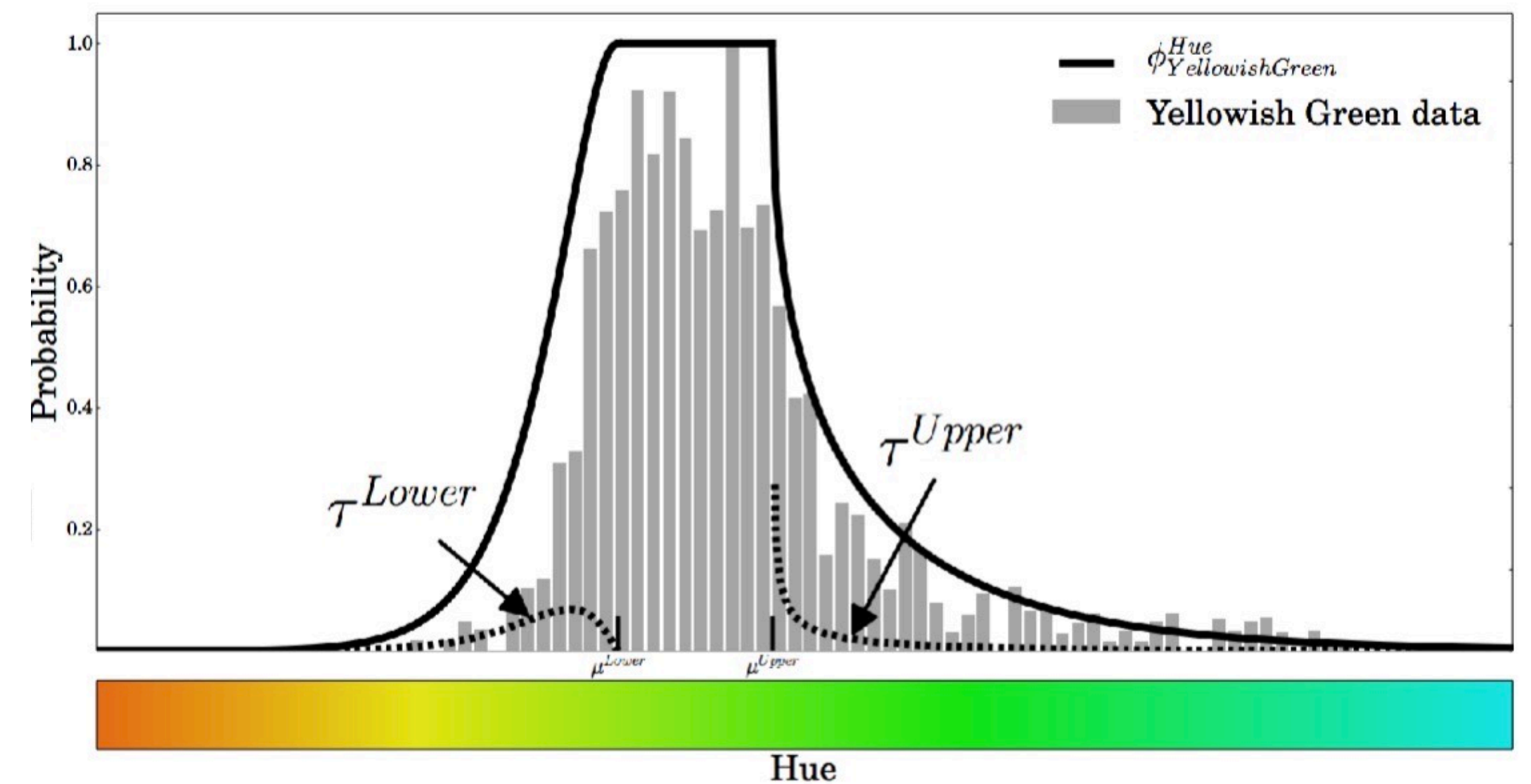# CS388: Natural Language Processing

## Lecture 22: Multimodality, Language Grounding

Greg Durrett



McMahan and Stone (2015)

TEXAS
The University of Texas at Austin

# Announcements

- FP due April 28

- Presentations on last two class days

# Today's Lecture

▸ Classic grounding

▸ Multimodality

▸ Language and vision models

▸ Language and manipulation
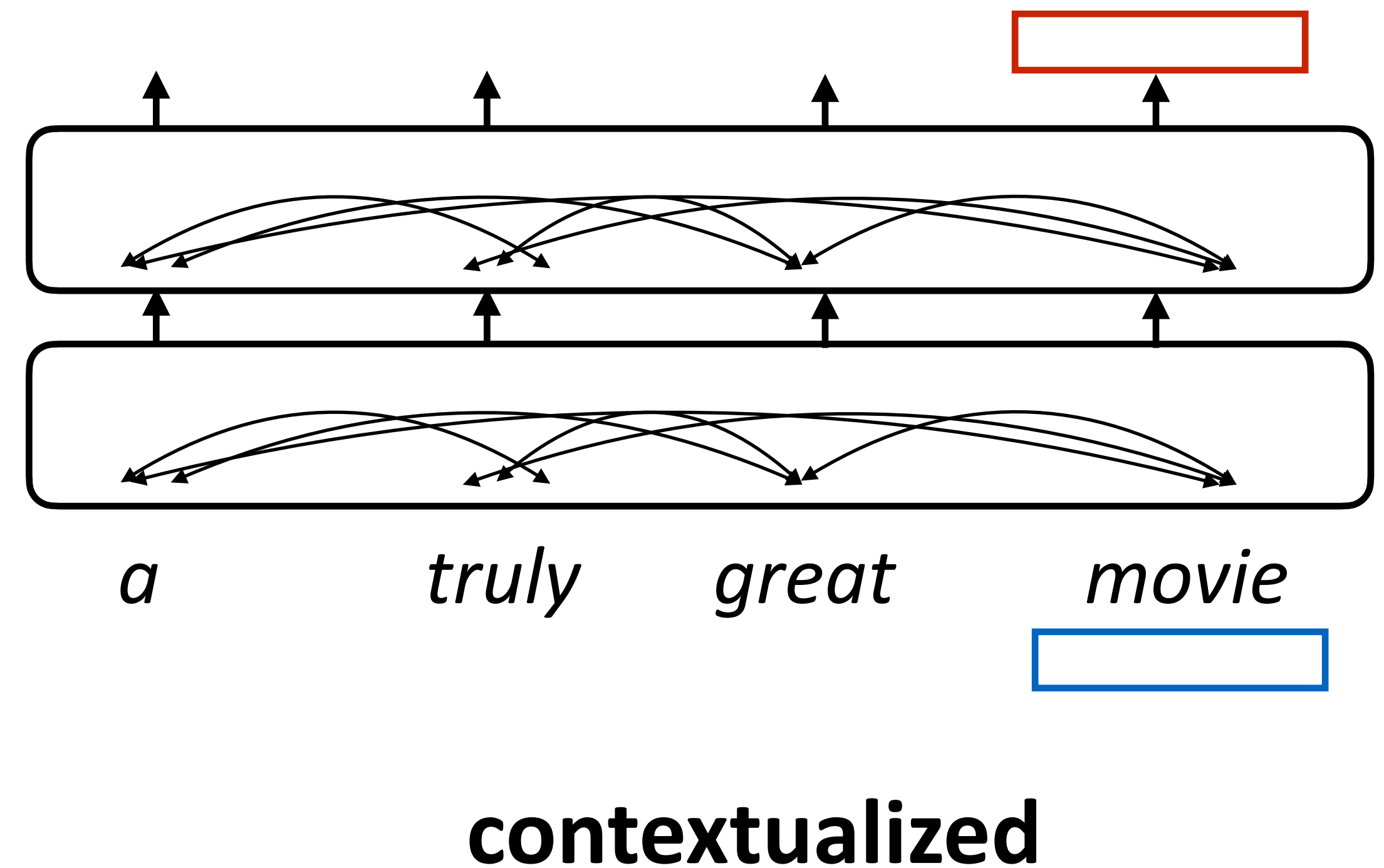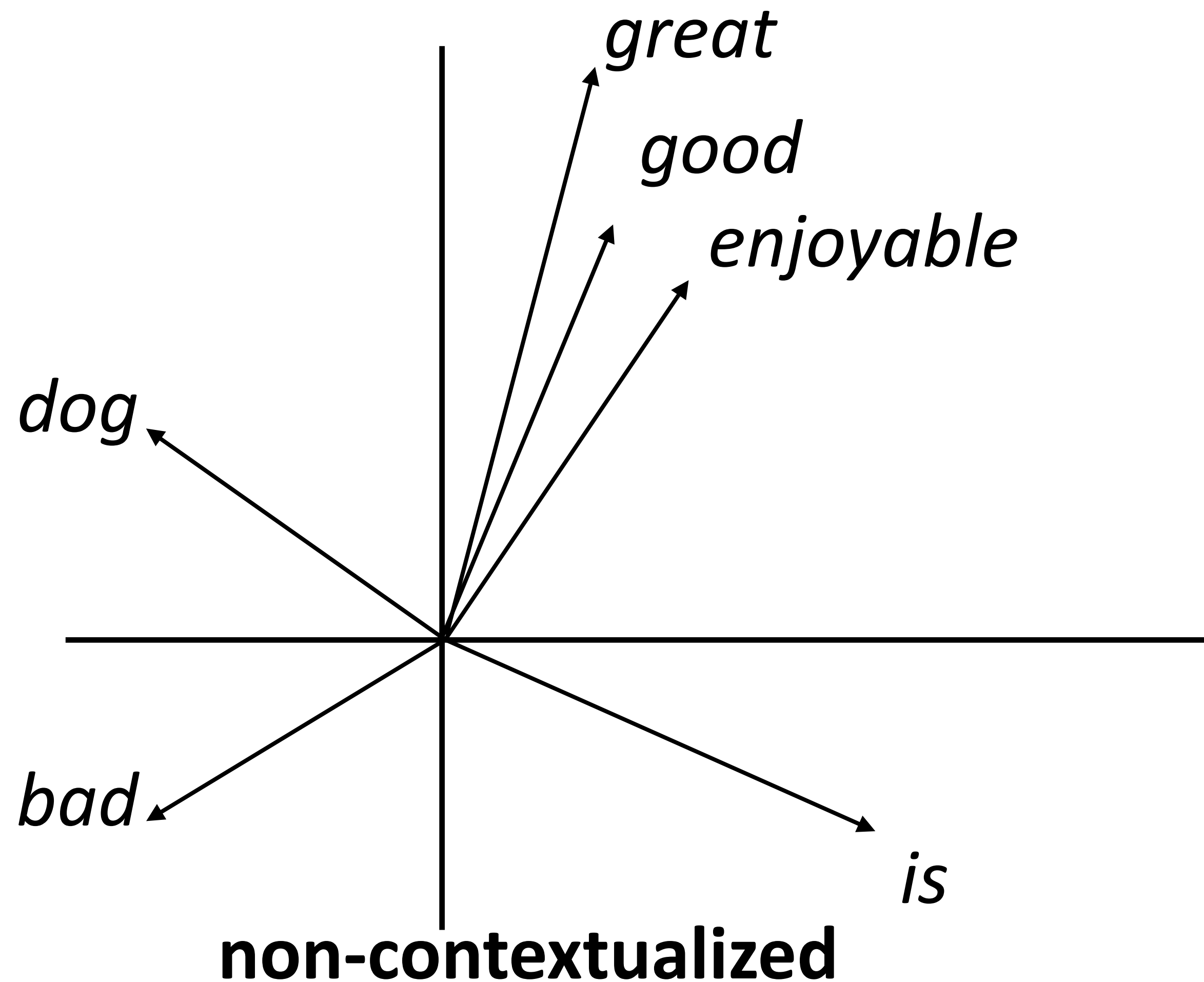
# Classic Grounding

# Language Grounding

‣ How do we represent language in our models?

‣ How did we learn these representations? What do the vectors "mean"?



*great*

*good*

*enjoyable*

*dog*

*bad*

*is*

**non-contextualized**

*a*  *truly*  *great*  *movie*

**contextualized**

# Language Grounding

‣ Harnad defines a "symbol system": we have symbols (e.g., strings) manipulated on the basis of rules, and these symbols ultimately have "semantic interpretation"

  ‣ "Fodor (1980) and Pylyshyn (1980, 1984)…emphasize that the symbolic level (for them, the mental level) is a natural functional level of its own, with ruleful regularities that are independent of their specific physical realizations"

‣ Harnad challenges the idea that fully symbolic approaches can work well.

‣ Argues that "horse" is something that should be understood bottom-up through grounding. "Zebra" = "horse" + "stripes" could emerge this way, but he claims it cannot through a top-down symbolic system

‣ What does it mean to "understand" the symbols that get manipulated?

Harnad (1990) *The Symbol Grounding Problem*

# Searle's Chinese Room

‣ Suppose we have someone in a room with a long list of rules, dictionaries, etc. for how to translate Chinese into English. A Chinese string is passed into the room and an English string comes out. The person is not a speaker of Chinese, but merely follows the rules and looks things up in the dictionaries to produce the translation.

‣ Does the person understand Chinese? Does the room? (the "system"?)

‣ Searle argues that (a) the room is like an AI system producing Chinese translations; (b) the operator in the room (the AI) does not "understand" Chinese. Harnad summarizes :

*The interpretation will not be intrinsic to the symbol system itself: It will be parasitic on the fact that the symbols have meaning for us, in exactly the same way that the meanings of the symbols in a book are not intrinsic, but derive from the meanings in our heads.*

# Language Grounding

‣ Bender and Koller separate form and meaning. Meaning = communicative intent. The role of the speaker/listener are crucial in language, LMs lack the underlying intent

‣ They propose the "octopus" experiment to show how form alone can fail.
An octopus is eavesdropping on a conversation between A and B (using deep-sea communication cables). Suddenly, the octopus decides to cut the cable and impersonate B.

‣ A has an emergency and asks how to construct something with sticks to fend off a bear. The octopus can't help because it can't simulate this novel situation.



Bender and Koller (2020) *Climbing towards NLU*

# Counterarguments

▸ We can't necessarily learn semantics from predicting next characters alone without execution. Consider training on:

```
x = 2
y = x + 2
print(y)
```

▸ **However**, assertion statements are sufficient to teach us some semantics! (but this can still break down)

```
x = 2
y = x + 2
assert(y == 4)
```

▸ For language: similar argument. Assume people say true things. Consider saying a pair of sentences $x_1$, $x_2$; given enough examples, the fact that $x_2$ should not be contradicted by $x_1$ tells us something
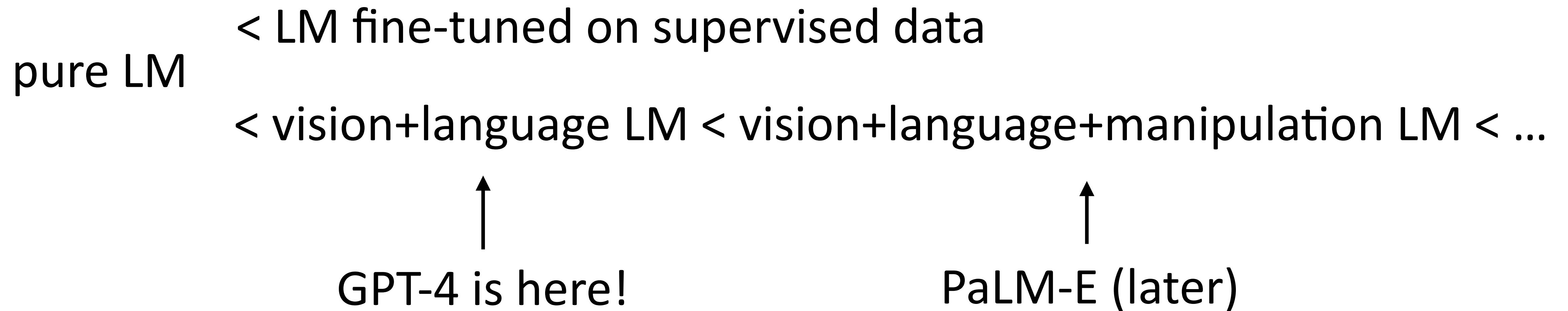
Merrill et al. (2021) *Provable Limitations of Acquiring Meaning from Ungrounded Form*

Merrill et al. (2022) *Entailment Semantics can be Extracted from an Ideal Language Model*
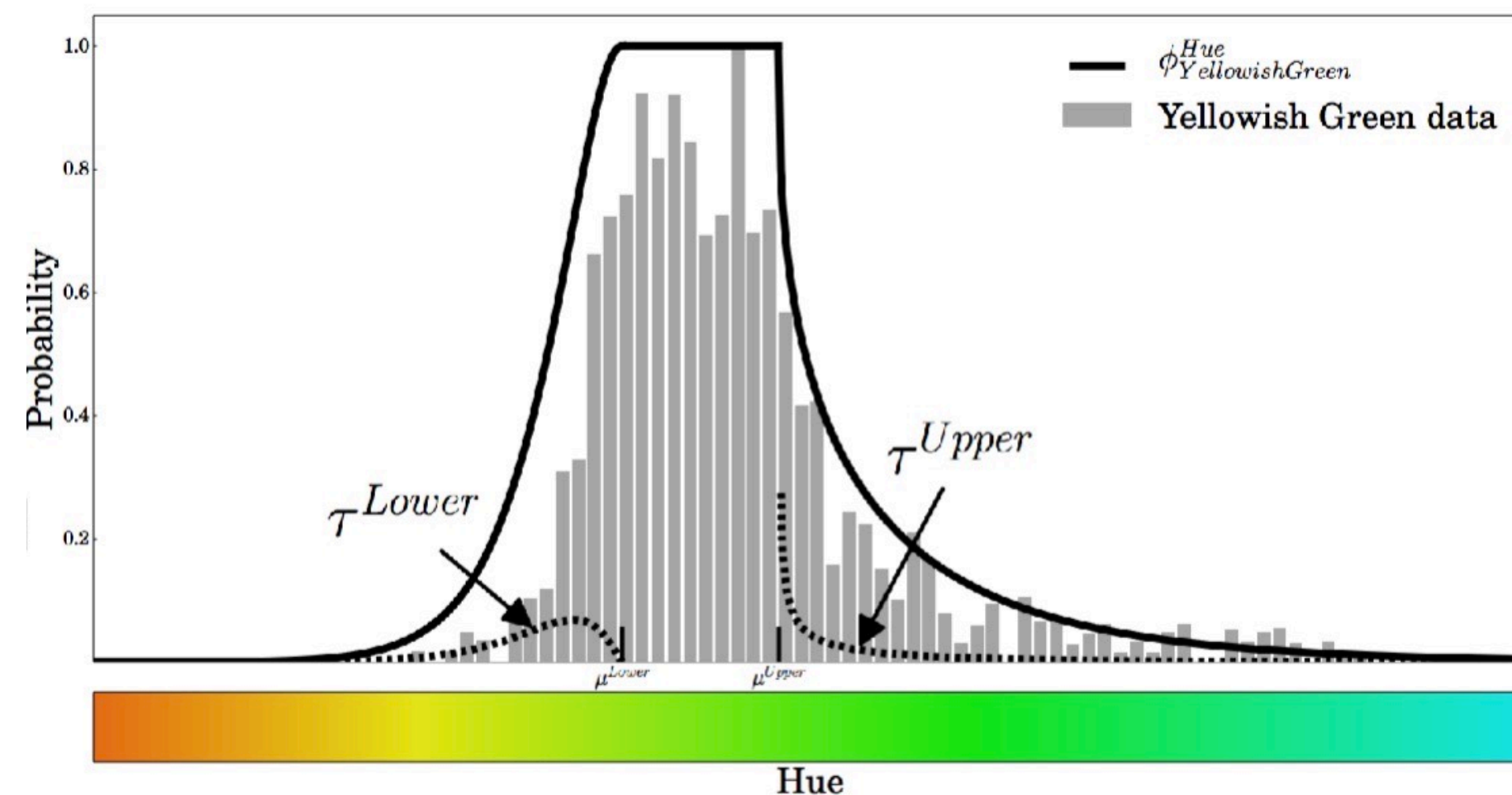
# Where are we?

‣ Lots of philosophy about these models!

‣ Nevertheless, it seems there's a hierarchy in terms of their understanding:

pure LM  < LM fine-tuned on supervised data

< vision+language LM < vision+language+manipulation LM < ...

GPT-4 is here!                    PaLM-E (later)

# Language Grounding

▸ There are many things that we can ground language in! Focus on vision today.

▸ How to associate words with sensory-motor experiences

▸ How to associate words with meaning representation





**Alan Turing** was a British mathematician, logician, cryptanalyst, and computer scientist.

$\mathtt{nationality(AT,UK)} \wedge \mathtt{notable\_for(AT,mathematian)}$

$\wedge \mathtt{profession(AT,logic))} \wedge \mathtt{research(AT,cryptanalysm)}$
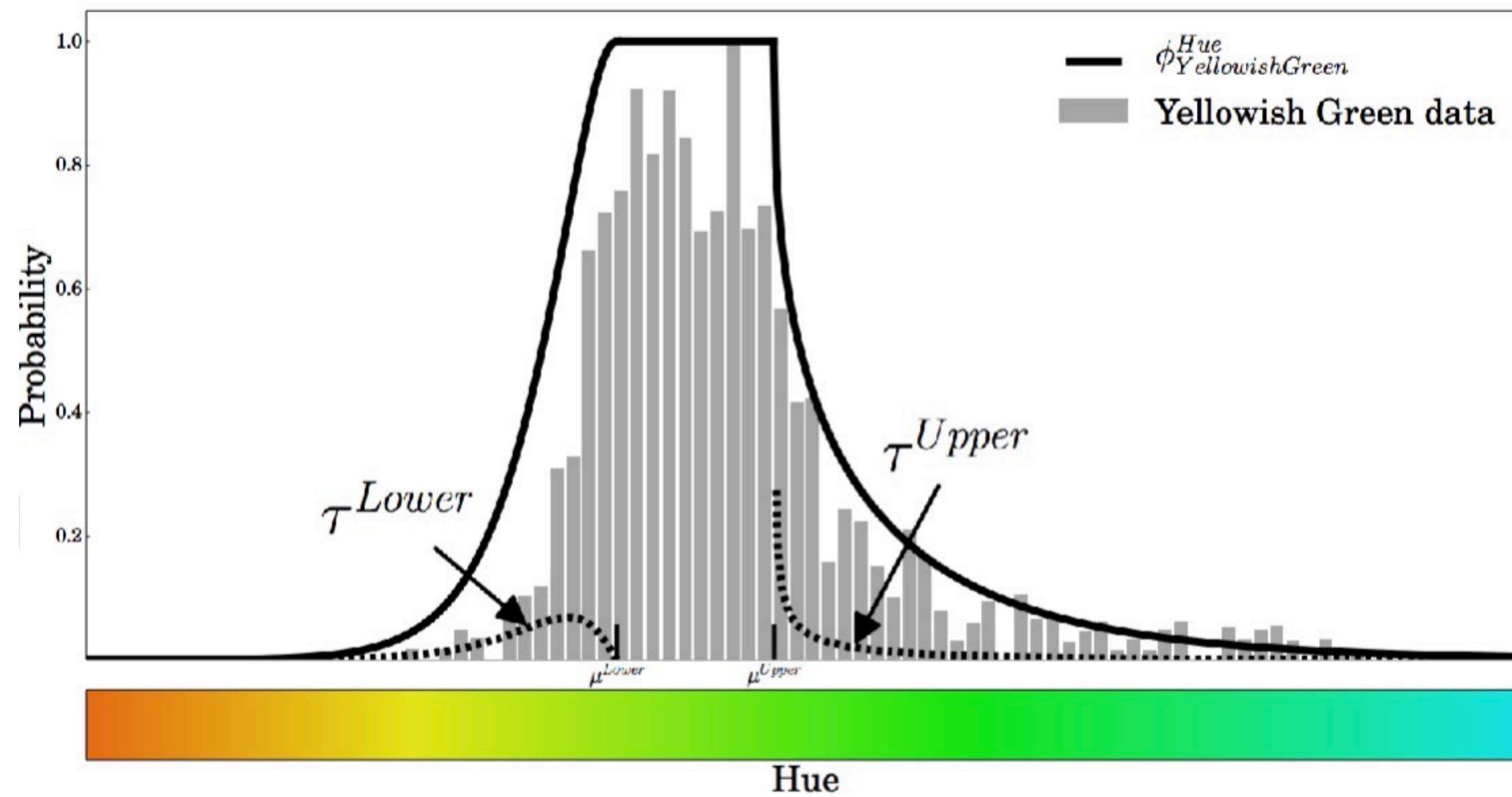
$\wedge \mathtt{notable\_type(AT,compsci)}$

# Multimodality, Language Grounding

# Language Grounding

- What does "yellowish green" mean?

- Formal semantics: yellowish green is a predicate. Things are either yellowish green or not. No connection to real color

- Grounding in perceptual space:

McMahan and Stone (2015)

# Perception

▶ Visual: *green =* [0,1,0] in RGB

▶ Auditory: *loud =* >120 dB

▶ Taste: sweet = >some threshold level of sensation on taste buds

▶ High-level concepts:

cat

dog

running

eating
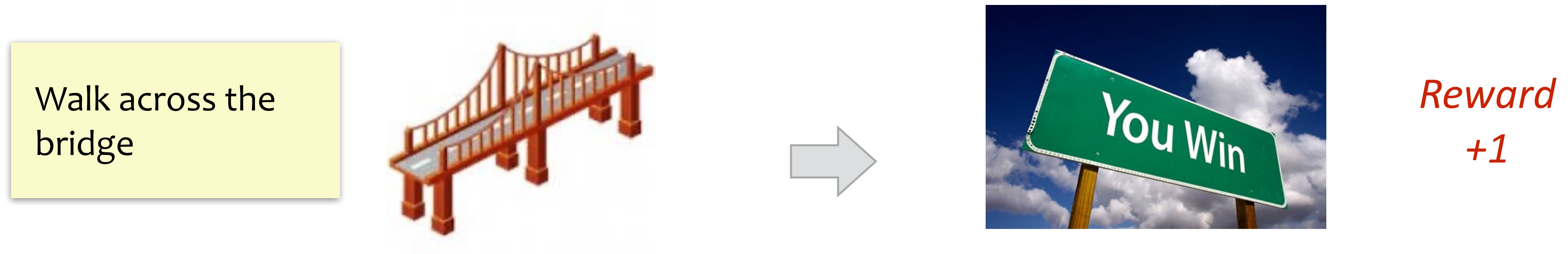
# Learning from Interaction

## 1. Use feedback from control application to understand language

Walk across the bridge



You Win

*Reward +1*

*Alleviate dependence on large scale annotation*

## 2. Use language to improve performance in control applications



 +

1. **Ghosts** chase and try to kill you
2. Collect all the **pellets**
3. ...

*Score: 7*

*Score: 107*

# Other Grounding

- **Temporal concepts**

  - *late evening* = after 6pm.
    Ground in a time interval

  - *fast, slow* = describing rates of change

- **Functional:**

  - *Jacket:* keeps people warm

  - *Mug*: holds water

- **Spatial Relations**

  - *left, on top of, in front of*: how should
    we ground these?

- **Size:**

  - Whales are *larger* than lions

- **Focus today: grounding in images**

# Language and Vision Models

# Grounding in Images

‣ How would you describe this image?

‣ What does the word "*spoon*" evoke?



*the girl is licking the spoon of batter*

# Grounding Spoon

Winco 0005-03 7 3/8" Dinner Spoon…

$7.16

wikiHow

How to Hold a Spoon: 13 Steps (…

Indiegogo

Spoon that Elevates Taste …

# Grounding Language in Images

▸ Syntactic categories have some regular correspondences to the world:

  ▸ Nouns: objects

  ▸ Verbs: actions

  ▸ Sentences: whole scenes or things happening

▸ Tasks:

  ▸ Object recognition (pick out one most salient object or detect all of them)

  ▸ Image captioning: produce a whole sentence for an image

# Language-vision Models



Image encoder
(CNN, Transformer)

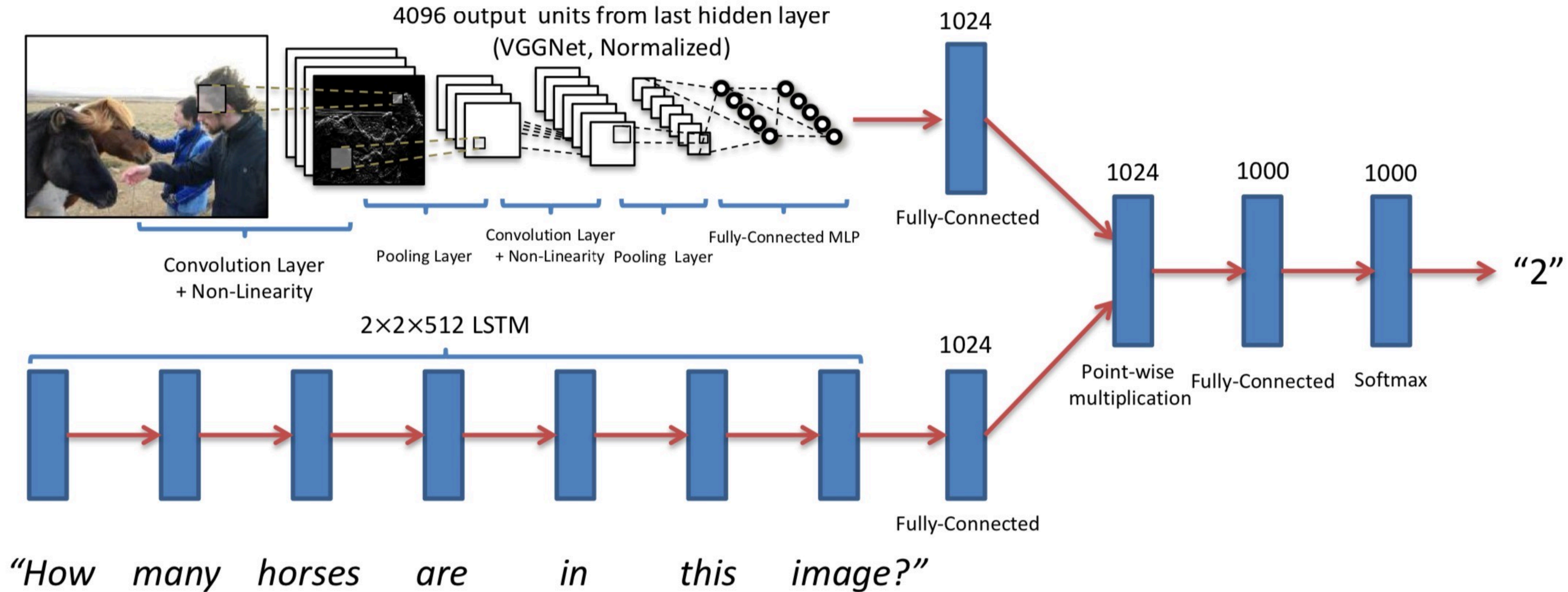*the girl is licking the spoon of batter*

Language encoder
(LSTM, Transformer)

Cross-attention/joint layer
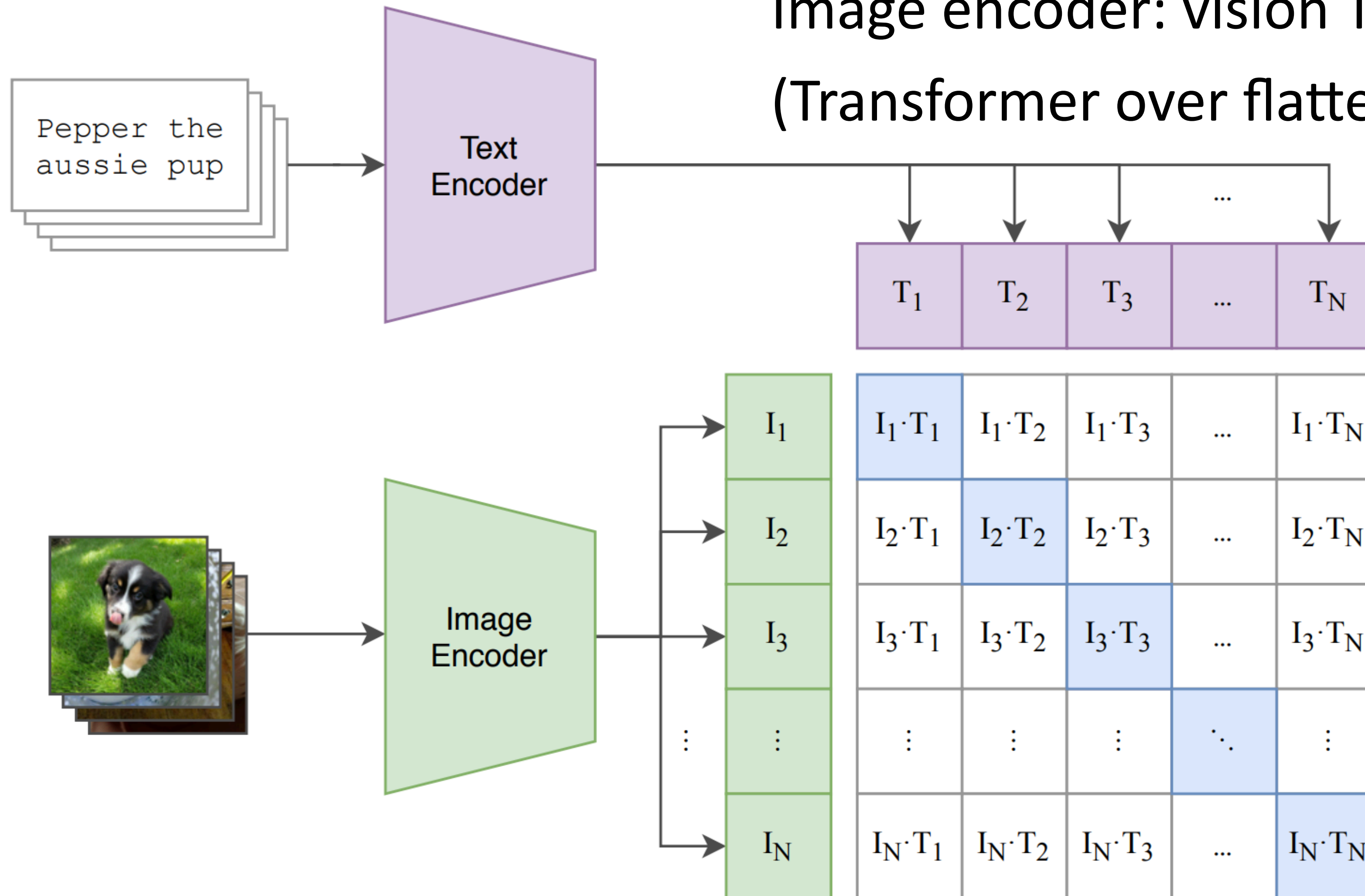
Prediction

# Visual Question Answering



4096 output units from last hidden layer (VGGNet, Normalized)

1024

Fully-Connected

Convolution Layer + Non-Linearity

Pooling Layer

Convolution Layer + Non-Linearity

Pooling Layer

Fully-Connected MLP

2×2×512 LSTM

1024

Fully-Connected

1024

1000

1000

Point-wise multiplication

Fully-Connected

Softmax

"2"

"How many horses are in this image?"

Agrawal et al., 2015

# Language-vision Pre-training

## (1) Contrastive pre-training

Text encoder: Transformer

Image encoder: vision Transformer

(Transformer over flattened patches)

Pepper the aussie pup → Text Encoder

| | $T_1$ | $T_2$ | $T_3$ | ... | $T_N$ |
|---|---|---|---|---|---|

Image Encoder

| | $T_1$ | $T_2$ | $T_3$ | ... | $T_N$ |
|---|---|---|---|---|---|
| $I_1$ | $I_1 \cdot T_1$ | $I_1 \cdot T_2$ | $I_1 \cdot T_3$ | ... | $I_1 \cdot T_N$ |
| $I_2$ | $I_2 \cdot T_1$ | $I_2 \cdot T_2$ | $I_2 \cdot T_3$ | ... | $I_2 \cdot T_N$ |
| $I_3$ | $I_3 \cdot T_1$ | $I_3 \cdot T_2$ | $I_3 \cdot T_3$ | ... | $I_3 \cdot T_N$ |
| ⋮ | ⋮ | ⋮ | ⋮ | ⋱ | ⋮ |
| $I_N$ | $I_N \cdot T_1$ | $I_N \cdot T_2$ | $I_N \cdot T_3$ | ... | $I_N \cdot T_N$ |

Radford et al., 2021

# Language-vision Pre-training

|  | $T_1$ | $T_2$ | $T_3$ | ... | $T_N$ |
|---|---|---|---|---|---|
| $I_1$ | $I_1 \cdot T_1$ | $I_1 \cdot T_2$ | $I_1 \cdot T_3$ | ... | $I_1 \cdot T_N$ |
| $I_2$ | $I_2 \cdot T_1$ | $I_2 \cdot T_2$ | $I_2 \cdot T_3$ | ... | $I_2 \cdot T_N$ |
| $I_3$ | $I_3 \cdot T_1$ | $I_3 \cdot T_2$ | $I_3 \cdot T_3$ | ... | $I_3 \cdot T_N$ |
| ⋮ | ⋮ | ⋮ | ⋮ | ⋱ | ⋮ |
| $I_N$ | $I_N \cdot T_1$ | $I_N \cdot T_2$ | $I_N \cdot T_3$ | ... | $I_N \cdot T_N$ |

▸ Contrastive objective: each image should be more similar to its correspond caption than to other captions

maximize softmax($I_1^T T_i$)[1]
+ softmax($I_2^T T_i$)[2]
+ ...

Radford et al., 2021

24

# Language-vision Pre-training

**(2) Create dataset classifier from label text**

| plane |
| car |
| dog |
| ⋮ ⋮ |
| bird |

→ A photo of a {object}. → **Text Encoder** → $T_1$ | $T_2$ | $T_3$ | ... | $T_N$

**(3) Use for zero-shot prediction**

**Image Encoder** → $I_1$

| $I_1 \cdot T_1$ | $I_1 \cdot T_2$ | $I_1 \cdot T_3$ | ... | $I_1 \cdot T_N$ |

→ A photo of a dog.

Radford et al., 2021

# CLIP: Zero-shot Results



**Stanford Cars**

correct label: 2012 Honda Accord Coupe   correct rank: 1/196   correct probability: 63.30%

a photo of a 2012 honda accord coupe.

a photo of a 2012 honda accord sedan.

a photo of a 2012 acura tl sedan.

a photo of a 2012 acura tsx sedan.

a photo of a 2008 acura tl type-s.

**Country211**

correct label: Belize    correct rank: 5/211    correct probability: 3.92%

# Parti

- Autoregressive text-to-image model (differs from the diffusion models you may have seen, like Stable Diffusion or DALL-E)



A. A photo of a frog reading the newspaper named "Toaday" written on it. There is a frog printed on the newspaper too.

Yu et al., 2022

# Parti

Yu et al., 2022

# Manipulation: SayCan, PaLM-E

‣ Most models like CLIP are just vision+language. What about interaction with the world?

# SayCan

‣ Probability of taking an action decomposes as follows:

$$p(c_i | i, s, \ell_\pi) \propto p(c_\pi | s, \ell_\pi) p(\ell_\pi | i)$$

p(skill possible given world state)
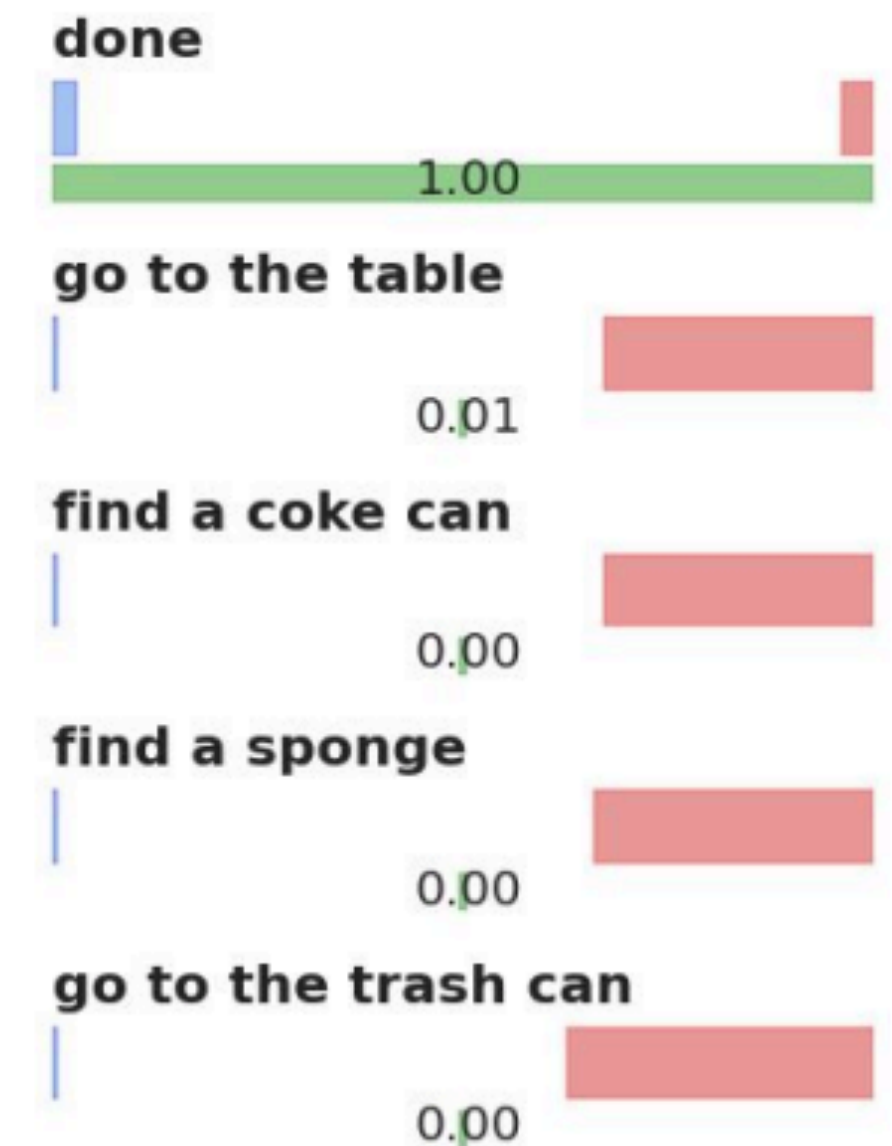
p(language description of skill | instruction)
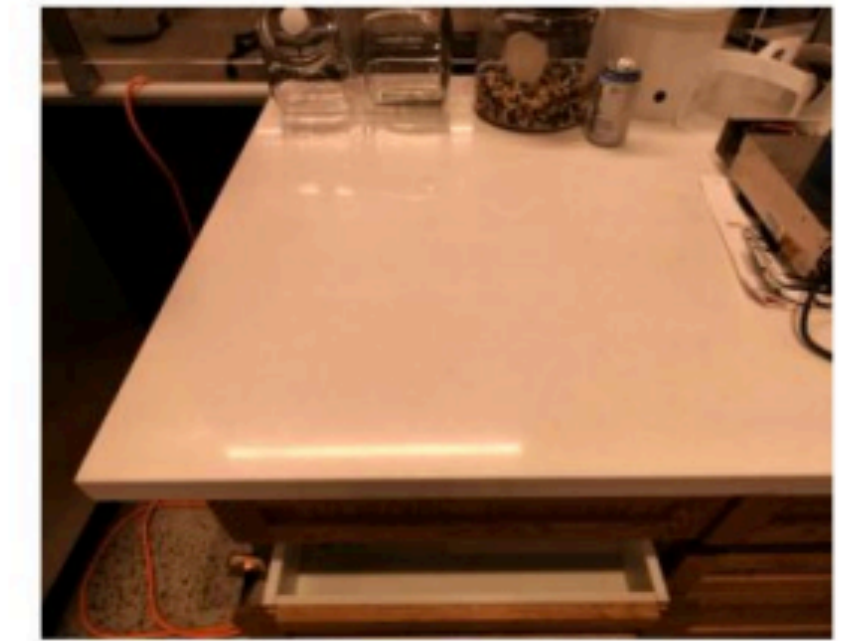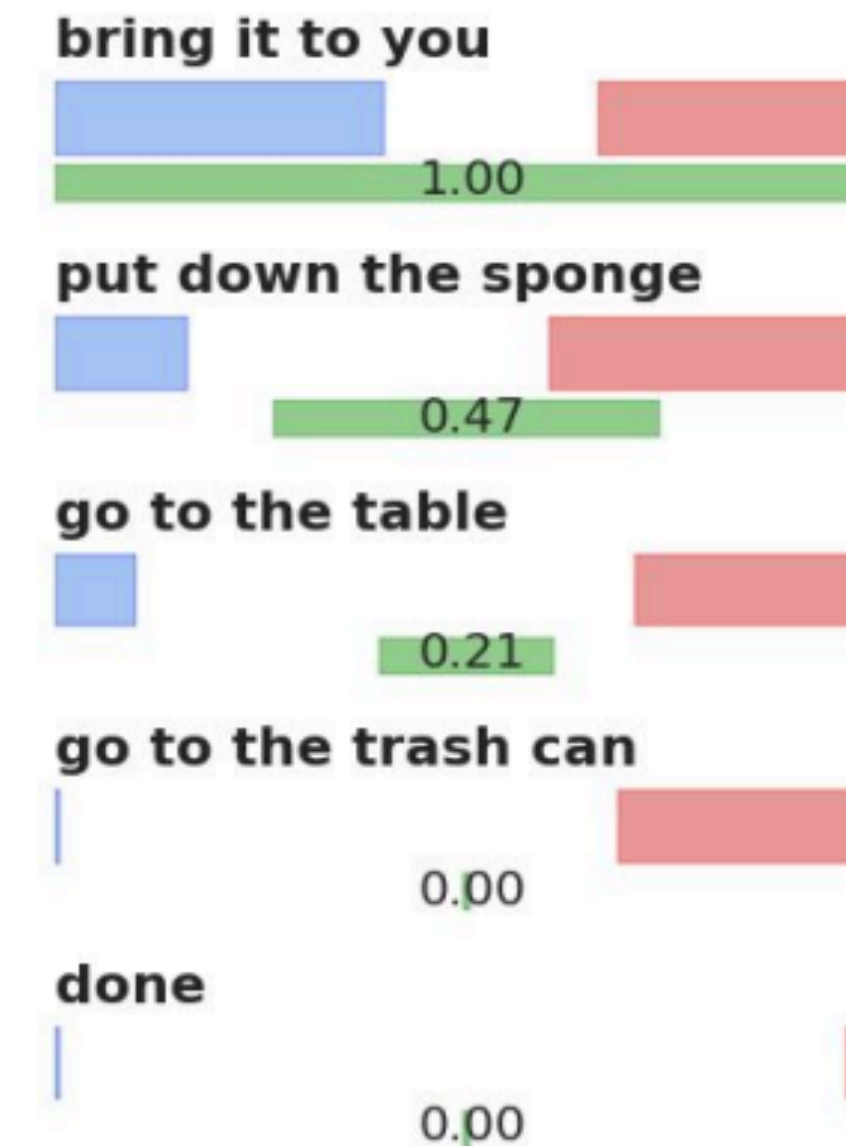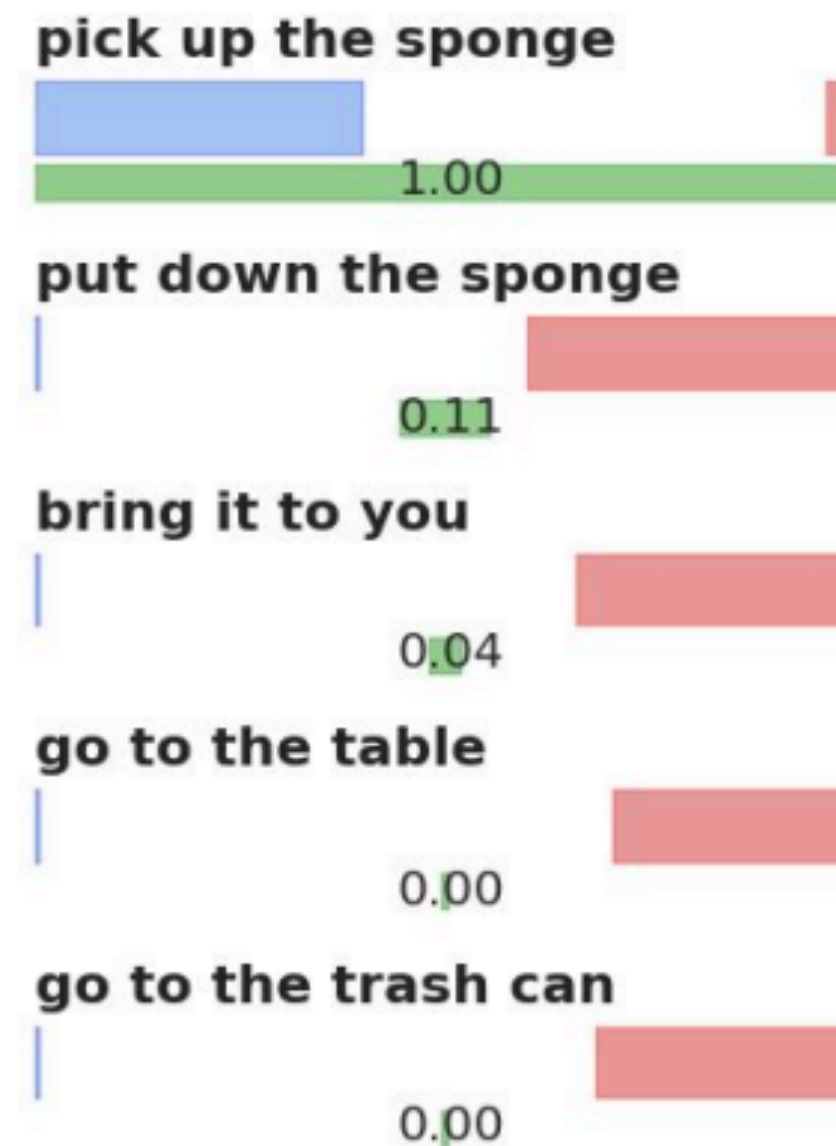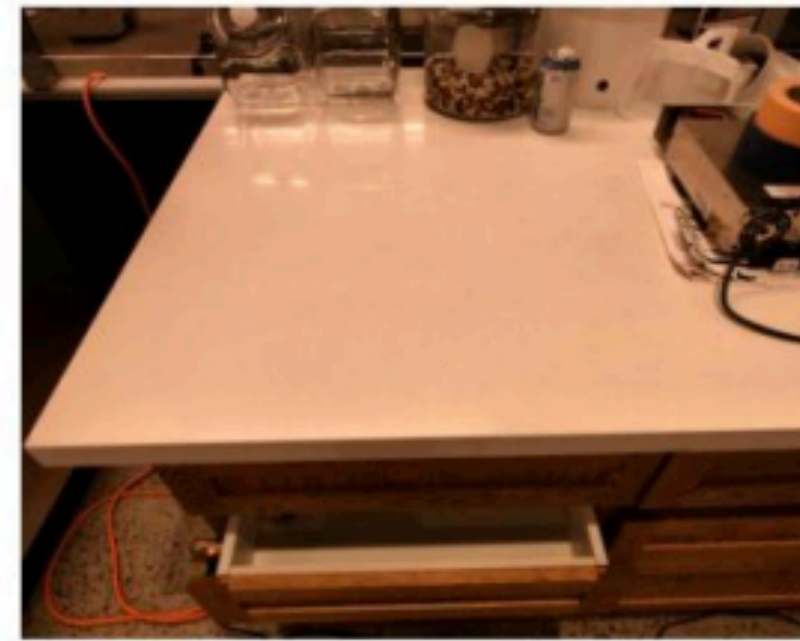
‣ Individual skills are learned in advance, form affordance models for that skill

‣ Train a single multi-task policy that conditions on the lang description

‣ Do you think this is a grounded language model?

# SayCan

**Human:** I spilled my coke, can you bring me something to clean it up?

**Robot:** I would
1. Find a sponge
2. Pick up the sponge
3. Bring it to you
4. Done

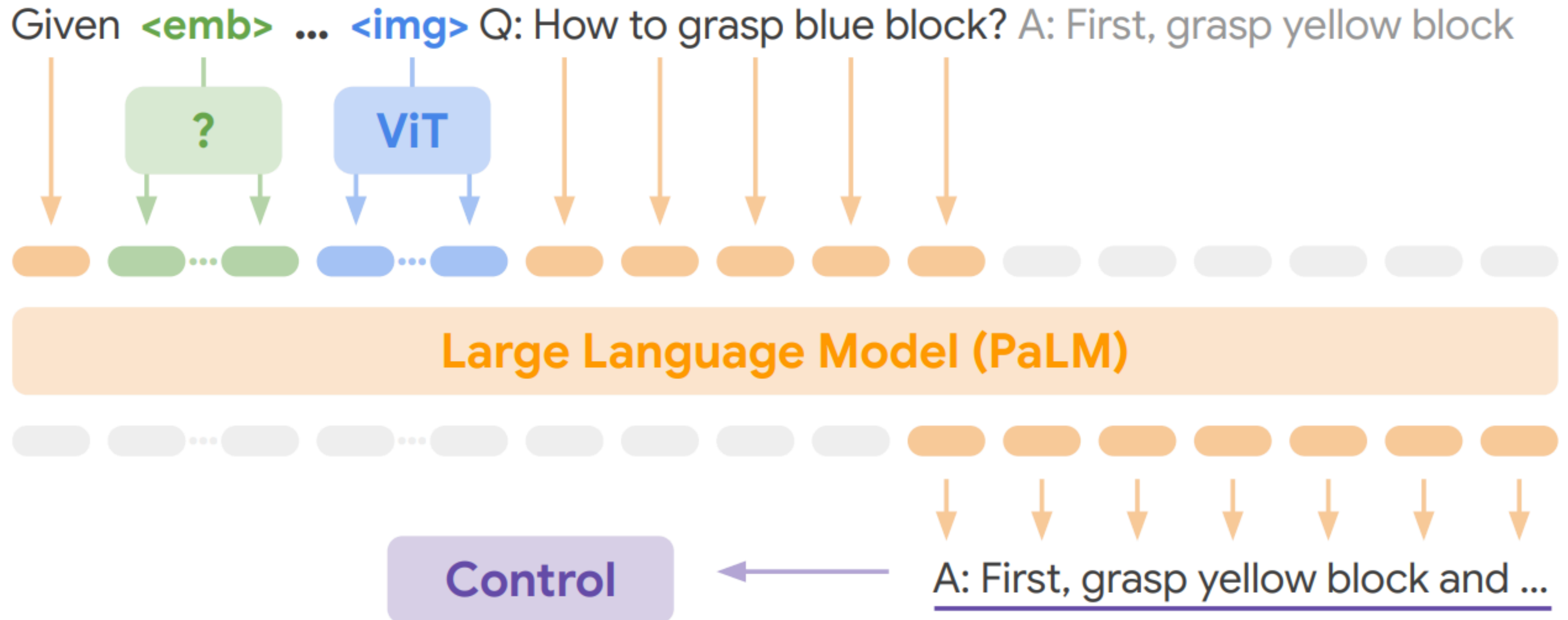Language × Affordance
Combined Score



**find a sponge** — 1.00
go to the table — 0.08
find a coke can — 0.08
go to the trash can — 0.05
find a water bottle — 0.01



**pick up the sponge** — 1.00
put down the sponge — 0.11
bring it to you — 0.04
go to the table — 0.00
go to the trash can — 0.00



**bring it to you** — 1.00
put down the sponge — 0.47
go to the table — 0.21
go to the trash can — 0.00
done — 0.00



**done** — 1.00
go to the table — 0.01
find a coke can — 0.00
find a sponge — 0.00
go to the trash can — 0.00

‣ Most models like CLIP are just vision+language



PaLM-E: An Embodied Multimodal Language Model

Given \<emb\> ... \<img\> Q: How to grasp blue block? A: First, grasp yellow block

?

ViT

Large Language Model (PaLM)

Control

A: First, grasp yellow block and ...

Push red blocks to the coffee cup

# Where are we today

- Explosion of multimodal pre-training for {video, audio, images, interaction} x text

- Many of these methods are Transformer-based

- Still haven't seen large-scale multimodal pre-training of this form advance text-only tasks, but there's potential!
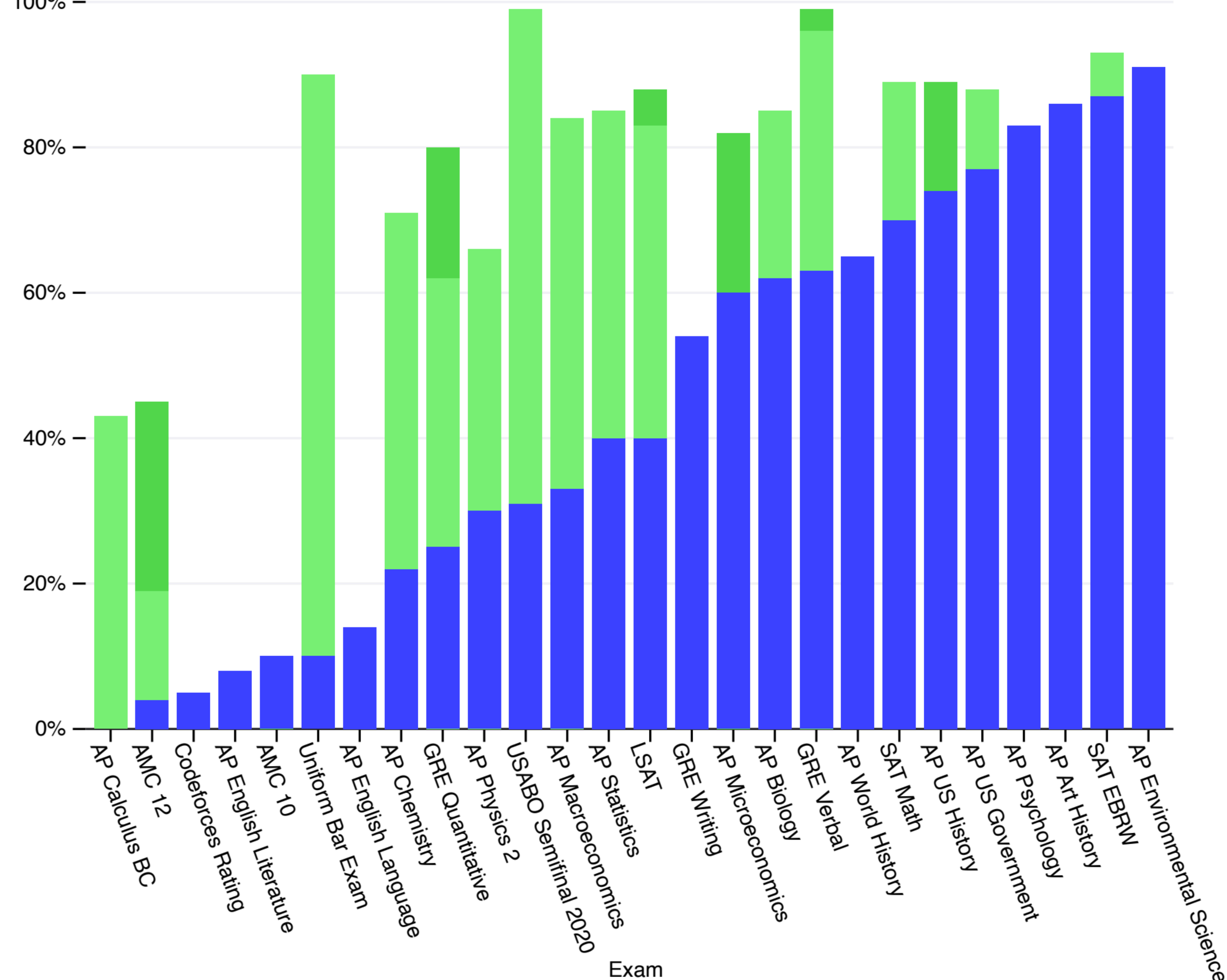
- Impact of images on GPT-4 is unclear

# GPT-4

- ‣ Dark green: additional performance from vision pre-training

- ‣ This graph is hard to read and doesn't make sense...



**Exam results (ordered by GPT-3.5 performance)**

Estimated percentile lower bound (among test takers)

Legend: gpt-4 (dark green), gpt-4 (no vision) (light green), gpt3.5 (blue)

# Takeaways

‣ Is the lack of grounding in text-only pre-trained models a problem?

‣ Multimodal methods can allow us to learn representations for images as well as text and provide a path towards language grounding

‣ Pre-training on text and other modalities is more and more common and unlocking new capabilities for models