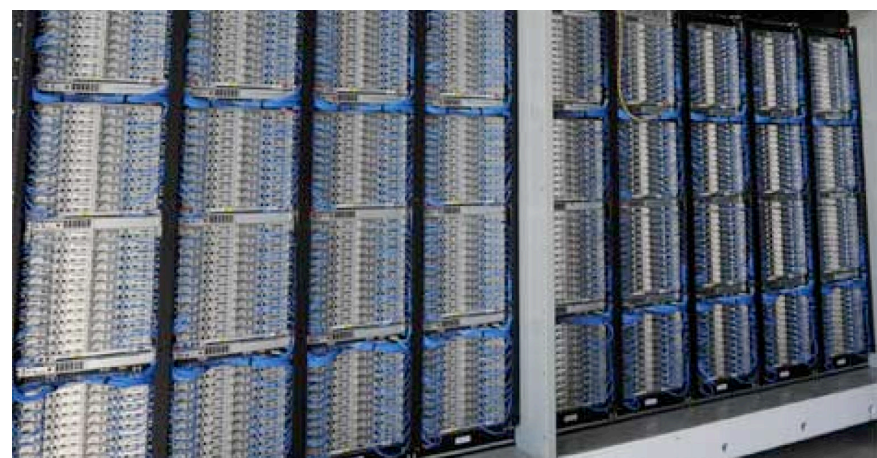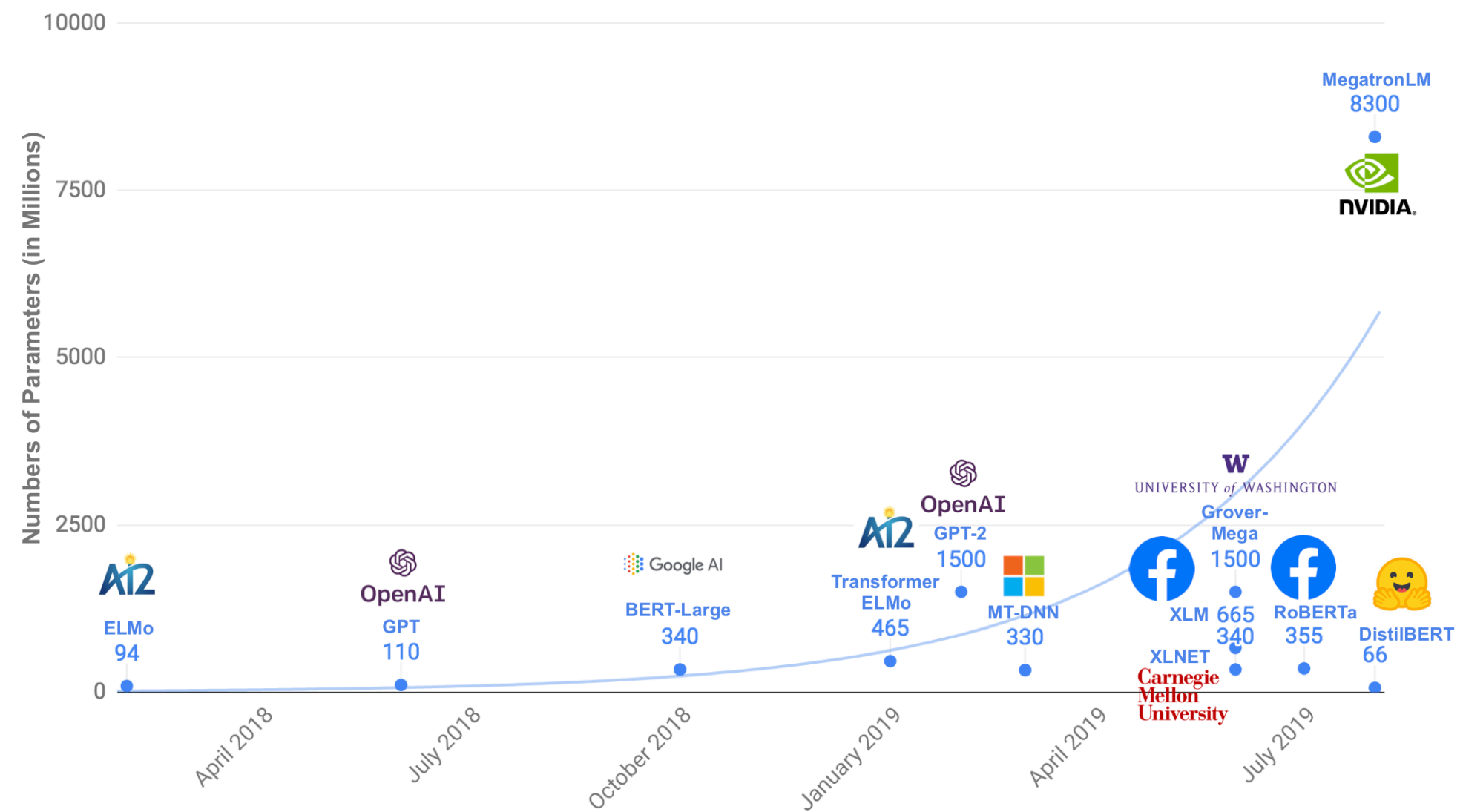# Learning Commonsense Understanding through Language and Vision

## Rowan Zellers

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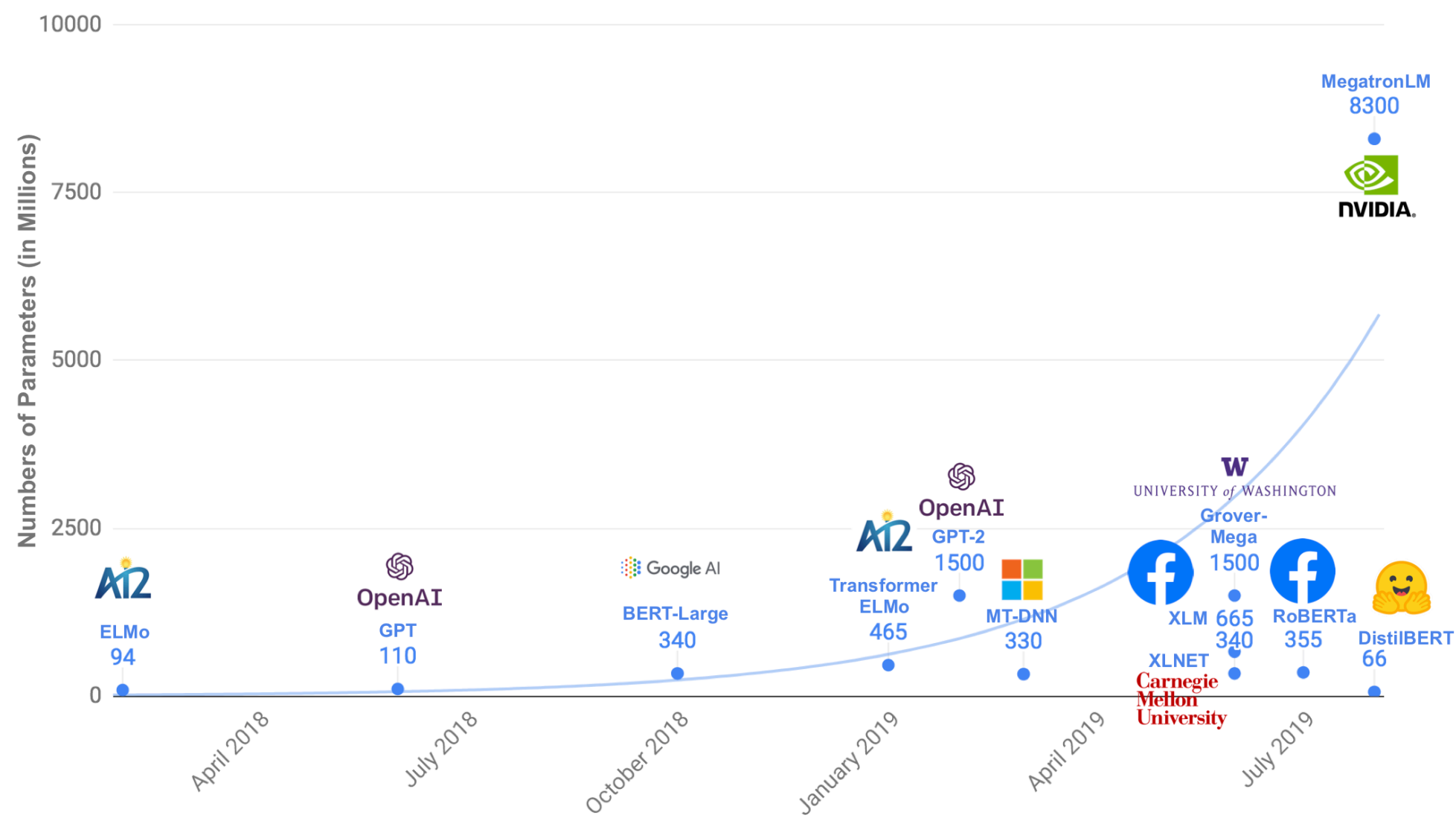University of Washington &
Allen Institute for Artificial Intelligence

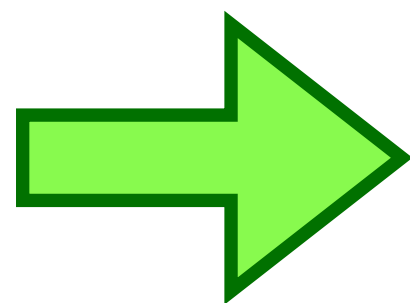# AI progress... vs humans

**exponential increase in model scale**
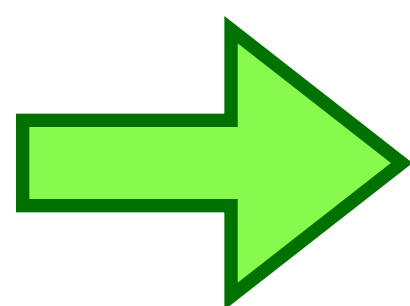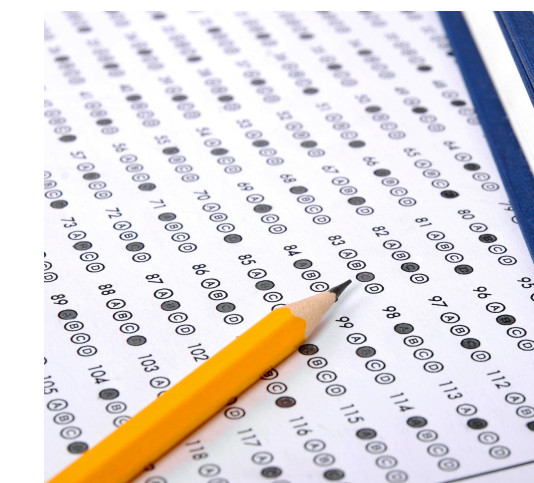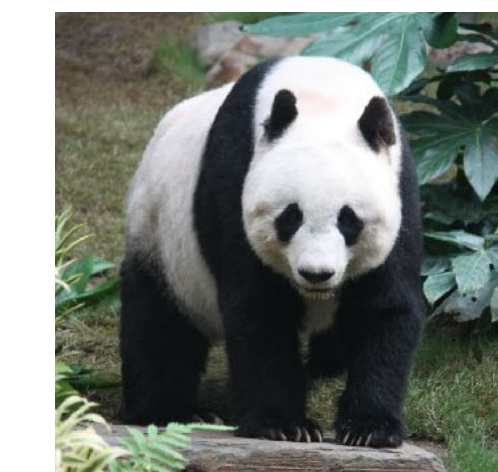
# exponential increase in model scale

**Text: multiple-choice QA**



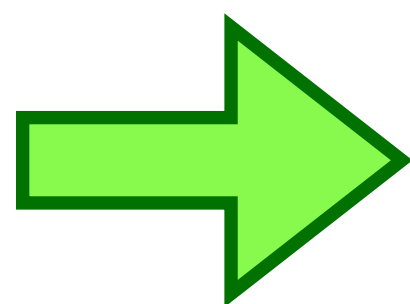*Raffel et al 2019, Brown et al 2020, inter alia*

**Vision: webly supervised classification, detection**



**Panda**

*Radford et al 2021, Kamath et al 2021, inter alia*
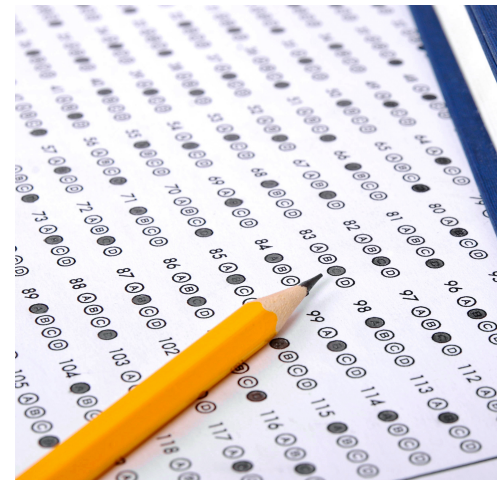
**Vision + Language: learning from captions**



**A train on the tracks**

*Chen et al 2019, Zhang et al 2021, inter alia*
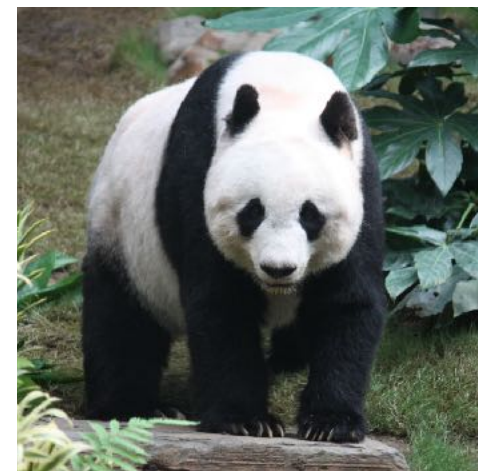
**Text: multiple-choice QA**

**Vision: webly supervised classification, detection**
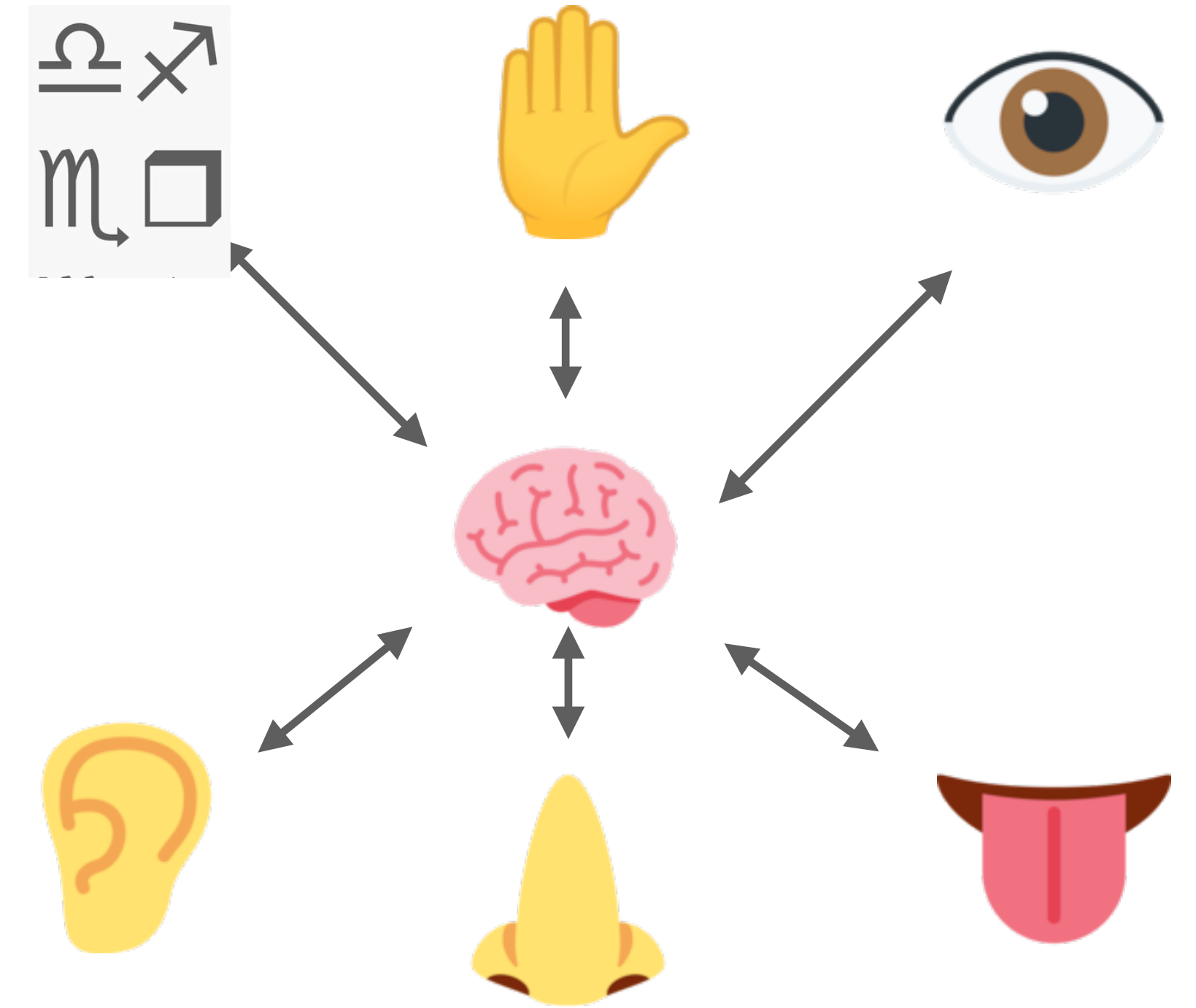
Panda
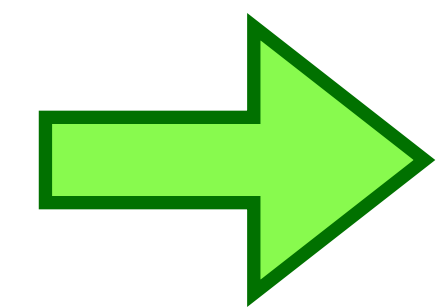
**Vision + Language: learning from captions**

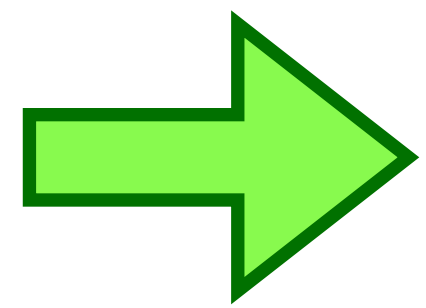A train on the tracks

*Humans....*

- *Integration of many modalities, learned from interaction*
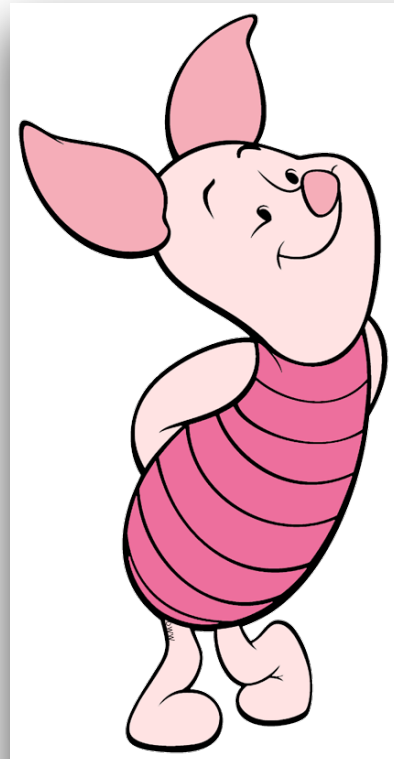
- *Grounded in events, and daily life*

# Today's talk

- *Integration of many modalities, learned from interaction*

- *Grounded in events, and daily life*

●*Integration of many modalities, learned from interaction*

PIGLeT

Language Grounding Through Neuro-Symbolic Interaction in a 3D World
(ACL 2021)

Me

Ari Holtzman

Matthew Peters
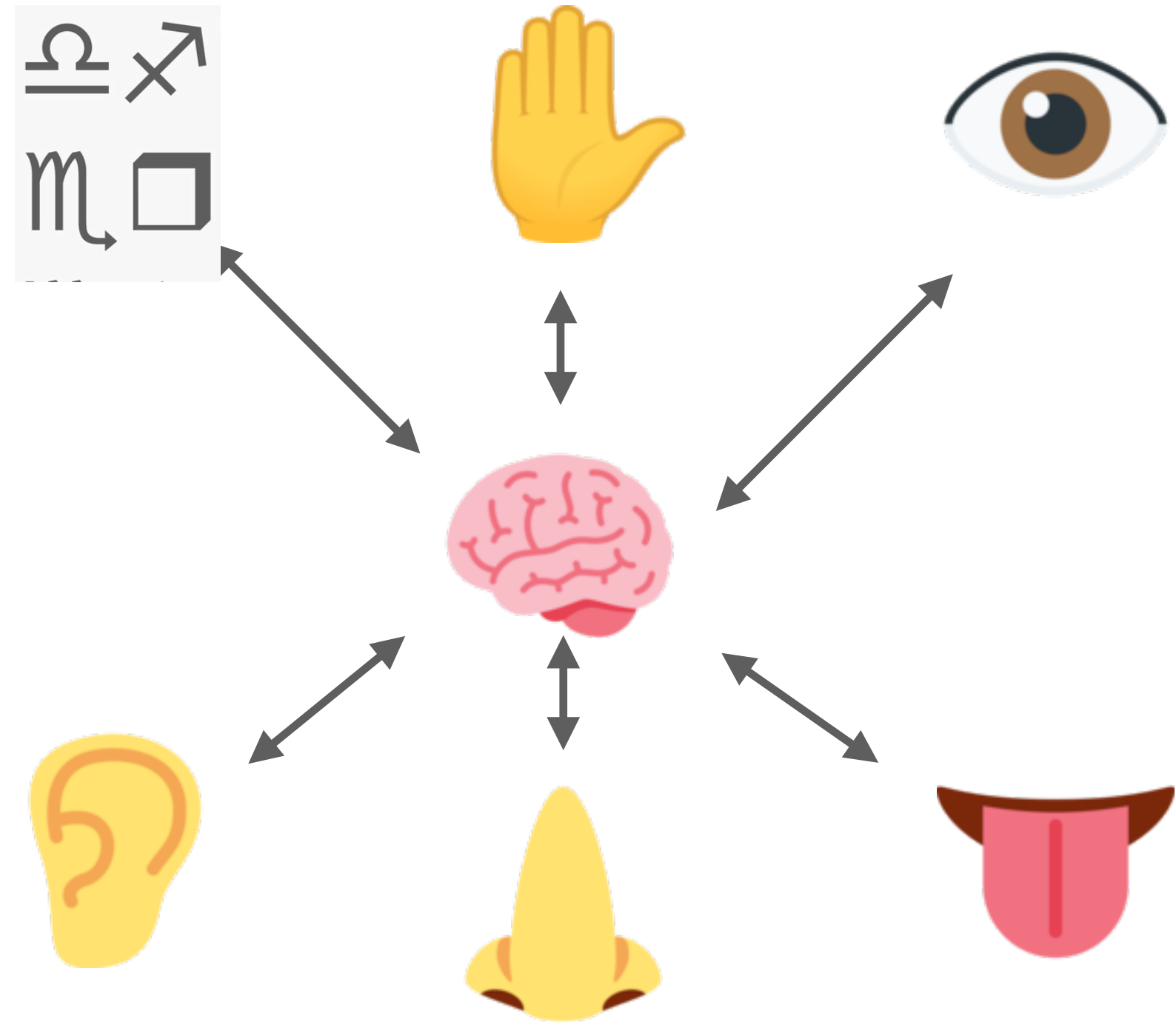
Roozbeh Mottaghi
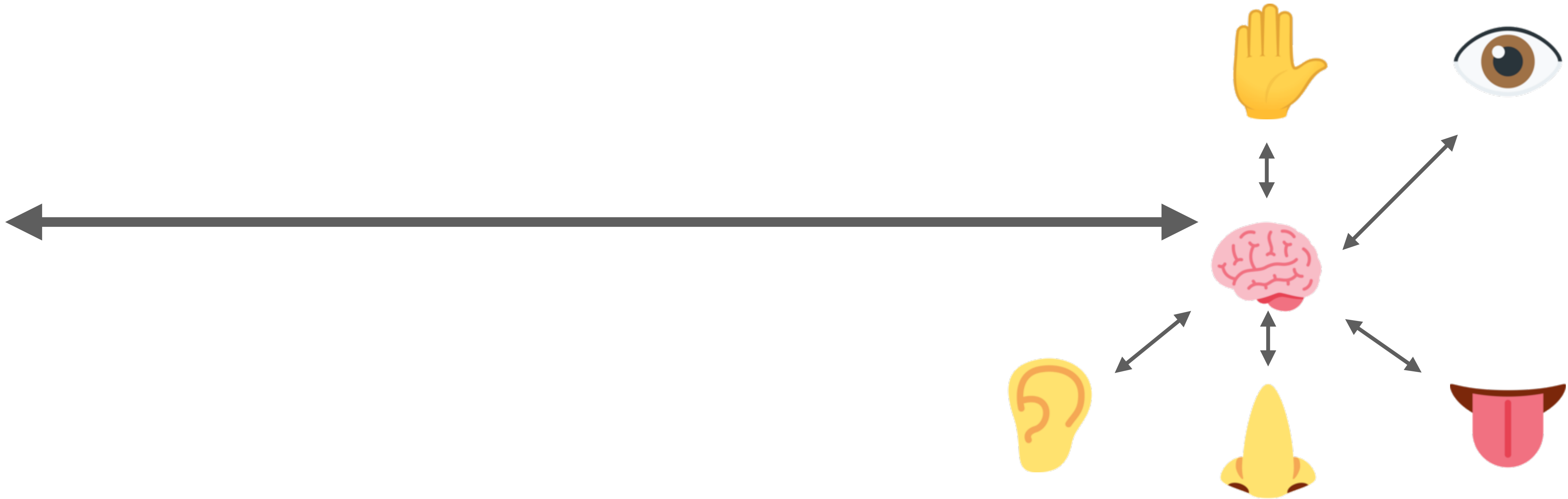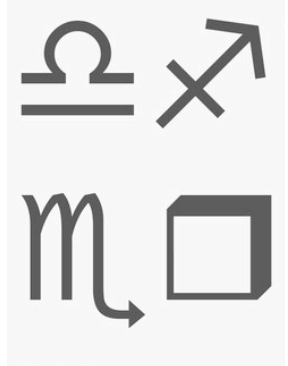
Aniruddha Kembhavi

Ali Farhadi

Yejin Choi

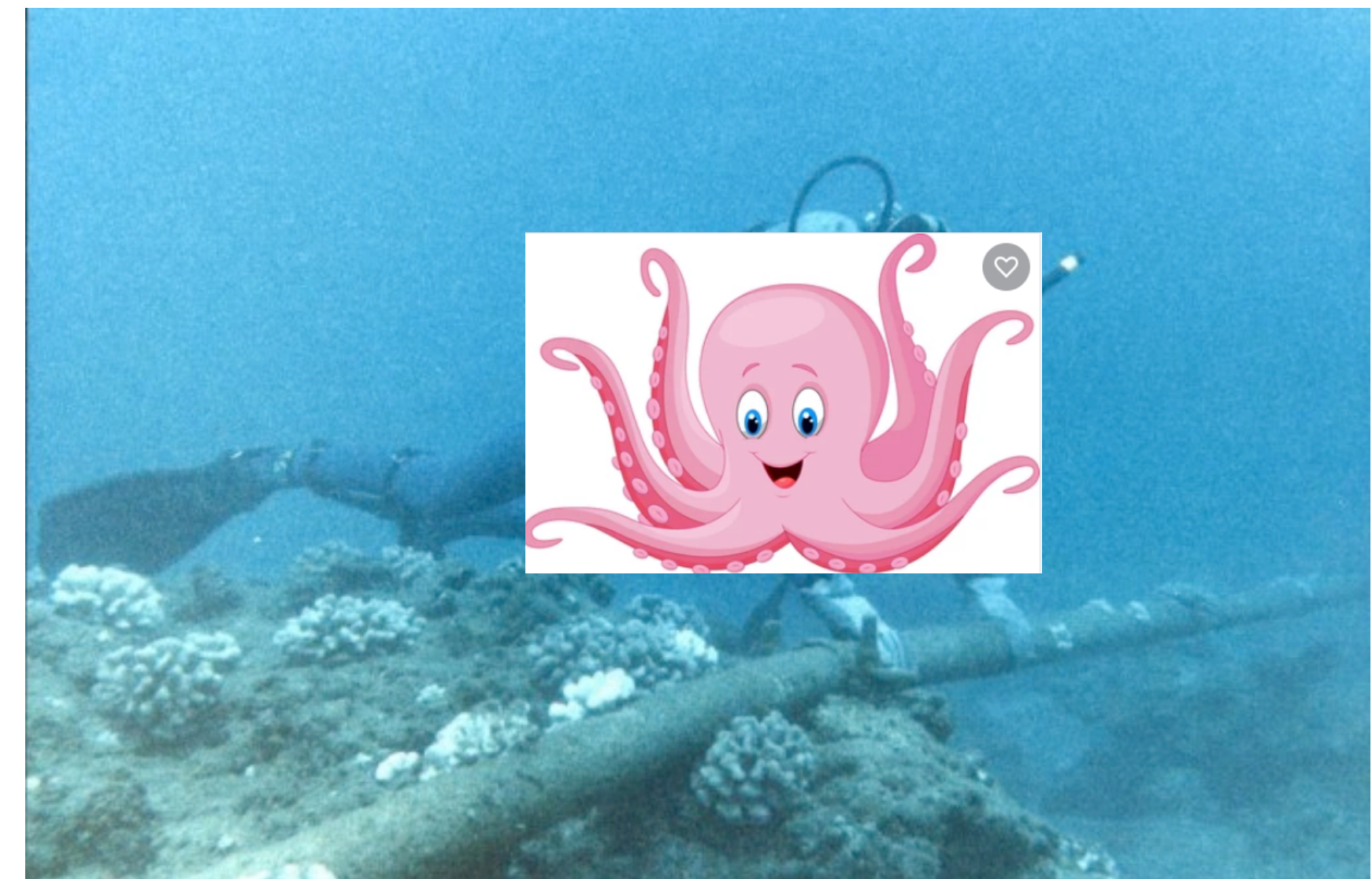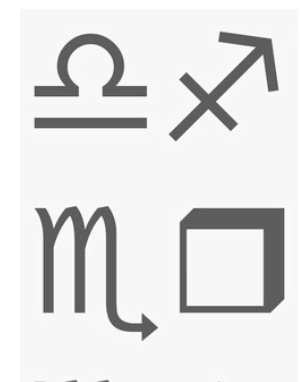# Problem: a gap between *language form* and ***commonsense grounded meaning***



**Written language**
(*symbols*)

**The world**
(*continuous, subjective experience*)

Harnad 1992, inter alia

# 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

| Name: | Mug |
|---|---|
| Temperature: | RoomTemp |
| isBreakable | True |
| isFilledWithLiquid | True |
| ... | |

**grounds**

"I accidentally dropped the mug and it broke"

"I filled up my mug with coffee"

"I'm holding that mug with my hand"

"Careful touching that mug, it's hot"

# Learning from THOR

- THOR: An interactive 3D environment with 20 actions, 125 object types
- Actions are contextual
- Objects have a state (expressed by 42 attributes)

Model: PIGLeT

- We'll predict *explicitly* "what happens next" to an object given an event written out in English

- Or, write an English sentence summarizing the state change.

# PIGLeT: **P**hysical **I**nteractions as **G**rounding for **L**anguag**e** **T**ransformers
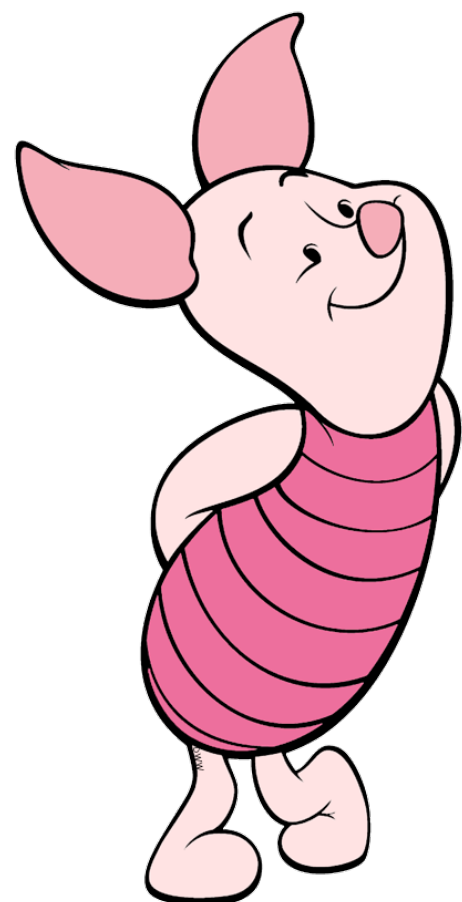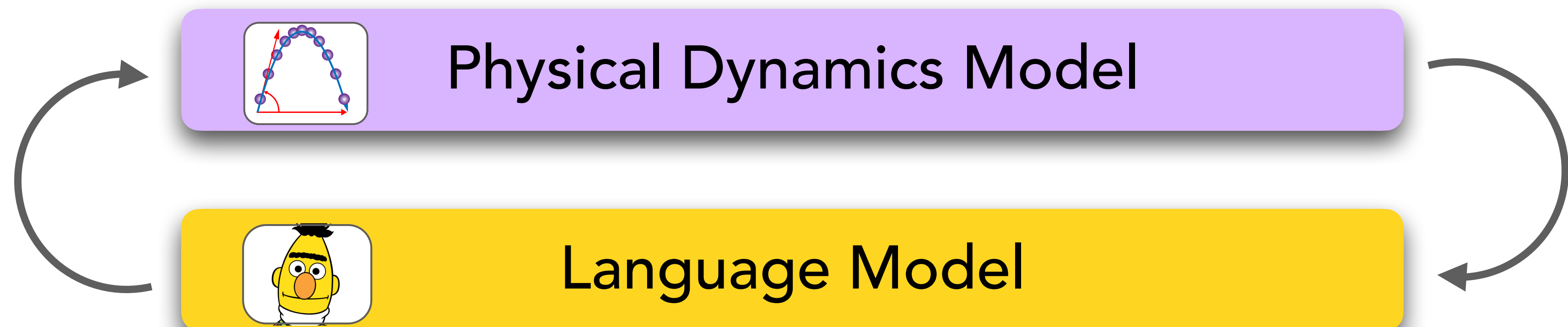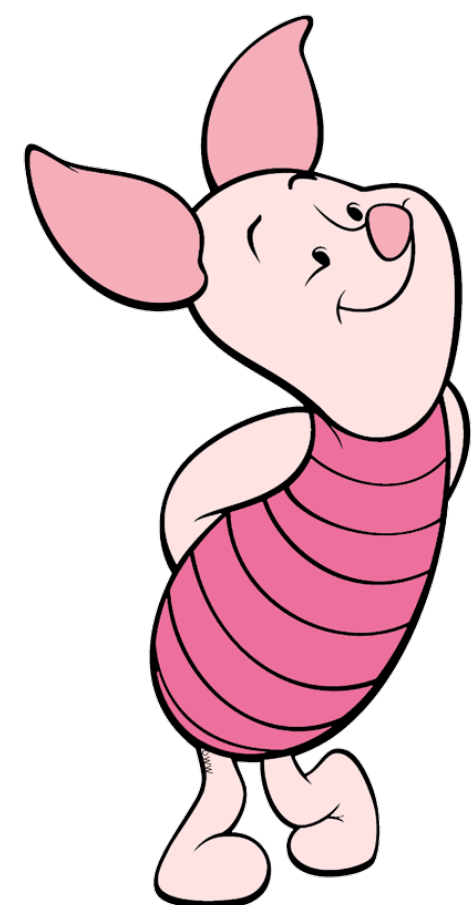


< RotateRight >

brew_coffee:
Try to brew coffee in [Mug] using [CoffeeMachine].

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

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

Physical Dynamics Model

Language Model

# 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 |

...

<heatUp, Pan>

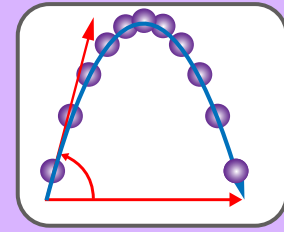| Name: | Egg |
|---|---|
| Temperature: | *Hot* |
| isCooked: | *True* |
| isBroken: | True |

...

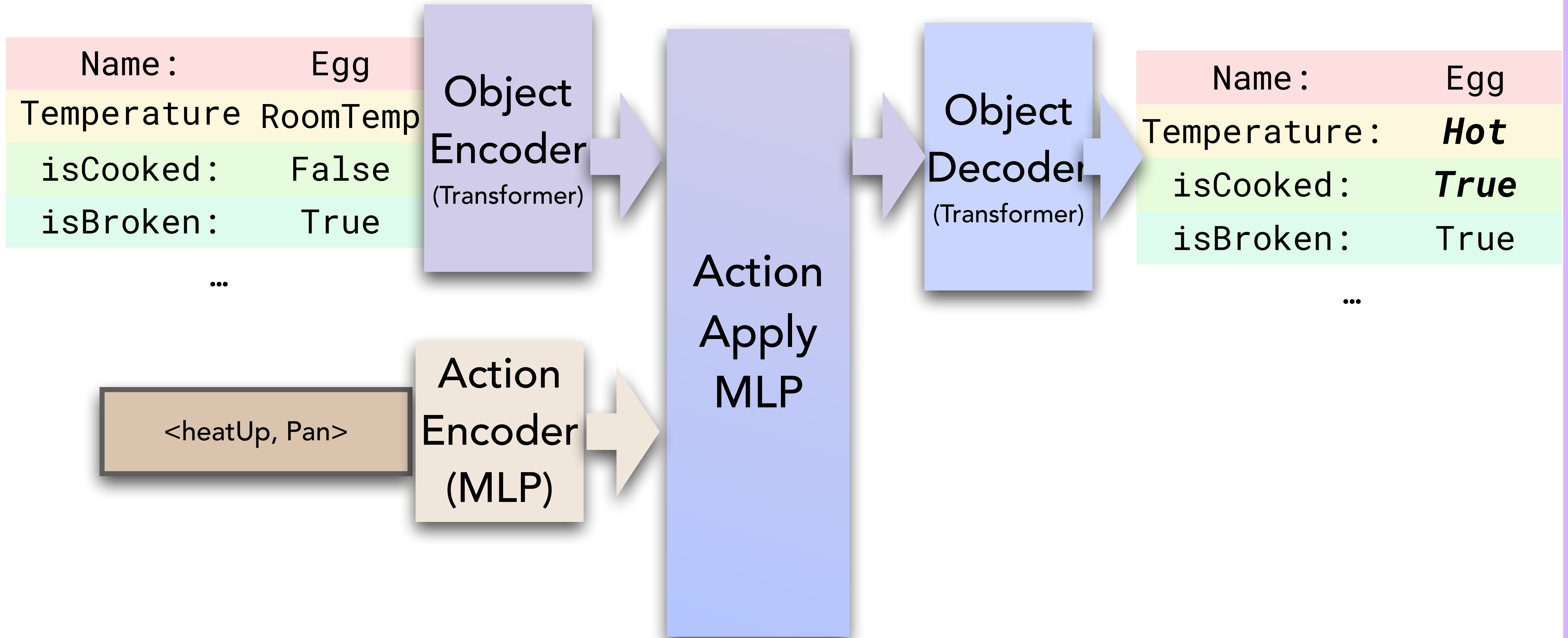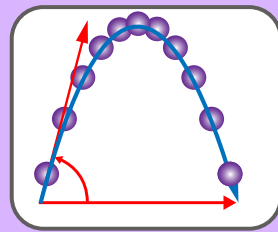**Physical Dynamics Model**

Language Model

Physical Dynamics Model

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

Object Encoder

Object Decoder

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

<heatUp, Pan>

Action Encoder

Action Apply

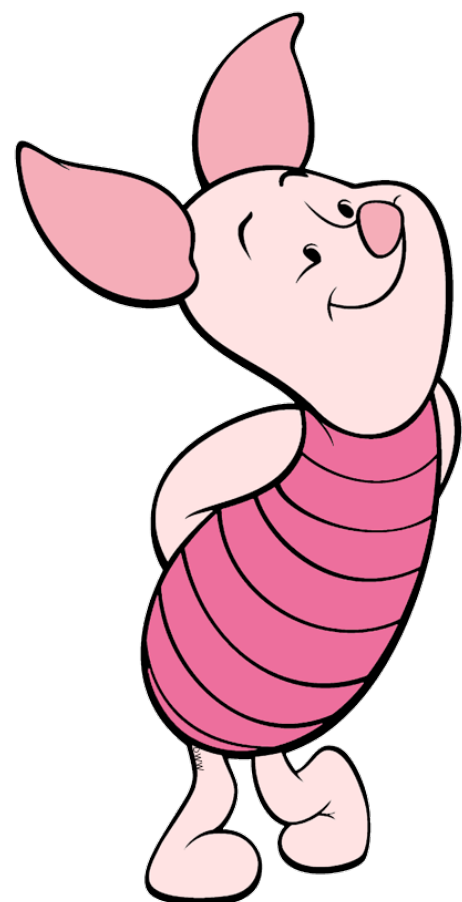The robot heats up the pan.

Language Model

Language Model

The pan becomes hot, and the egg gets cooked.

Model: PIGLeT

- We'll predict *explicitly* "what happens next" to an object given an event written out in English

- Or, write an English sentence summarizing the state change.

- predict "what happens next" to an object given an event written out in English



| Name: | Sink |
|---|---|
| filledWith Liquid | True |

| Name: | Mug |
|---|---|
| filledWith Liquid | True |
| isPickedUp | True |

The robot empties the mug.

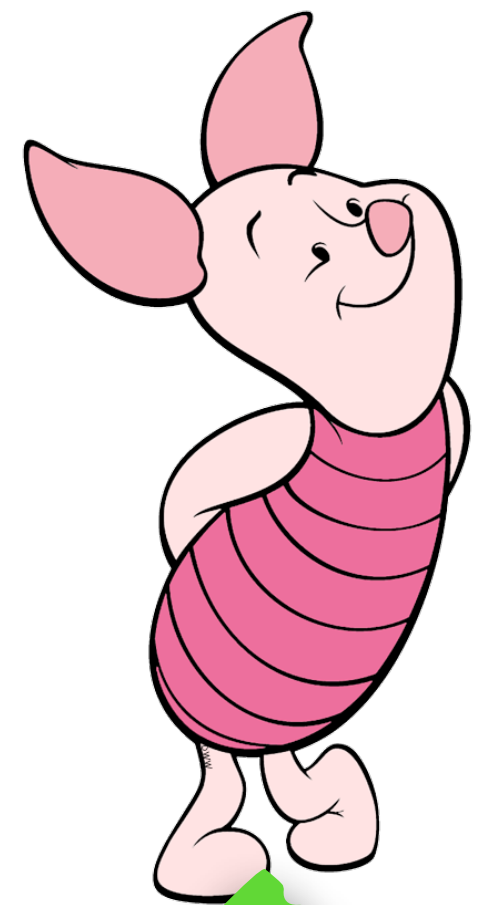| Name: | Sink |
|---|---|
| filledWith Liquid | True |

| Name: | Mug |
|---|---|
| filledWith Liquid | *False* |
| isPickedUp | True |

*Evaluation: Accuracy (of getting all attributes right)*

# Results (accuracy of getting all attributes right)

PIGLET —- learned through physical experiences in the THOR 3D env
—- outperforms a 100x larger model (T5-11B) by over 10%

Accuracy (%)

| | |
|---|---|
| 100 | |
| 75 | |
| 50 | |
| 25 | |
| 0 | |

No change — 26
GPT3 — 22
T5 — 64
Ours — 81

# Qualitative Example

t

| Name: | Sink |
|---|---|
| filledWith Liquid | True |

| Name: | Mug |
|---|---|
| filledWith Liquid | True |
| isPickedUp | True |

The robot empties the mug.

| Name: | Sink |
|---|---|
| filledWith Liquid | True |

✓

| Name: | Mug |
|---|---|
| filledWith Liquid | *False* |
| isPickedUp | True |

✓

# Qualitative Example



Name: Sink

filledWith Liquid    True

Name: Mug

filledWith Liquid    True

isPickedUp    True

The robot empties the mug.

Name: Sink

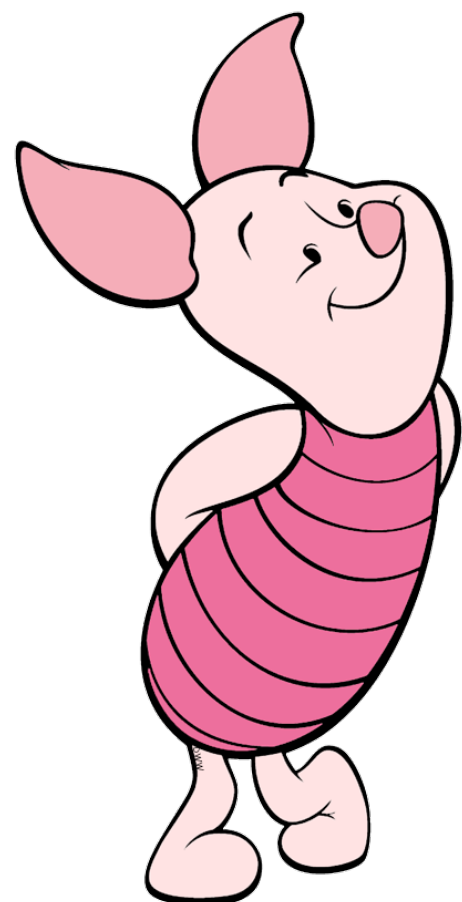filledWith Liquid    True ✅

Name: Sink

filledWith Liquid    False ❌

T5-11B, through text, learns "emptying liquid from an object" makes all objects in the room empty 🤔

Model: PIGLeT

- We'll predict *explicitly* "what happens next" to an object given an event written out in English

- Or, write an English sentence summarizing the state change.

- and summarize this prediction in English

| Name: | Sink |
|---|---|
| filledWith Liquid | True |

| Name: | Mug |
|---|---|
| filledWith Liquid | True |
| isPickedUp | True |

```
<empty,
Mug>
```

*The mug is no longer filled with water.*
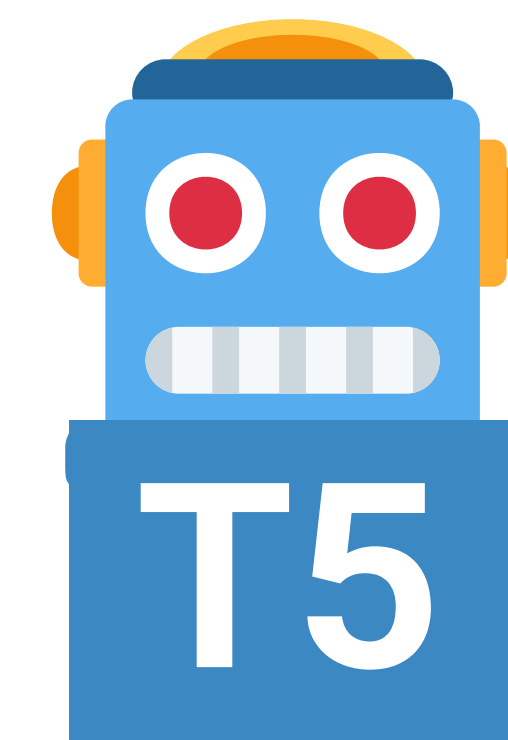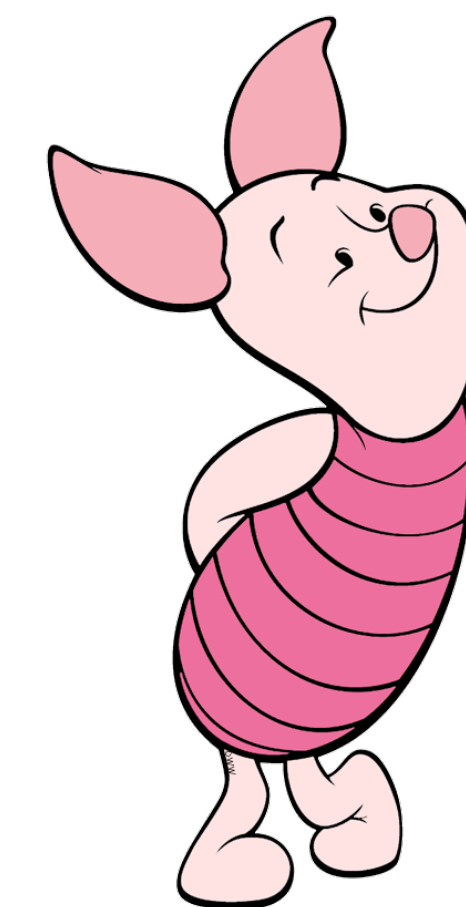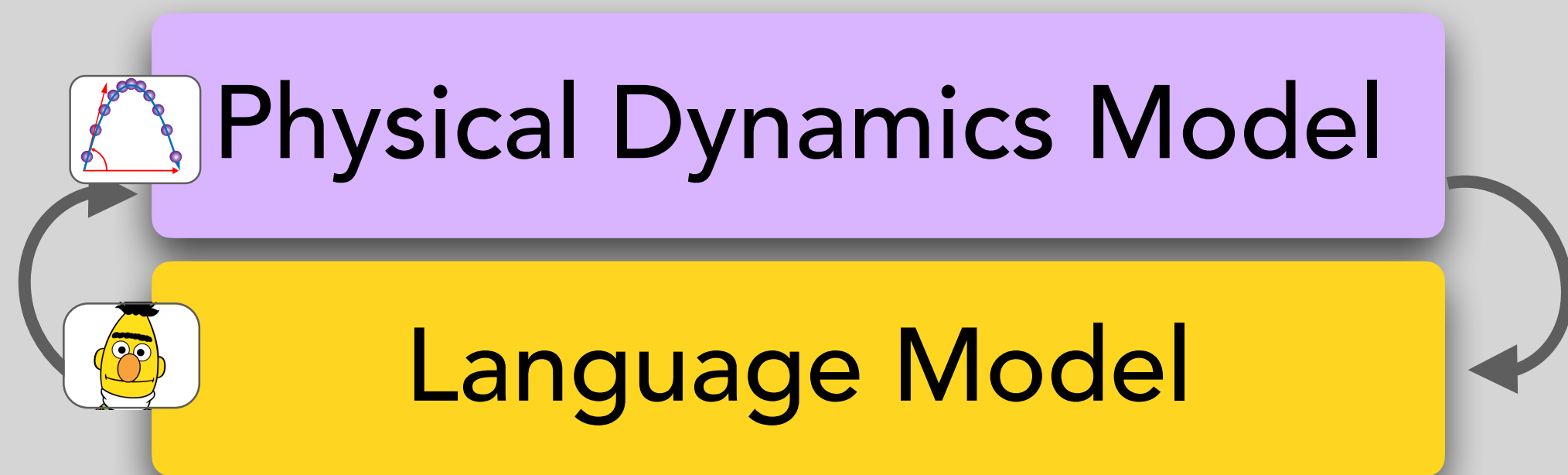
*Evaluation: human, BLEU, BERTScore*

# PIGLeT's generations

# PIGLeT's generations



| Name: | Sink |
|---|---|
| filledWith Liquid | True |

| Name: | Mug |
|---|---|
| filledWith Liquid | True |
| isPickedUp | True |

<empty, Mug>

*The mug is now empty.*

*The sink is now empty.*

# _PiGLET_ vs text-only learning

**Physical Dynamics Model**

**Language Model**

GPT-N

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

Learn a lightweight factorized world model
for predicting _what might happen next_

A single, heavyweight, entangled model

**Paper-only bonus!!**

Can generalize to new concepts like
"Dax" without words

Limited generalization to new concepts

# Today's talk

➡ ● *Integration of many modalities, learned from interaction*

● *Grounded in events, and daily life*

●*Grounded in events, and daily life*

# MERLOT: Multimodal Neural Script Knowledge Models

arxiv 2021

Rowan
Zellers*

Ximing
Lu*
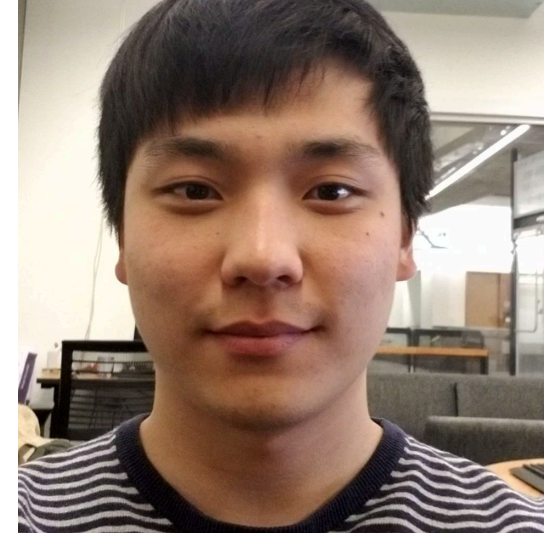
Jack
Hessel*

Youngjae
Yu

Jae Sung
(James) Park

Jize
Cao

Ali
Farhadi

Yejin
Choi

Previously on my slide deck...

**Why is he pointing?**

<object: syrup bottle>

scene: a diner

# *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 Pro-
metheus*, 1818)

## Abstract

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

zation of knowledge can result in a real under-
standing system in the not too distant future.  We
expect that programs based on the theory we out-
line 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 under-
stander.  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



```
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


    ...
```

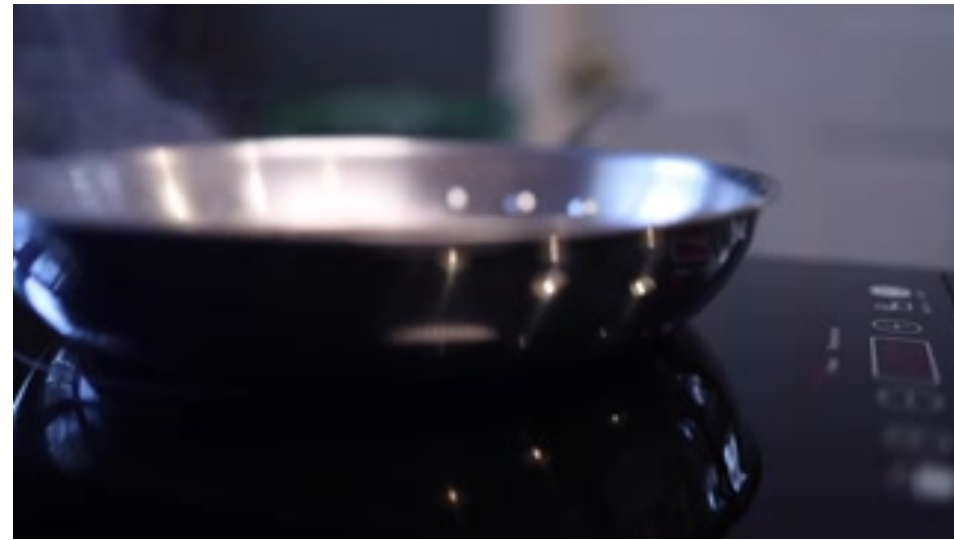# *Multimodal   Script Knowledge*

*(Neural)*

# Multimodal   Script Knowledge

## (Neural)

From 6M youtube videos, we'll learn:

# From 6M youtube videos, we'll learn:



**Recognition-level Knowledge**
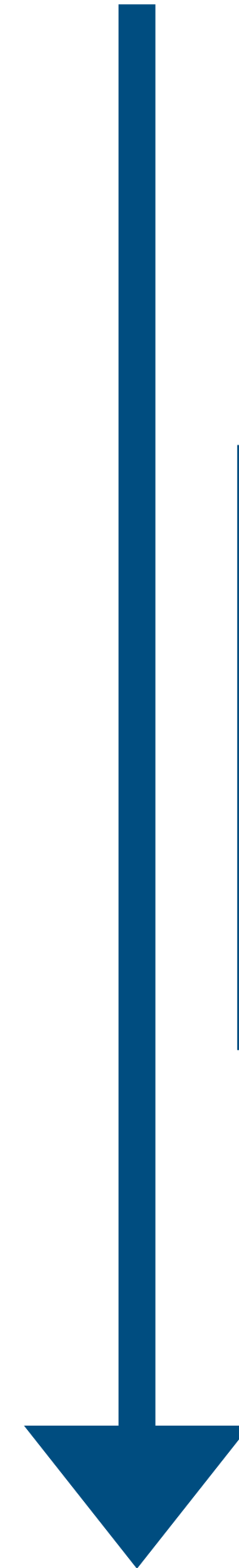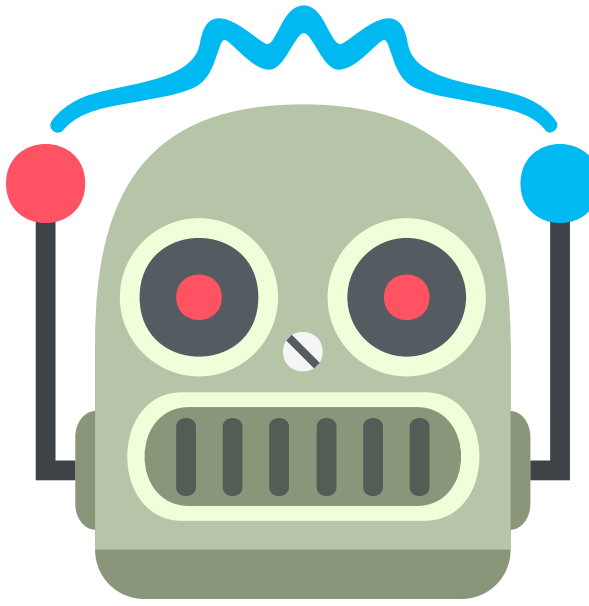
person

pan    Burner

stopwatch
water pitcher

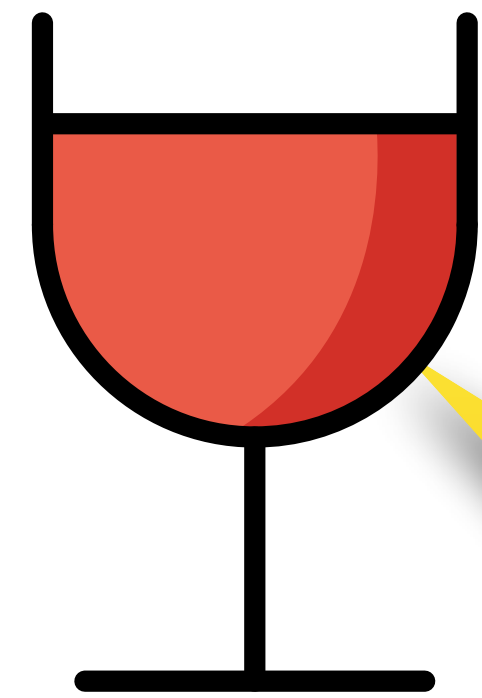thermometer

**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

**M**ultimodal **E**vent **R**epresentation **L**earning **O**ver **T**ime

The result:

- Trained fully from scratch, we get...

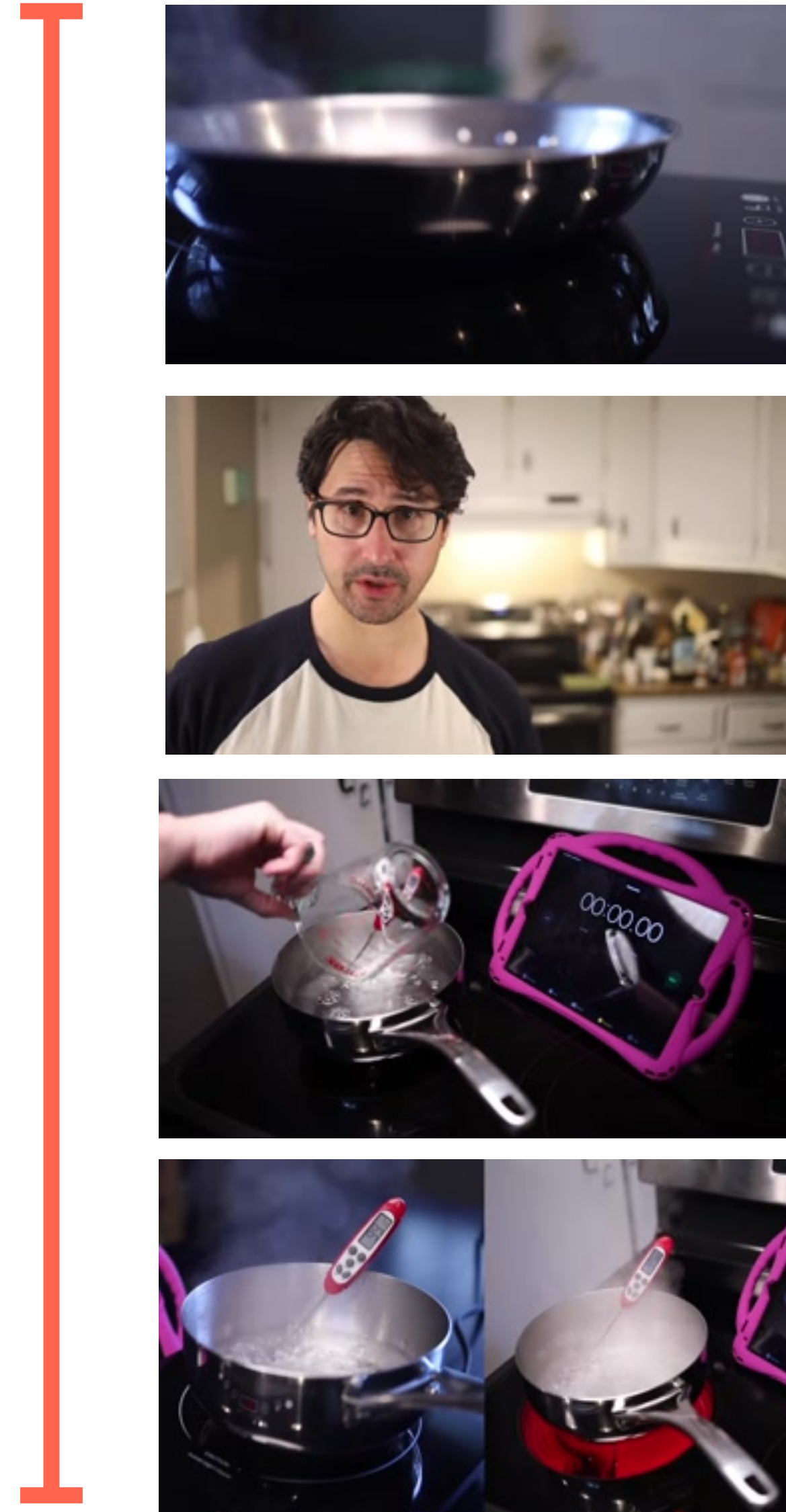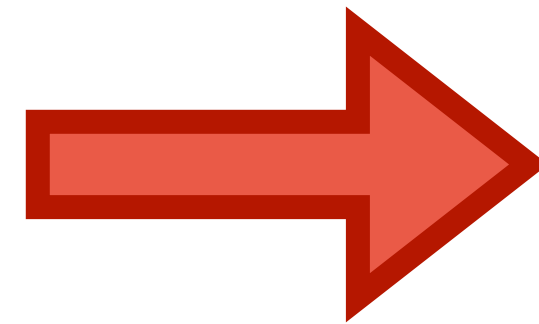- zero-shot temporal commonsense,

- Fine-tuned SOTA on 13 tasks

**M**ultimodal **E**vent **R**epresentation **L**earning **O**ver **T**ime

- Pretraining Strategy + Objectives

- Evaluation

# Setup: Videos and Transcripts



"In this video I'm ..."

# 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...."
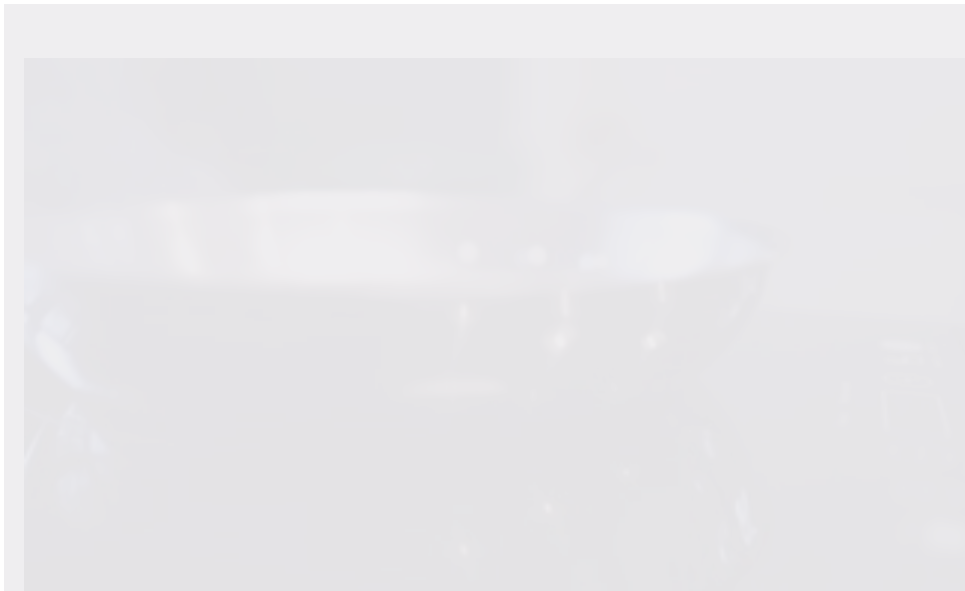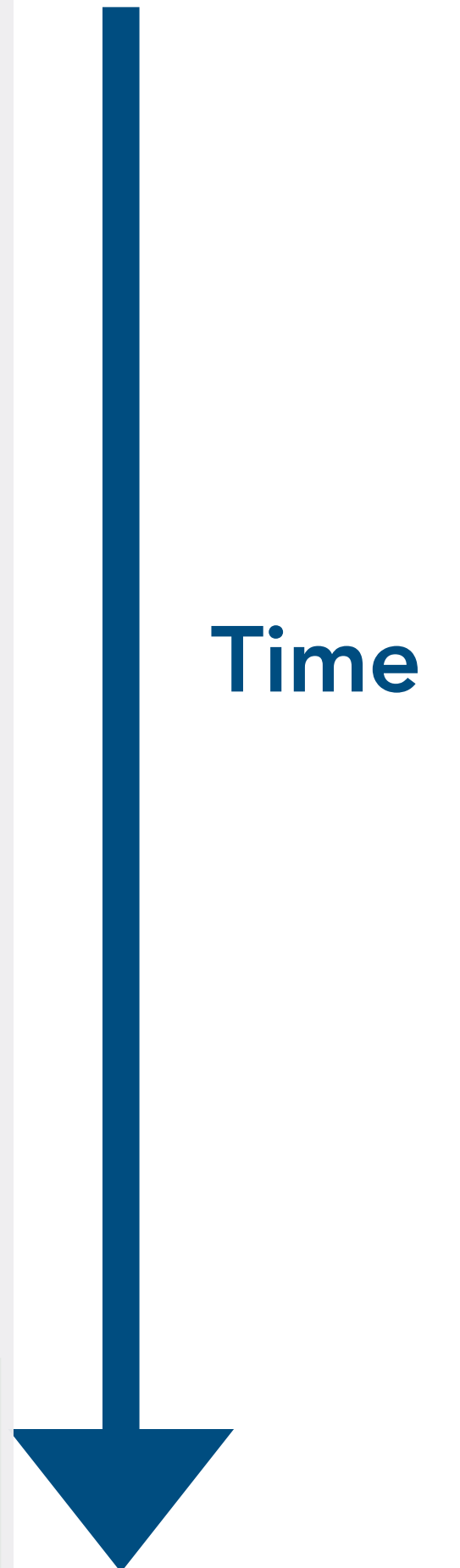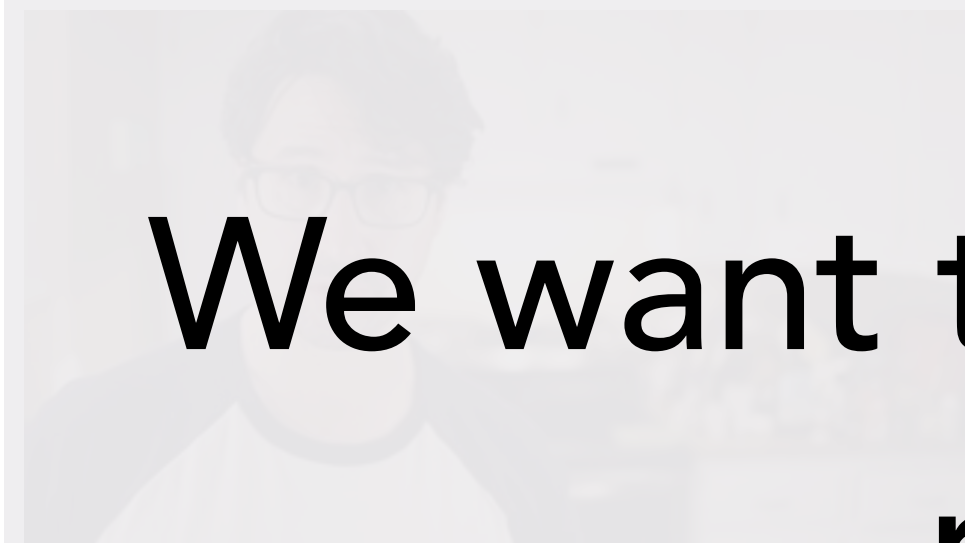
"In goes the cold water..."

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

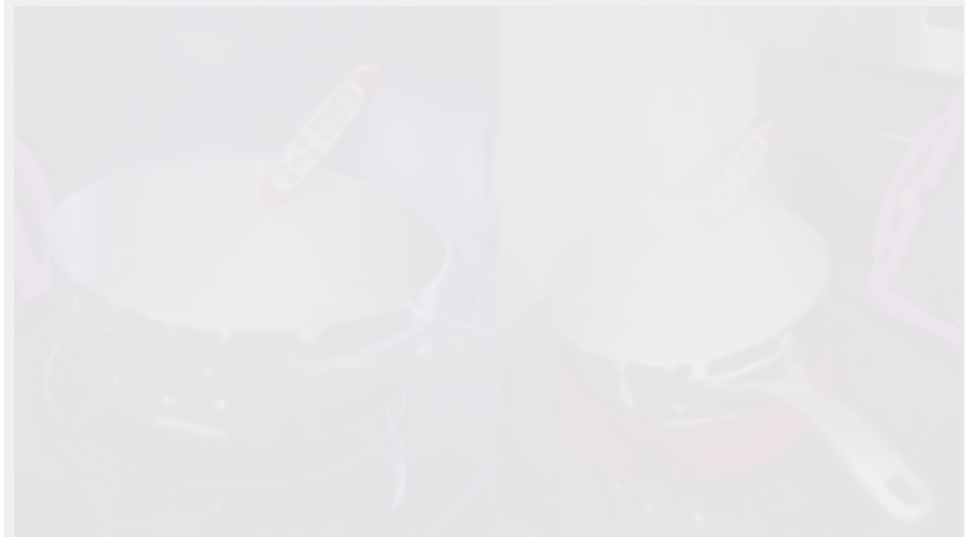Time

"I'm going to compare electric and induction stoves…"

"I'll use a stopwatch to time how fast my electric stove boils water…"

We want to use this (dynamic) data to first learn recognition-level reasoning…
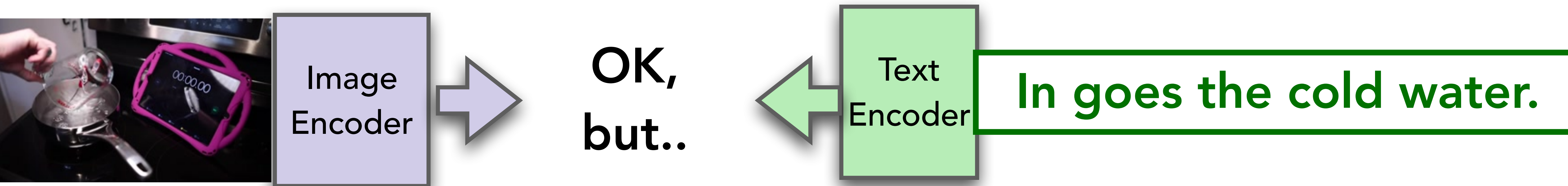without training on manually labeled data

"It took 4 and a half minutes to reach full boil…"
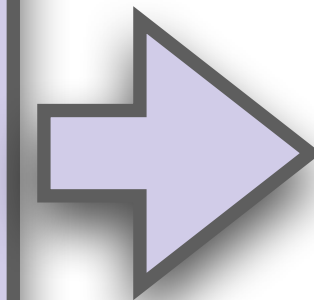
Time

# Recognition-level learning



Image Encoder

OK, but..

Text Encoder

**In goes the cold water.**

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

# Recognition-level learning



"I'll use a stopwatch to time how fast my electric stove boils water."

OK, but..

In goes the cold water.

Image Encoder

Text Encoder

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

# Recognition-level learning



[CLS] "I'm going to compare electric and induction stoves."

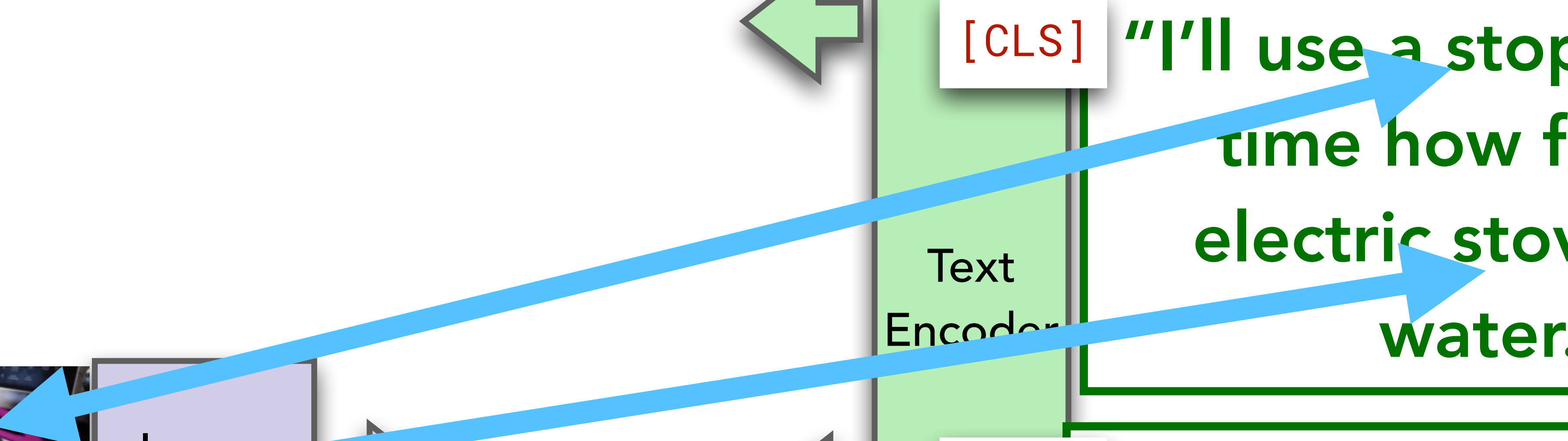[CLS] "I'll use a stopwatch to time how fast my electric stove boils water."
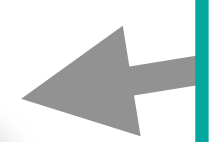
[CLS] In goes the cold water.

[CLS] "It took 4 and a half minutes to reach full

Text Encoder

Image Encoder
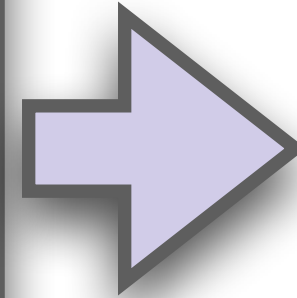
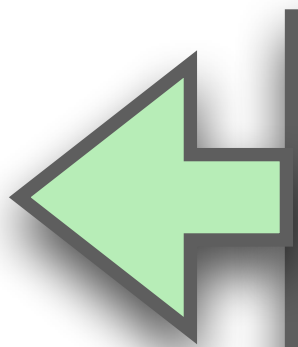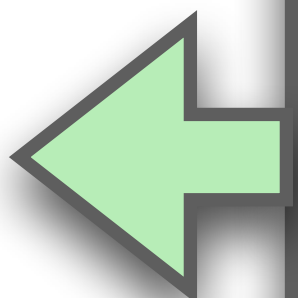Better!

# Recognition-level learning



Objective 1: maximize similarity between contextualized language and individual frames
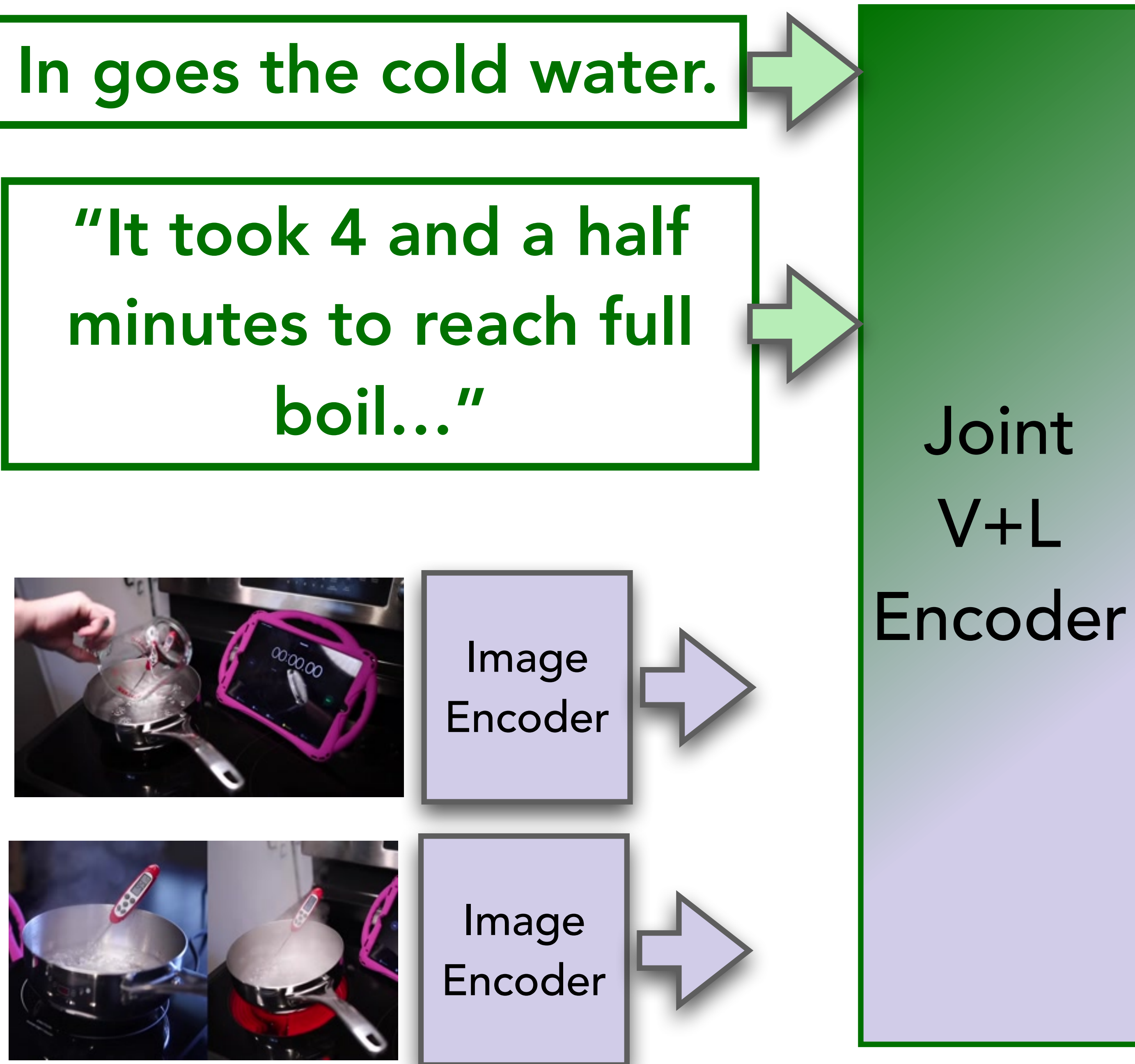
Image Encoder

Image Encoder

Text Encoder

Objective1:
Contextual Frame-
Text Matching

Objective 2:
Mask LM

Objective 3:
Unshuffle frames

Using a 12-layer 'base' Transformer, train everything E2E on 6M videos

Image Encoder

Joint V+L Encoder

Text Encoder

# **M**ultimodal **E**vent **R**epresentation **L**earning **O**ver **T**ime

- Pretraining Strategy + Objectives

- 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.



(3)

(1)

(4)

(5)

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. |

(1) (2) (3) (4) (5)

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
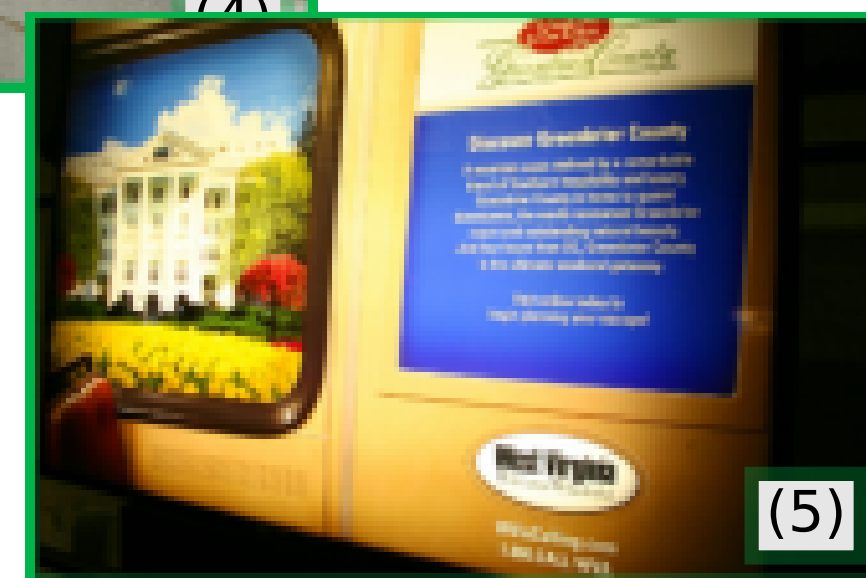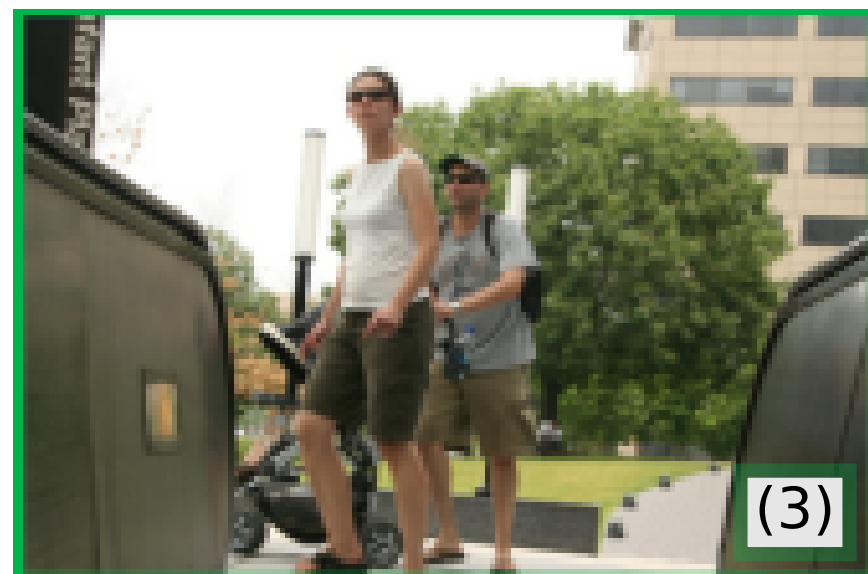
The old man was riding the escalator.

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.

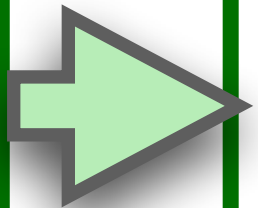Visual Coref over time!

 (1)
 (2)
 (3)
 (4)
 (5)

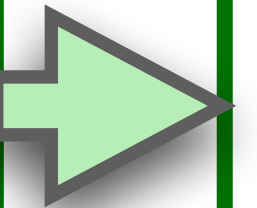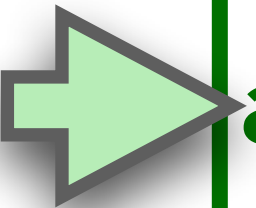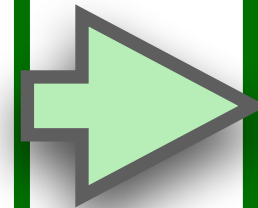The old man was riding the escalator. → He was almost to the top. → His kids were already at the top. → At the top was a train station. → They then got on the train.



CLIP   (Radford et al 2021)

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



0.49     0.63     0.74

CLIP     UNITER

(Chen et al 2019)

# 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. |


(3)


(2)


(1)


(5)


(4)

# 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

■ MERLOT ■ Prev SOTA

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



MERLOT: 65.1, UNITER: 58.2, VILLA: 60.6, ERNIE-ViL: 62.1

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

# Analysis (on TVQA+)



Legend: Ours (red), No contrastive V+L Loss (green)

- 80
- 76
- 75
- 70
- 68
- 65
- 60

# Analysis (on TVQA+)

# Analysis (on TVQA+)

■ Ours ■ Trained on HowTo100M ■ Trained on Captions

# Performance increases with # epochs

# Discussion

- **Simulation pros:**
  - **Learning to *act*, not just see/write**
  - **Future work: Models guiding the training loop, maybe based on curiosity**
- **Cons:**
  - **Limited vocabulary in simulation**
  - **Hard to learn human behavior**

- **Web video pros:**
  - **Super wide vocabulary**
  - **Learning human norms, behavior, events**
- **Cons:**
  - **Can't participate in the video**
  - **Privacy**

# Privacy (and other negative societal implications of training on multimodal Web Data)

- Things we did for MERLOT

  - data curation focused on big channels, not randos

    - on a public platform that people expect is public (Kang et al 2015)

    - … at a scale so that people are "in public without being public" (Marwick and boyd 2011)

  - distributing links, not the videos, for the "right to be forgotten"

  - Encouraging future work into these foundation models — not advocating for product use right now

# Privacy (and other negative societal implications of training on multimodal Web Data)

- data curation focused on big channels, not randos

  - on a public platform that people expect is public (Kang et al 2015)

  - … at a scale so that people are "in public without being public" (Marwick and boyd 2011)

- distributing links, not the videos, for the "right to be forgotten"

- Encouraging future work into these foundation models — not advocating for product use right now

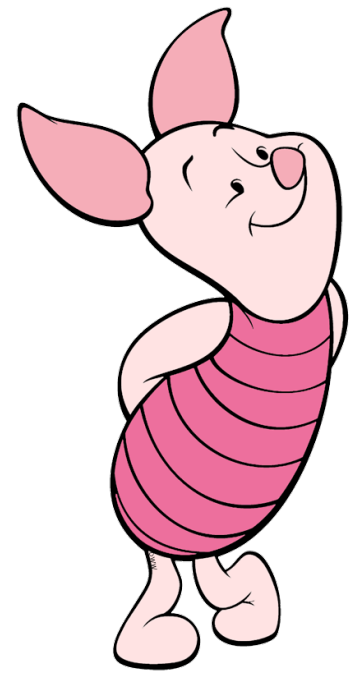Lots of local news… which has bias issues (Gilliam Jr et al 1996)

Inherent bias with training on data that encodes a "view from nowhere" (Haraway et al 1988, Waseem et al 2021)

… bias that is amplified by culture and the "YouTube Algorithm" (Strangelove et al 2020)

# Privacy (and other negative societal implications of training on multimodal Web Data)

- Future work: studying privacy, bias, and dual use,

- … exploring possibly a mix of technical and non-technical fixes here

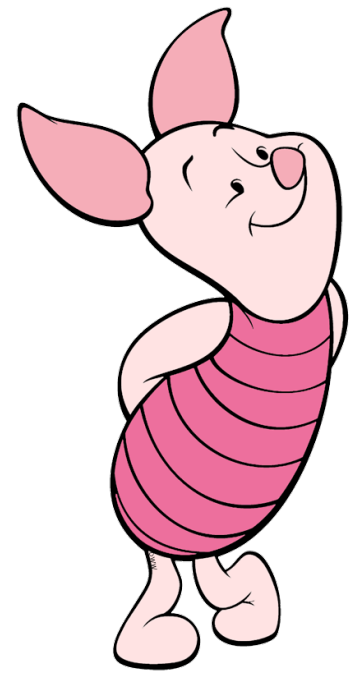- Hopefully the beginning, not the end, of this key conversation

# Questions?

# Thanks!!