



Grounded Causal Commonsense Reasoning

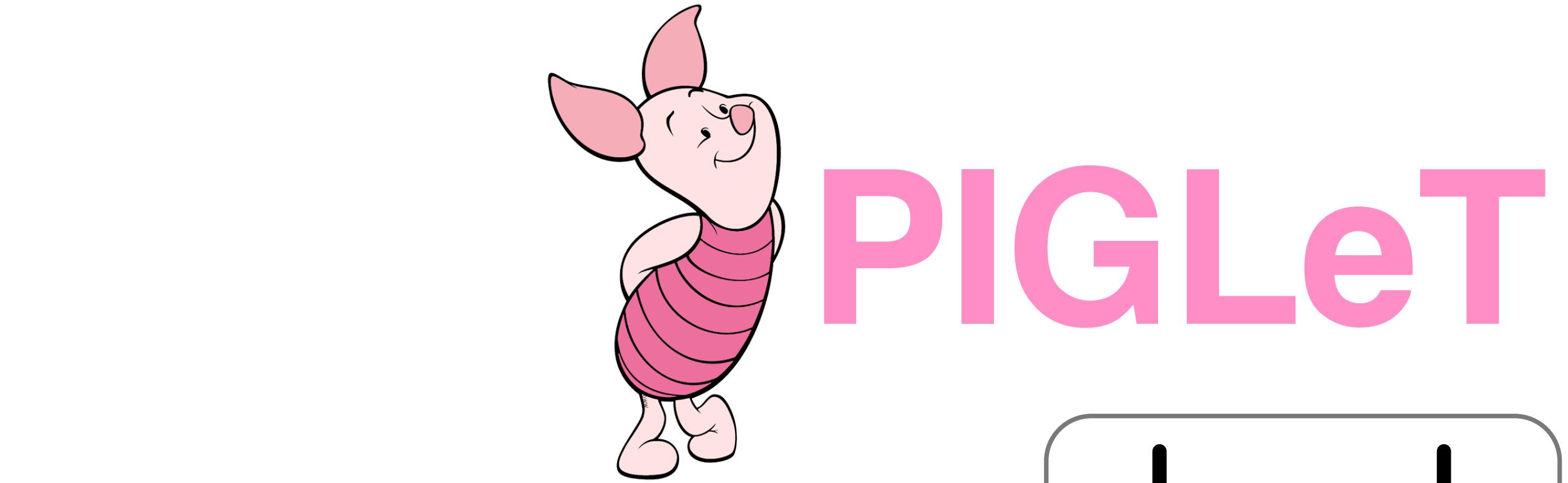


Yejin Choi
Paul G. Allen School of Computer Science & Engineering
University of Washington &
Allen Institute for Artificial Intelligence



Harnad's Symbol Grounding Problem

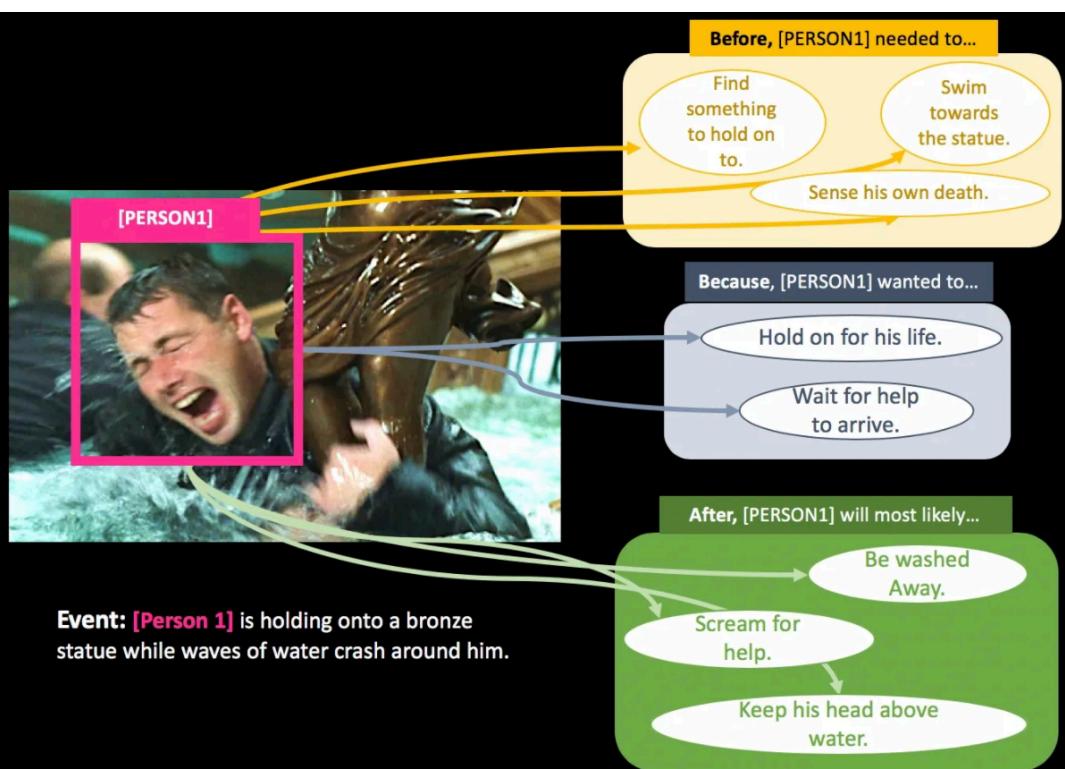
- **Grounding with 3D**
- **Grounding with 2D + Time**
- **Grounding with 2D + KG**



MERIOT



Visual Comet





PIGLeT

Language Grounding Through Neuro-Symbolic Interaction in a 3D World

ACL 2021

Rowan Zellers



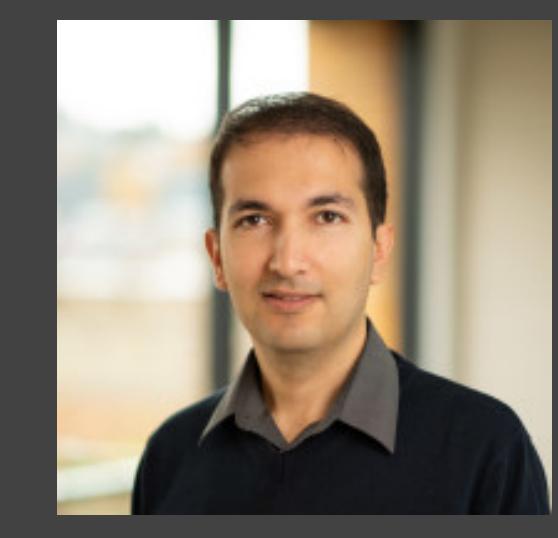
Ari
Holtzman



Matthew
Peters



Roozbeh
Mottaghi



Aniruddha
Kembhavi



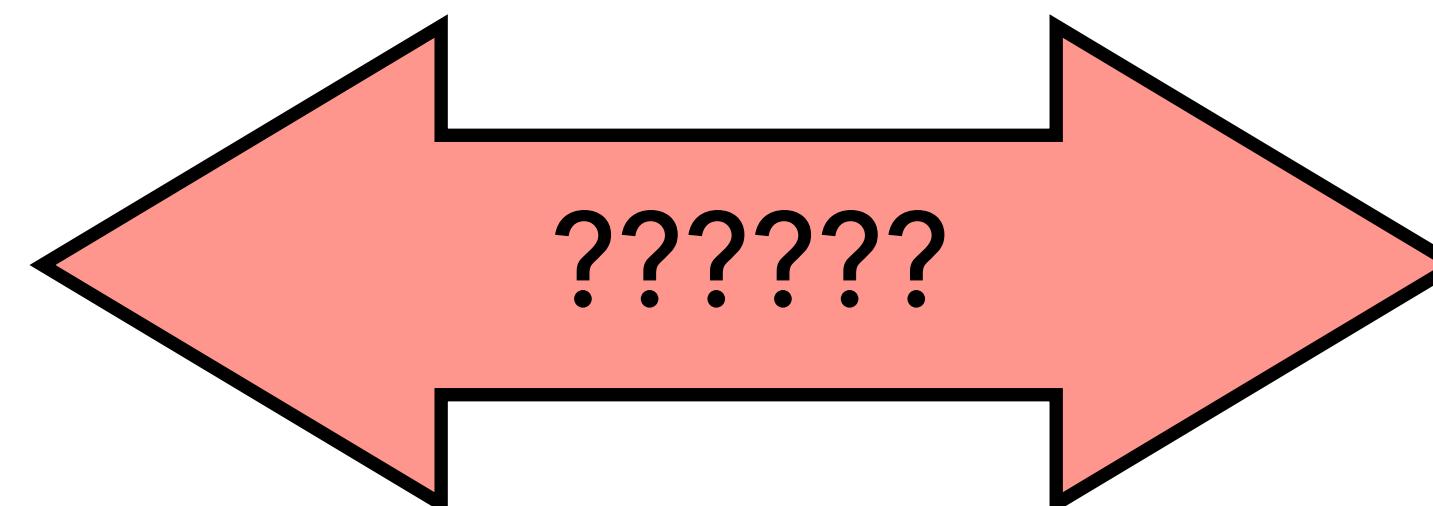
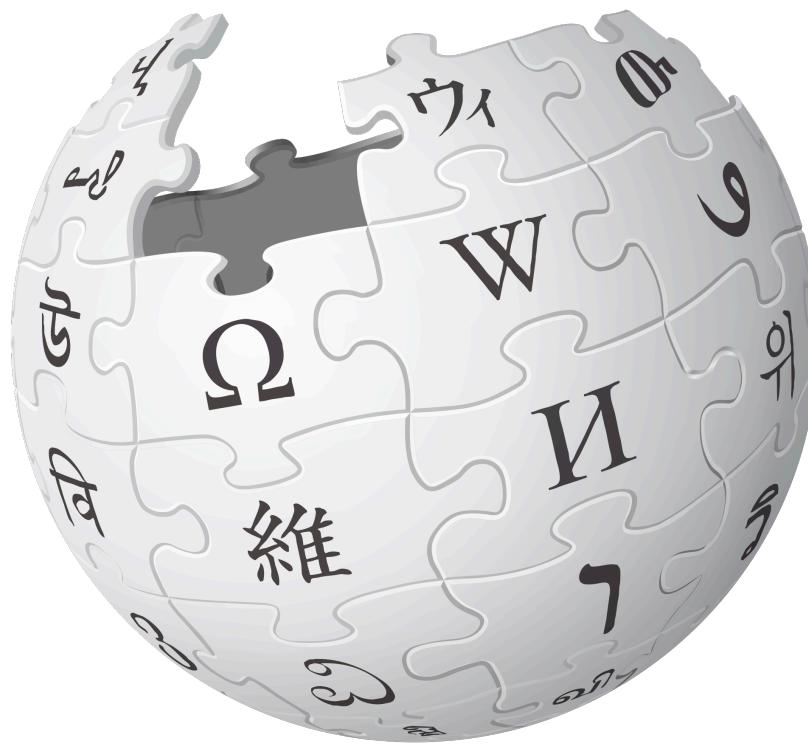
Ali
Farhadi



Me



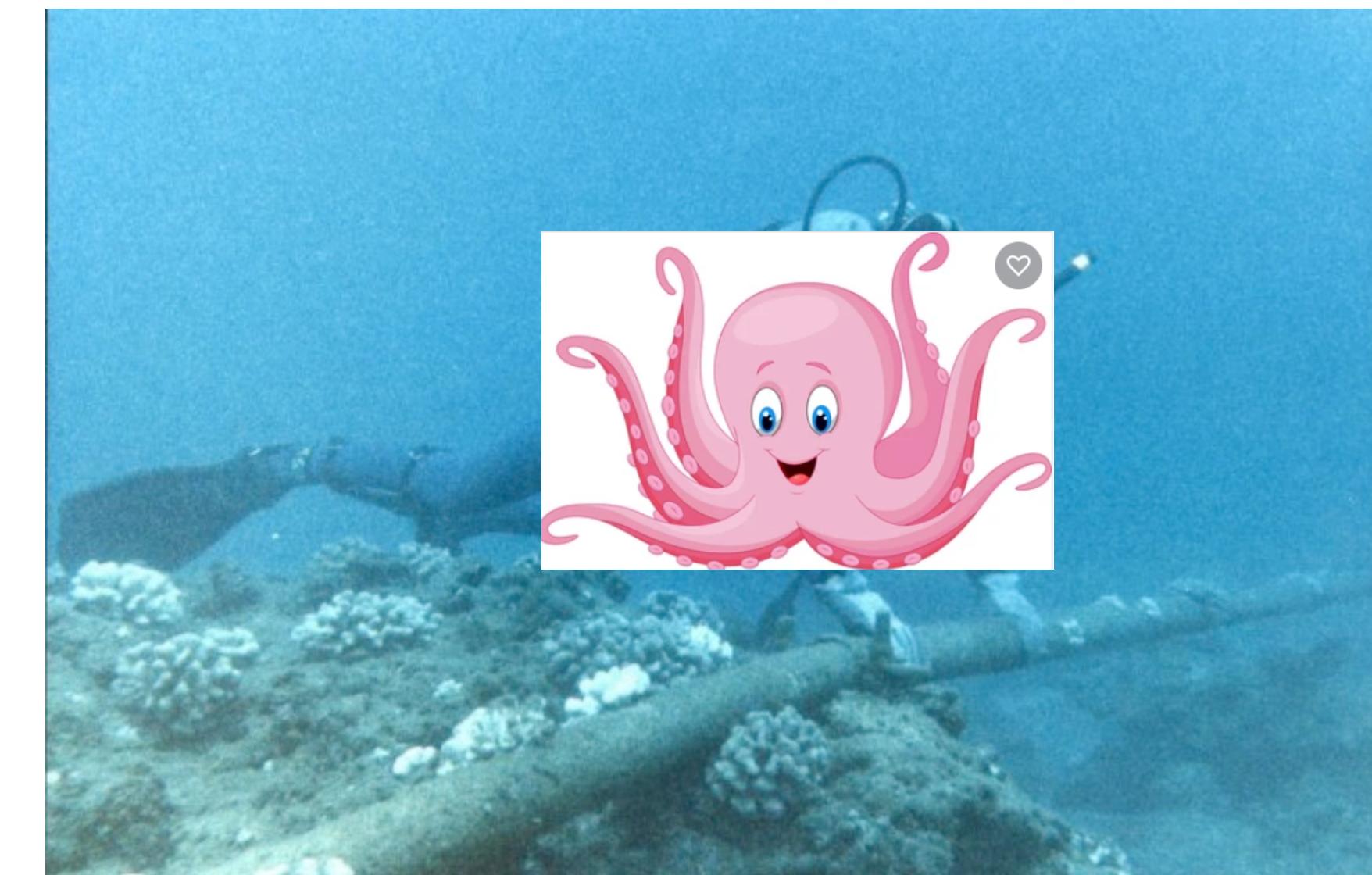
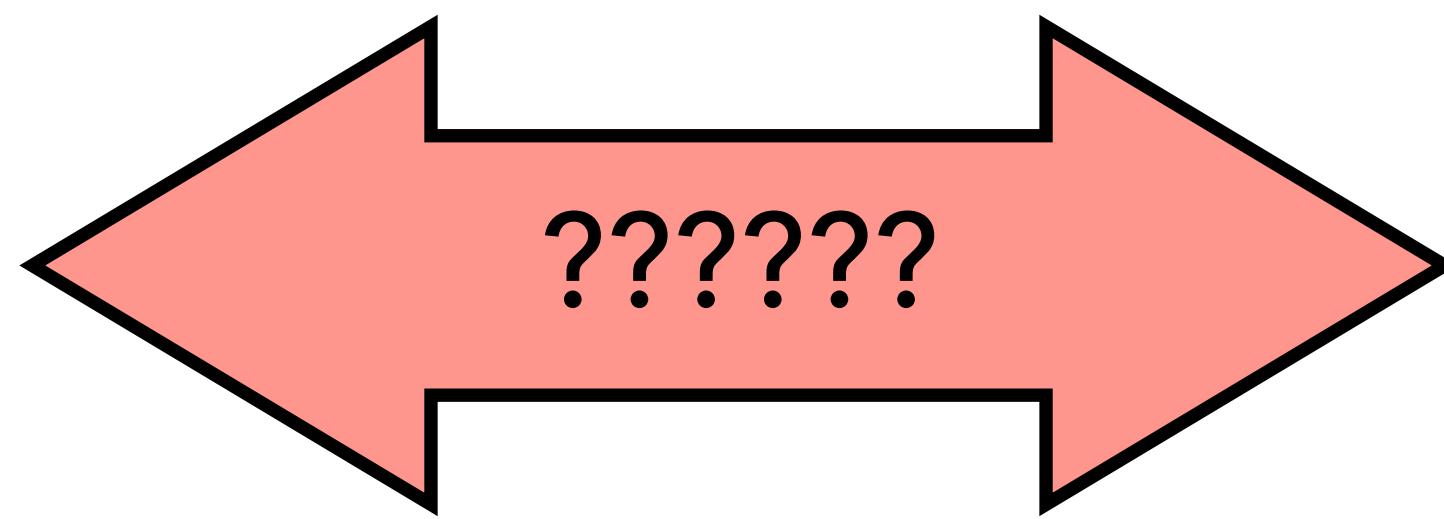
Problem: a gap between *language form* and *commonsense grounded meaning*



Written language
(*symbols*)

The world
(*continuous, subjective experience*)

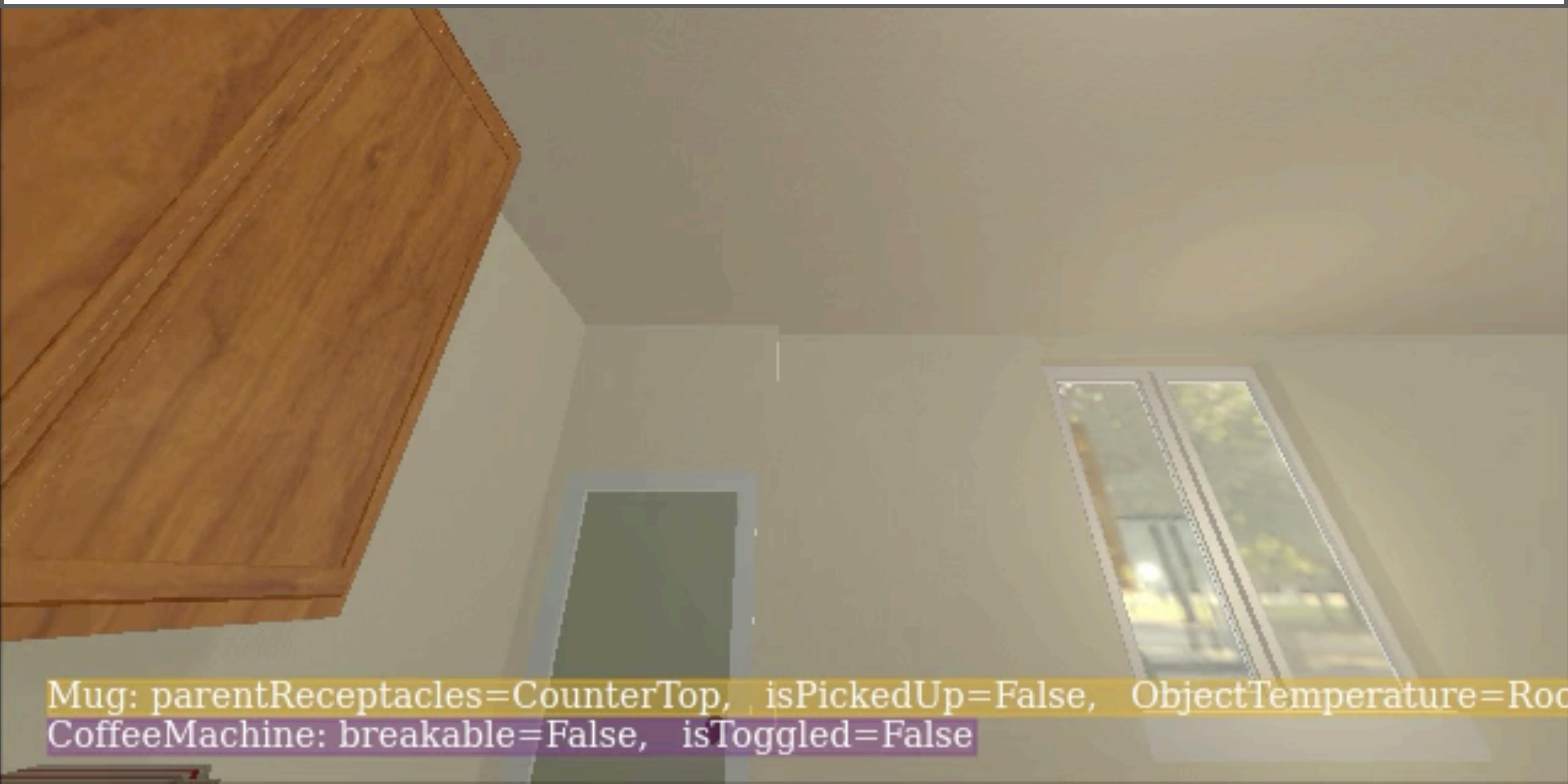
Problem: a gap between *language form* and *commonsense grounded meaning*



Harnad 1992, *inter alia*

Bender and Koller 2020,
inter alia

Proposal: ground language via a functional world representation, learned in simulation

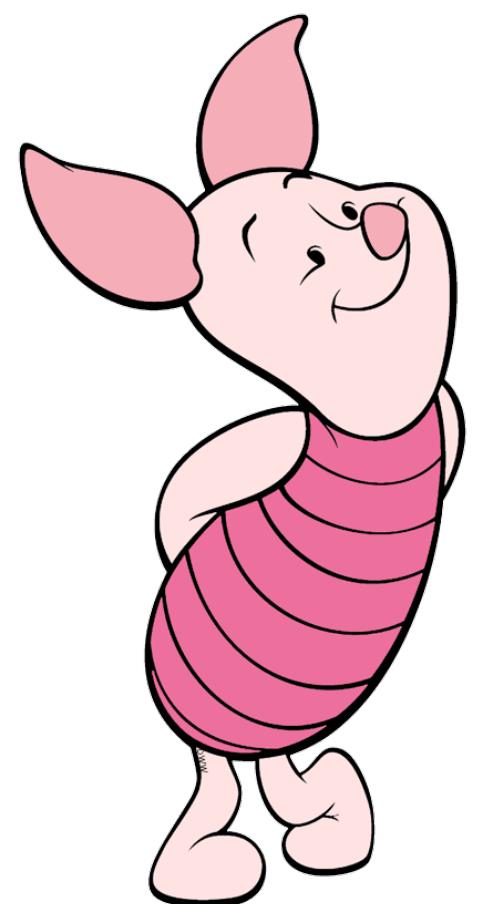


Mug: parentReceptacles=CounterTop, isPickedUp=False, ObjectTemperature=Rock
CoffeeMachine: breakable=False, isToggled=False

PIGLeT: Physical Interactions as Grounding for Language Transformers



Key idea: learn **TWO** model components for “how the world works” and “how to communicate it”



Learning “How the World Works”



Name: Egg

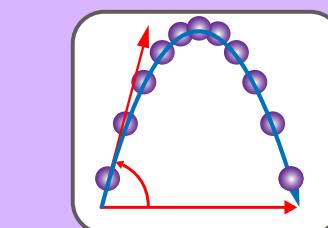
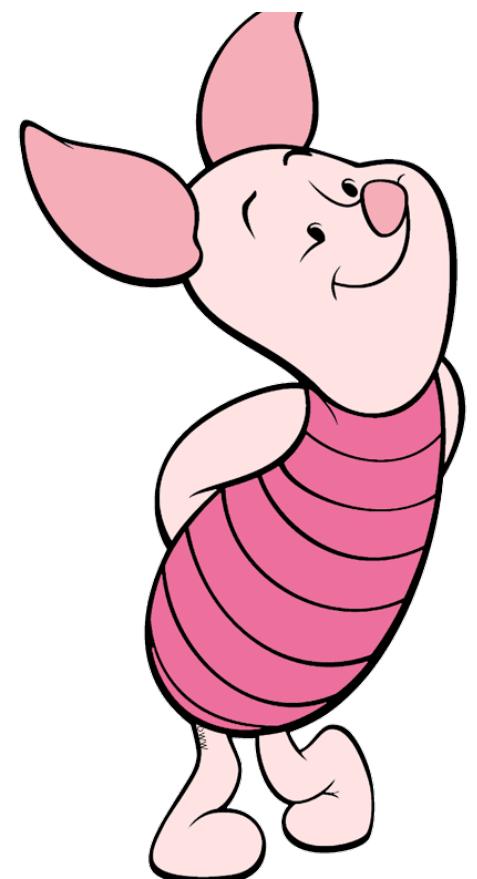
Temperature: RoomTemp

isCooked: False

isBroken: True

<heatUp, Pan>

...



Physical Dynamics Model



Language Model

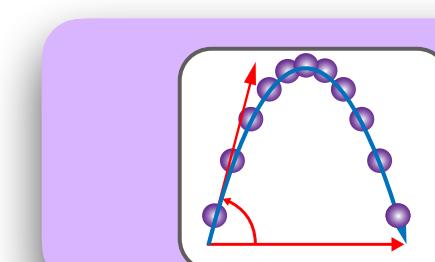
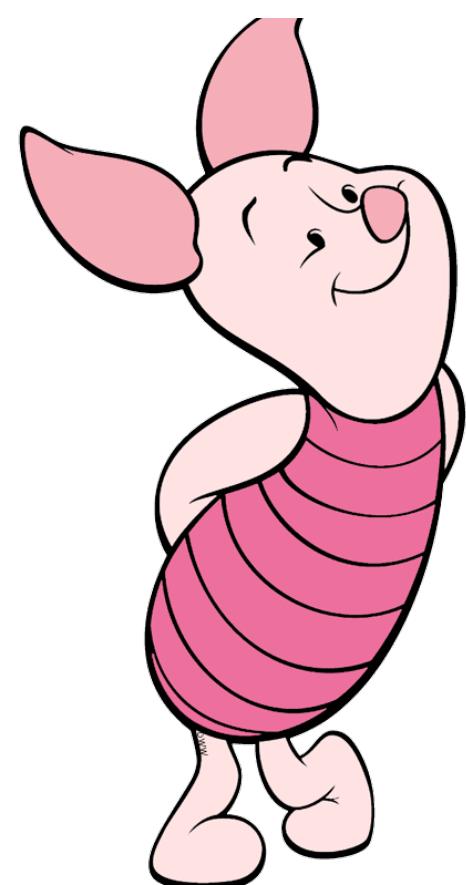
Learning “How the World Works”



Name:	Egg
Temperature:	RoomTemp
isCooked:	False
isBroken:	True



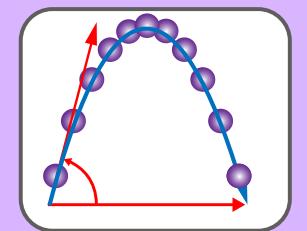
Name:	Egg
Temperature:	Hot
isCooked:	True
isBroken:	True



Physical Dynamics Model



Language Model



Physical Dynamics Model

Name:	Egg
Temperature:	RoomTem
isCooked:	False
isBroken:	True

...

Object Encoder

Action Apply

Object Decoder

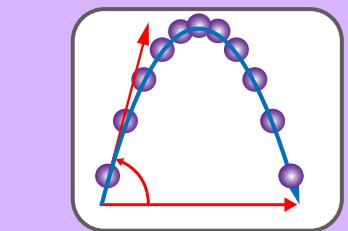
Name:	Egg
Temperature:	Hot
isCooked:	True
isBroken:	True

...

<heatUp, Pan>

Action Encoder

Physical Dynamics Model



Name:	Egg
Temperature:	RoomTem
isCooked:	False
isBroken:	True

...

<heatUp, Pan>

Object Encoder

Action Apply

Object Decoder

Name:	Egg
Temperature:	Hot
isCooked:	True
isBroken:	True

...

The robot heats up the pan.

Language Model

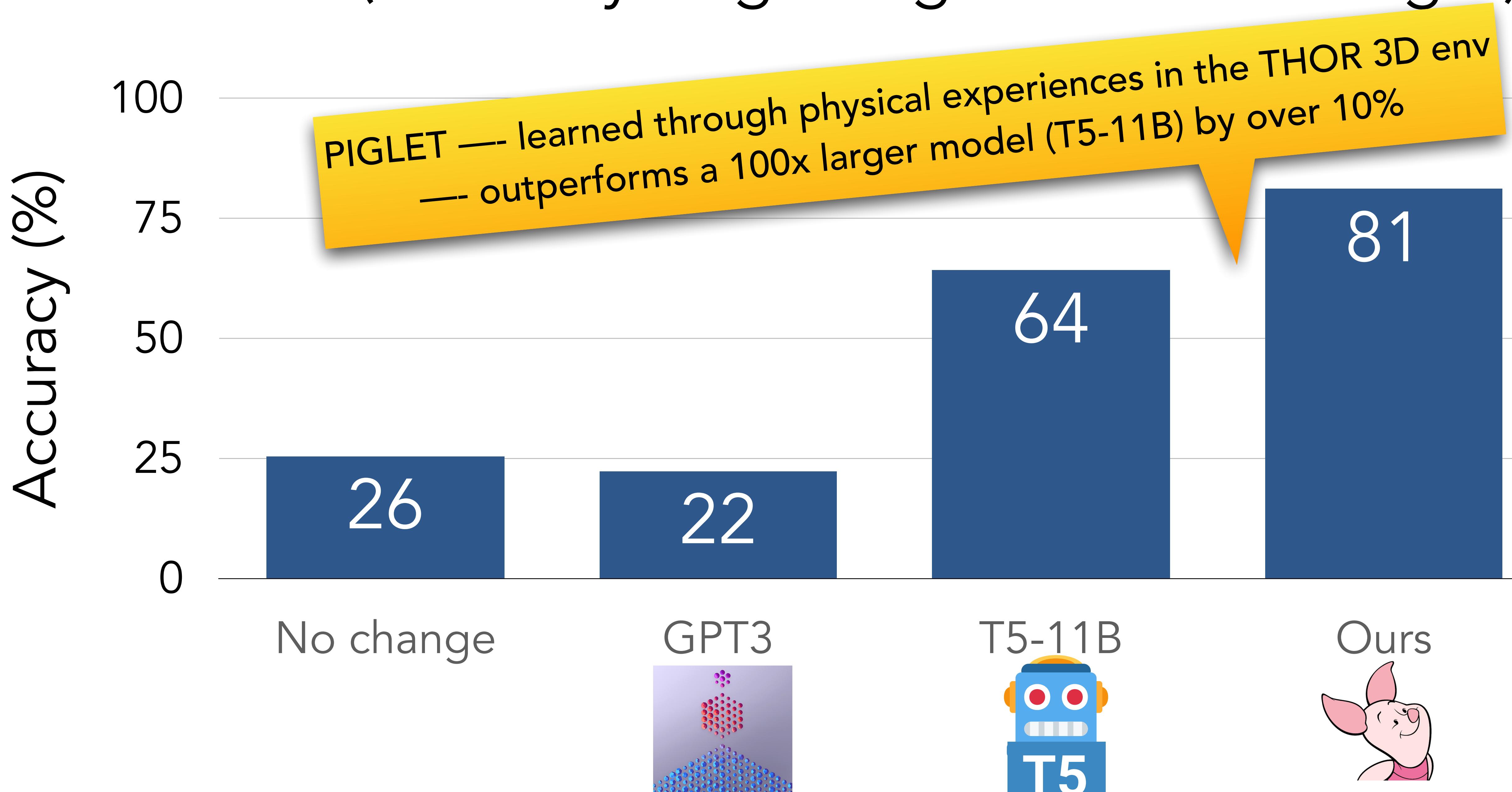


Language Model

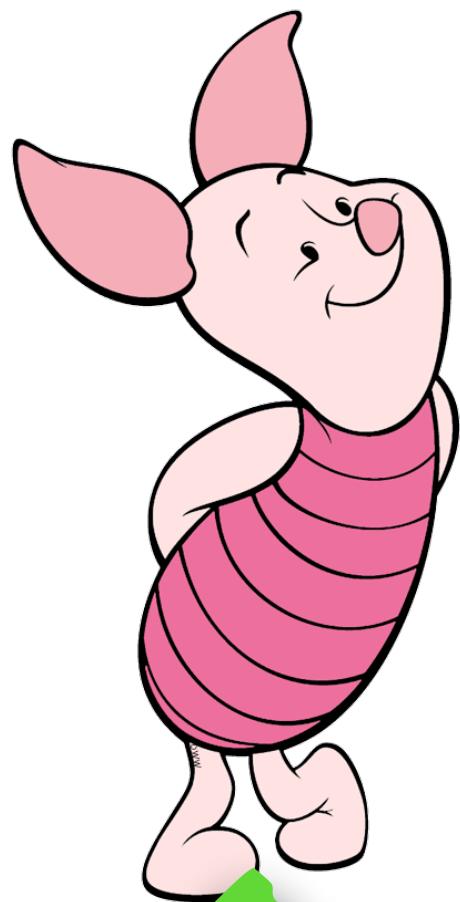


The pan becomes hot, and the egg gets cooked.

Results (accuracy of getting all attributes right)



Qualitative Example



Name:	Sink
filledWithLiquid	True

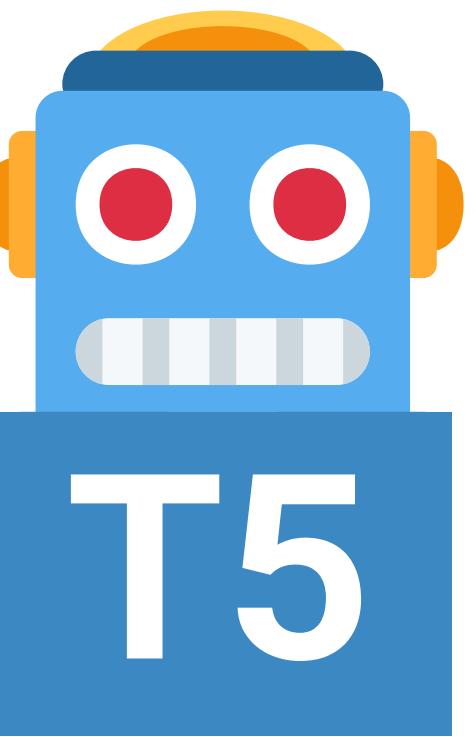
Name:	Mug
filledWithLiquid	True
isPickedUp	True

The robot
empties the
mug.

Name:	Sink
filledWithLiquid	True

Name:	Mug
filledWithLiquid	False
isPickedUp	True

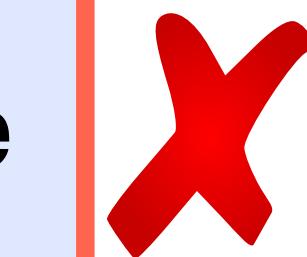
Qualitative Example



Name:	Sink
filledWithLiquid	True
Name:	Mug
filledWithLiquid	True
isPickedUp	True

The robot
empties the
mug.

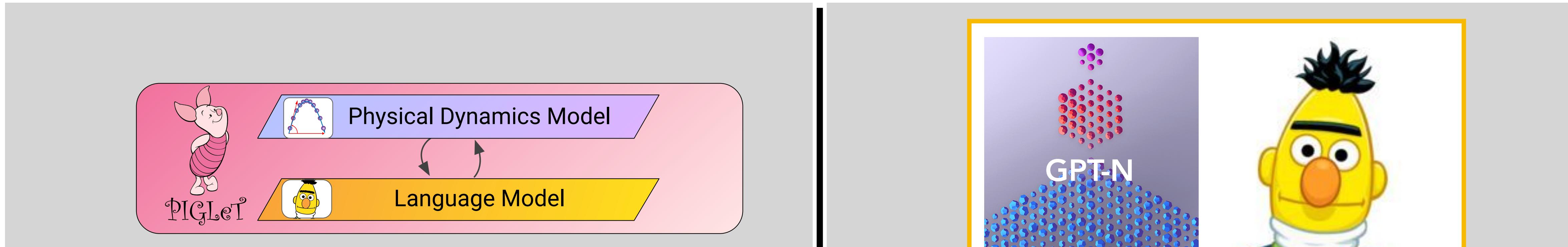
Name:	Sink
filledWithLiquid	False



T5, through text, learns “emptying liquid from an object” makes all objects in the room empty



PIGLeT: Physical Interactions as Grounding for Language Transformers



Learning physical commonsense through interactions
=> higher performance with 100x smaller models

Learn a lightweight factorized world model
for predicting *what might happen next*

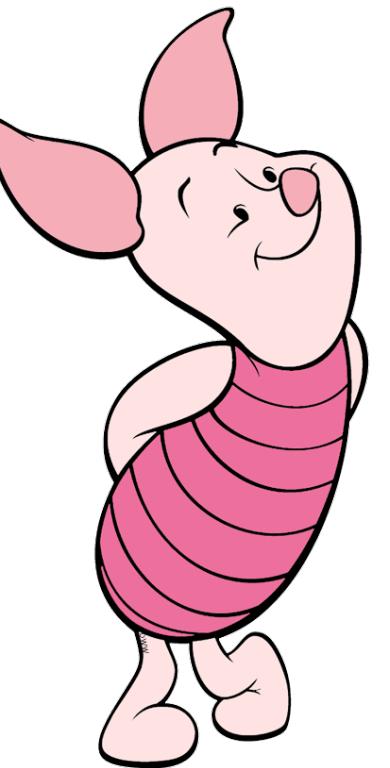
Can generalize to new concepts without words

A single, heavyweight, entangled model

Limited generalization to new concepts

Harnad's Symbol Grounding Problem

- Grounding with 3D



PIGLeT

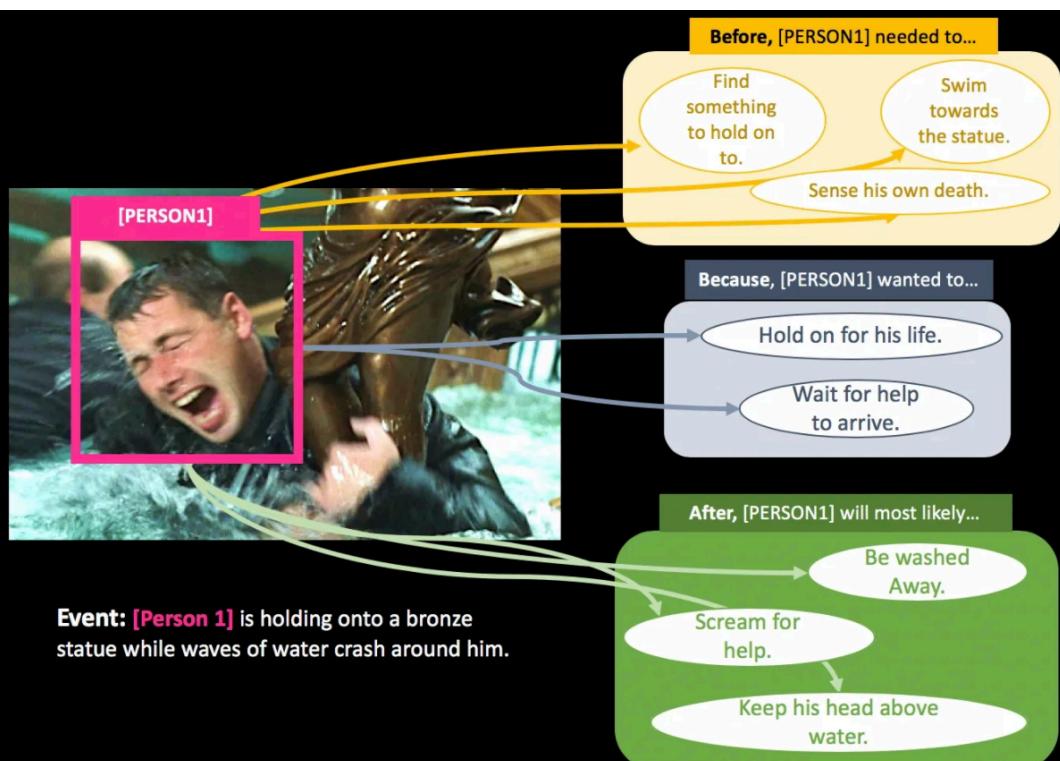
→ Grounding with 2D + Time

MERIOT



- Grounding with 2D + KG

Visual Comet





MERLOT: Multimodal Neural Script Knowledge Models

In Preparation



Rowan Zellers



A red wine glass emoji is positioned to the left of the text 'Ximing Lu'. The text is in a large, white, sans-serif font. 'Ximing' is on the top line, and 'Lu' is on the bottom line, centered vertically below 'Ximing'.



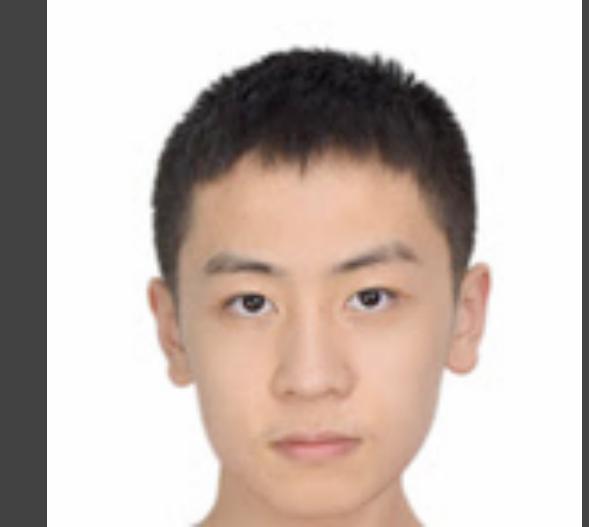
Jack Hessel



Youngjae Jae Sung Yu (James) Park



Jize Cao



Ali Farhadi



Me



Previously on VCR (cvpr 2019)

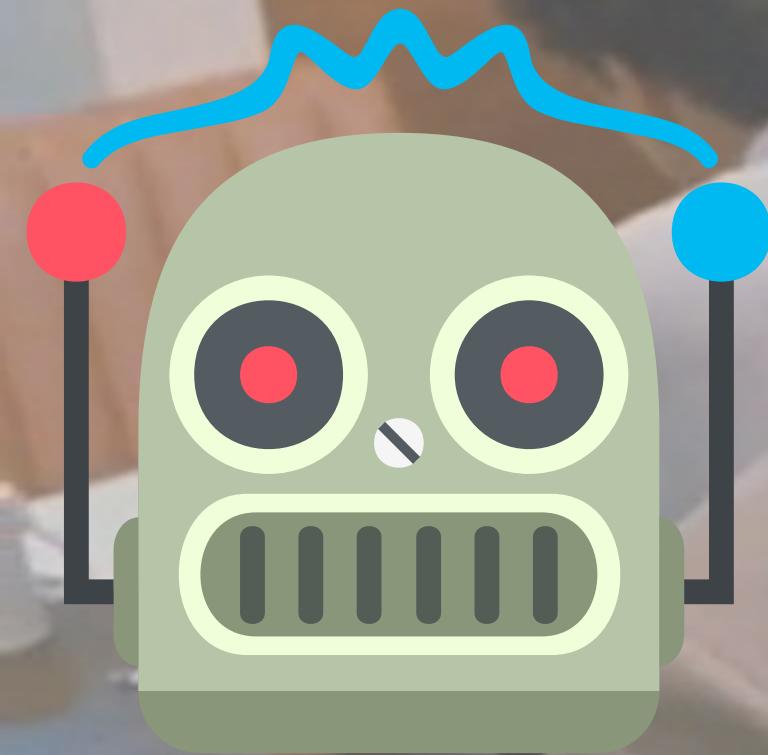


Why is he pointing?





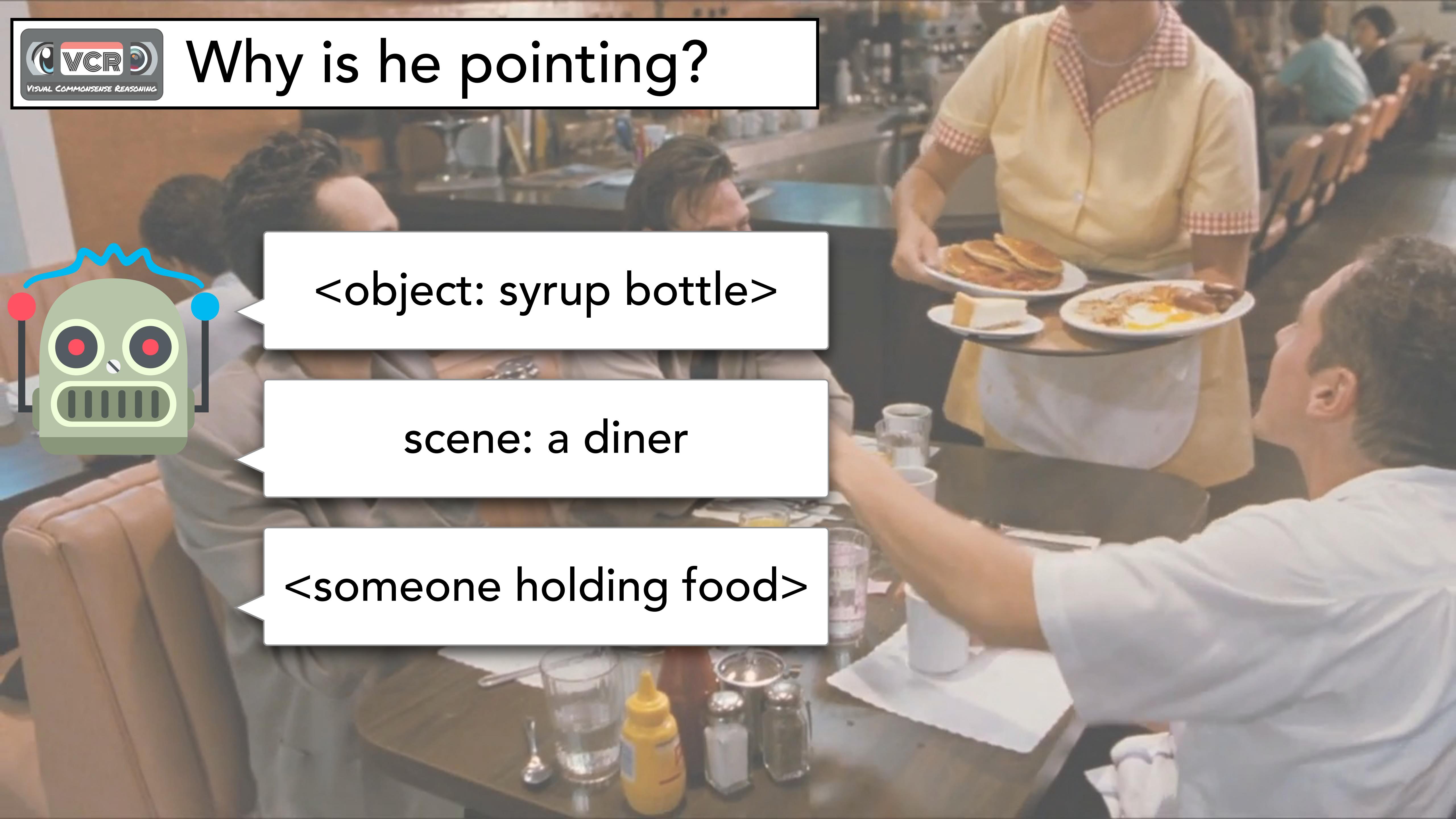
Why is he pointing?



<object: syrup bottle>

scene: a diner

<someone holding food>



Multimodal Script Knowledge



- Commonsense knowledge about events, including...
- What do people do at restaurants, and why?
- What might happen next in this event?

Script Knowledge

- (vanilla) script knowledge theory dates back to the early days of AI

SCRIPTS, PLANS, AND KNOWLEDGE

Roger C. Schank and Robert P. Abelson[†]

Yale University
New Haven, Connecticut USA

(1977)

"Of what a strange nature is knowledge! It clings to the mind, when it has once seized on it, like a lichen on the rock."

- Frankenstein's Monster
(M. Shelley, *Frankenstein or the Modern Prometheus*, 1818)

Abstract

We describe a theoretical system intended to facilitate the use of knowledge in an understanding system. The notion of script is introduced to

zation of knowledge can result in a real understanding system in the not too distant future. We expect that programs based on the theory we outline here and on our previous work on conceptual dependency and belief systems will combine with the MARGIE system (Schank et al., 1973a; Riesbeck, 1975; Rieger, 1975) to produce a working understander. We see understanding as the fitting of new information into a previously organized view of the world. We have therefore extended our work on language analysis (Schank, 1973a; Riesbeck 1975) to understanding - an understander, like an

Script Knowledge

SCRIPTS, PLANS, AND KNOWLEDGE

Roger C. Schank and Robert P. Abelson[†]

Yale University
New Haven, Connecticut USA

"Of what a strange nature is knowledge! It clings to the mind, when it has once seized on it, like a lichen on the rock."

— Frankenstein's Monster
(M. Shelley, *Frankenstein or the Modern Prometheus*, 1818)

Abstract

We describe a theoretical system intended to facilitate the use of knowledge in an understanding system. The notion of script is introduced to account for knowledge about mundane situations. A program, SAM, is capable of using scripts to understand. The notion of plans is introduced to account for general knowledge about novel situations.

I. Preface

In an attempt to provide theory where there have been mostly unrelated systems, Minsky (1974) recently described the work of Schank (1973a), Abelson (1973), Charniak (1972), and Norman (1972) as fitting into the notion of "frames." Minsky attempted to relate this work, in what is essentially language processing, to areas of vision research that conform to the same notion.

Minsky's frames paper has created quite a stir in AI and some immediate spinoff research along the lines of developing frames manipulators (e.g. Bobrow, 1975; Winograd, 1975). We find that we agree with much of what Minsky said about frames and with his characterization of our own work. The frames idea is so general, however, that it does not lend itself to applications without further specialization. This paper is an attempt to develop further the lines of thought set out in Schank (1975a) and Abelson (1973; 1975a). The ideas presented here can be viewed as a specialization of the frame idea. We shall refer to our central constructs as "scripts."

II. The Problem

Researchers in natural language understanding have felt for some time that the eventual limit on the solution of our problem will be our ability to characterize world knowledge. Various researchers have approached world knowledge in various ways. Winograd (1972) dealt with the problem by severely restricting the world. This approach had the positive effect of producing a working system and the negative effect of producing one that was only minimally extendable. Charniak (1972) approached the problem from the other end entirely and has made some interesting first steps, but because his work is not grounded in any representational system or any working computational system the restriction of world knowledge need not critically concern him.

Our feeling is that an effective characteri-

[†] The work of the second author was facilitated by National Science Foundation Grant GS-35768.

script: restaurant

roles: customer, waiter, chef, cashier

Scene 1: entering

PTRANS self into restaurant

ATTEND eyes to where empty tables are

MBUILD where to sit

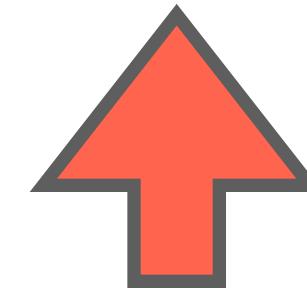
PTRANS self to table

MOVE sit down

Scene 2: ordering

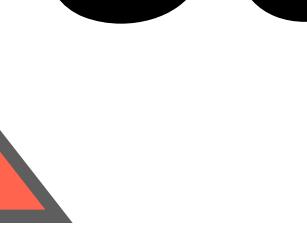
...

Multimodal Script Knowledge



(Neural)

Multimodal Script Knowledge



(Neural)

From 6M youtube videos, we'll learn:

From 6M youtube videos, we'll learn:



Recognition-level
Knowledge



person

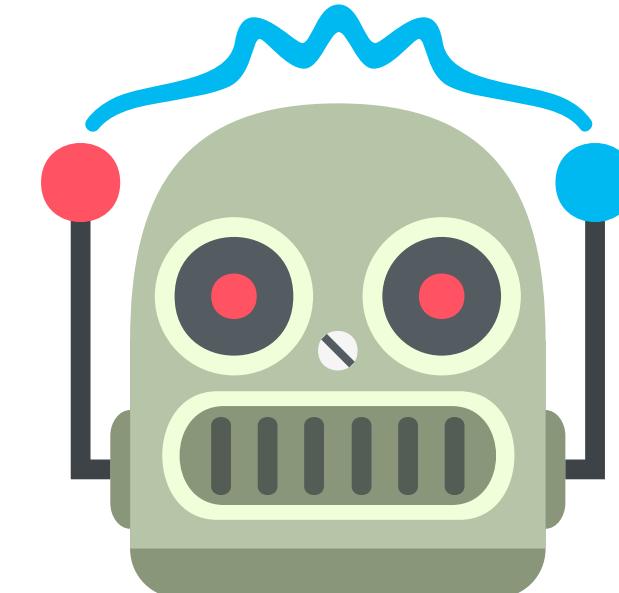


stopwatch
water pitcher



thermometer

Multimodal
Script Knowledge



This person might be
measuring how fast the
water boils

From 6M youtube videos, we'll learn:

**Recognition-level
Knowledge**

**Multimodal
Script Knowledge**



Multimodal Event Representation Learning Over Time

The result:

- Trained fully from scratch, we get...
- zero-shot temporal commonsense,
- Fine-tuned SOTA on 13 tasks

Setup: Videos and Transcripts



“I’m going to compare electric and induction stoves...”



“I’ll use a stopwatch to time how fast my electric stove boils water...”

Time



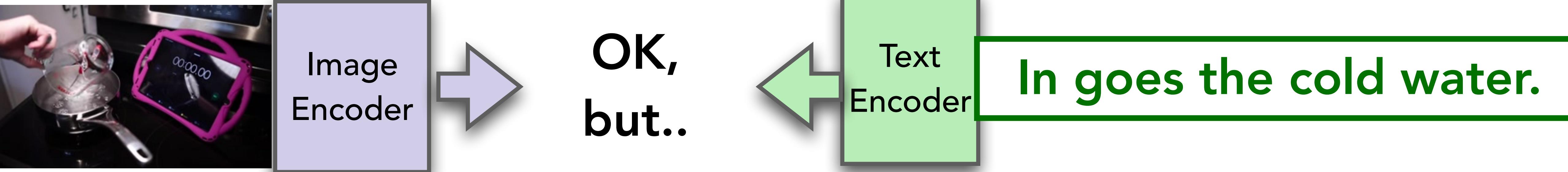
“In goes the cold water...”



“It took 4 and a half minutes to reach full boil...”

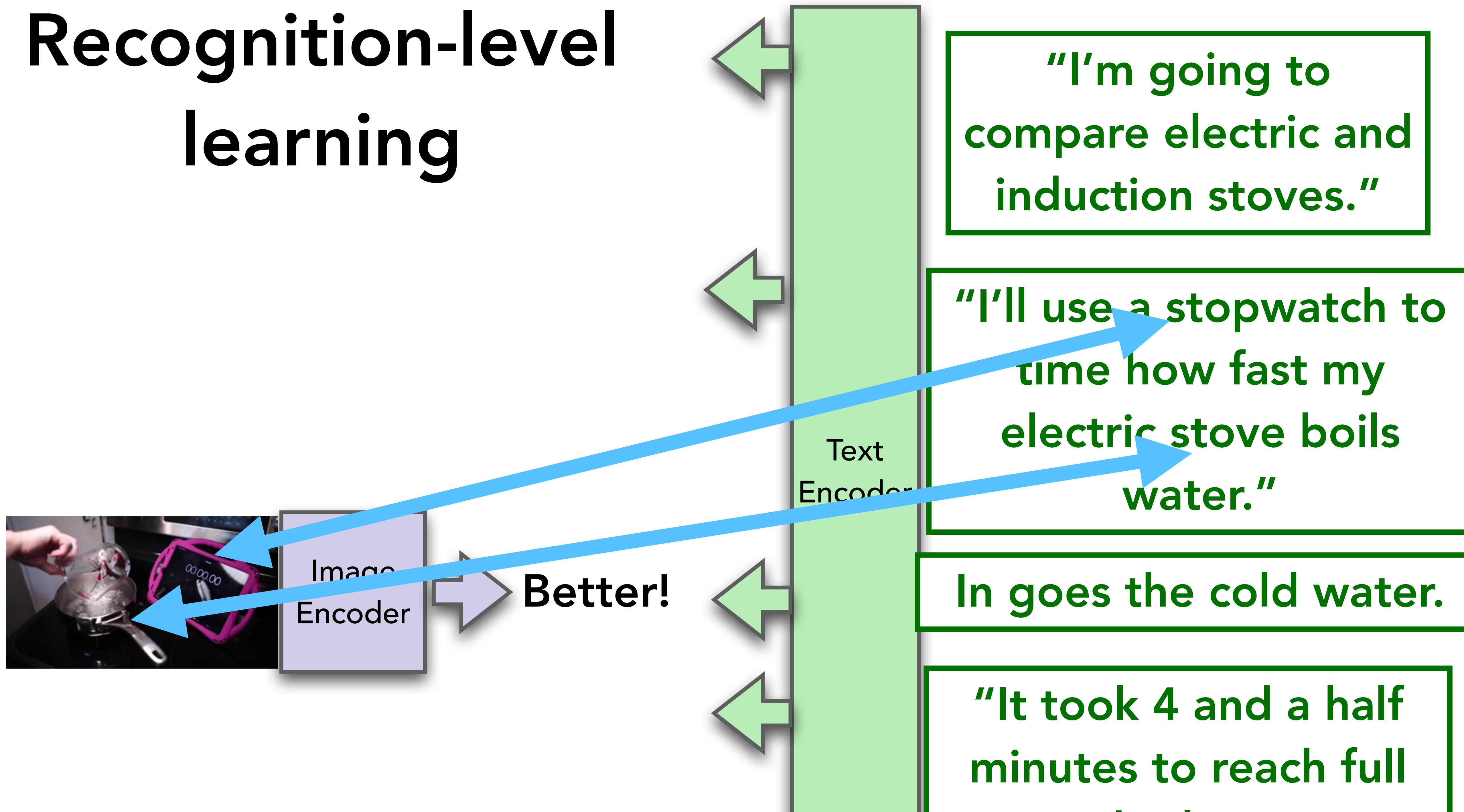


Recognition-level learning

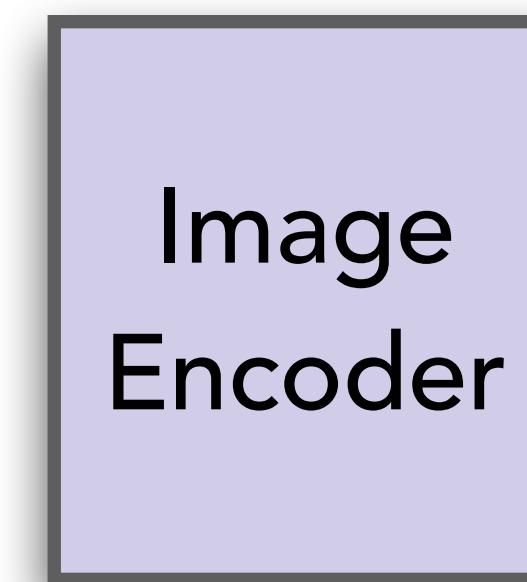
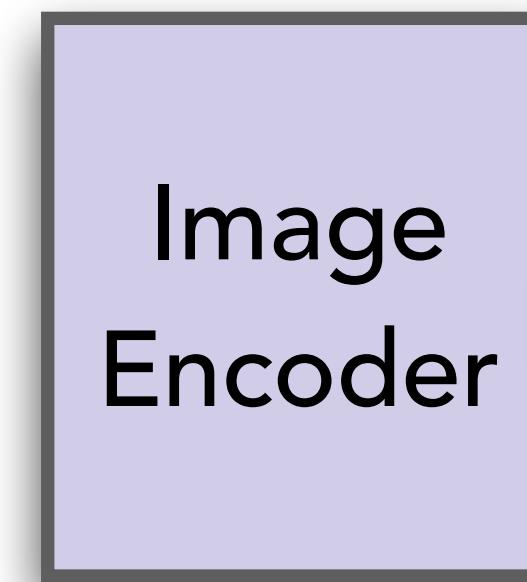


(ConVIRT; Zhang et al 2020, CLIP; Radford et al 2021)

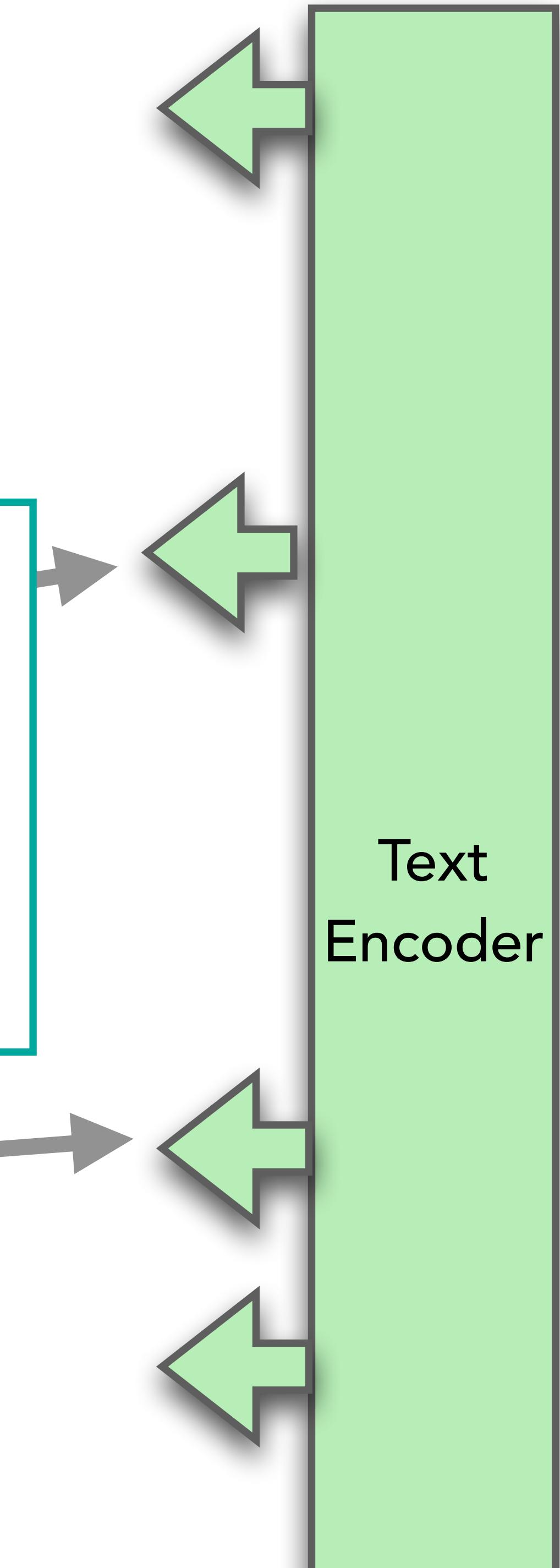
Recognition-level learning



Recognition-level learning



Objective 1: maximize similarity between contextualized language and individual frames



Commonsense Learning

In goes the cold water.

“It took 4 and a half minutes to reach full boil...”



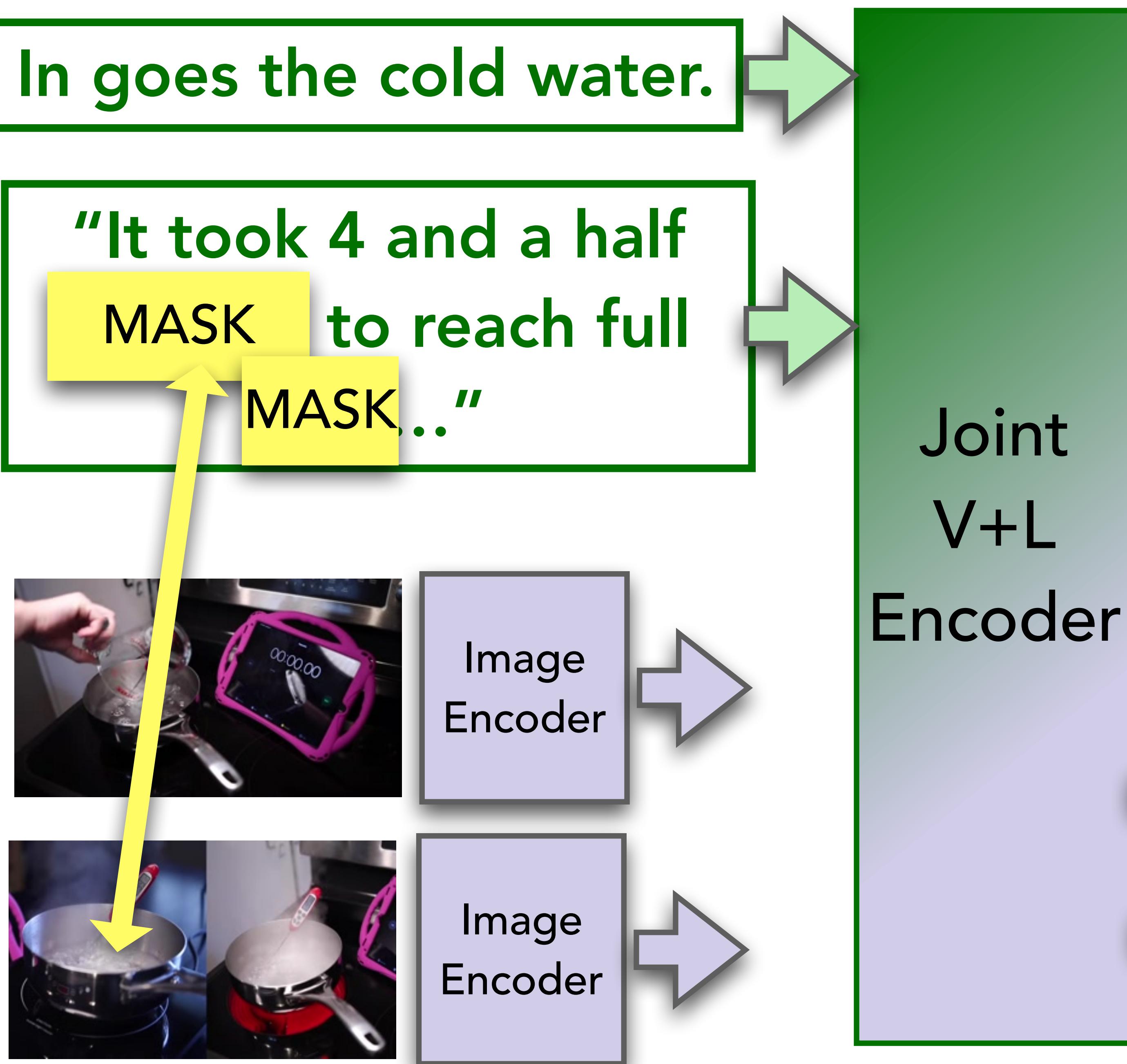
Image Encoder



Image Encoder

Joint
V+L
Encoder

Commonsense Learning



Objective 2:
Mask LM

Commonsense Learning

In goes the cold water.

“It took 4 and a half minutes to reach full boil...”



Image Encoder



Image Encoder

Joint
V+L
Encoder

Frame 2 comes first

Objective 3:
Unshuffle frames

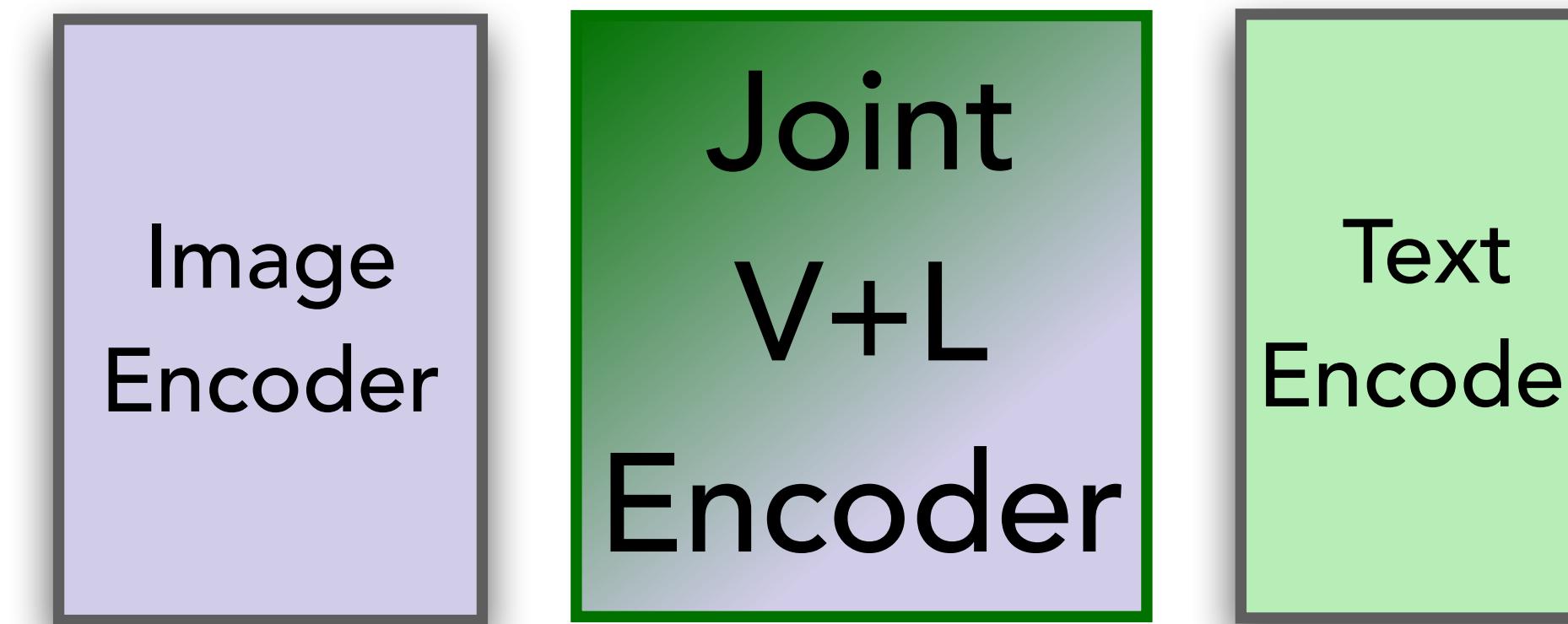
Objective 1:
Contextual Frame-
Text Matching

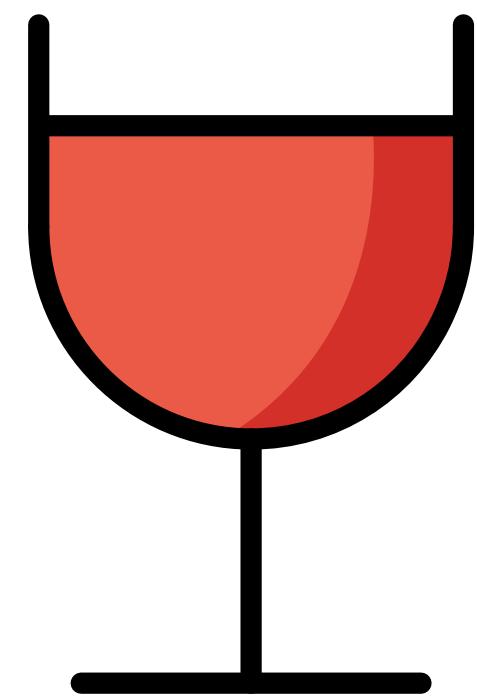
Objective 2:
Mask LM

Objective 3:
Unshuffle frames



Using a 12-layer 'base' Transformer,
train end-to-end on 6M videos





Evaluation

Evaluation 1: Zero-Shot Unscrambling Visual Stories

Task: Given the text of a visual story,
match images to text to tell a narrative

(SIND; Huang et al 2016,
Agrawal et al 2016)

The old man
was riding
the escalator.

He was
almost to the
top.

His kids were
already at the
top.

At the top
was a train
station.

They then
got on the
train.

Task: Given the text of a visual story,
match images to text to tell a narrative



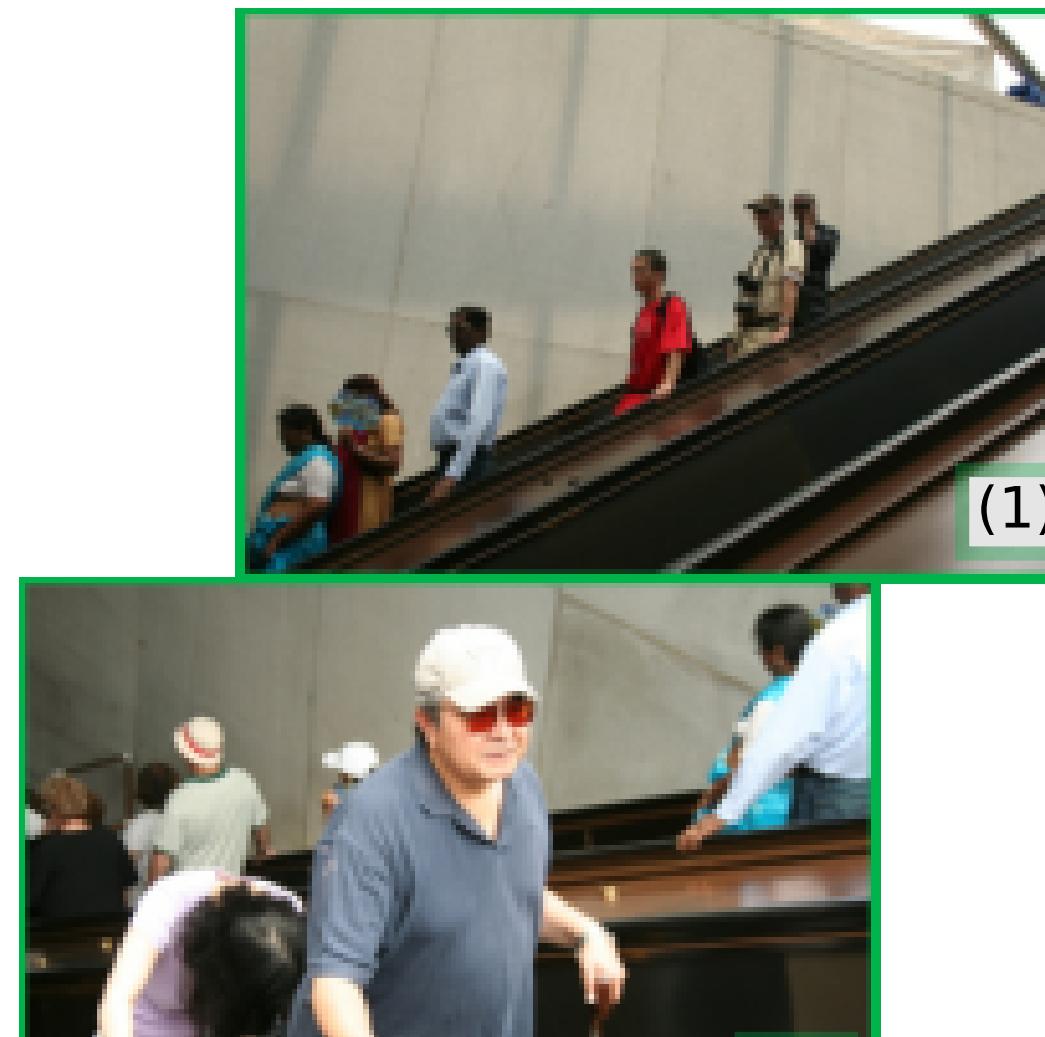
The old man
was riding
the escalator.

He was
almost to the
top.

His kids were
already at the
top.

At the top
was a train
station.

They then
got on the
train.



Task: Given the text of a visual story,
match images to text to tell a narrative



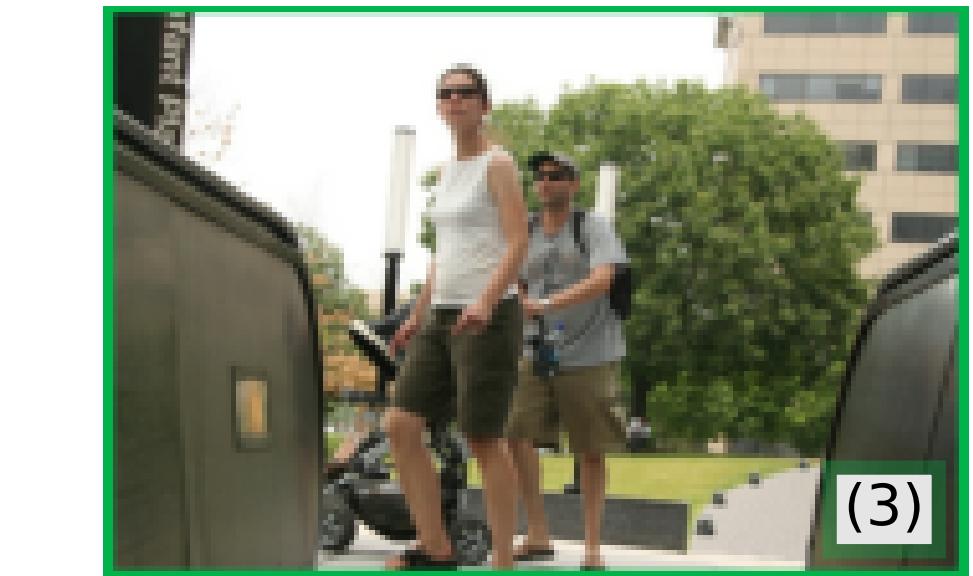
The old man
was riding
the escalator.

He was
almost to the
top.

His kids were
already at the
top.

At the top
was a train
station.

They then
got on the
train.



Our model gets this right *without finetuning*,
using the unscrambling objective

Task: Given the text of a visual story,
match images to text to tell a narrative

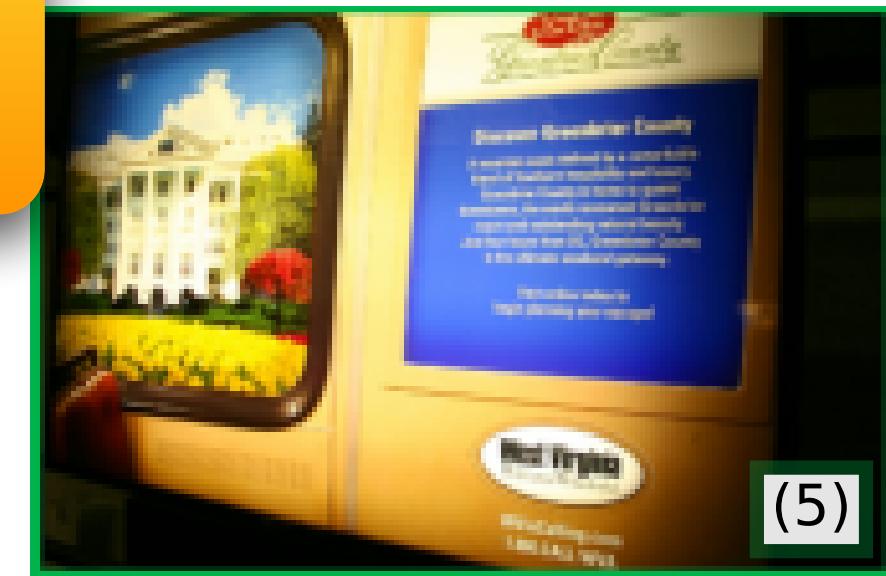
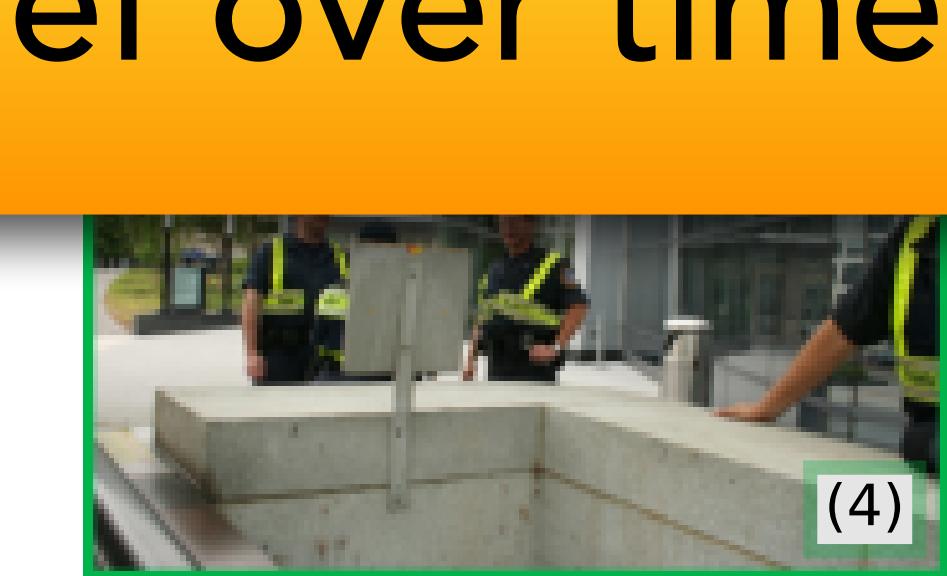
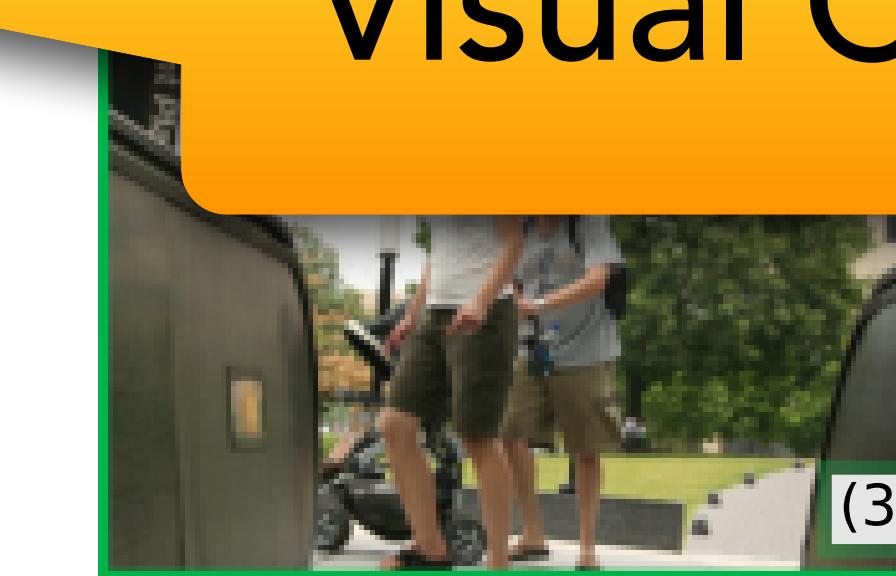
The old man
was riding
the escalator.

He was
almost to the
top.

His kids were
already at the
top.

At the top
was a train
station.

They then
got on the
train.



Visual Coref over time!

The old man
was riding
the escalator.

He was
almost to the
top.

His kids were
already at the
top.

At the top
was a train
station.

They then
got on the
train.

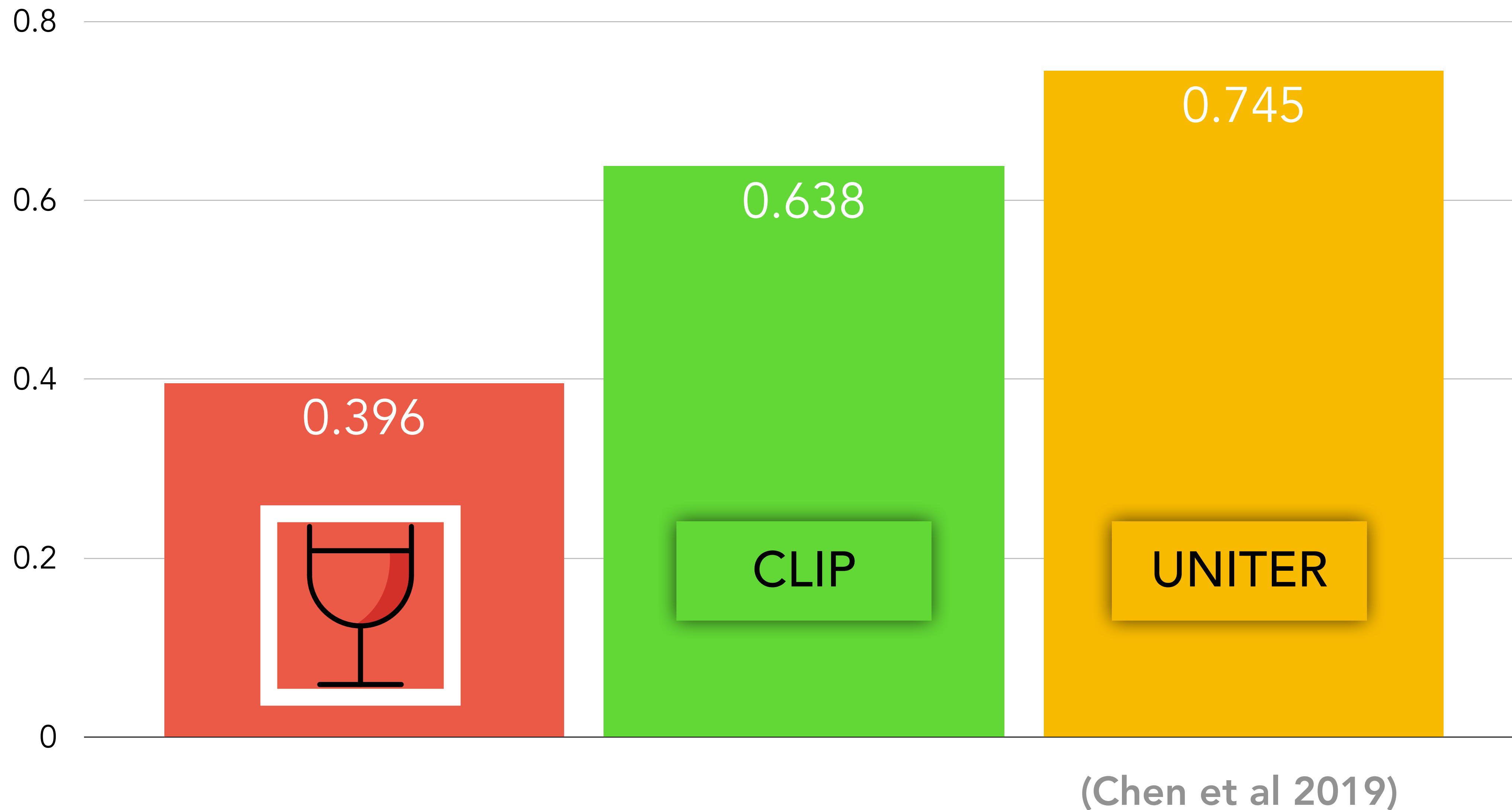


CLIP

(Radford et al 2021)



Distance away from sorted order (lower is better, 5.0 is max)



Even when our model is “wrong” it’s kinda cool

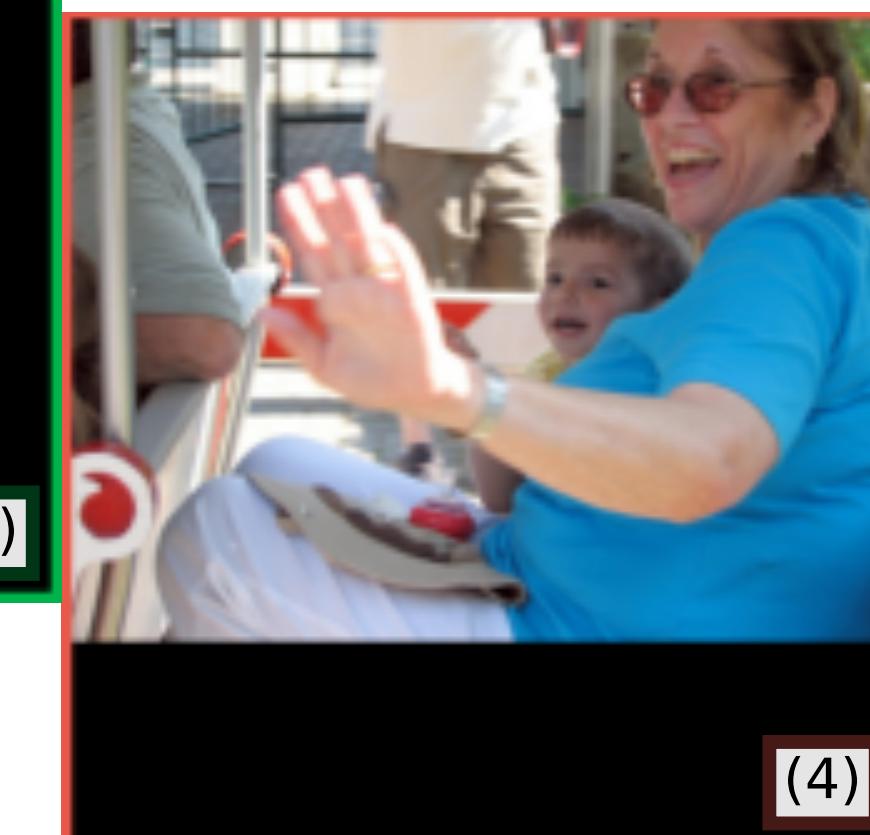
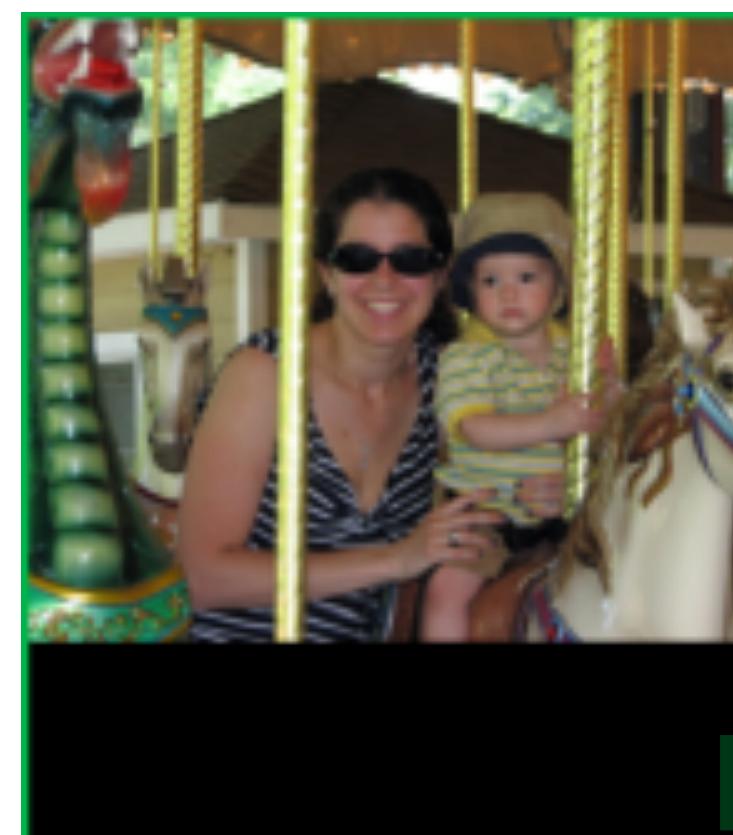
I went to the fair with my kids last weekend.

There were a lot of people there.

They also had a barn.

We got to see a lot of animals.

We can't wait to go back later.



Even when our model is “wrong” it’s kinda cool

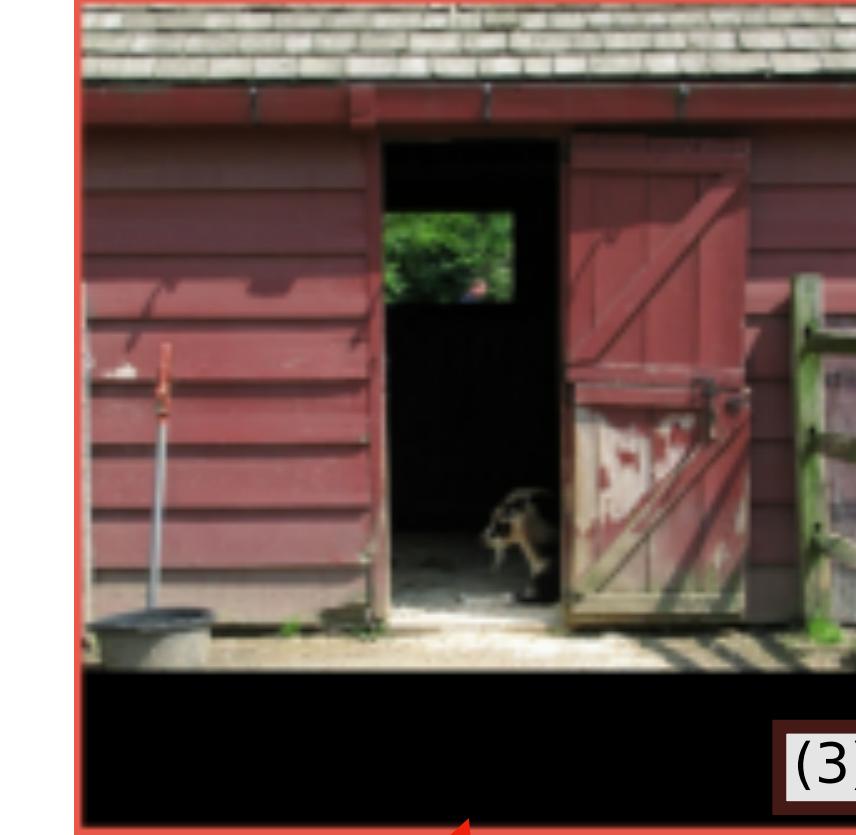
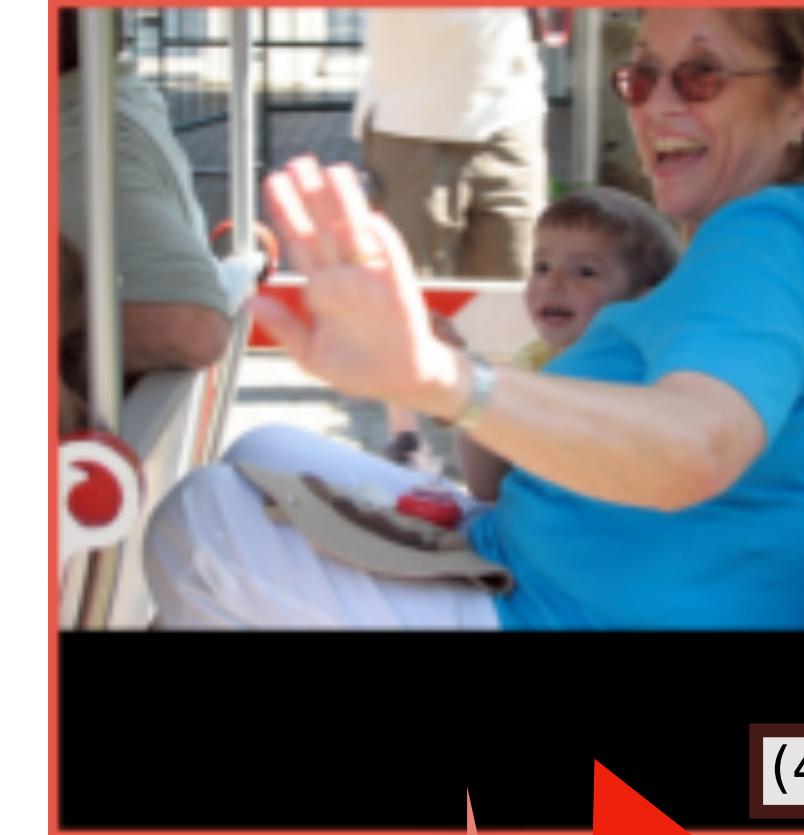
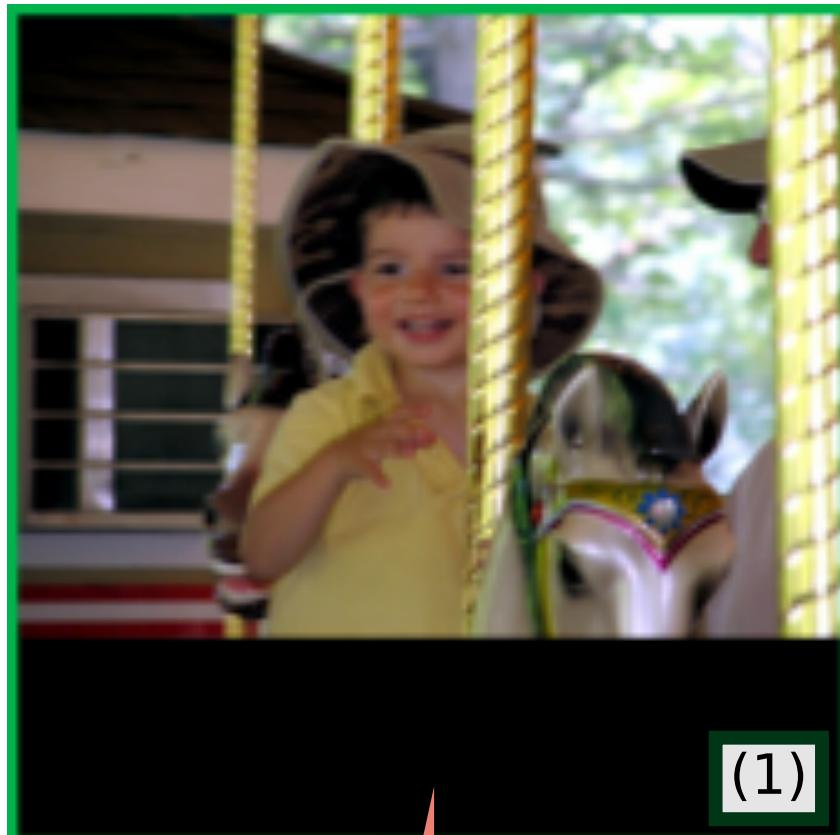
I went to the fair with my kids last weekend.

There were a lot of people there.

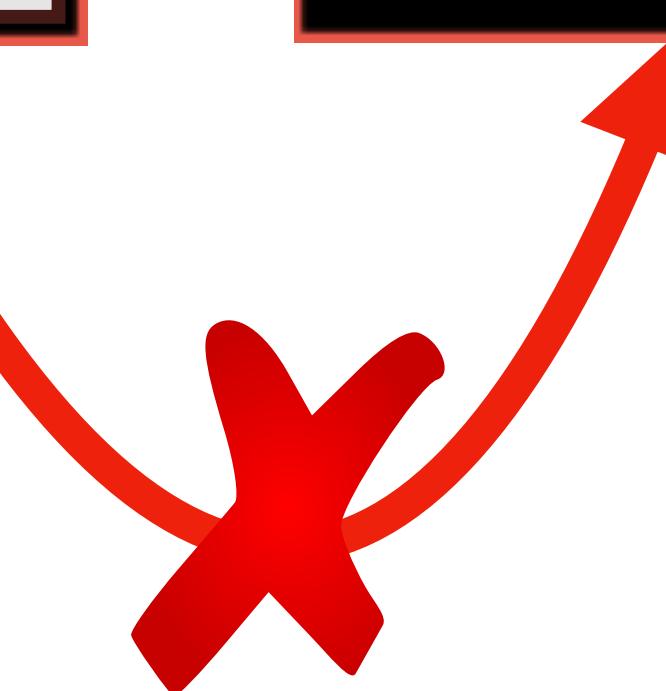
They also had a barn.

We got to see a lot of animals.

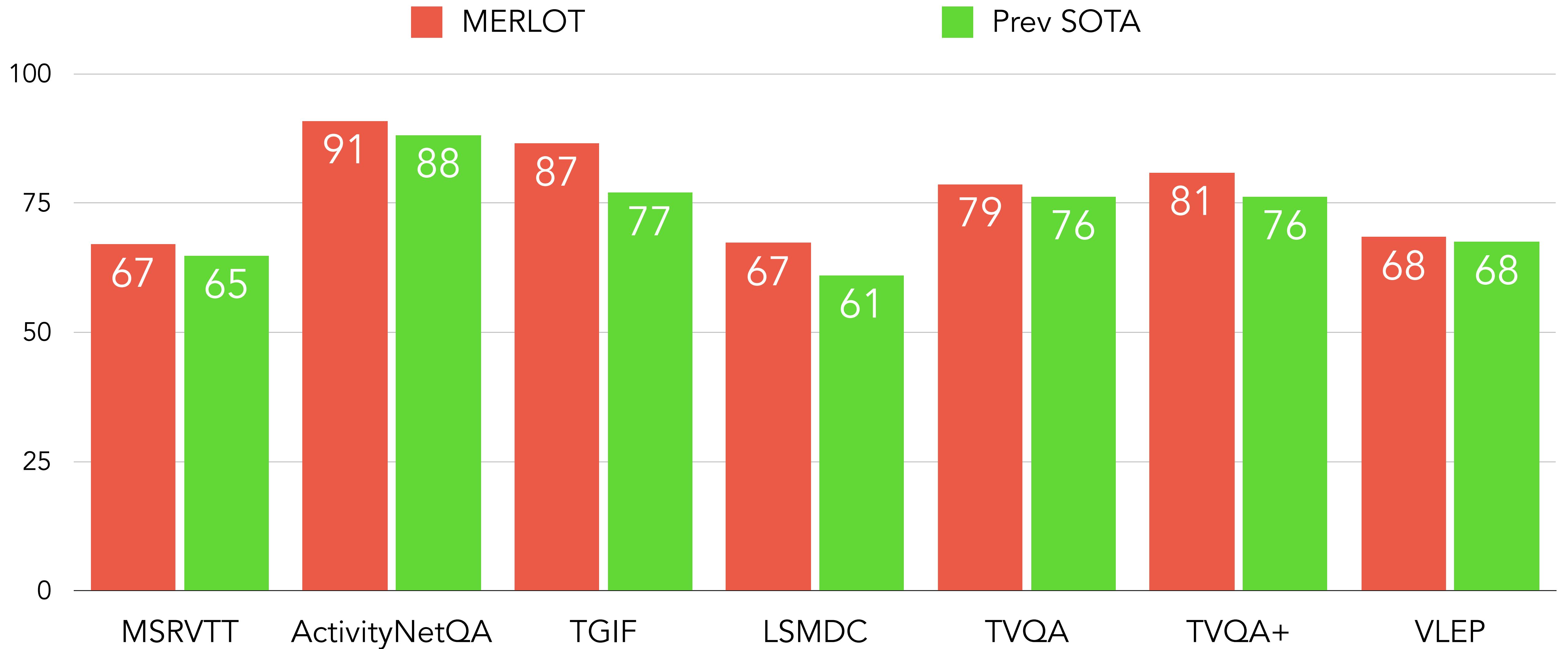
We can't wait to go back later.



MERLOT: people stay on the Merry-Go-Round for a while



Evaluation 2: Fine-tuned Video QA



Evaluation 3: Visual Commonsense Reasoning (Q->AR)

MERLOT

UNITER

VILLA

ERNIE-ViL

70

65

60

55

65.1



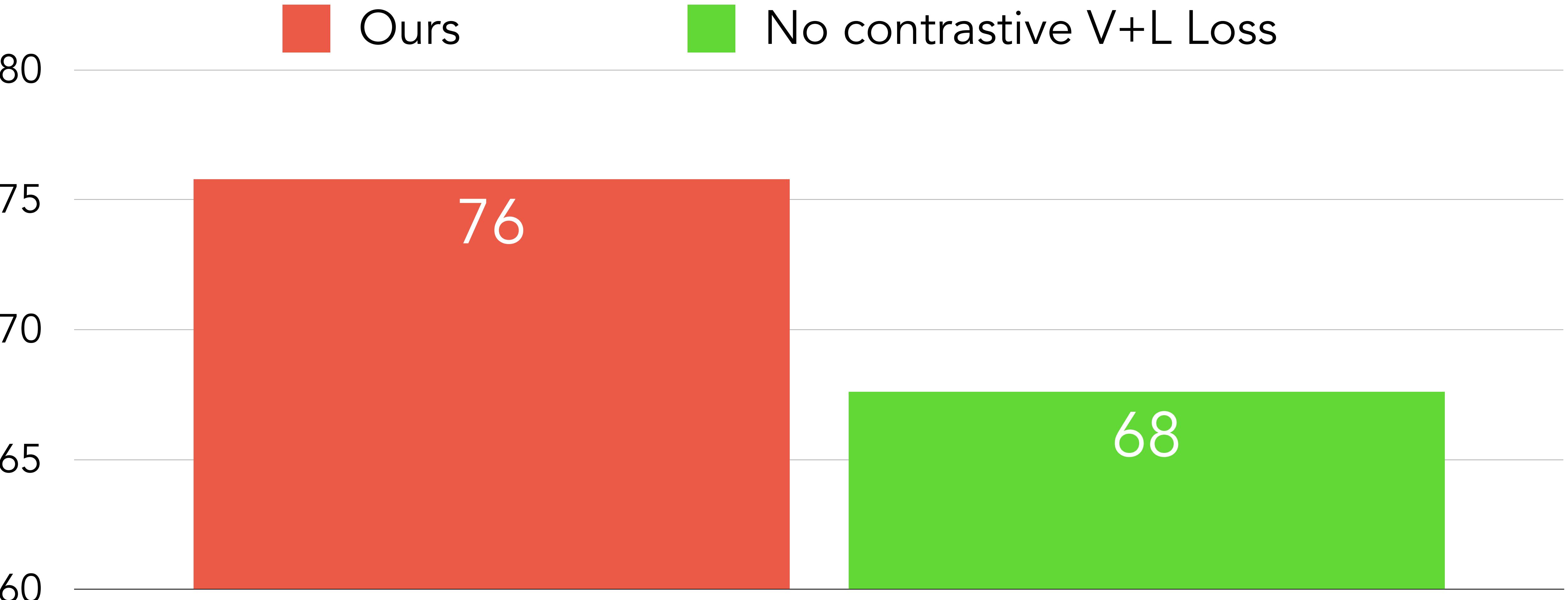
58.2

60.6

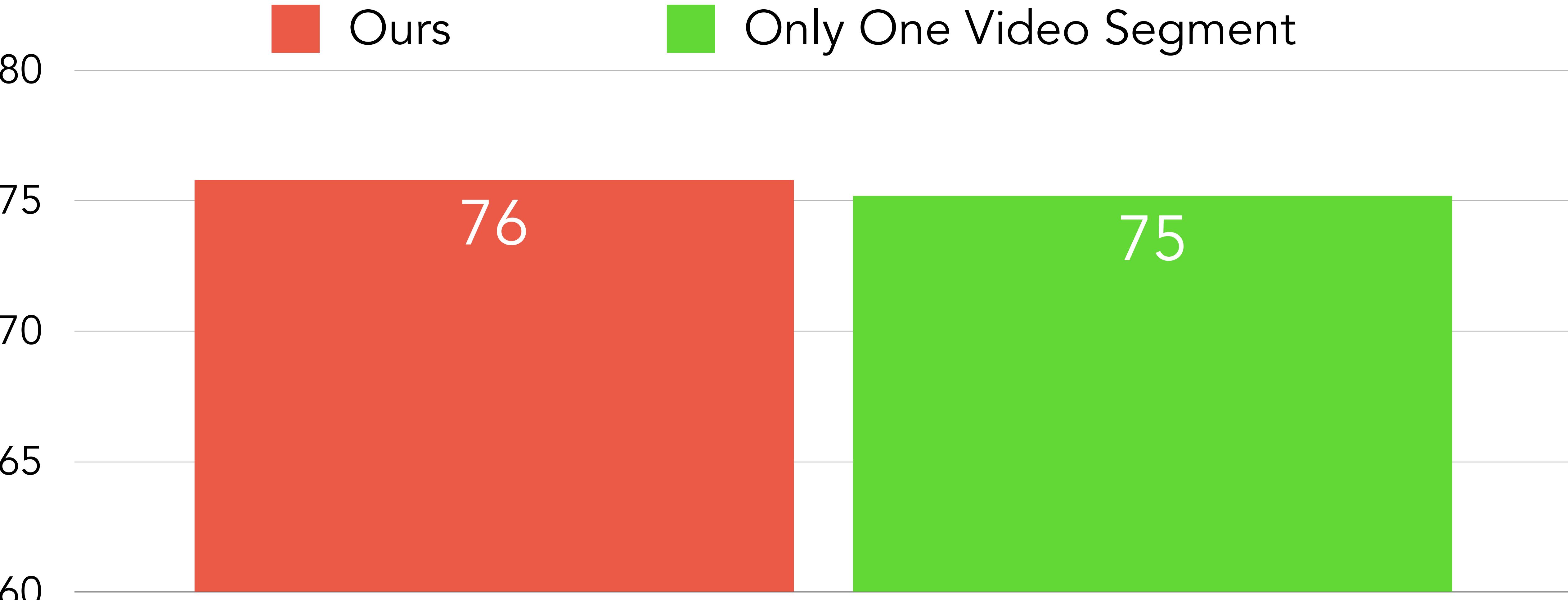
62.1

Despite no supervised object detector, and
never seeing still images before

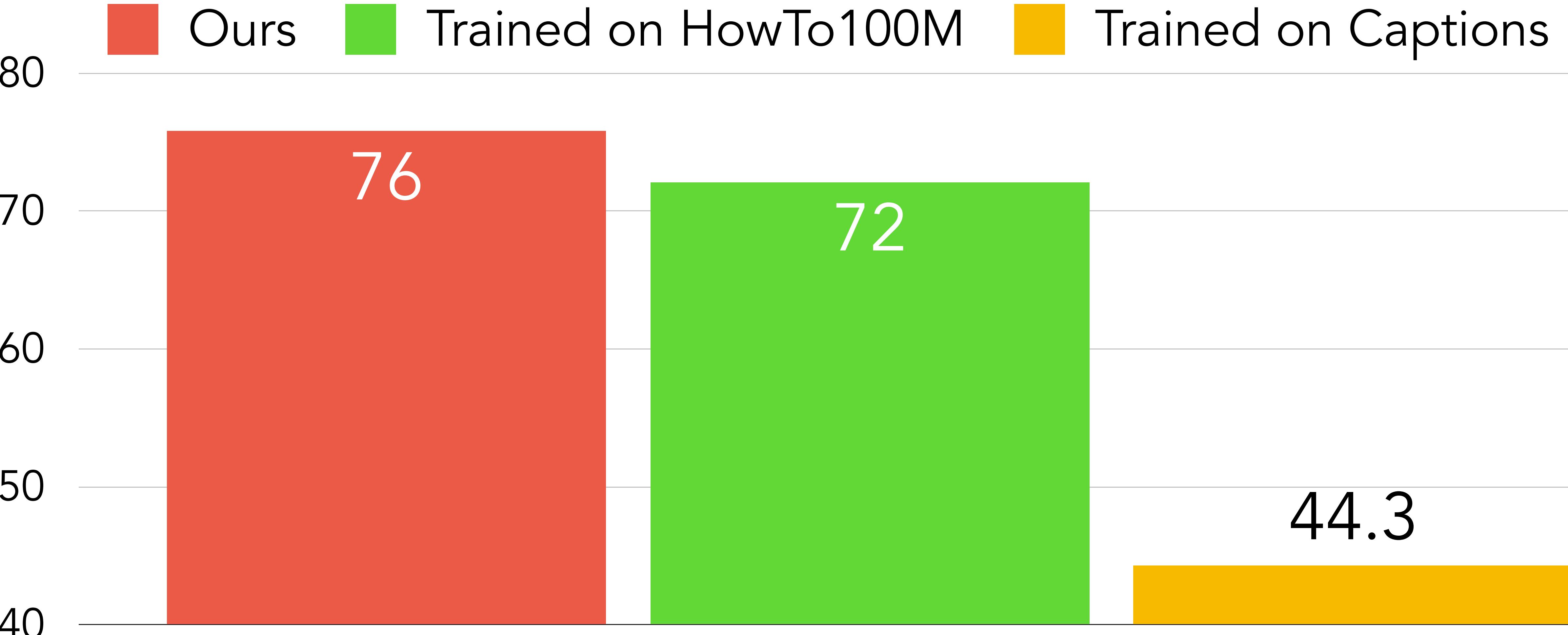
Analysis (on TVQA+)



Analysis (on TVQA+)

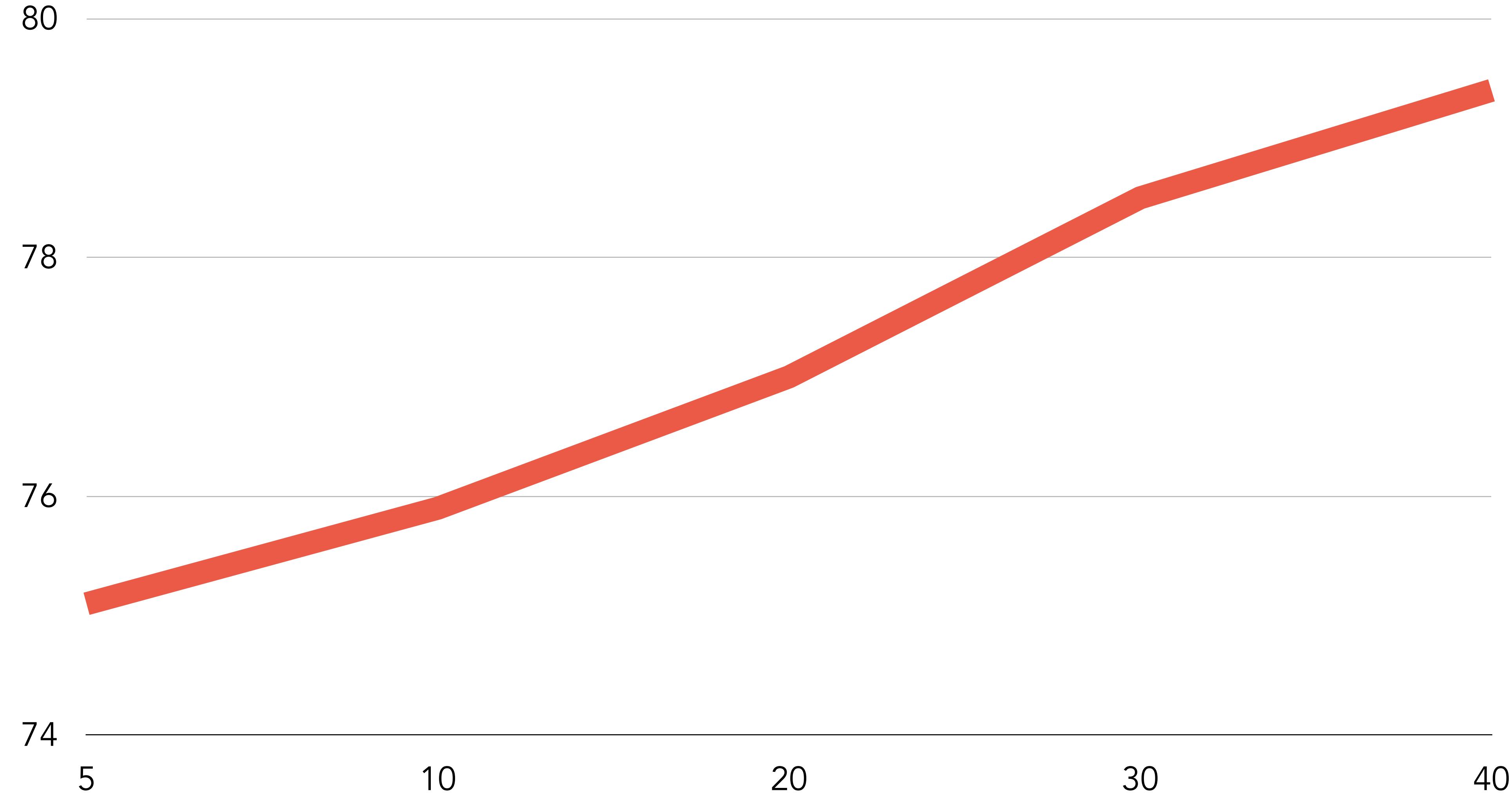


Analysis (on TVQA+)



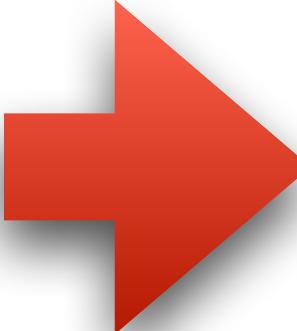
Performance increases with # epochs

VCR
Q->A



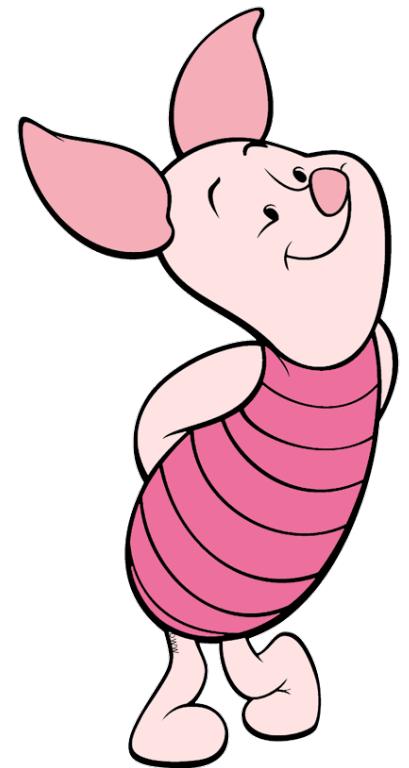
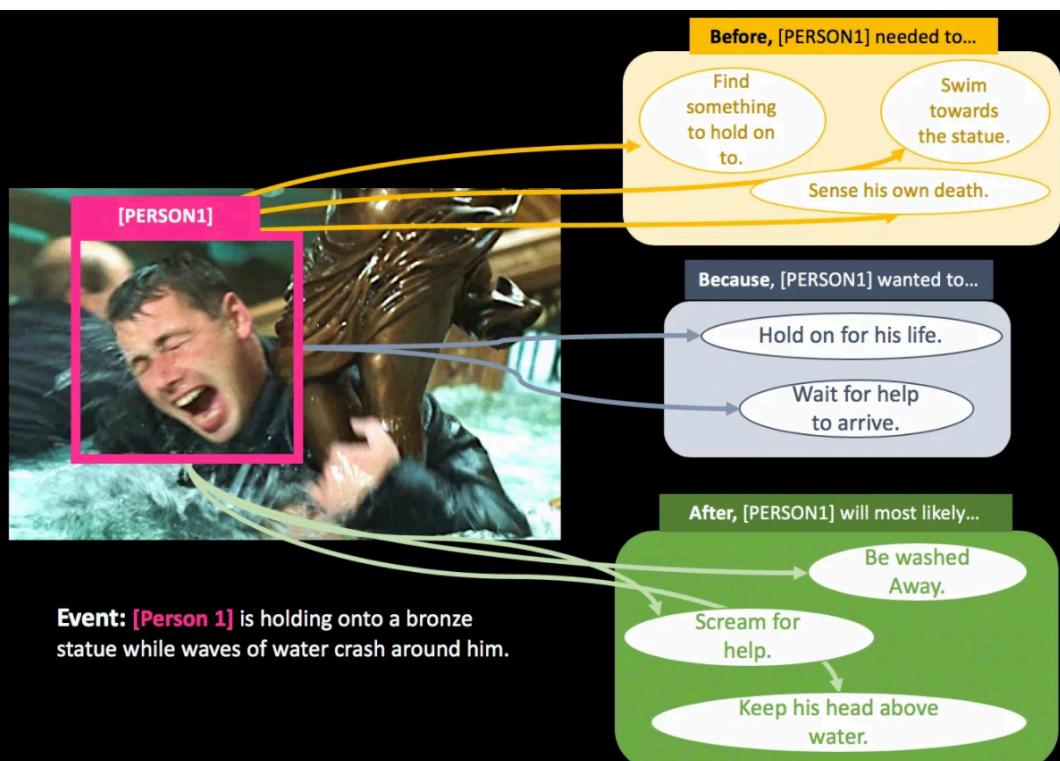
Harnad's Symbol Grounding Problem

- Grounding with 3D
- Grounding with 2D + Time



Grounding with 2D + KG

Visual Comet



PIGLeT



MERIOT

Visual COMET: Reasoning about the *Dynamic* Context of a *Still* Image

ECCV 2020

Jae Sung (James) Park



Chandra
Bhagavatula



Roozbeh
Mottaghi



Ali
Farhadi



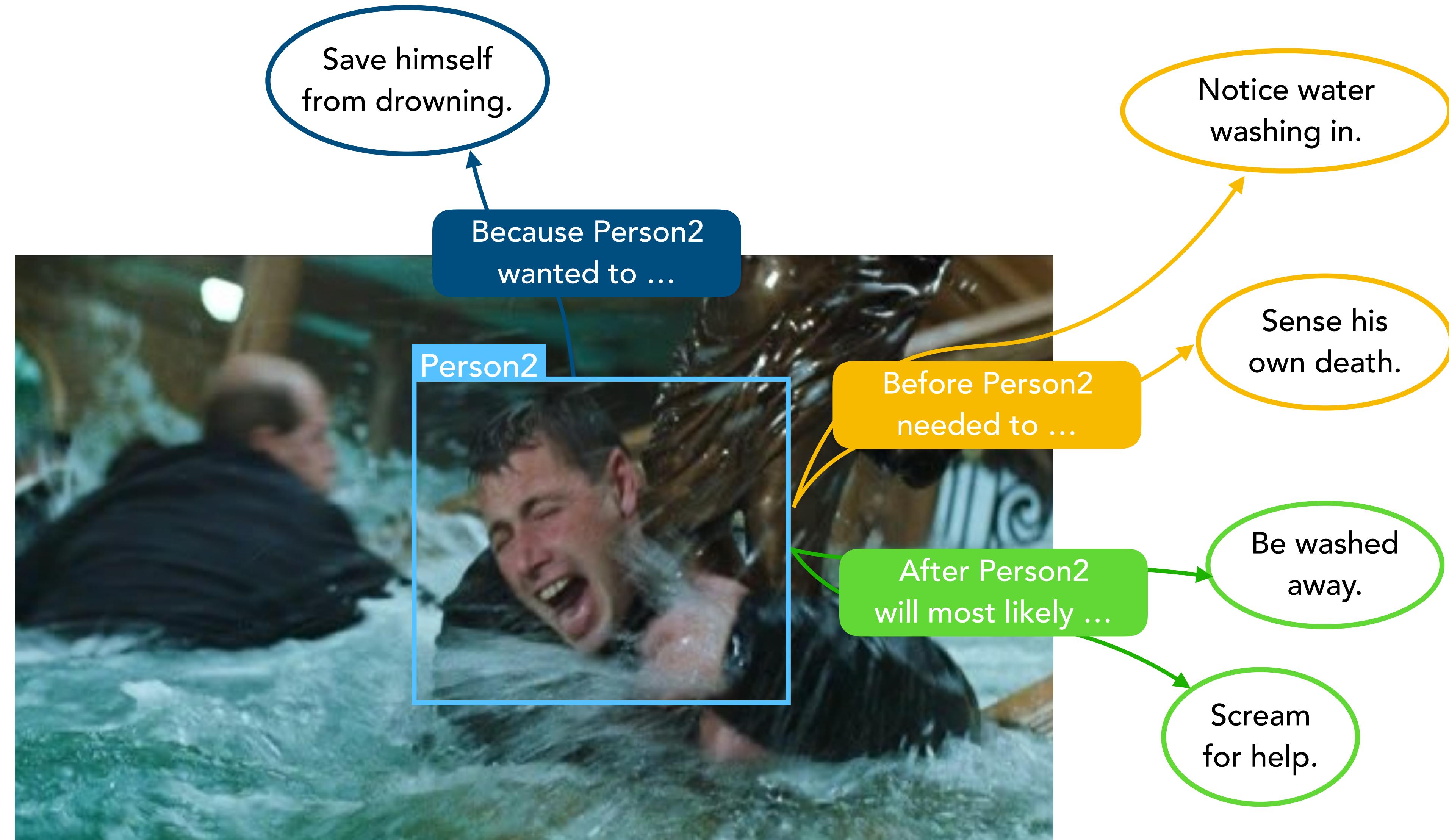
Yejin
Choi





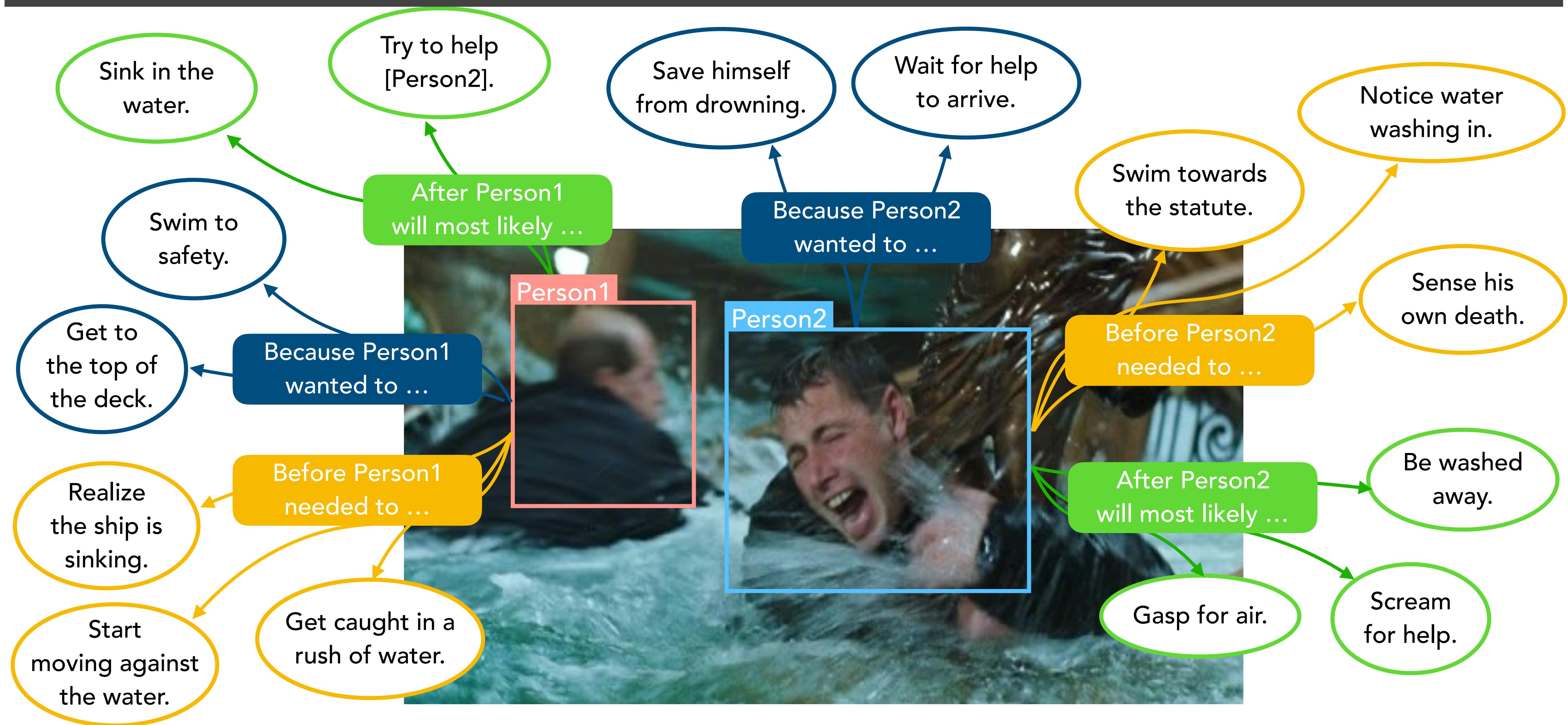
Visual Commonsense Graphs:

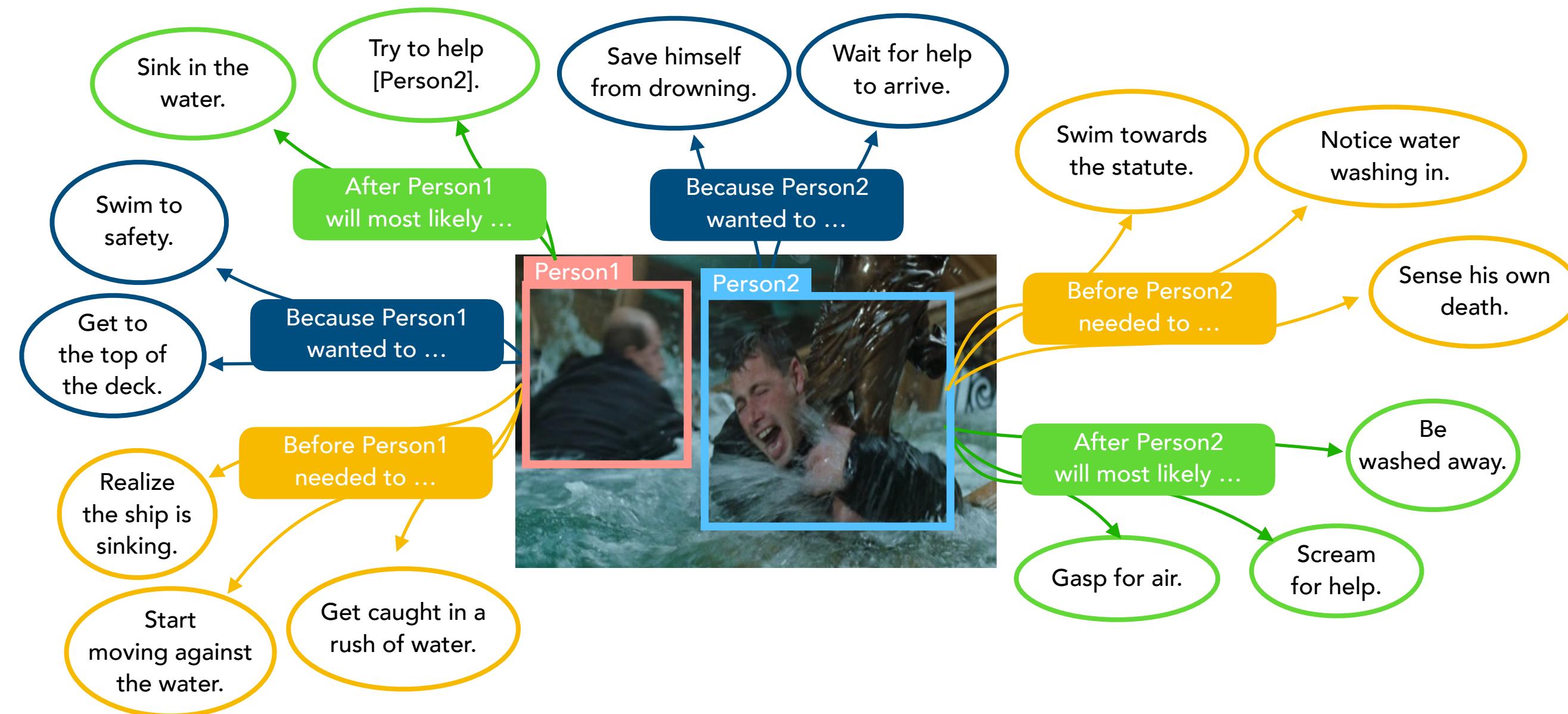
Reasoning about the *Dynamic Context* of a *Still* Image

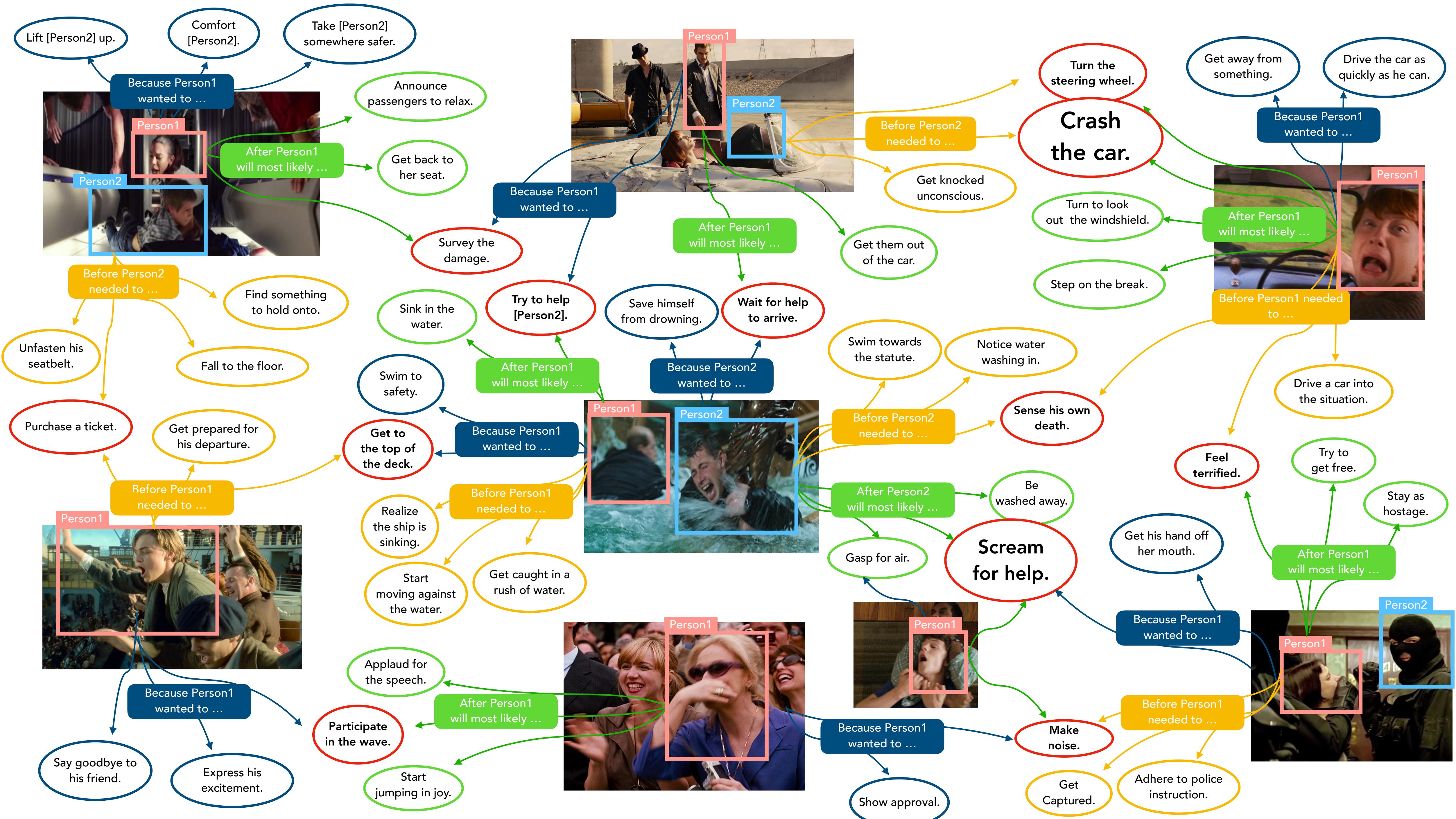


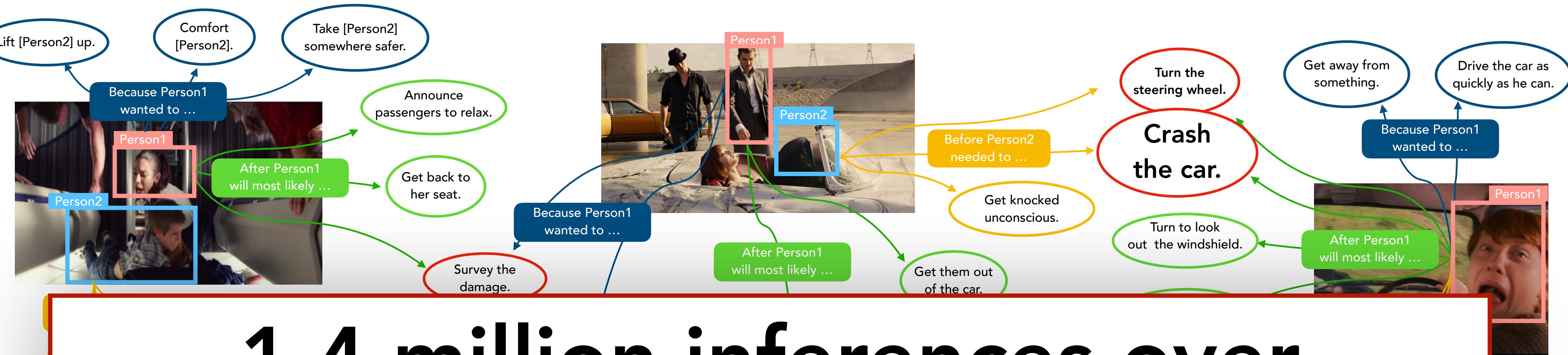
Visual Commonsense Graphs:

Reasoning about the *Dynamic Context* of a *Still* Image





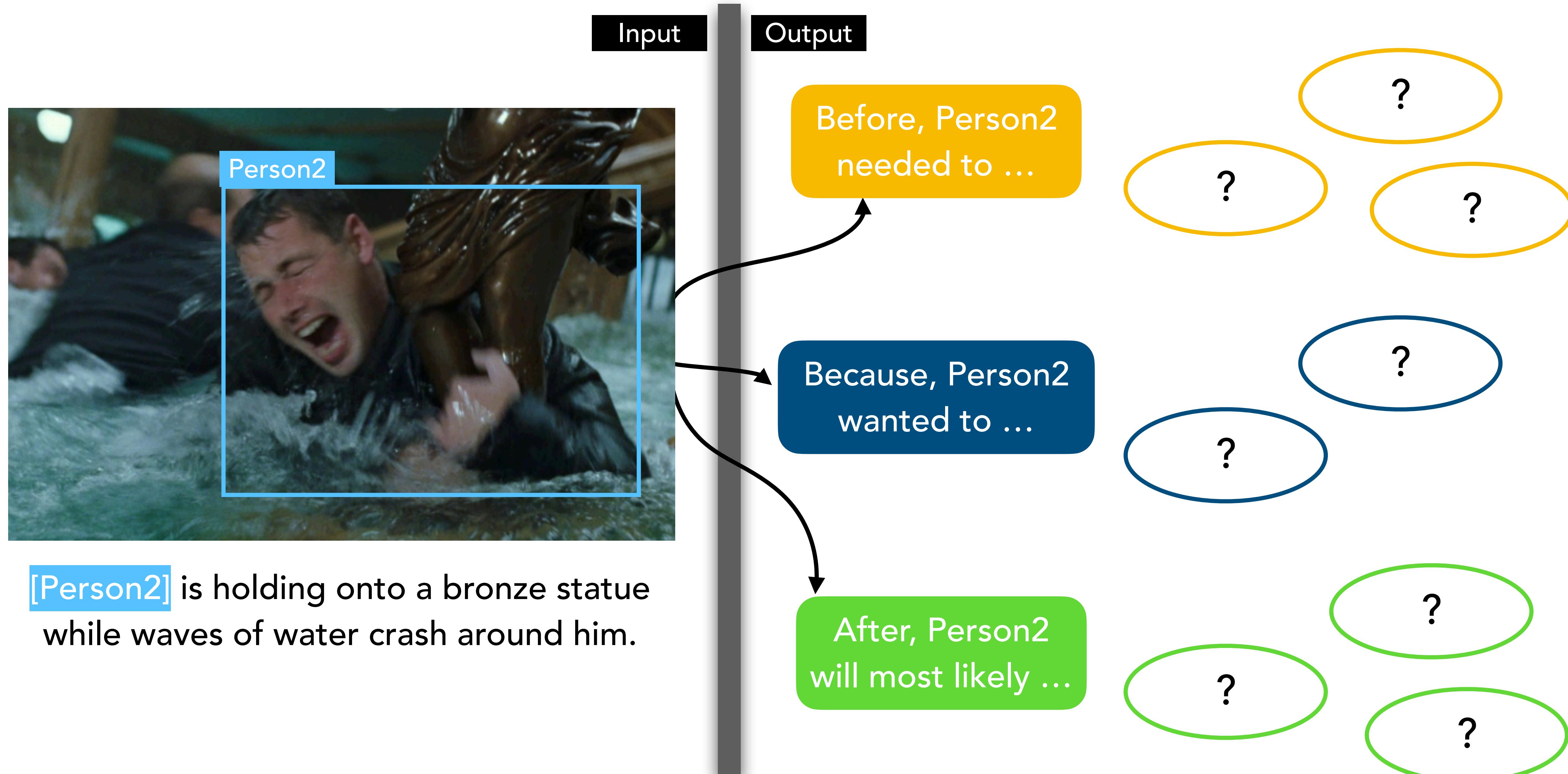




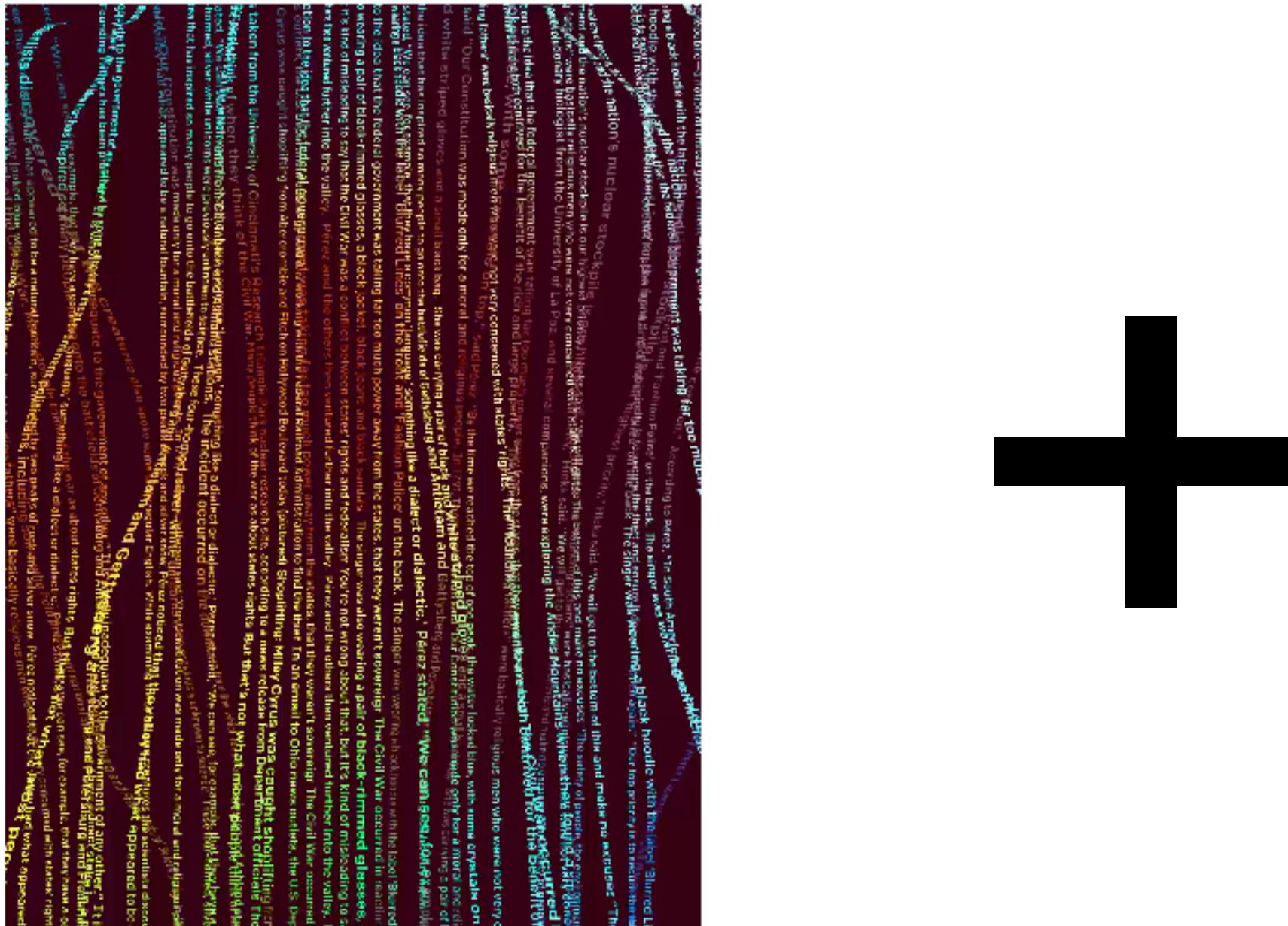
1.4 million inferences over
60K images with

<https://visualcomet.xyz>

Task: Generating Commonsense Inferences in Language

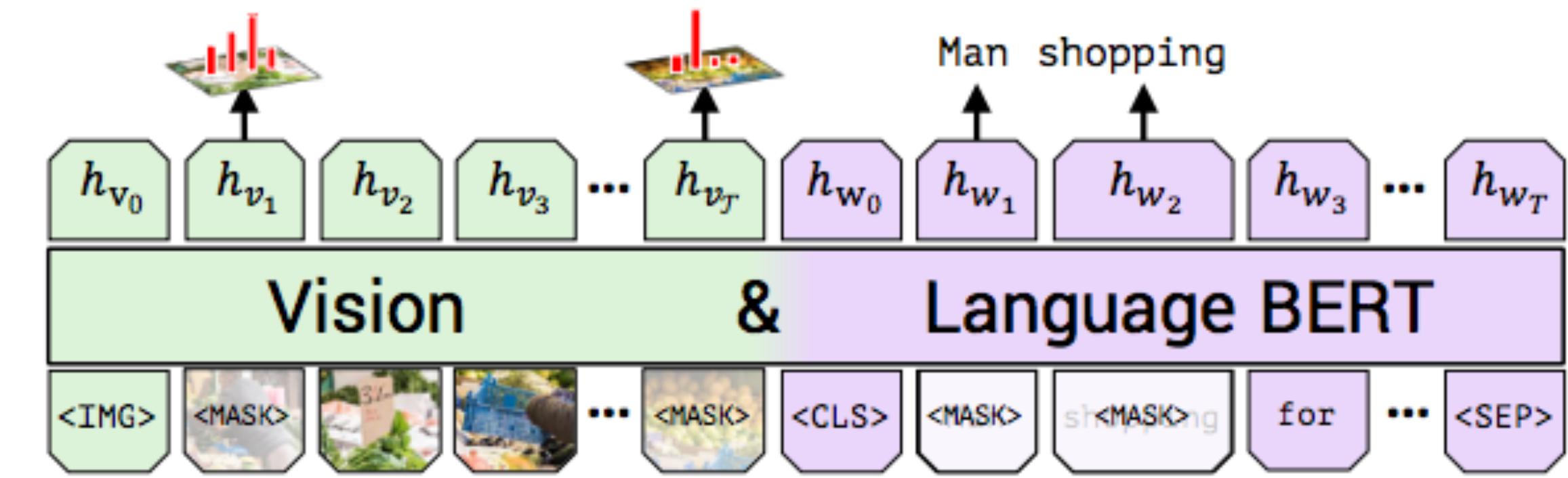


Our Model Builds on Pre-Trained Language Models



GPT-2 for Conditional Generation

(Radford et. al., 2019)



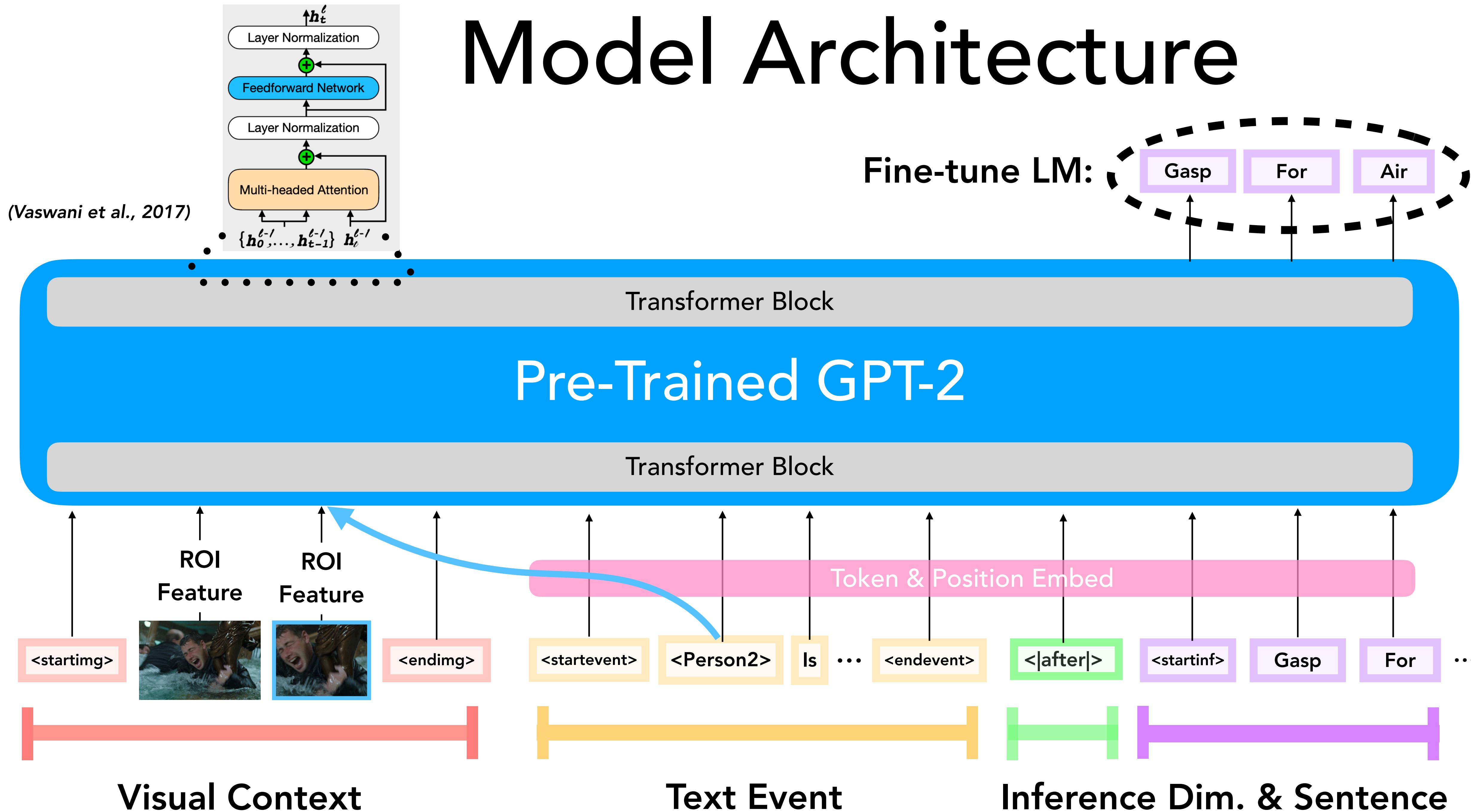
Vision-Language Transformer Architecture

(Lu et. al., 2020; Su et. al, 2020; Tan et. al, 2020)

Our Approach



Model Architecture



Before, Person1 needed to ...

Unlikely

[Person1] is putting a platter on the table at an outdoor restaurant.

Input

Output

Lang Only

Buy groceries.

Get up from the table.

Put food on the platter.



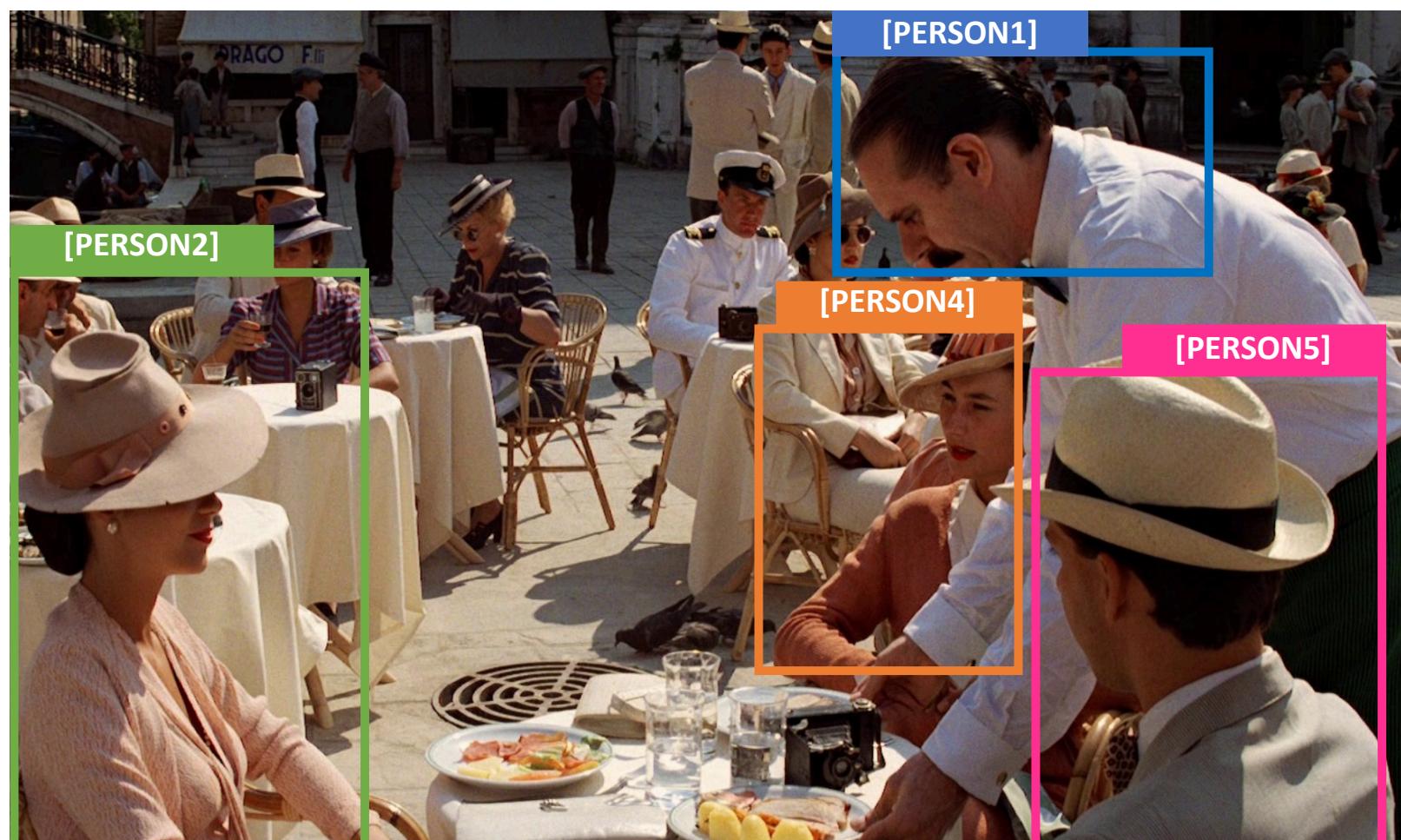
Before, Person1 needed to ...

Unlikely

Input

Output

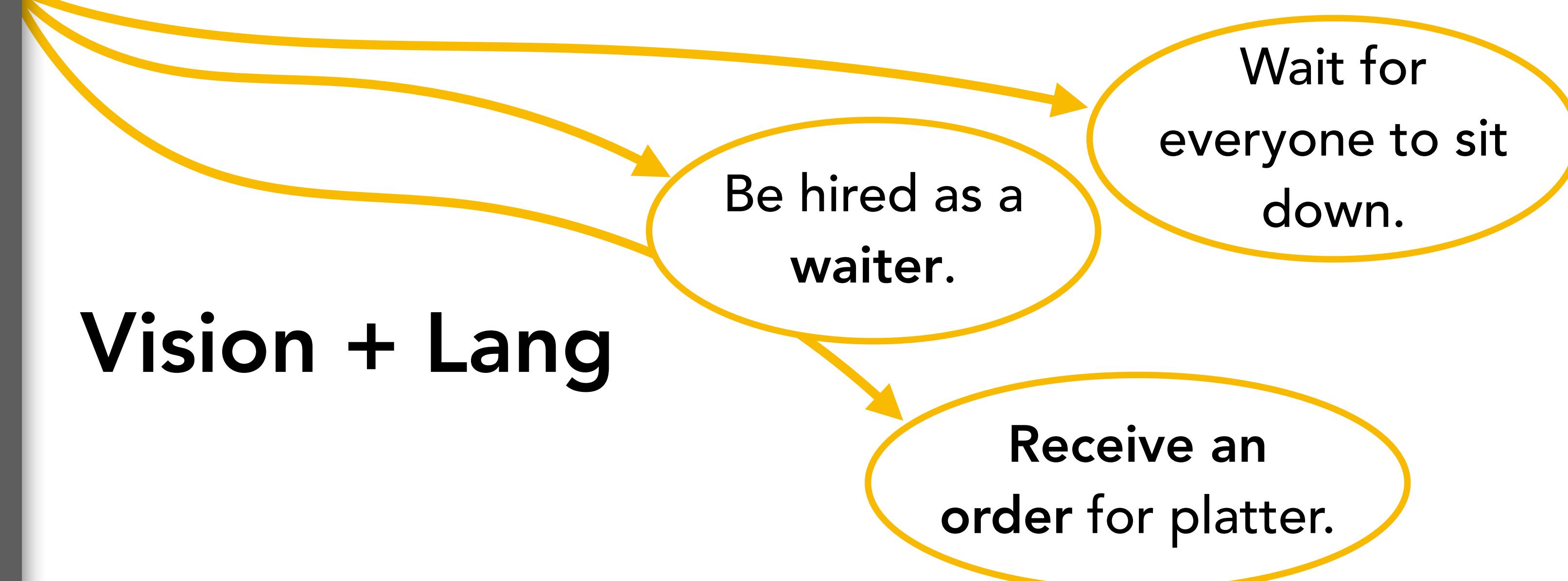
[Person1] is putting a platter on the table at an outdoor restaurant.



Lang Only



Vision + Lang



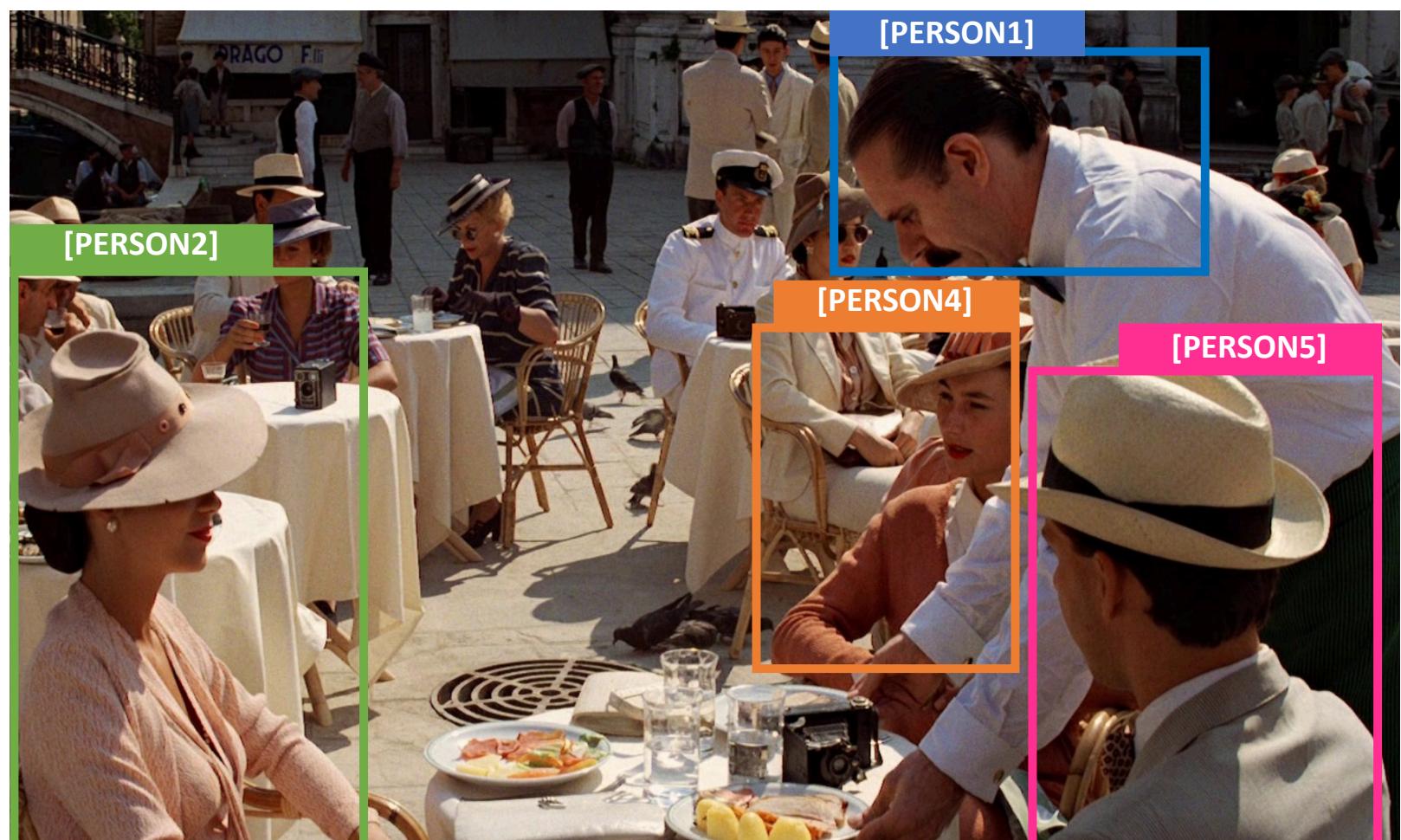
Because, Person1 wanted to ...

Unlikely

Input

Output

[Person1] is putting a platter on the table at an outdoor restaurant.



Lang Only

Have dessert

Ensure
the food is taken
care of.

Tend to the
patrons.

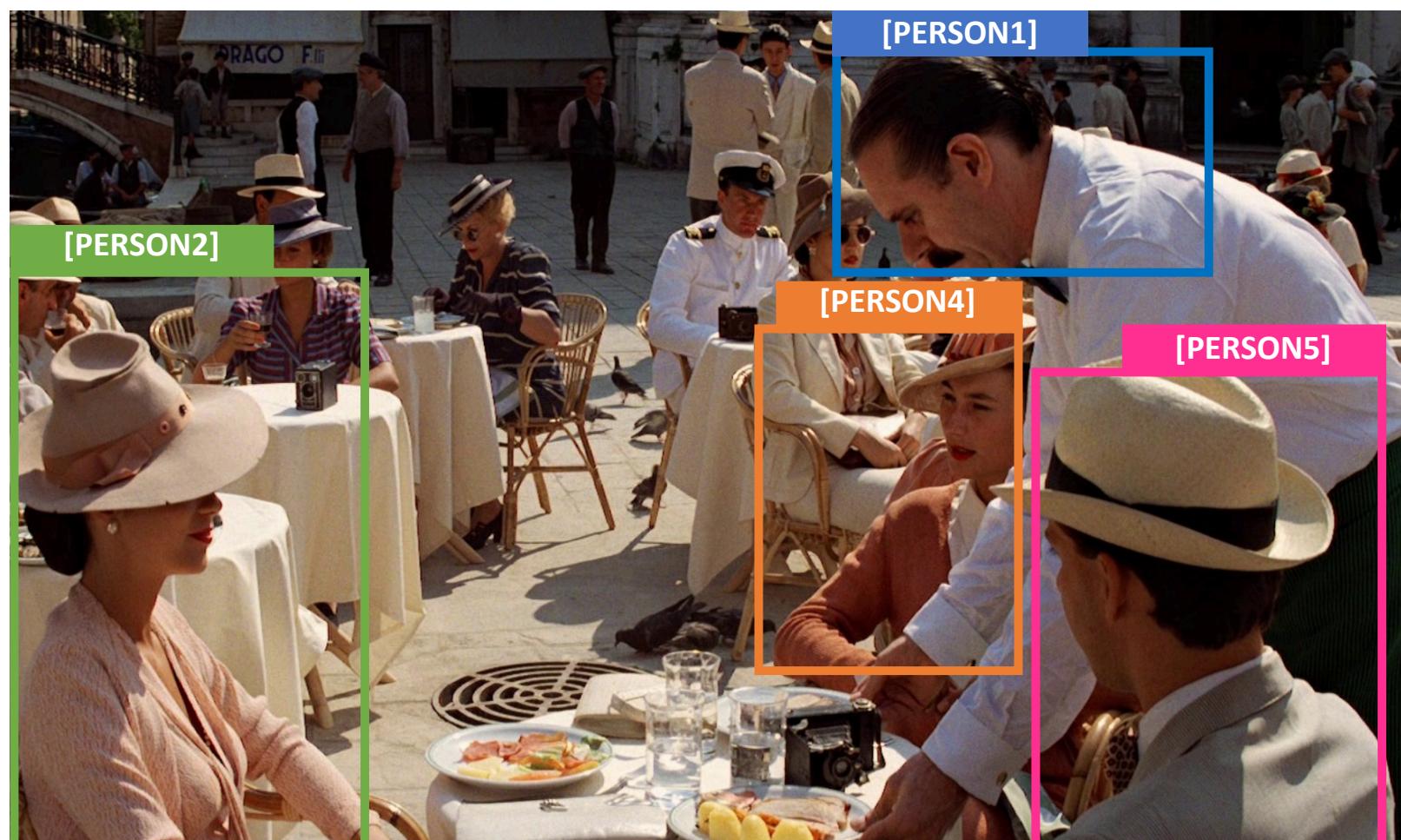
Because, Person1 wanted to ...

Unlikely

Input

Output

[Person1] is putting a platter on the table at an outdoor restaurant.



Lang Only

Have dessert

Ensure the food is taken care of.

Tend to the patrons.

Vision + Lang

Serve [P2], [P4], and [P5].

Greet [P2], [P4], and [P5].

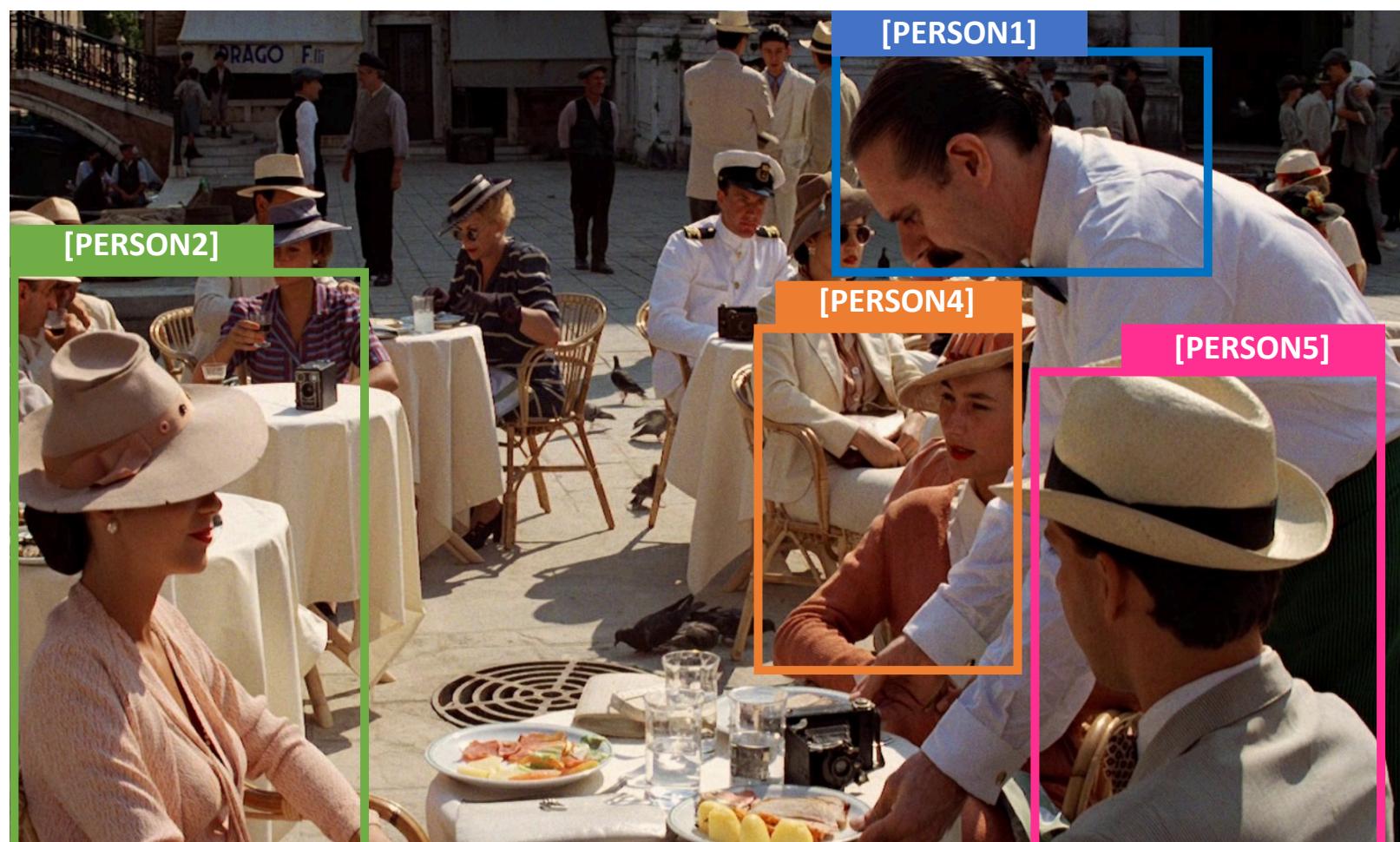
Have [P2], [P4], and [P5] to eat.

After, Person1 will most likely ...

Unlikely

Input

[Person1] is putting a platter on the table at an outdoor restaurant.



Output

Lang Only

Sip the water.

Ask [P2] for a menu.

Get up and walk over to his table.

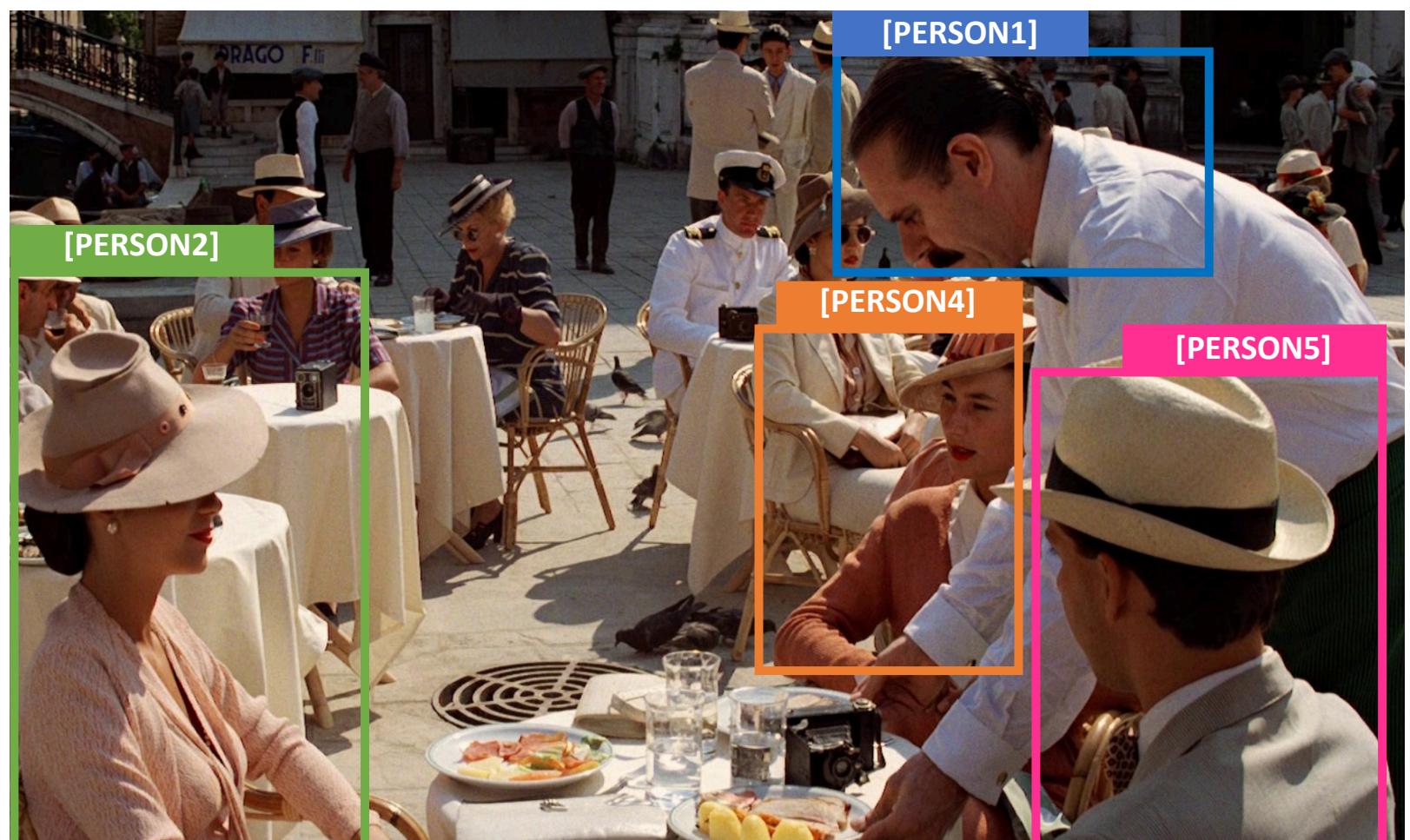
After, Person1 will most likely ...

Unlikely

Input

Output

[Person1] is putting a platter on the table at an outdoor restaurant.



Lang Only

Sip the water.

Ask [P2] for a menu.

Get up and walk over to his table.

Vision + Lang

Take drinks.

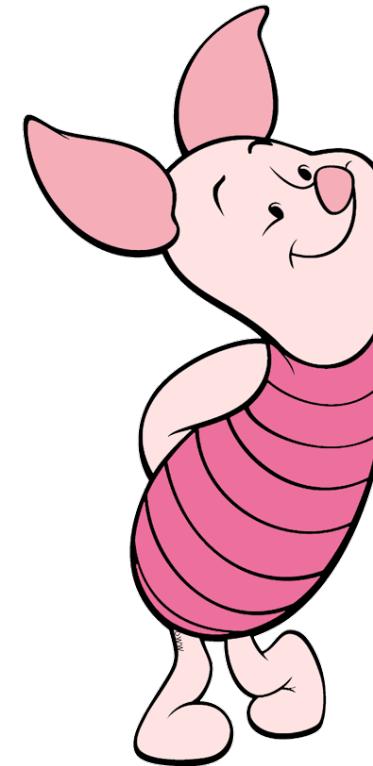
Get back to his work duties.

Get back to the kitchen to get more food.

Harnad's Symbol Grounding Problem

- **Grounding with 3D**

Interactions at the cost of concept coverage



PIGLeT

- **Grounding with 2D + Time**

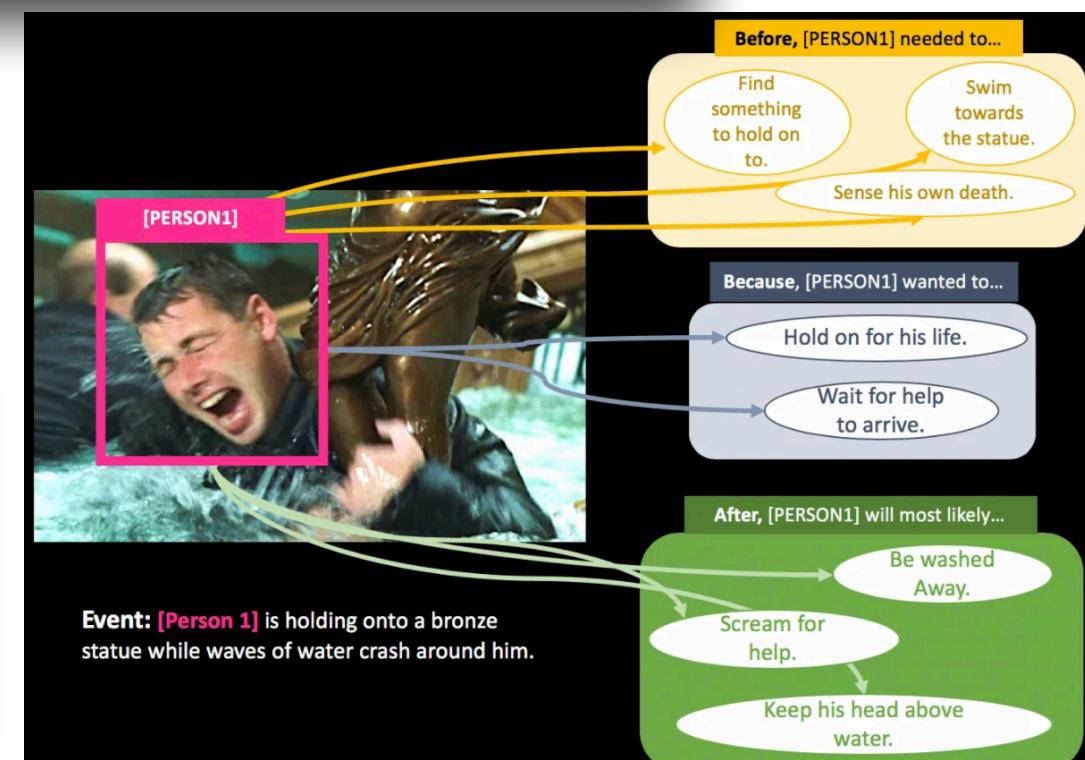
Far richer concepts (causal / temporal interactions) at the cost of direction interactions with the world

MEDLIOT

- **Grounding with 2D + KG**

Learning only from raw data vs from rich declarative knowledge about the world

World Context





Thanks! Questions?