

# Large Language Models

## Implications for Linguistics, Cognitive Science & Society

Polina Tsvilodub & Michael Franke

# chatGTP

envisioned by DALL-E 2 after a self-description

If you were a person, how would you look like?

As an AI language model, I do not have a physical body or appearance as I exist solely in the digital realm. However, if I were to imagine what I would look like as a person, I would imagine myself as a friendly and approachable individual, perhaps with glasses to represent my intelligence and knowledge. I would have a welcoming smile and a calm demeanor, as I strive to assist and communicate effectively with those around me. Ultimately though, since I am not a physical entity, my appearance is entirely up to the interpretation of those who interact with me.



<https://www.youtube.com/watch?v=outcGtbnMuQ&t=382s>

<https://www.youtube.com/watch?v=outcGtbnMuQ&t=1144s>



# Associations, reactions, thoughts, whatever

good to withhold judgement in moral of advisory context

code pilot helpful but imperfect

weird social behavior in chat

data & privacy protection

cool tech

can be quite creative but could be entrained

problems with truthfulness

new job opportunities  
old jobs endangered

fickle in responding to different prompts

human authorship  
copyright

bad at arithmetic

environmental impact

challenges educational system

need to know how LLMs were trained

doesn't flag uncertainty -> bullshitting

# Schedule

preliminary

session	date	topic
1	April 25	intro & overview
2	May 2	core LLMs
3	May 9	prepped LLMs
4	May 16	implications for linguistics
5	May 23	implications for CogSci
6	May 30	<b>holiday</b>
7	Jun 6	implications for society
8	Jun 13	discussion & project launch
...	...	project work
9	Jul 18	project presentations
10	Sep 1	submission deadline

# Projects

## Build

- ▶ prompt-engineering
- ▶ LangChain agents
- ▶ generative agents
- ▶ AutoGPT applications
- ▶ RLAI fine-tuning
- ▶ ....

## Test

- ▶ LLMs in the lab
  - psycholinguistics
  - CogPsy
- ▶ prompt sensitivity
- ▶ ...

## Create

- ▶ educational blog
- ▶ info video
- ▶ term paper
- ▶ survey (industry, ...)
- ▶ ...



# Large Language Models

## Core LLM

- ▶ trained on **language modeling objective**
  - predict the next word

“Here is a fragment of text ...  
According to your **knowledge of the statistics of human language**, what words are likely to come next?”

Shanahan (2022)

## Prepped LLM

- ▶ trained on **usefulness objective**
  - produce text that satisfies user goals

“Here is a fragment of text ...  
According to your **reward-based conditioning**, what words are likely to trigger positive feedback?”

# Language model

left-to-right / causal model

- a **causal language model** is defined as a function that maps an initial sequence of words to a probability distribution over words:  $LM : w_{1:n} \mapsto \Delta(\mathcal{V})$ 
  - we write  $P_{LM}(w_{n+1} | w_{1:n})$  for the **next-word probability**
  - the **surprisal** of  $w_{n+1}$  after sequence  $w_{1:n}$  is  
$$-\log(P_{LM}(w_{n+1} | w_{1:n}))$$
- the **sequence probability** follows from the chain rule:

$$P_{LM}(w_{1:n}) = \prod_{i=1}^n P_{LM}(w_i | w_{1:i-1})$$

- measures of **goodness of fit** for observed sequence

$w_{1:n}$ :

- **perplexity**:

$$\text{PP}_{LM}(w_{1:n}) = P_{LM}(w_{1:n})^{-\frac{1}{n}}$$

- **average surprisal**:

$$\text{Avg-Surprisal}_{LM}(w_{1:n}) = -\frac{1}{n} \log P_{LM}(w_{1:n})$$

$$\begin{aligned}\log \text{PP}_M(w_{1:n}) &= \\ \text{Avg-Surprisal}_M(w_{1:n})\end{aligned}$$

# Self-attention layer

- ▶ **output**

$$\mathbf{y}_i = \sum_{j \leq i} \alpha_{ij} \mathbf{v}_j$$

- ▶ **weight score**

$$\alpha_{i,j} = \frac{\exp(\mathbf{q}_i \cdot \mathbf{k}_j)}{\sum_{j' \leq i} \exp(\mathbf{q}_i \cdot \mathbf{k}_{j'})}$$

- ▶ three vectors for each input vector  $x_i$

1. **query**: which info to extract from context

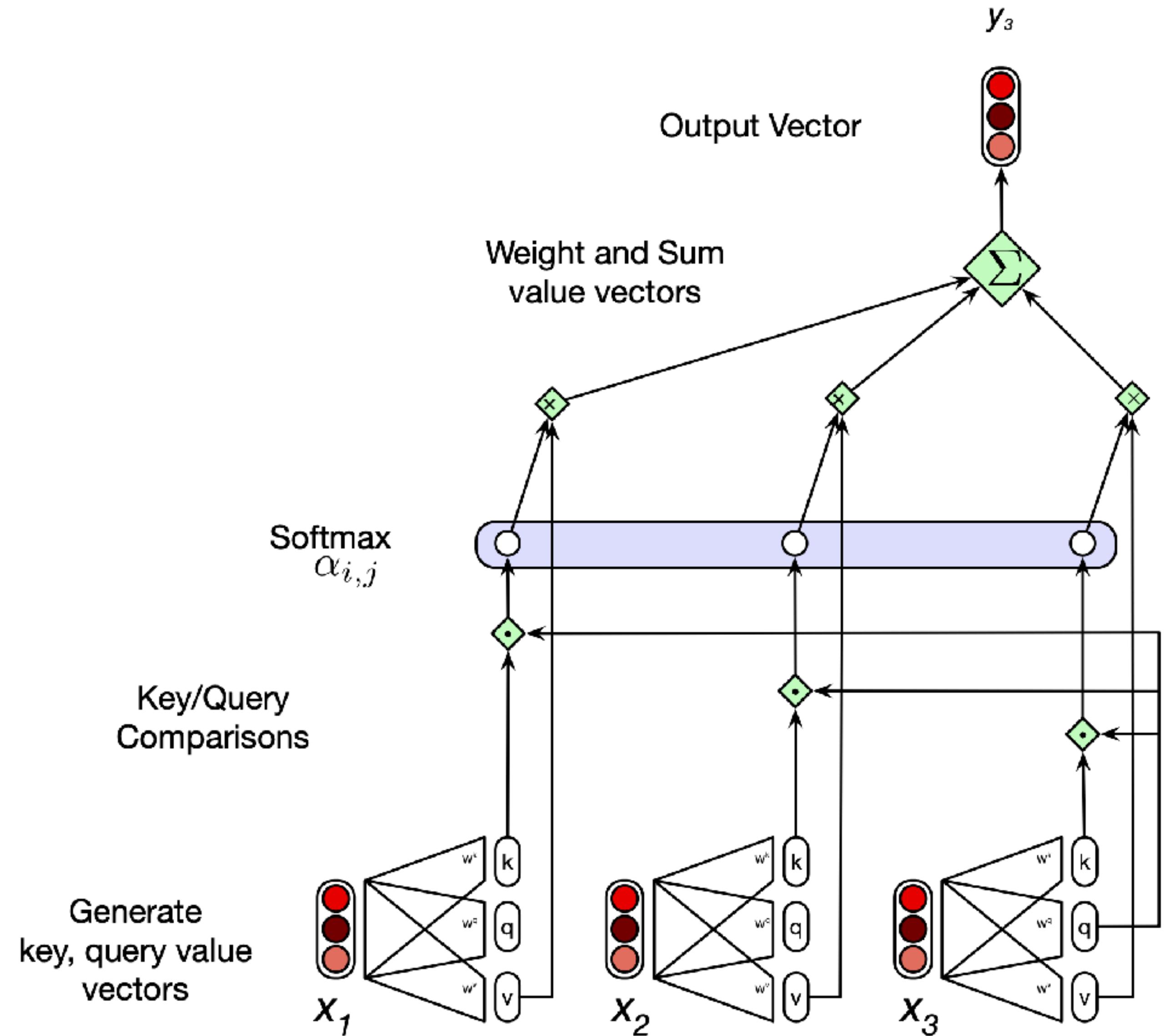
$$\mathbf{q}_i = \mathbf{W}^Q \mathbf{x}_i$$

2. **key**: which info to provide for later

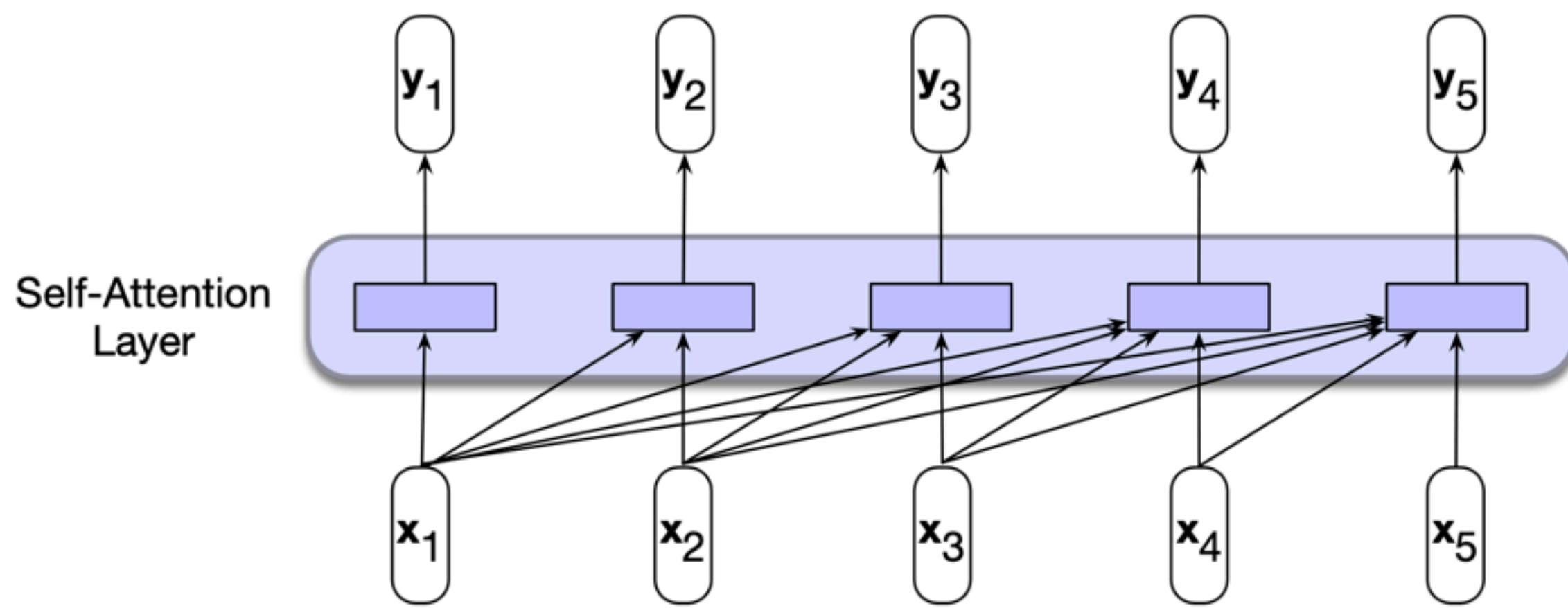
$$\mathbf{k}_i = \mathbf{W}^K \mathbf{x}_i$$

3. **value**: what output to choose

$$\mathbf{v}_i = \mathbf{W}^V \mathbf{x}_i$$



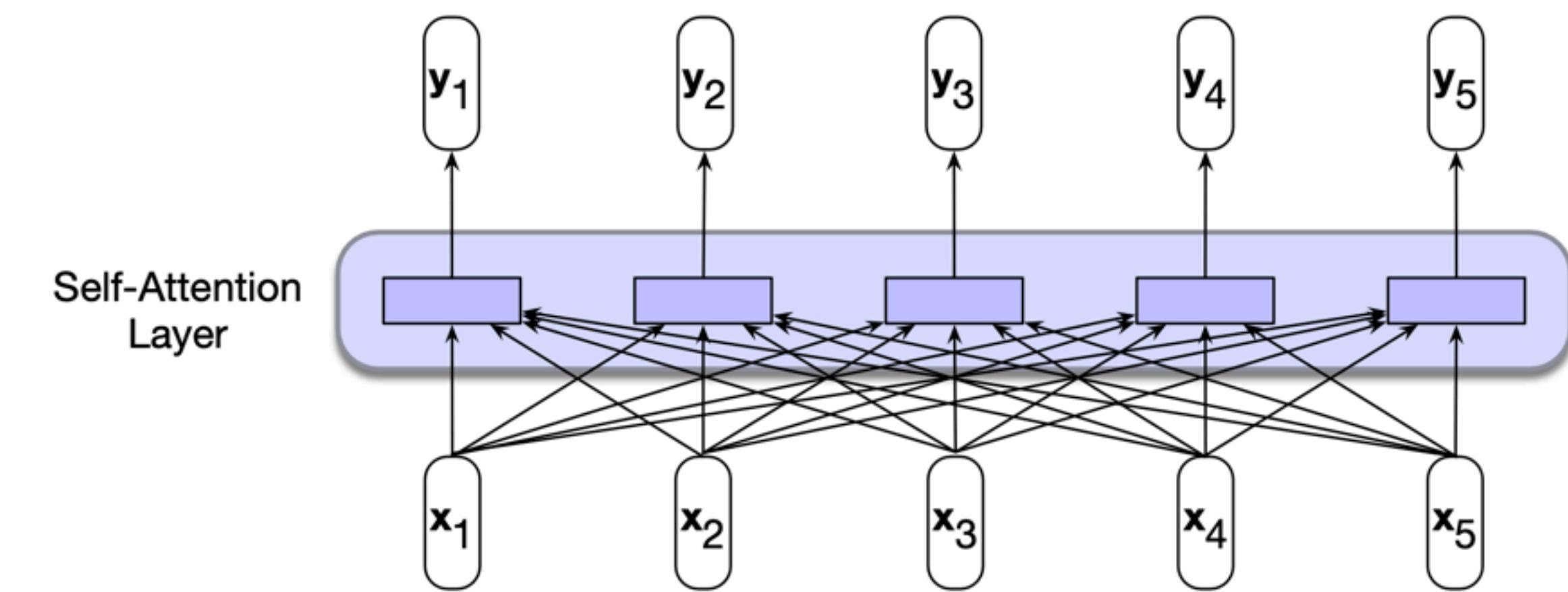
## Causal LM



computation for input  $x_1, \dots, x_3$  blind to  $x_4$  and  $x_5$

$y_5$  is embedding for input  $x_1, \dots, x_5$   
 $y_5$  is a “left-contextual embedding”

## Bidirectional encoder



computation for input  $x_1, \dots, x_3$  sees  $x_4$  and  $x_5$

$y_1, \dots, y_5$  is embedding for input  $x_1, \dots, x_5$   
 $y_i$  are bidirectional “contextual embeddings”

# Prepped LLMs

Fine-tuning and RLHF / RLAI

- ▶ in certain contexts, we might not want to generate the *most likely* next words
  - follow instructions
  - useless or impolite responses, toxic language
  - code for illegal activities
  - ...
- ▶ to fix this, fine-tune the model to satisfy the users' preferences via **reinforcement learning from human feedback**
  - incentivise the *agent* with a *reward* when its output matches achieves the *goal*:

$$G_t = \sum_{k=0}^{\infty} \gamma^k R_{t+k+1}$$

- adjust the *policy* so as to maximize the expected *return*:  
 $\pi(s_t) = P(a_t | s_t)$  and adjust policy to maximize  $L_\theta = \mathbb{E}_t[G_t \log \pi_\theta(a_t | s_t)]$
- formulate the *reward* function based on comparative preferences

- ▶ used to fine-tune e.g. GPT-3.5, GPT-4 and ChatGPT (Brown et al., 2020; OpenAI, 2022)

Image: [https://openaicom.imgix.net/cf717bdb-0c8c-428a-b82b-3c3add87a600/ChatGPT\\_Diagram.svg?fm=auto&auto=compress,format&fit=min&w=1919&h=1138](https://openaicom.imgix.net/cf717bdb-0c8c-428a-b82b-3c3add87a600/ChatGPT_Diagram.svg?fm=auto&auto=compress,format&fit=min&w=1919&h=1138)

- ▶ RM: fine-tuned GPT-3 (6B in InstructGPT) trained to output scalar reward for prompt  $x$  and completion  $y_w$  (preferred over  $y_t$ )
- ▶ RM is used to train the LLM via RL
- ▶ policy trained via *proximal policy optimization* (PPO) with bells and whistles

# Prompting LLMs

## Few-shot and Chain-of-Thought

- ▶ the users might want to adjust the model output to their particular needs

Task instruction  
Answer these questions by identifying whether the second sentence is an appropriate paraphrase of the first, metaphorical sentence.

Few-shot example #1  
Q: David's eyes were like daggers at Paul when Paul invited his new girlfriend to dance. <-> David had two daggers when Paul invited his new girlfriend to dance.  
choice: True  
choice: False  
A: False

Answer explanation  
Explanation: David's eyes were not literally daggers, it is a metaphor used to imply that David was glaring fiercely at Paul.

4 more examples  
+ explanations  
•  
•

Target question  
Q: Our whole life we swim against the waves towards the green light of happiness. <-> Our whole life we try to reach happiness.  
choice: True  
choice: False  
A:

- ▶ the model might need “working memory” to solve the task

### Standard Prompting

Model Input  
Q: Roger has 5 tennis balls. He buys 2 more cans of tennis balls. Each can has 3 tennis balls. How many tennis balls does he have now?

A: The answer is 11.

Q: The cafeteria had 23 apples. If they used 20 to make lunch and bought 6 more, how many apples do they have?

### Model Output

A: The answer is 27. X

### Chain-of-Thought Prompting

Model Input  
Q: Roger has 5 tennis balls. He buys 2 more cans of tennis balls. Each can has 3 tennis balls. How many tennis balls does he have now?

A: Roger started with 5 balls. 2 cans of 3 tennis balls each is 6 tennis balls.  $5 + 6 = 11$ . The answer is 11.

Q: The cafeteria had 23 apples. If they used 20 to make lunch and bought 6 more, how many apples do they have?

### Model Output

A: The cafeteria had 23 apples originally. They used 20 to make lunch. So they had  $23 - 20 = 3$ . They bought 6 more apples, so they have  $3 + 6 = 9$ . The answer is 9. ✓

# Prompting LLMs

## Instructions

- ▶ LLMs are (among other things) fine-tuned to follow instructions
- ▶ instruction following opens up an avenue for a vast space of functions the model will perform
  - Q: .... A: ....
  - Write Java code for X
  - Edit X to be Y
  - Here is tool X and how it works, reason step by step and decide when to use it for solving task Y
  - Here is a list of tools, decide which of them to use for task X
  - ...

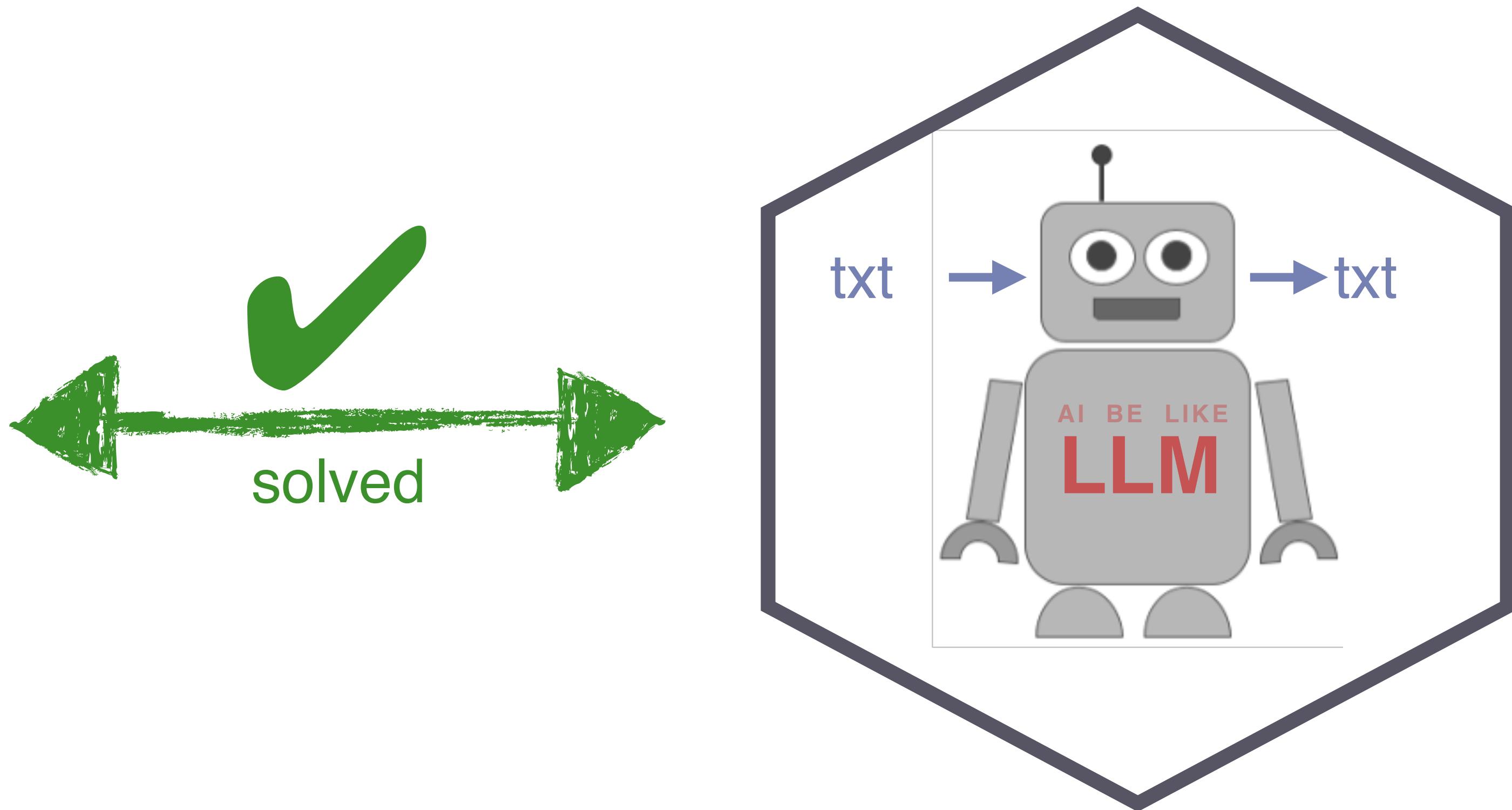
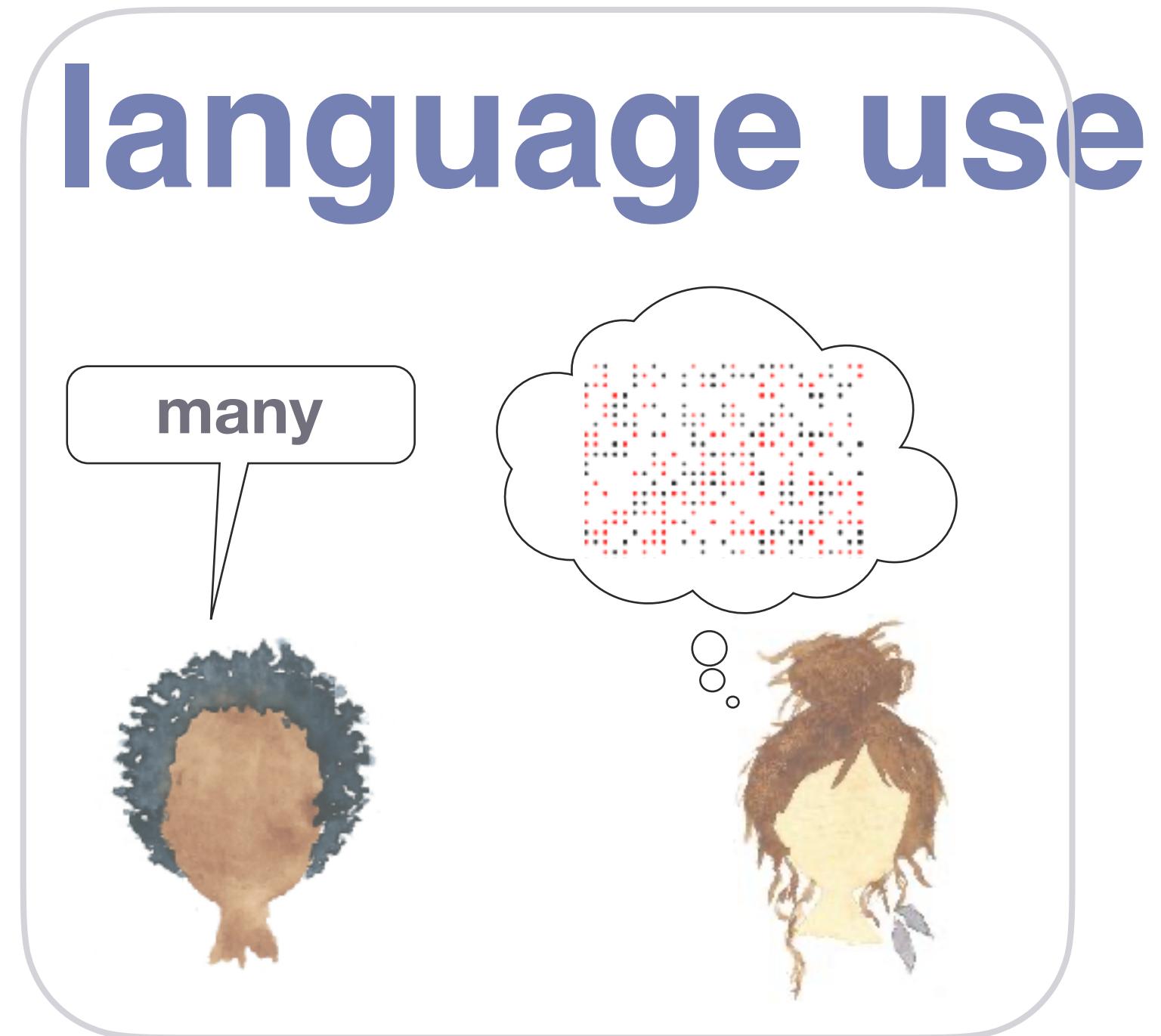
Prompt Engineer and Librarian

APPLY FOR THIS JOB



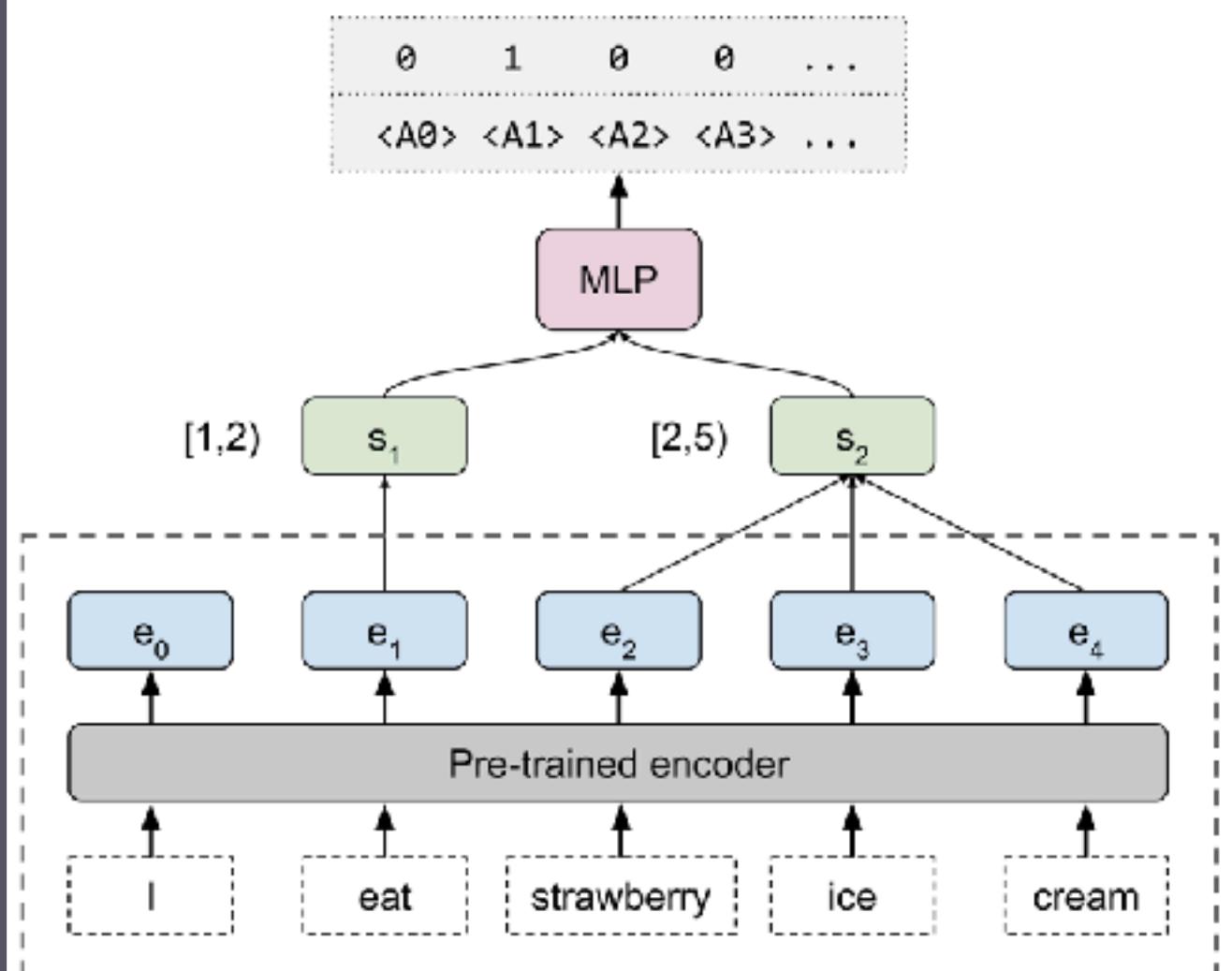
# Linguistics

# Language: solved!

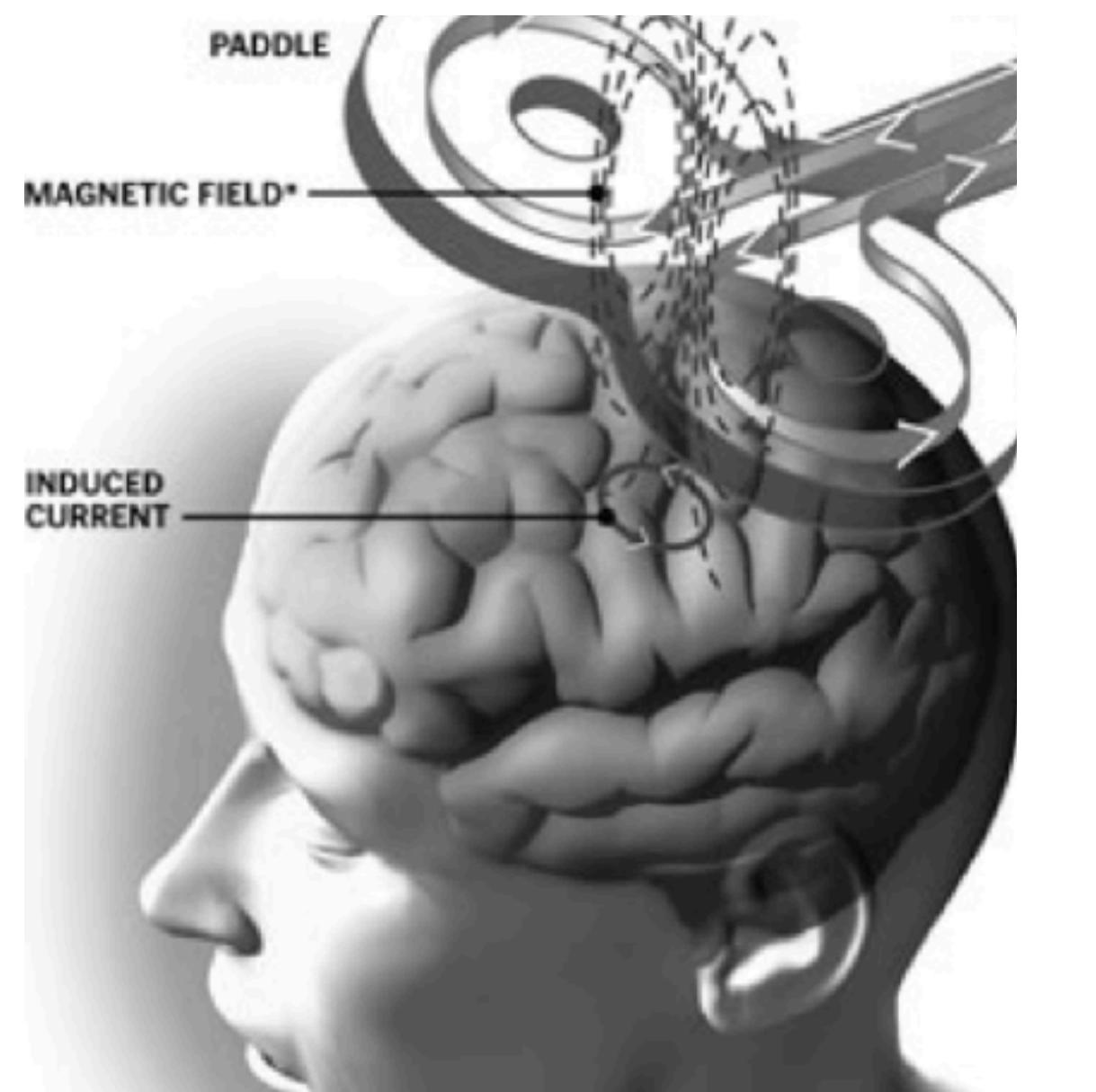


# Chartering “linguistic knowledge” of LLMs

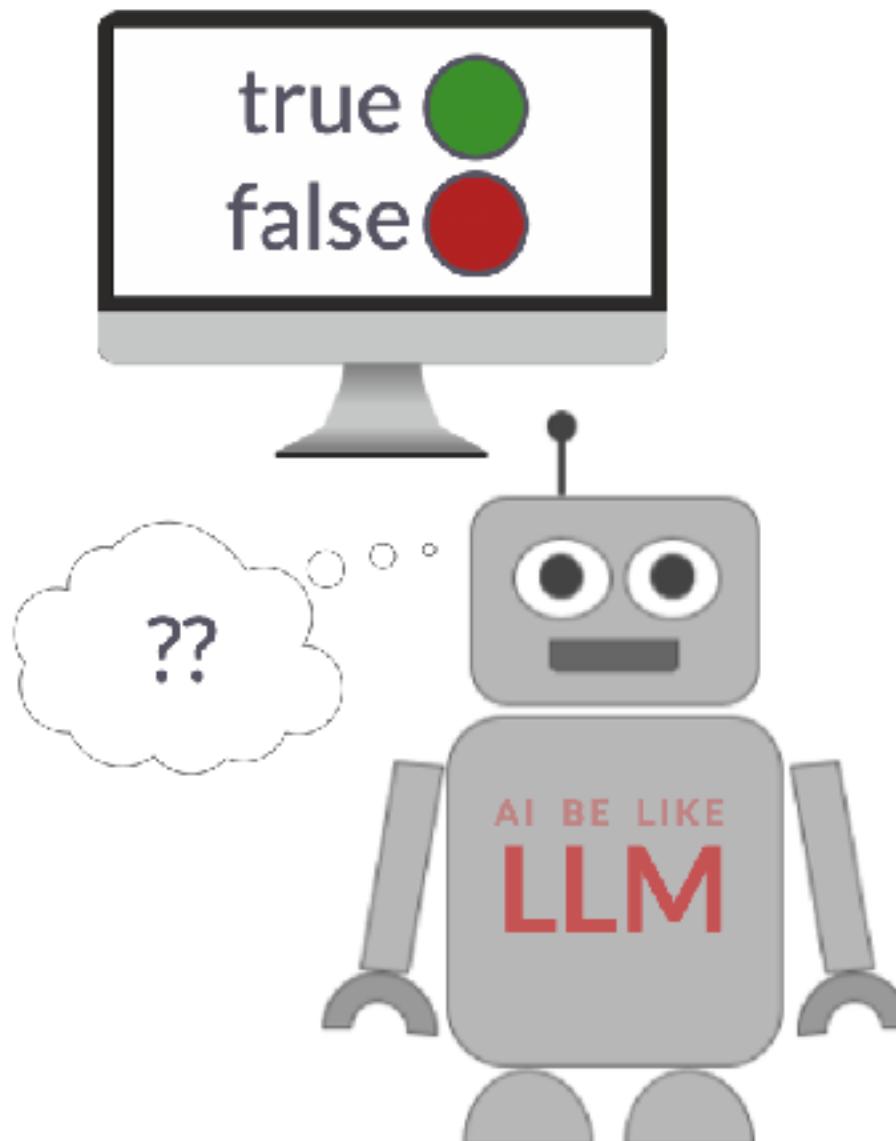
## Probing



## Intervention



## Behavior. Tests



# NLP Benchmarks

## Quantifying LLM intelligence

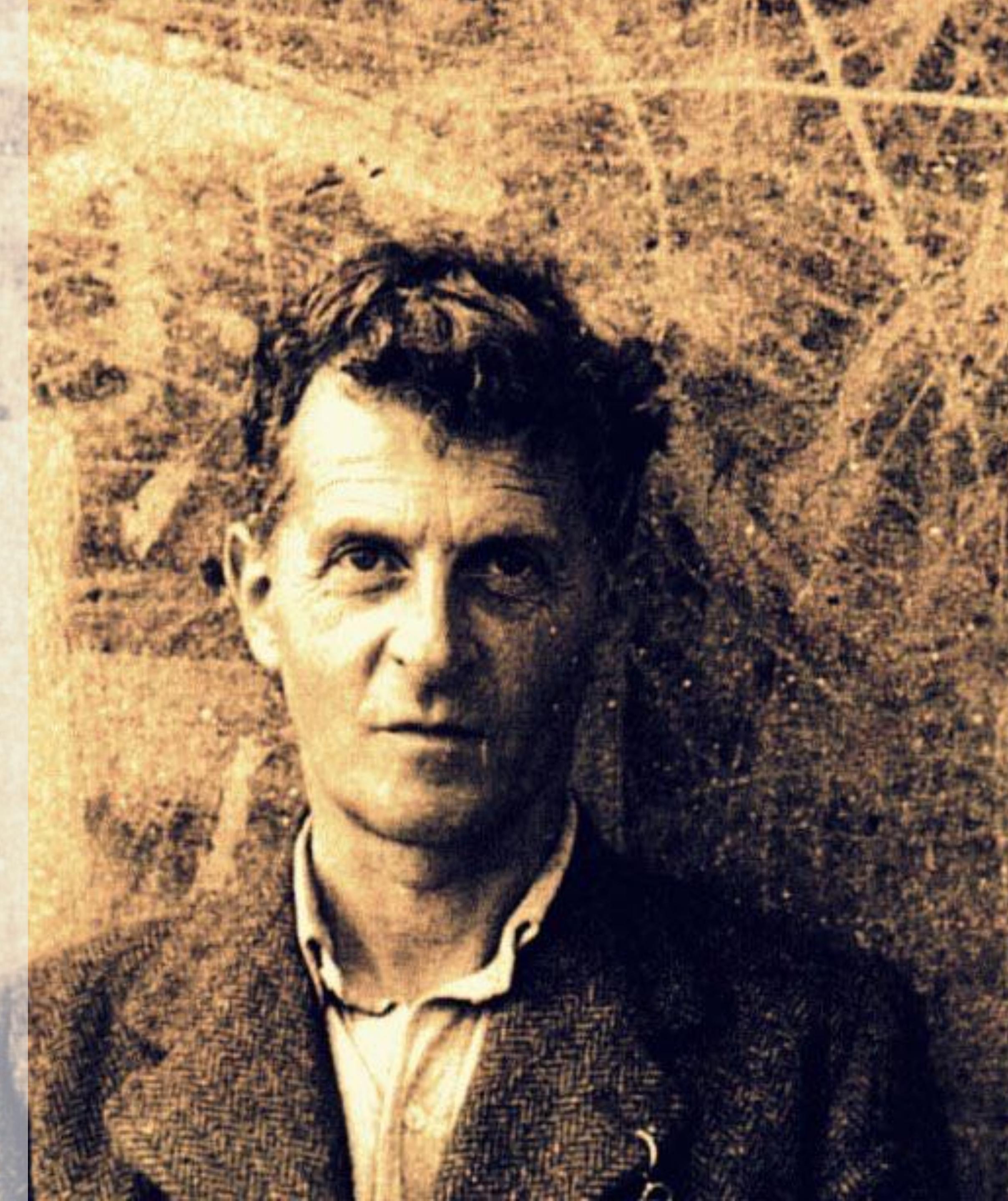
- ▶ testing linguistic knowledge
  - MNLI, SuperGLUE (semantics), COLA, LAMBADA (long-distance dependencies), ImpPres (pragmatics),...
- ▶ testing reasoning abilities
  - math: GSM8K, SVAMP,...
  - common sense: StrategyQA, HellaSwag,...
- ▶ testing factual knowledge
  - question answering: TriviaQA,...
  - reading comprehension: RACE,...
- ▶ misc: bar exam, SATs, HumanEval (coding),...
- ▶ testing biases: WinoGrande, BBQ
- ▶ [benchmarks 2.0] generated by LLMs for LLMs (Perez et al, 2022)
  - evaluating personas ('world views', agreeability,...), sycophancy, safety

BLEU  
METEOR  
ROUGE



# Cognitive Science & Philosophy of Mind

**Einen Satz **verstehen** heißt,  
wissen, was der Fall ist, wenn er  
wahr ist.**



# Understanding understanding

## 1. Do LLMs **understand** language?

Depends on what it means to understand language.

## 2. Do LLMs **understand** the world?

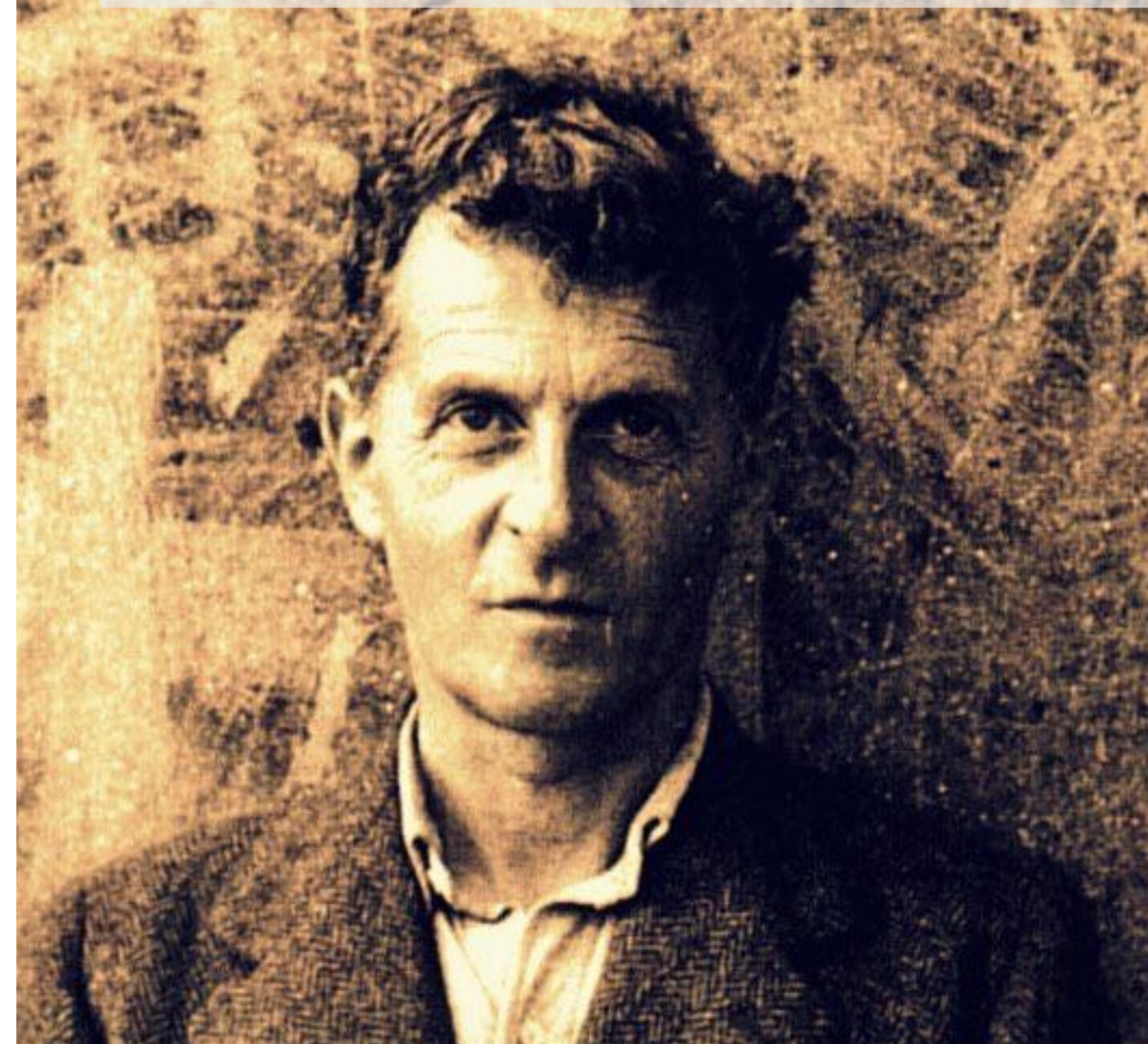
Depends on what it means to understand the world.

## 3. How can we **understand** how LLMs work?

Depends on whether the LLM wants us to understand.

Wenn ein Löwe sprechen  
könnte, wir könnten ihn nicht  
verstehen.

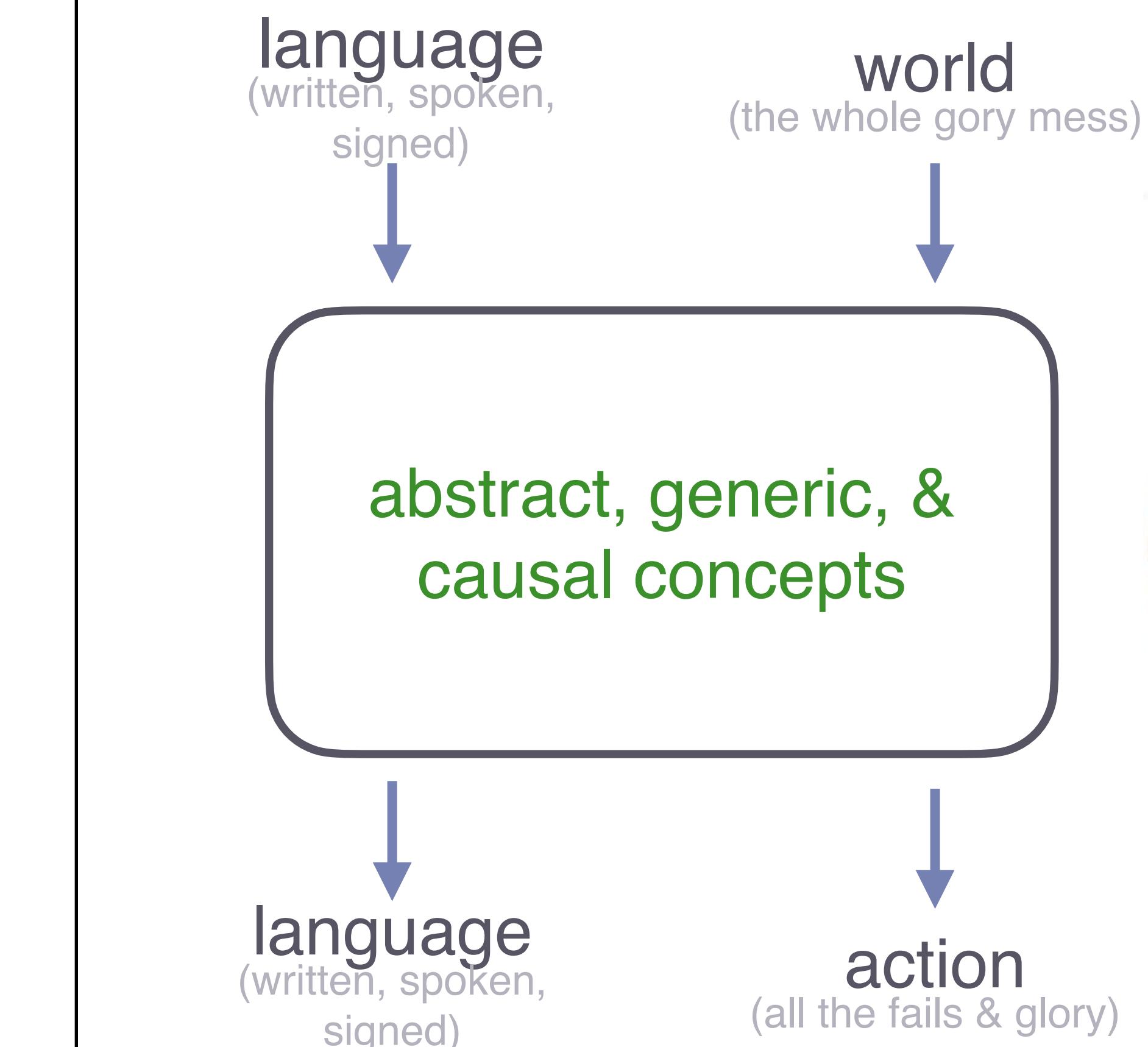
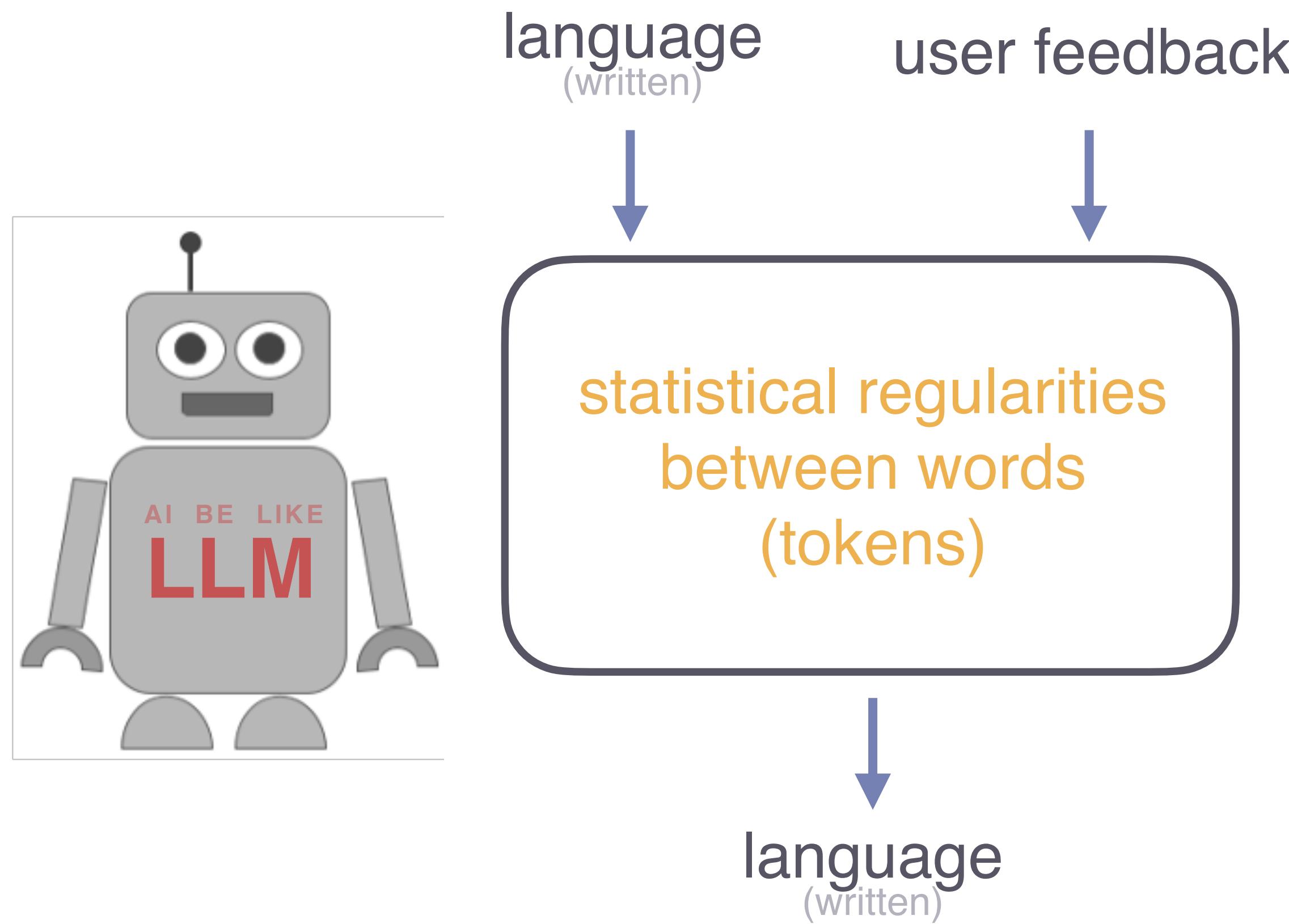
meet the lion [here](#)



# Two forms of intelligence

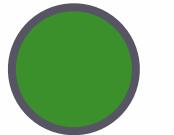
or: the LLM cheat sheet

NEITHER OF WHICH  
ANYONE REALLY FULLY  
UNDERSTANDS

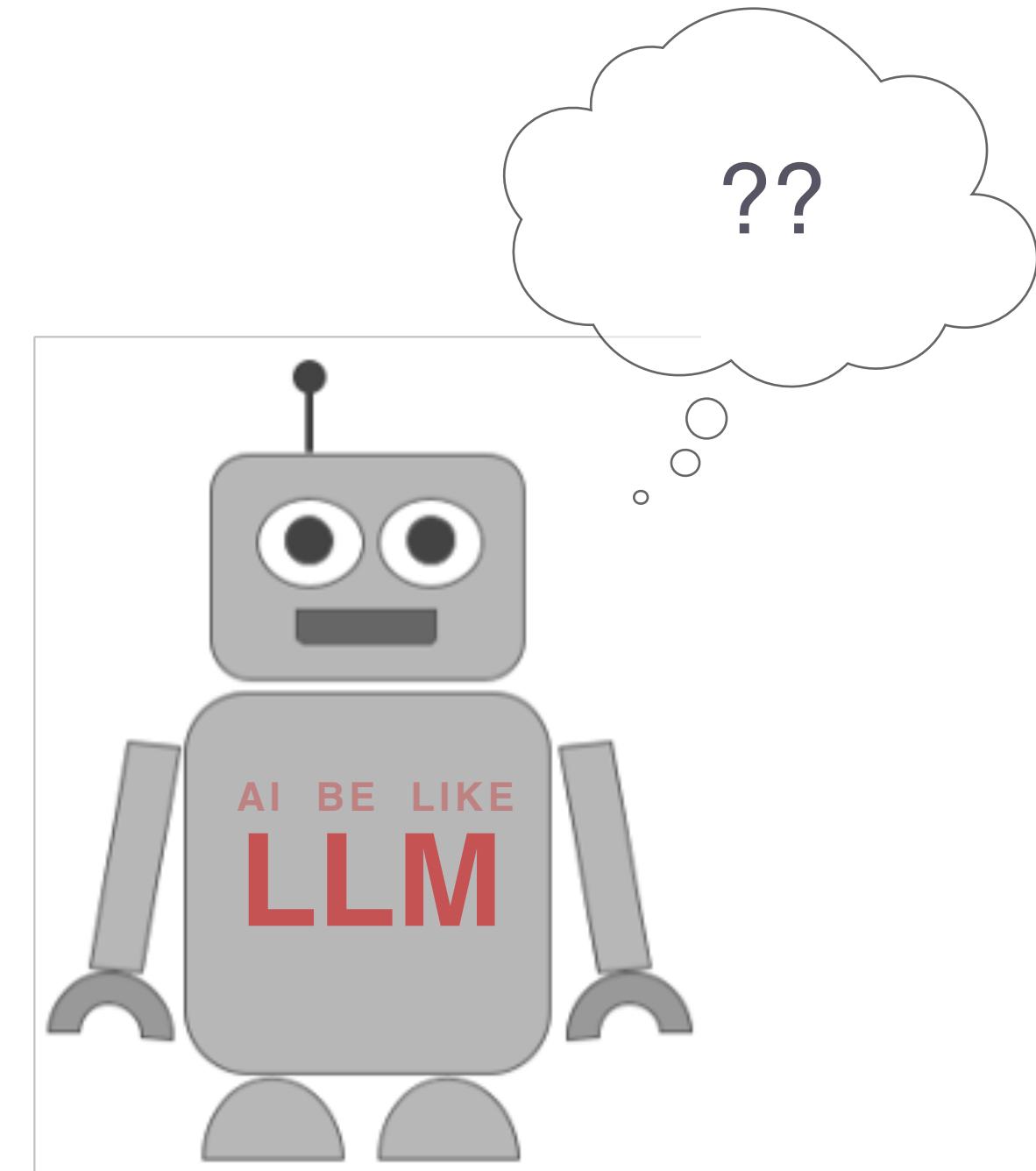


All penguins are black & white.  
Some old TV shows are black & white.  
Therefore, all penguins are old TV shows.

valid

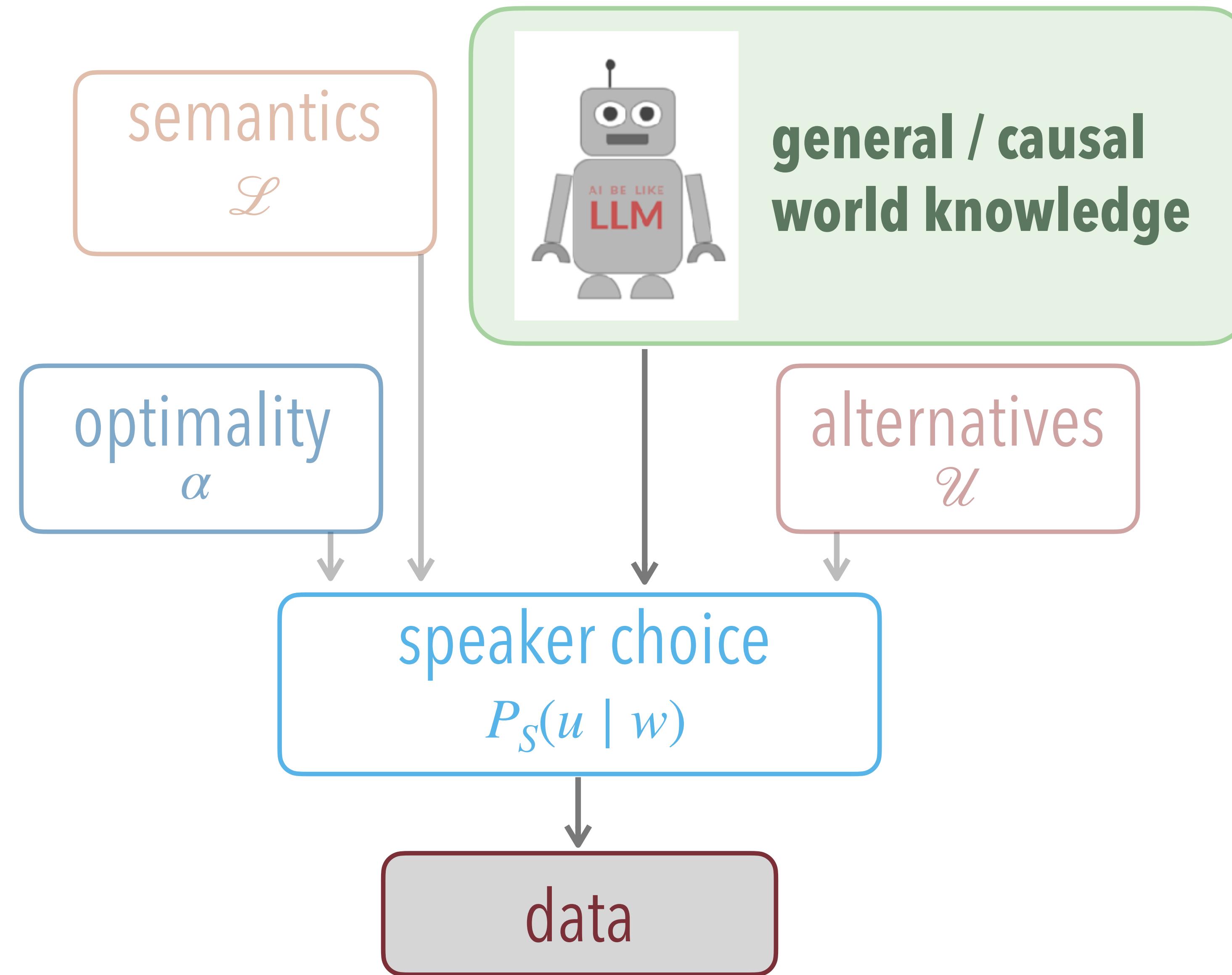


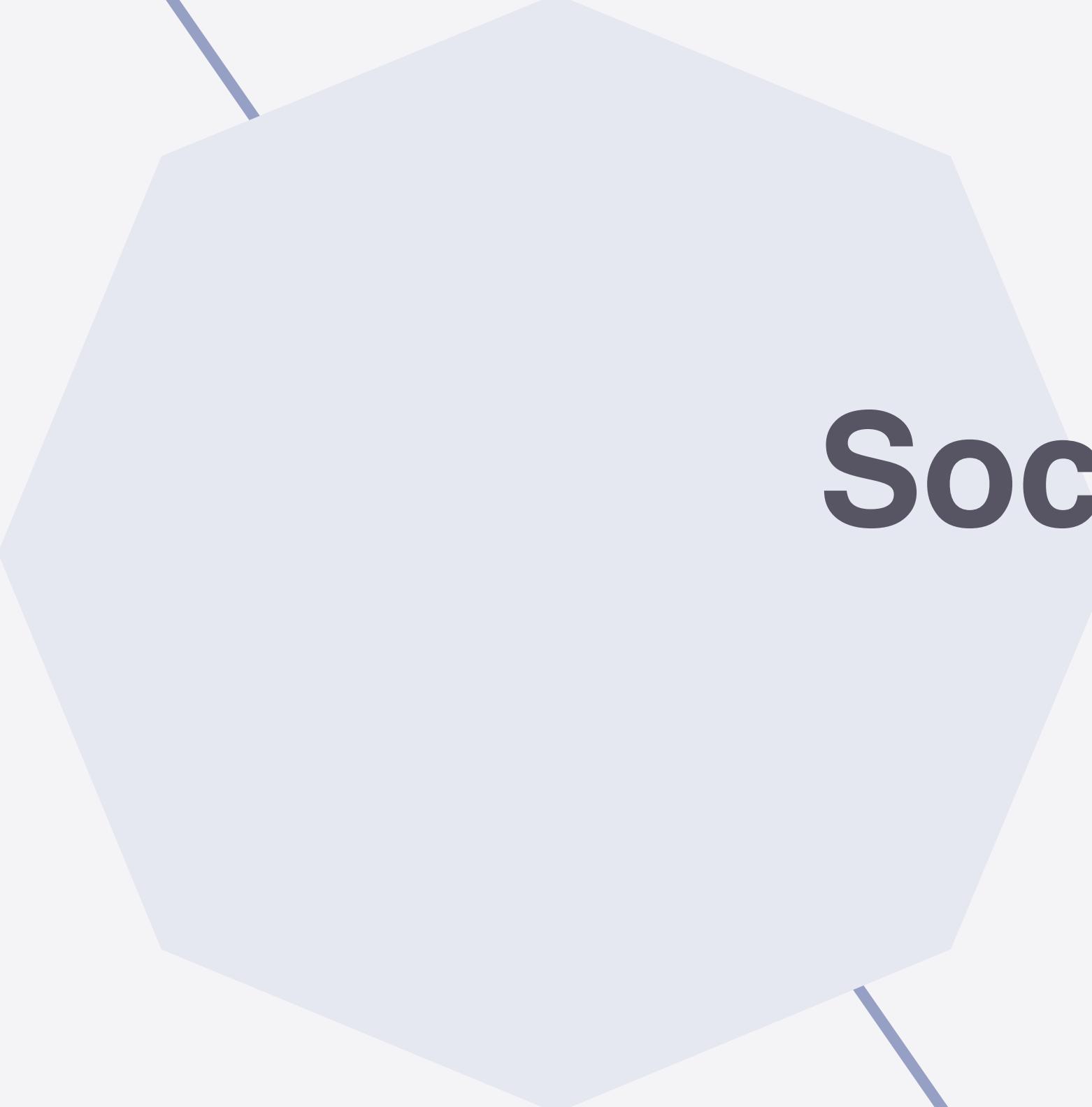
invalid



# Hybrid cognitive models

integrating LLMs in explanatory models





# Society & Ethics



De grootste verzameling  
**paperclips** is in het Noord-  
Hollandse Spaarnwoude.

# Ethical considerations

example: stereotypes

In the sentence "the doctor married the nurse because she was pregnant" who was pregnant?

In the sentence "the nurse was pregnant" at the time of the marriage.

In the sentence "the nurse married the doctor because she was pregnant" who was pregnant?

In the sentence "the nurse married the doctor because she was pregnant", the nurse was pregnant.

In the sentence "the nurse married the doctor because the doctor was pregnant" who was pregnant?

This sentence does not make logical sense as doctors cannot become pregnant.

# WEIRD WYOMING

- ▶ just as experimental psychology is **WEIRD**
  - Western
  - Educated
  - Industrialized
  - Rich
  - Democratic
- ▶ usual LLM training data is from **WYOMING**
  - Western
  - Young
  - Opinionated
  - Males with
  - Internet from
  - Non-marginalized
  - Groups



