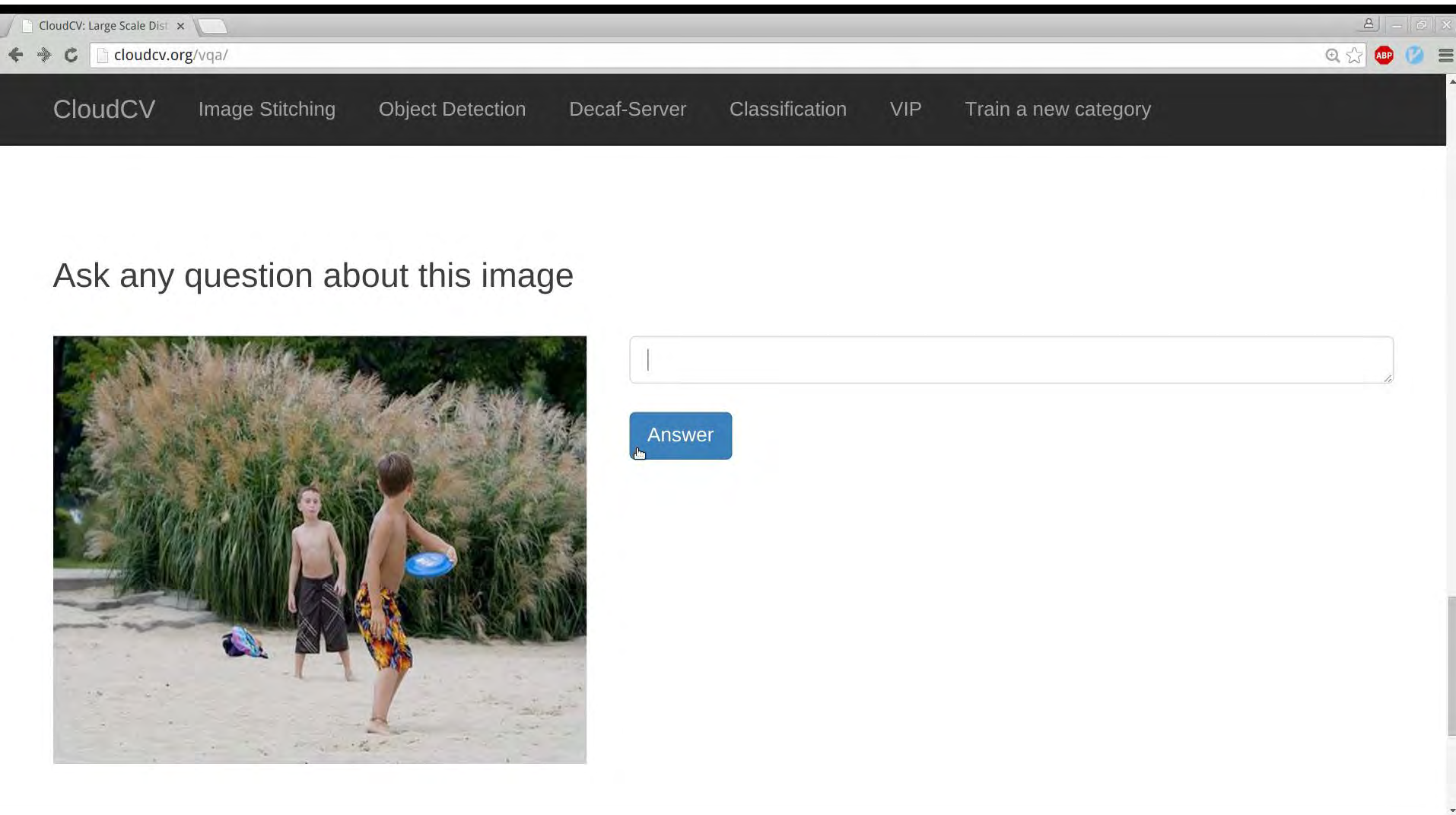# Human Accuracy, Inter-Human Agreement

Human agreement: 83%

First model (Summer 2015): 54%

State-of-the-art machine accuracy: 68%

# Visual Question Answering (VQA)



www.visualqa.org

# Papers using VQA

**Ask Me Anything: Free-form Visual Question Answering
Based on Knowledge from External Sources**

**Simple Baseline for Visual Question Answering**

## Academia, industry, start ups
## Dataset, Code

**ABC-CNN: An Attention Based Convolutional Neural Network for Visual Question Answering**

Kan Chen
University of Southern California
kanchen@usc.edu

Jiang Wang
Baidu Research - IDL
wangjiang03@baidu.com

Liang-Chieh Chen
UCLA
lcchen@cs.ucla.edu

Haoyuan Gao
Baidu Research - IDL
gaohaoyuan@baidu.com

Wei Xu
Baidu Research - IDL
wei.xu@baidu.com

Ram Nevatia
University of Southern California
nevatia@usc.edu

**Stacked Attention Networks for Image Question Answering**

Zichao Yang[1], Xiaodong He[2], Jianfeng Gao[2], Li Deng[2], Alex Smola[1]
[1]Carnegie Mellon University, [2]Microsoft Research, Redmond, WA 98052, USA
oy@cs.cmu.edu, {xiaohe, jfgao, deng}@microsoft.com, alex@smola.org

Slide credit: Devi Parikh

35

# VQA Challenge @ CVPR16

# VQA Challenge @ CVPR16

| | By Answer Type | | | Overall |
|---|---|---|---|---|
| | Yes/No | Number | Other | |
| UC Berkeley & Sony[14] | 83.24 | 39.47 | 58 | 66.47 |
| Naver Labs[10] | 83.31 | 38.7 | 54.62 | 64.79 |
| DLAIT[5] | 83.25 | 40.07 | 52.09 | 63.68 |
| snubi-naverlabs[25] | 83.16 | 39.14 | 51.33 | 63.18 |
| POSTECH[11] | 81.67 | 38.16 | 52.79 | 63.17 |
| Brandeis[3] | 82.11 | 37.73 | 51.91 | 62.88 |
| VTComputerVison[19] | 79.95 | 38.22 | 51.95 | 62.06 |
| MIL-UT[7] | 81.98 | 37.56 | 49.75 | 61.77 |
| klab[23] | 81.53 | 39.27 | 49.61 | 61.69 |
| SHB_1026[13] | 82.07 | 36.81 | 47.77 | 60.76 |
| MMCX[8] | 80.43 | 36.82 | 48.33 | 60.36 |
| VT_CV_Jiasen[20] | 80.56 | 38.14 | 47.87 | 60.33 |
| LV-NUS[6] | 81.34 | 35.67 | 46.1 | 59.54 |
| ACVT_Adelaide[1] | 81.07 | 37.12 | 45.83 | 59.44 |
| UC Berkeley (DNMN)[15] | 80.98 | 37.48 | 45.81 | 59.44 |
| CNNAtt[4] | 81.04 | 36.44 | 45.76 | 59.33 |
| san[24] | 79.11 | 36.41 | 46.42 | 58.85 |
| UC Berkeley (NMN)[16] | 81.16 | 37.7 | 44.01 | 58.66 |
| global_vision[22] | 78.24 | 36.27 | 46.32 | 58.43 |
| vqateam-deeperLSTM_NormlizeCNN[27] | 80.56 | 36.53 | 43.73 | 58.16 |
| Mujtaba hasan[9] | 80.28 | 36.92 | 42.24 | 57.36 |
| RIT[12] | 78.82 | 35.97 | 42.13 | 56.61 |
| Bolei[2] | 76.76 | 34.98 | 42.62 | 55.89 |
| UPV_UB[18] | 78.88 | 36.33 | 40.27 | 55.77 |
| att[21] | 78.1 | 35.3 | 40.27 | 55.34 |
| vqateam-lstm_cnn[28] | 79.01 | 35.55 | 36.8 | 54.06 |
| UPC[17] | 78.05 | 35.53 | 36.7 | 53.62 |
| vqateam-nearest_neighbor[29] | 71.73 | 24.31 | 22 | 42.73 |
| vqateam-prior_per_qtype[30] | 71.17 | 35.63 | 9.32 | 37.55 |
| vqateam-all_yes[26] | 70.53 | 0.43 | 1.26 | 29.72 |

~ 30 teams

# What such a model can't do



How many vegetarian slices are left in the pizza box?

# It can't count...



How many vegetarian slices are left in the pizza box?

# It doesn't have commonsense / knowledge...



How many vegetarian slices are left in the pizza box?

# It can't reason…



How many vegetarian slices are left in the pizza box?

# It doesn't leverage compositionality…



# How many vegetarian slices are left in the pizza box?

# It lacks consistency…



How many vegetarian slices are left in the pizza box?

# Visual Dialog

# VisDial Dataset

## Live Two-Person Chat on Amazon Mechanical Turk

Questioner ←——————————→ Answerer

# VisDial Dataset
## Live Two-Person Chat on Amazon Mechanical Turk

# VisDial Dataset

- ~100k images

- 1 dialog per image

- 2 people

- 10 rounds per person

# Qualitative Results

Slide credit: Dhruv Batra

# Applications

# Looking forward:
# Visual Dialog for Action

# Applications

## Natural language instructions to an agent

# Describer sees a reference image
# Actor does now

53

Actor: What is going on in the image overall?

Describer: Two people are playing

Actor: Is Mike to Jenny's left?

Describer: No, on the right

Actor: What is going on in the image overall?

Describer: Two people are playing

Actor: Is Mike to Jenny's left?

Describer: No, on the right
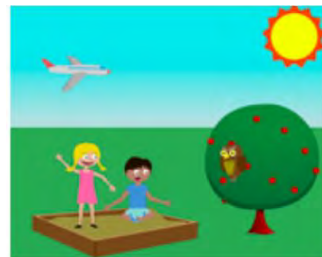Describer: There is an airplane

Actor: What is going on in the image overall?

Describer: Two people are playing

Actor: Is Mike to Jenny's left?

Describer: No, on the right
Describer: There is an airplane

Actor: What else is there?

Describer: Sun is on the top right corner
Describer: And a tree on the right side with fruits

Actor: Is the plane going towards the sun?

Describer: Yes

59

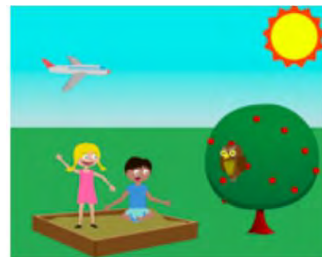Actor: What is going on in the image overall?

Describer: Two people

Actor: Is Mike to Jenny's left?

Describer: No, on the right
Describer: There is an airplane

Actor: What else is there?

Describer: Sun is on the top right corner
Describer: And a tree on the right side with fruits

Actor: Is the plane going towards the sun?

Describer: Yes
Describer: Jenny is standing
Describer: They are both standing in a sandbox

Actor: What expressions do they have?

Describer: Jenny is smiling while Mike is surprised

Actor: So there are just 6 objects in the image?

Describer: and there is an owl sitting on the tree
Describer: 7 objects
Describer: The plane is farther away from the sun. The whole tree is visible

# Common Sense

Man in blue wetsuit is surfing on wave
Karpathy and Fei-Fei (Stanford) 2015


A group of young people playing a game of Frisbee
Vinyals et al. (Google) 2015
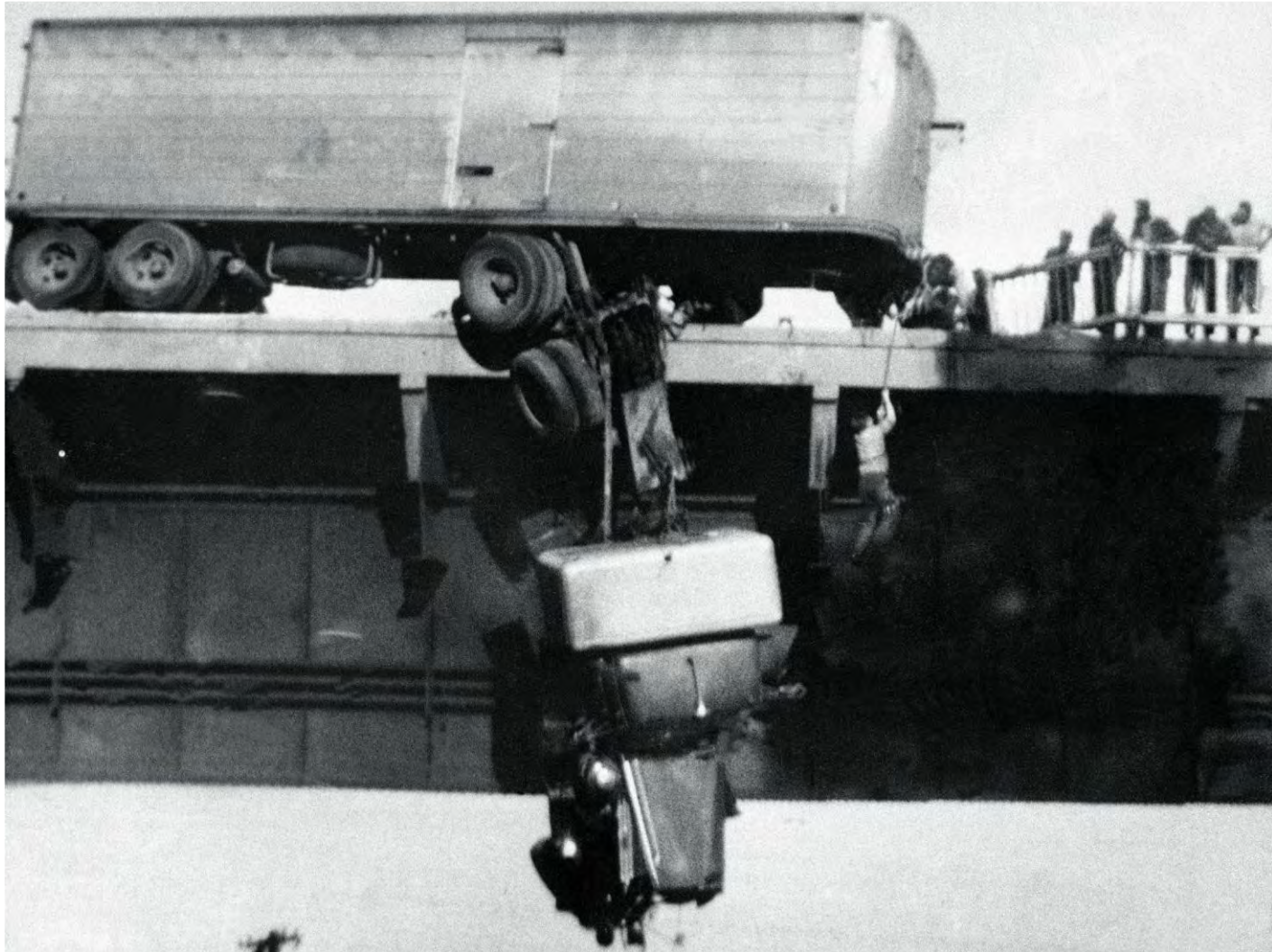

A car is parked in the middle of nowhere
Kiros et al. (University of Toronto) 2015


A pot of broccoli on a stove.
Fang et al. (Microsoft Research) 2015

Slide credit: Devi Parikh

# A man is rescued from his truck that is hanging dangerously from a bridge.

# A man is *rescued* from his truck that is hanging *dangerously* from a bridge.

# Learning Common Sense

- Text
  - Reporting bias

65

# Reporting bias in text

| Word | Teraword | Knext | Word | Teraword | Knext |
|---|---|---|---|---|---|
| spoke | 11,577,917 | 244,458 | hugged | 610,040 | 10,378 |
| laughed | 3,904,519 | 169,347 | blinked | 390,692 | 20,624 |
| murdered | 2,843,529 | 11,284 | was late | 368,922 | 31,168 |
| inhaled | 984,613 | 4,412 | exhaled | 168,985 | 3,490 |
| breathed | 725,034 | 34,912 | was punctual | 5,045 | 511 |

[Gordon et al. 2013]

# Reporting bias in text

| Word | Teraword | Knext | Word | Teraword | Knext |
|------|----------|-------|------|----------|-------|
| spoke | 11,577,917 | 244,458 | hugged | 610,040 | 10,378 |
| laughed | 3,904,519 | | | 390,692 | 20,624 |
| murdered | 2,843,529 | 11,284 | was late | 368,922 | 31,168 |
| inhaled | 984,613 | 4,412 | exhaled | 168,985 | 3,490 |
| breathed | 725,034 | 34,912 | was punctual | 5,045 | 511 |

inhale:exhale = 6:1

[Gordon et al. 2013]

# Reporting bias in text

| Word | Teraword | Knext | Word | Teraword | Knext |
|------|----------|-------|------|----------|-------|
| spoke | 11,577,917 | 244,458 | hugged | 610,040 | 10,378 |
| laughed | 3,904,519 | 169,347 | blinked | 390,692 | 20,624 |
| murdered | 2,843,529 | 11,284 | was late | 368,922 | 31,168 |
| inhaled | 984,613 | 4,412 | exhaled | 168,985 | 3,490 |
| breathed | 725,034 | 34,912 | was punctual | 5,045 | 511 |

murder:exhale = 17:1

[Gordon et al. 2013]

# Reporting bias in text

| Body Part | Teraword | Knext | Body Part | Teraword | Knext |
|-----------|----------|-------|-----------|----------|-------|
| Head | 18,907,427 | 1,332,154 | Liver | 246,937 | 10,474 |
| Eye(s) | 18,455,030 | 1,090,640 | Kidney(s) | 183,973 | 5,014 |
| Arm(s) | 6,345,039 | 458,018 | Spleen | 47,216 | 1,414 |
| Ear(s) | 3,543,711 | 230,367 | Pancreas | 24,230 | 1,140 |
| Brain | 3,277,326 | 260,863 | Gallbladder | 17,419 | 1,556 |

[Gordon et al. 2013]

# Reporting bias in text

| Body Part | Teraword | Knext | | Body Part | Teraword | Knext |
|---|---|---|---|---|---|---|
| Head | 18,907,427 | 1,332,154 | | Liver | 246,937 | 10,474 |
| Eye(s) | 18,455,030 | 1,090,640 | | Kidney(s) | 183,973 | 5,014 |
| Arm(s) | | | | | | 414 |
| Ear(s) | 3,543,711 | 230,367 | | Pancreas | 24,230 | 1,140 |
| Brain | 3,277,326 | 260,863 | | Gallbladder | 17,419 | 1,556 |

People have heads:gallbladders = 1085:1

[Gordon et al. 2013]
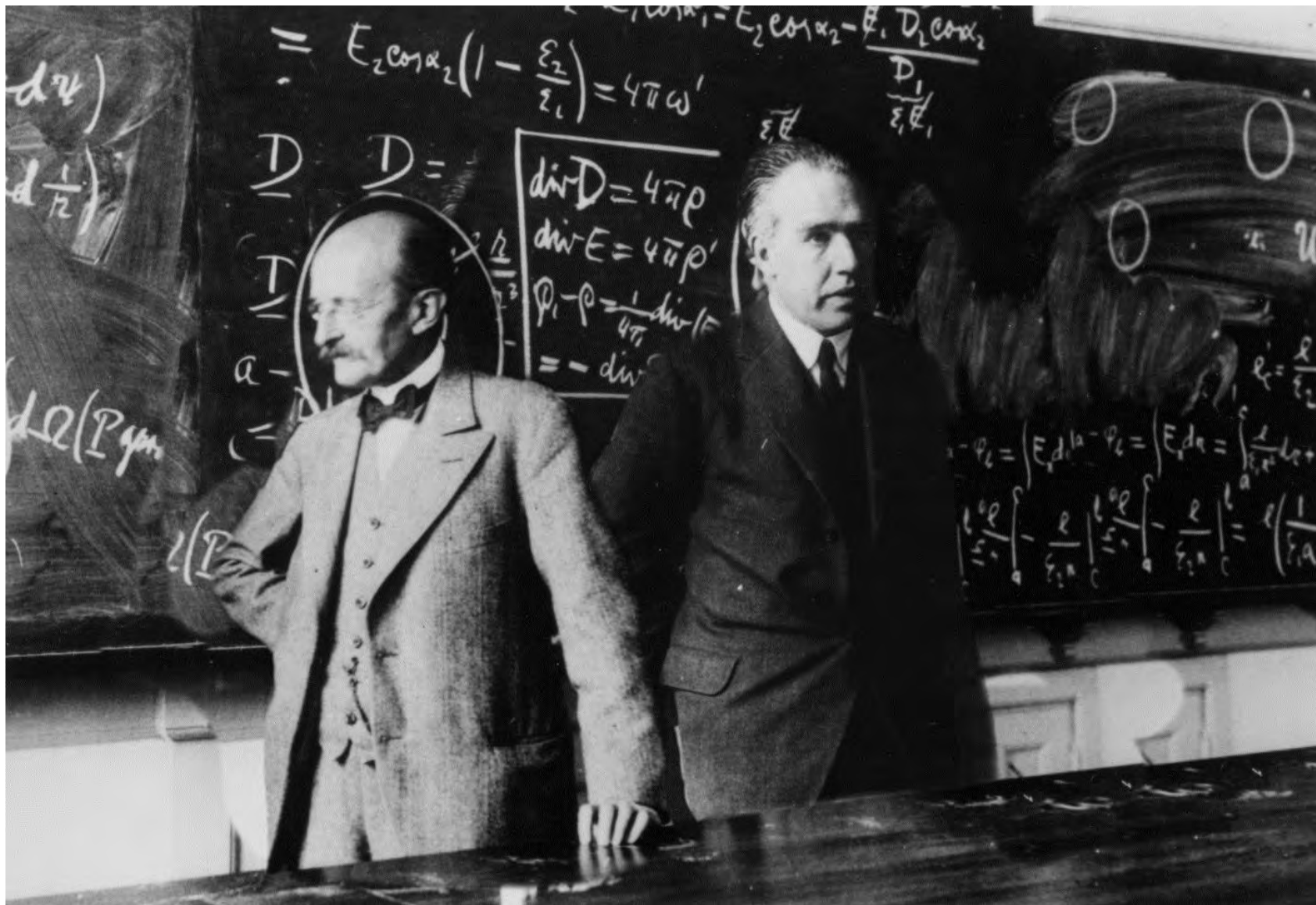
Slide credit: Devi Parikh

# Learning Common Sense

- Text
  - Reporting bias

- From structure in our visual world?

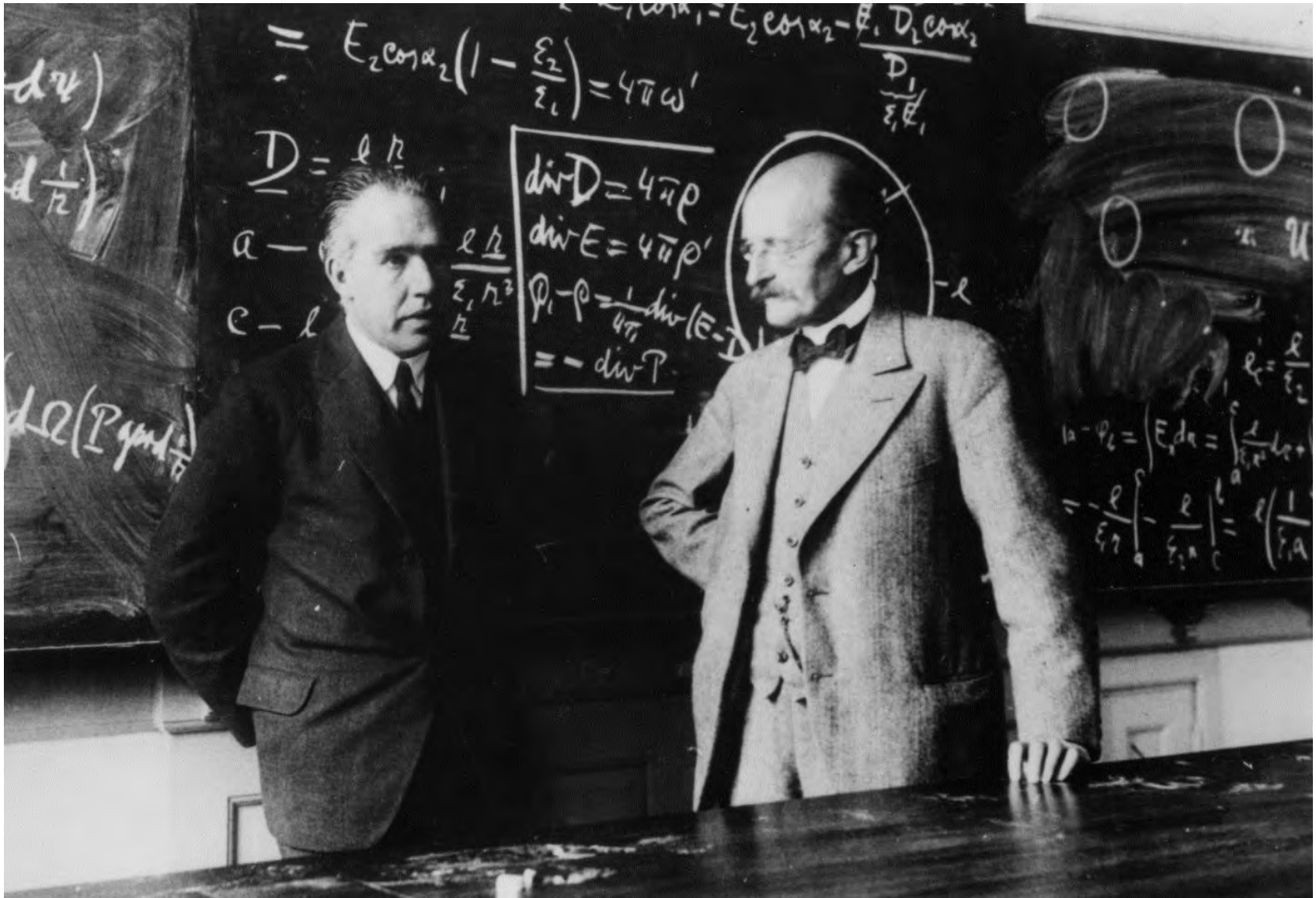# Two professors converse in front of a blackboard.

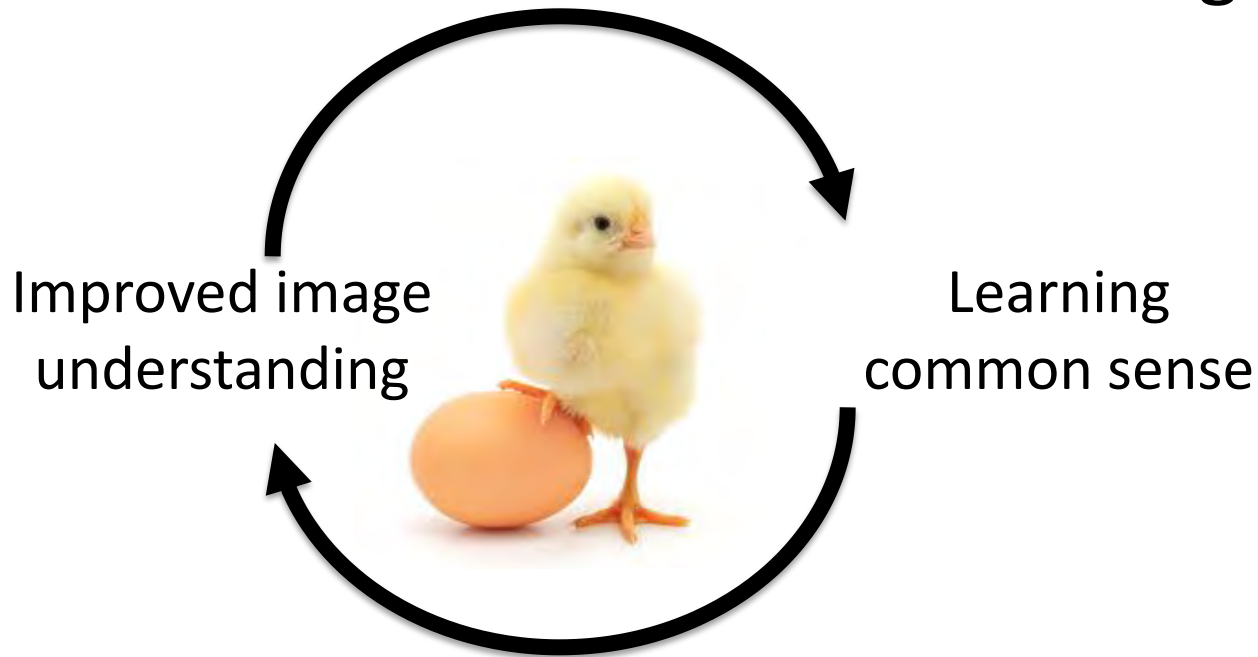# Two professors stand in front of a blackboard.

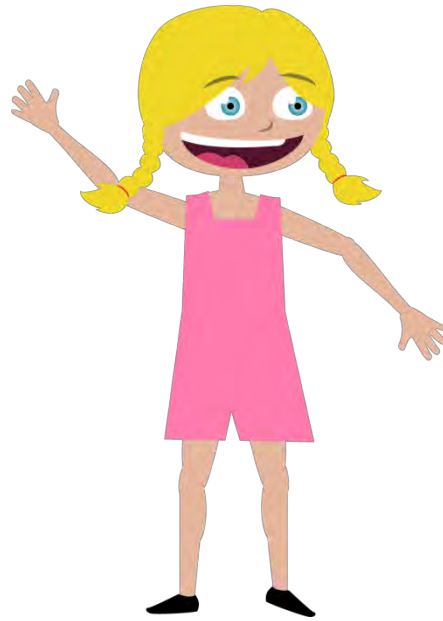# Two professors converse in front of a blackboard.

# Challenges

- Lacking visual density
- Annotations are expensive
- Computer vision doesn't work well enough

Improved image understanding

Learning common sense

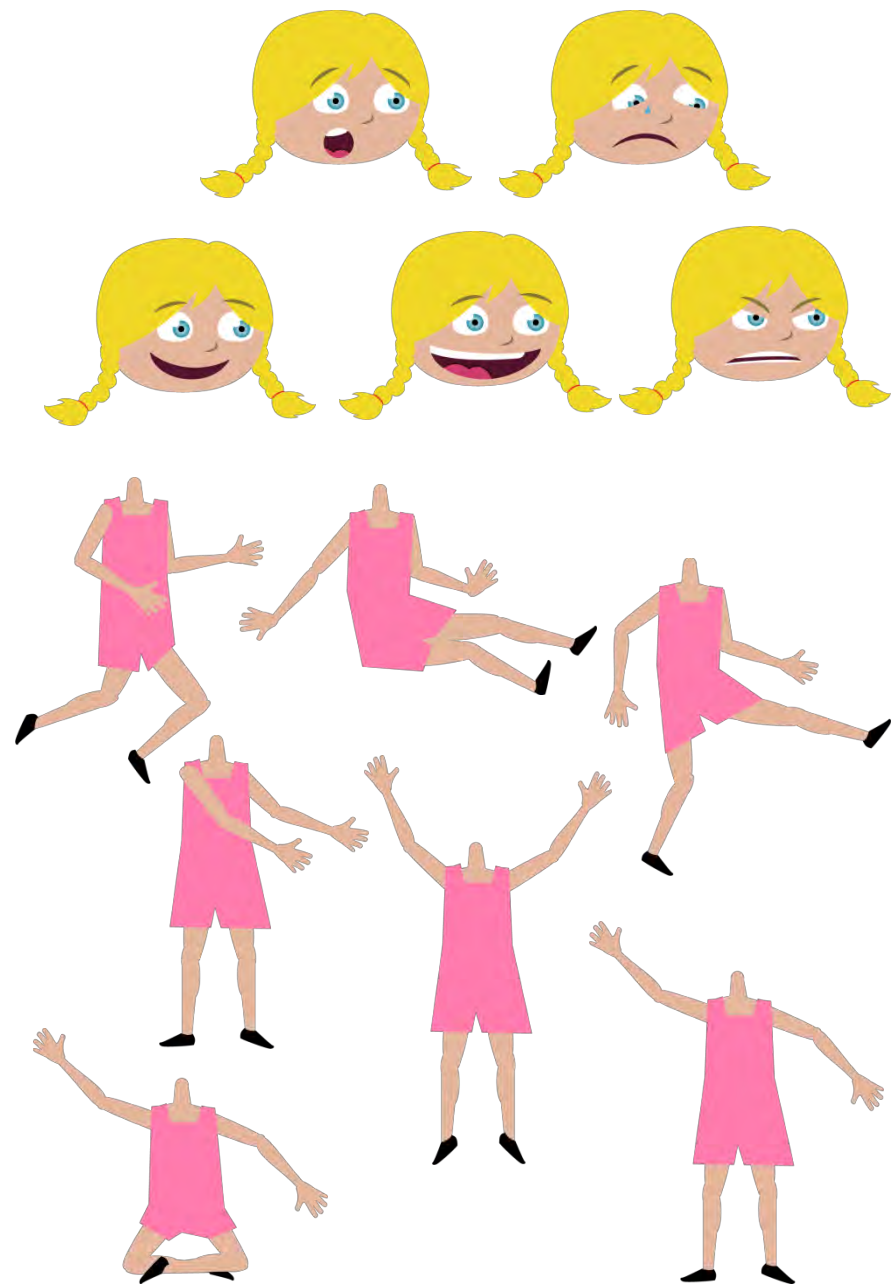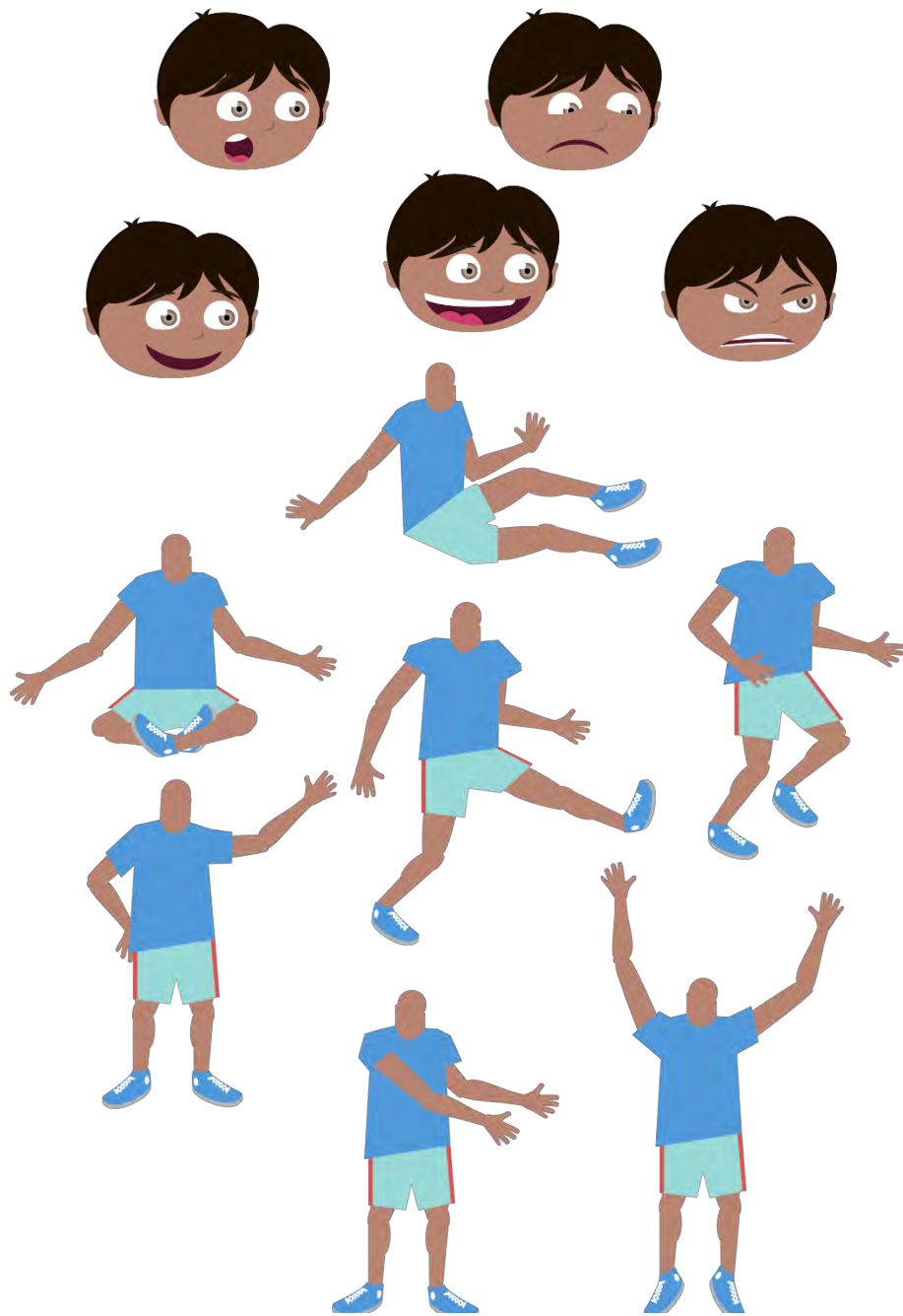# Is photorealism necessary?

Jenny          Mike

# Mike fights off a bear by giving him a hotdog while Jenny runs away.

# Commonsense Tasks

- Text-based tasks

Fill-in-the-blank:

Mike is having lunch when he sees a bear.

_____.

A. Mike orders a pizza.
B. Mike hugs the bear.
C. Bears are mammals.
D. Mike tries to hide.

# Key idea

- Imagine the scene behind the text
- Reason about the visual interpretation of the text, not just the text alone

Slide credit: Devi Parikh

Visual Paraphrasing:
Are these two descriptions describing the same scene?

1. Jenny was going to throw her pie at Mike.

2. Jenny is very angry.
   Jenny is holding a pie.

[Lin and Parikh, CVPR 2015]

84

# Approach: Imagination

_____.
Mike is wearing a blue cap.
Mike is telling Jenny to get off the swing.

A. There is a tree near a table.

B. The brown dog is standing next to Mike.

C. The sun is in the sky.

D. Jenny is standing dangerously on the swing.

# Approach: Imagination

There is a tree near a table.
Mike is wearing a blue cap.
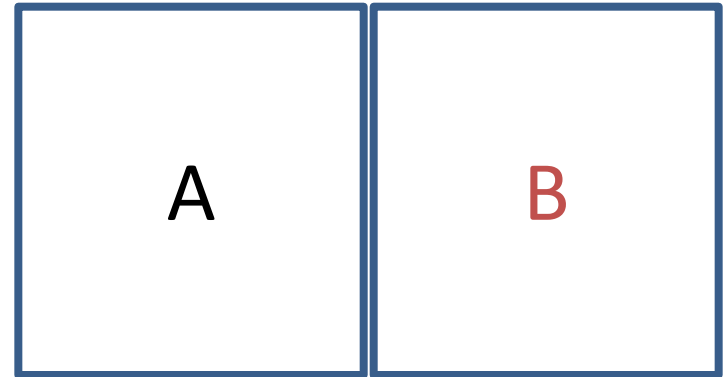Mike is telling Jenny to get off the swing.

A

A. There is a tree near a table.

B. The brown dog is standing next to Mike.

C. The sun is in the sky.

D. Jenny is standing dangerously on the swing.

# Approach: Imagination

The brown dog is standing next to Mike.
Mike is wearing a blue cap.
Mike is telling Jenny to get off the swing.

| A | B |
|---|---|

A. There is a tree near a table.

B. The brown dog is standing next to Mike.

C. The sun is in the sky.

D. Jenny is standing dangerously on the swing.

# Approach: Imagination

The sun is in the sky.
Mike is wearing a blue cap.
Mike is telling Jenny to get off the swing.

A. There is a tree near a table.

B. The brown dog is standing next to Mike.

C. The sun is in the sky.

D. Jenny is standing dangerously on the swing.

A

B

C

# Approach: Imagination

Jenny is standing dangerously on the swing.
Mike is wearing a blue cap.
Mike is telling Jenny to get off the swing.

A. There is a tree near a table.

B. The brown dog is standing next to Mike.

C. The sun is in the sky.

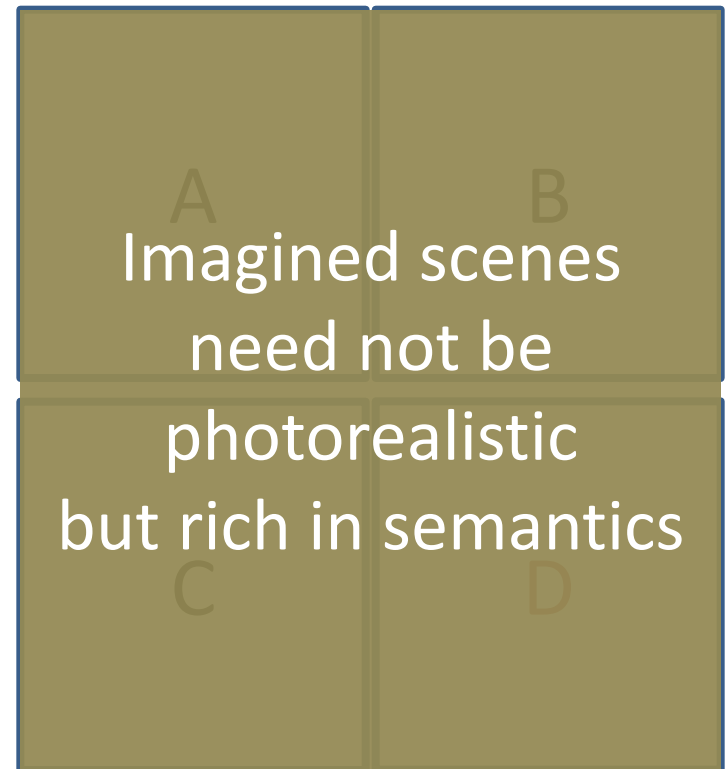D. Jenny is standing dangerously on the swing.

A    B

Imagined scenes
need not be
photorealistic
but rich in semantics

C    D

# Approach: Imagination

_____.
Mike is wearing a blue cap.
Mike is telling Jenny to get off the swing.

A. There is a tree near a table.

B. The brown dog is standing next to Mike.

C. The sun is in the sky.
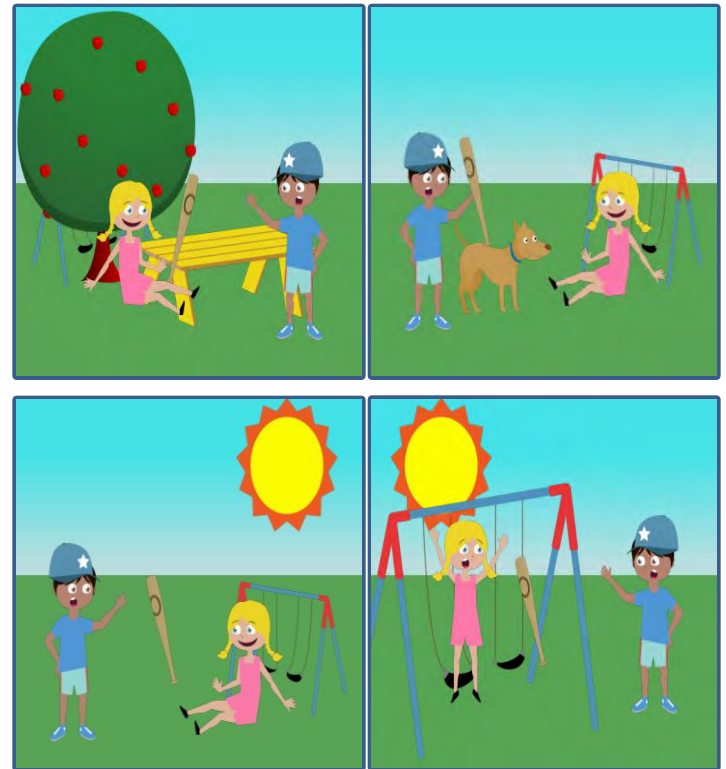
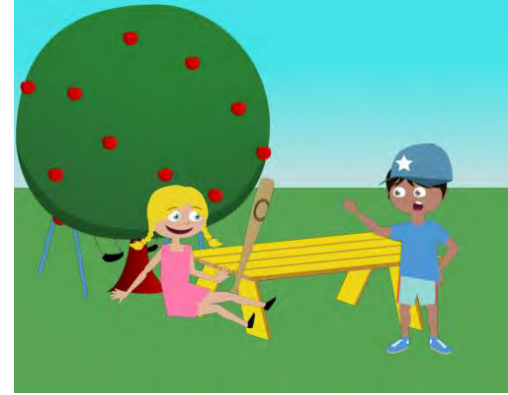D. Jenny is standing dangerously on the swing.

# Approach: Joint Text + Visual Reasoning



Jenny is standing dangerously on the swing. Mike is wearing a blue cap. Mike is telling Jenny to get off the

$\geq$

There is a tree near a table. Mike is wearing a blue cap. Mike is telling Jenny to get off the swing.

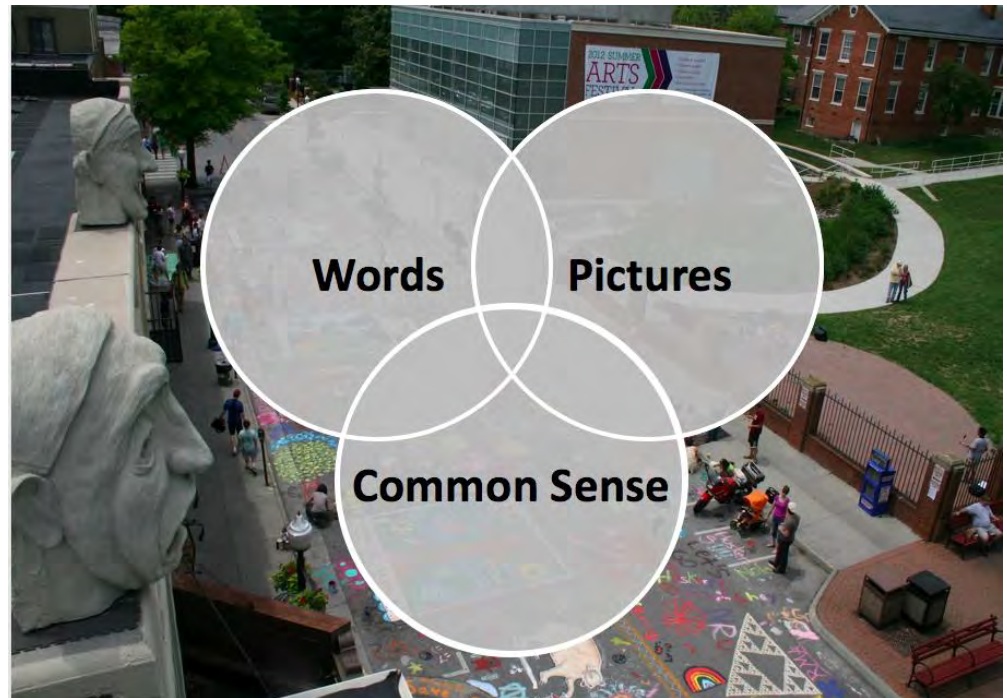# "This terrified woman's home is being invaded by mice as the cat sleeps."

# "The man is about to trip on his child's car and spill wine on his wife."

# Summary: Visual Dialog

- ## Applications
  - Today's chatbots are blind!

- ## AI
  - Vision
  - Language
  - Attention
  - Reasoning
  - External knowledge
  - Common Sense
  - Action, Manipulation

# Thank you.