# Bridging NLP & Survey Research: On LLMs, Language and Trust - An NLP researcher's perspective

## Barbara Plank

LMU Munich

& IT University of Copenhagen

October 7, 2024

**SurvAI Day: NLP meets Survey Science**
University of Maryland

# LLMs: A Swiss Knife for Science?

# NLP: The beauty & challenge of working with LANGUAGE

"Asking a Question Can Be a Science "
Frauke Kreuter

# Language is ambiguous

You said you were looking for some mixed nuts?

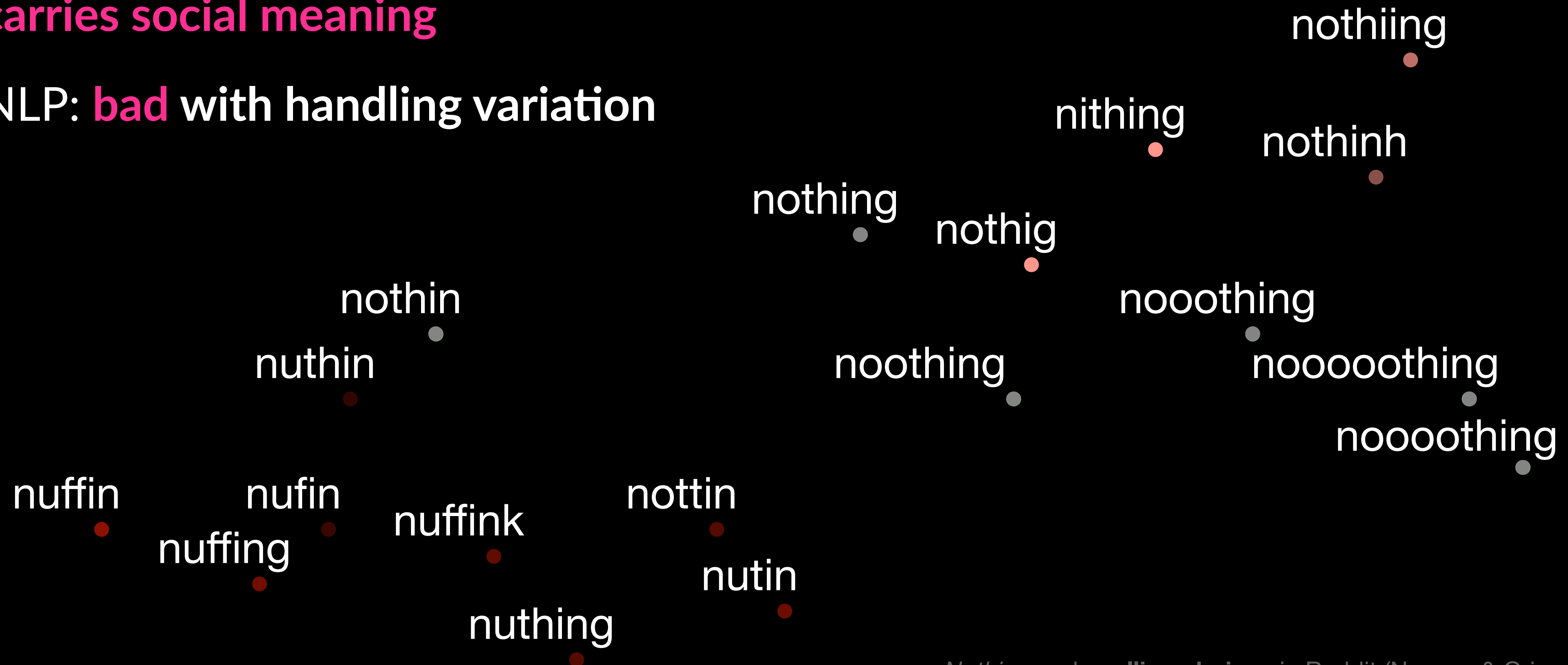# Language is full of variation

‣ The way we express a message **carries social meaning**

‣ NLP: **bad with handling variation**

You said nothing?

nothiing

nithing

nothinh

nothing

nothig

nothin

nooothing

nuthin

noothing

noooooothing

noooothing

nuffin      nufin

nottin

nuffink

nuffing

nutin

nuthing

*Nothing* and **spelling choices** in Reddit (Nguyen & Grieve, 2020)

# Language is dynamic and constantly changing

How to sunny-day Saturday in Seattle:
- ✅ pop out of bed and fling open the drapes
- ✅ brew coffee ☕ and grab your **go-cup**
- ✅ get outside asap
- ✅ dog walk, hike, run, bike, kayak, sail
- ✅ 🍺 soak up the sun in your favorite beer garden
- ✅ 😃👍✅

💬 2          🔁          ♡ 22          ✉
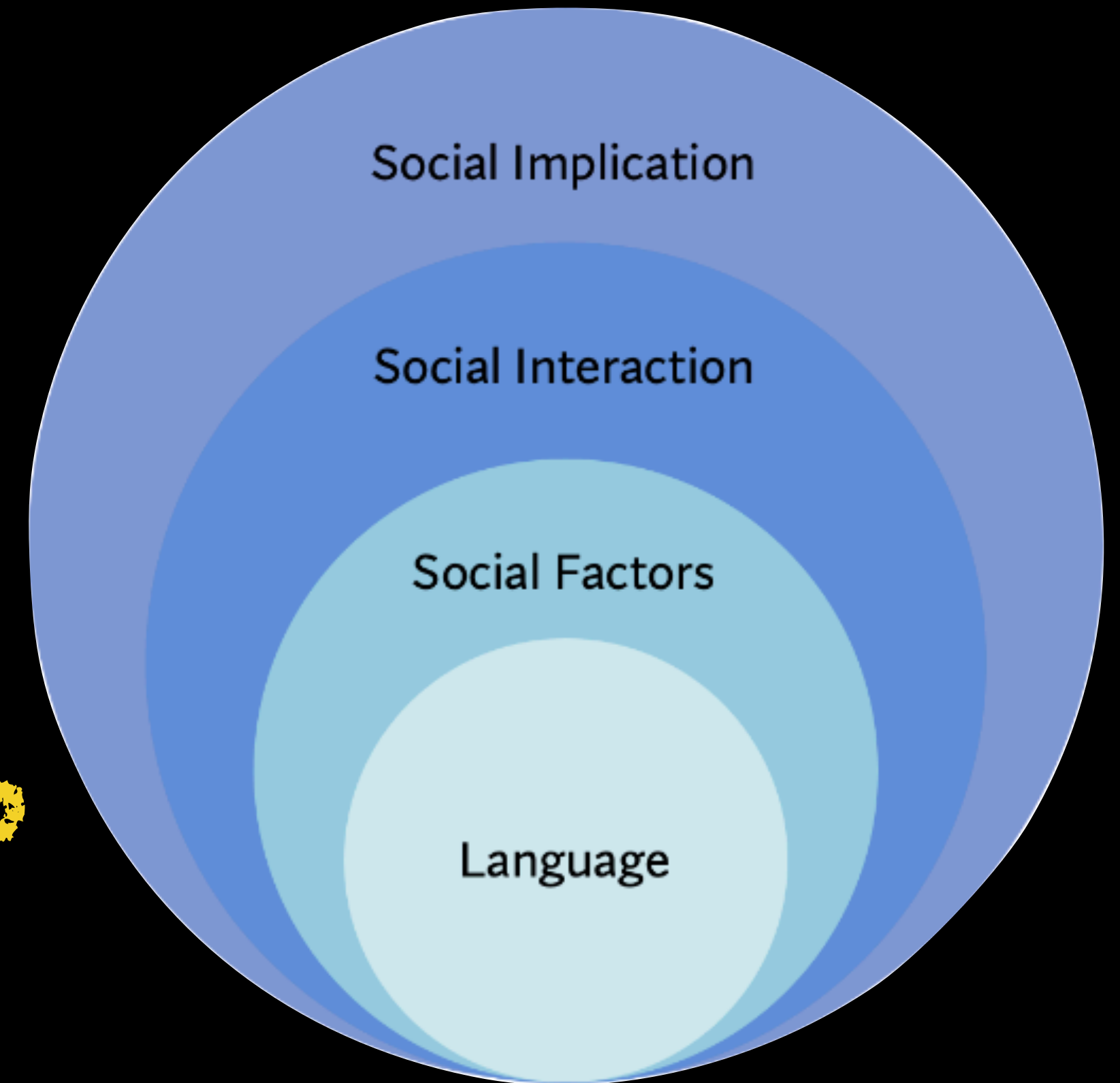
Cambridge Dictionary has revealed its word of the year for 2023 is 'hallucinate', as the term got a new additional definition relating to artificial intelligence (AI) producing false information.

# Language is for and by people

"The common misconception [is] that language use has primarily to do with words and what they mean. It doesn't. It has primarily to do with people and what they mean.

*Socially aware NLP*



Social Implication

Social Interaction

Social Factors

Language

## The Call for Socially Aware Language Technologies

**Diyi Yang**
Stanford University
diyiy@stanford.edu

**Dirk Hovy**
Bocconi University
mail@dirkhovy.com

**David Jurgens**
University of Michigan
jurgens@umich.edu
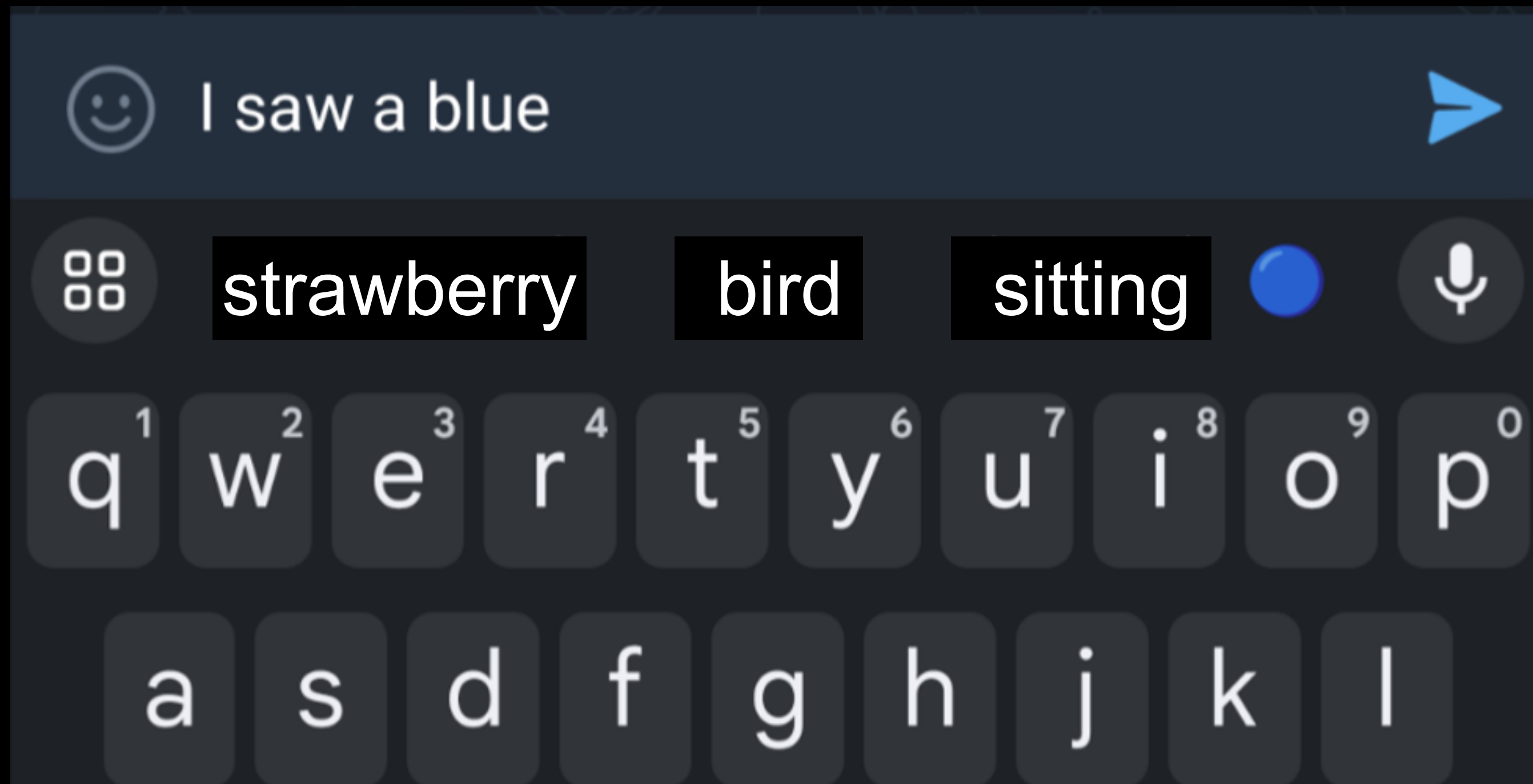
**Barbara Plank**
LMU Munich
bplank@cis.lmu.de

Slide credits: Diyi Yang

# What Can We Learn From Each Other?

# Roadmap

**1** Past: LLMs & Trust - How Did We Get There?

**2** Present: Trust Issues with LLMs

**3** Future: Trustworthy Human-Facing NLP
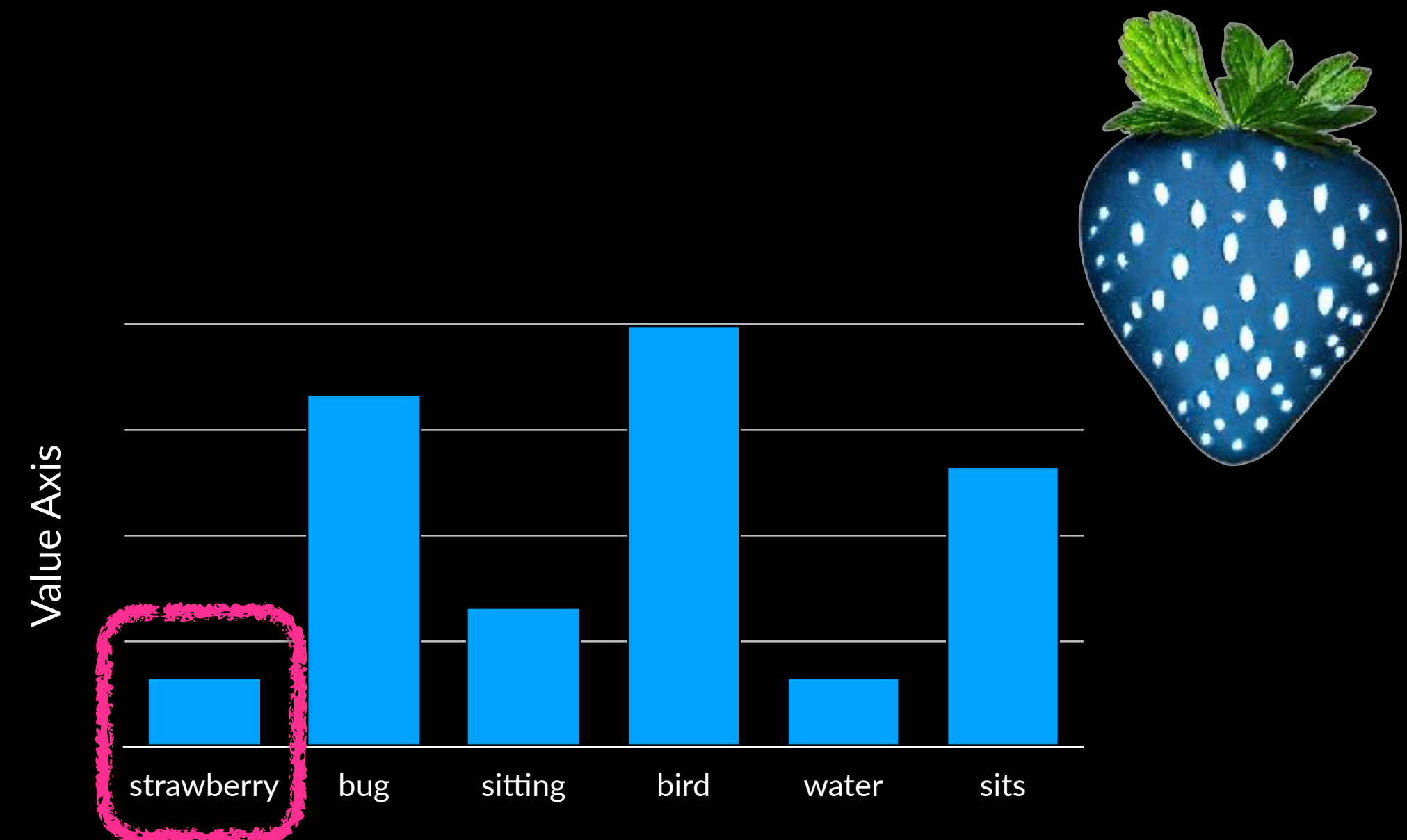
# A Language Model - The most likely text completion

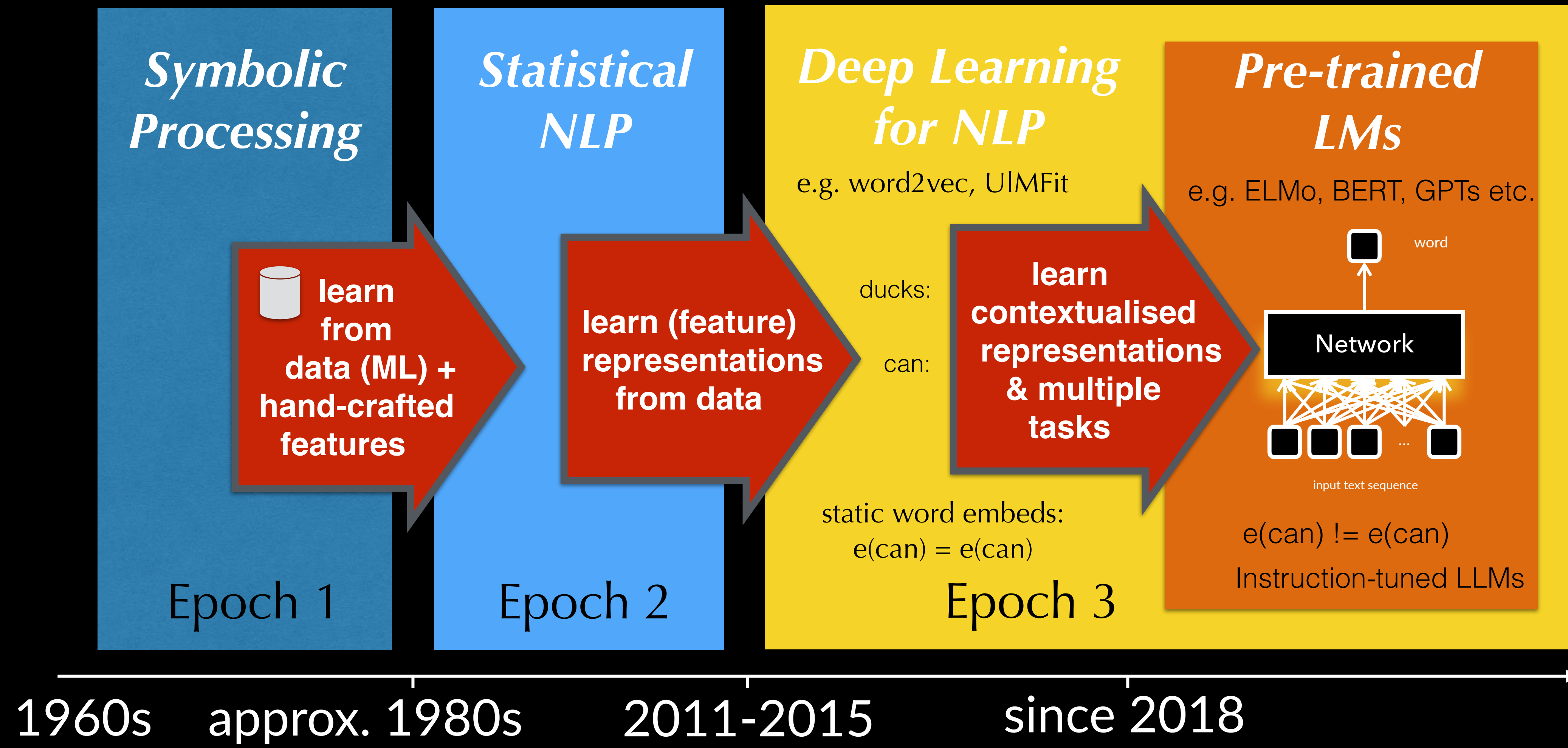‣ A LM computes the probability for a word given its previous words (=context)

strawberry: https://huggingface.co/spaces/stabilityai/stable-diffusion

# A Look Back - How Did We Get There?

# A Short NLP History



**Symbolic Processing**

learn from data (ML) + hand-crafted features

**Statistical NLP**

learn (feature) representations from data

**Deep Learning for NLP**

e.g. word2vec, UlMFit

ducks:

can:

learn contextualised representations & multiple tasks

static word embeds:
e(can) = e(can)

**Pre-trained LMs**

e.g. ELMo, BERT, GPTs etc.

word

Network

input text sequence

e(can) != e(can)

Instruction-tuned LLMs

Epoch 1          Epoch 2          Epoch 3

1960s    approx. 1980s    2011-2015    since 2018

# Gained Power - At What Cost?

2022-today: 💥 Explosion of LLMs

Power

LLM

Feature Engineering

Representation Learning

Trust

Knowledge about Model Input

⚠ Output:
Issues with factuality, bias, robustness, explainability

1990   2000   2010   2020

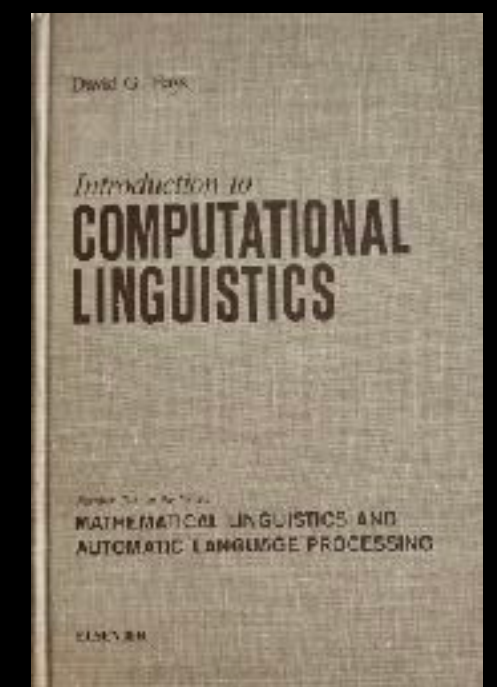Epoch 2: Statistical Processing

Epoch 3: Deep Learning (DL) for NLP

14

# What is trust?

# Trustworthiness - Working Definition

"Trust arises from **knowledge of origin** as well as from **knowledge of functional capacity**."
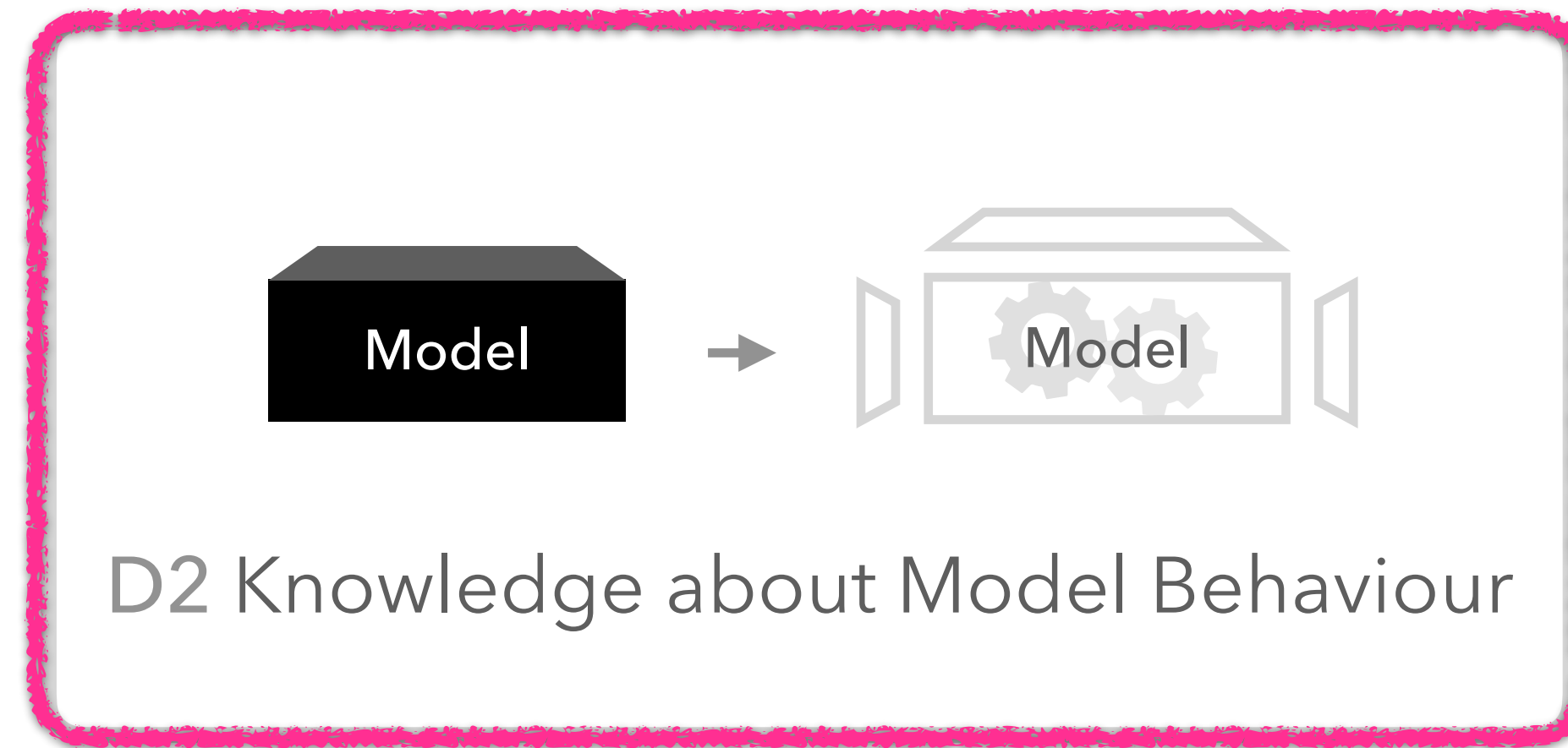


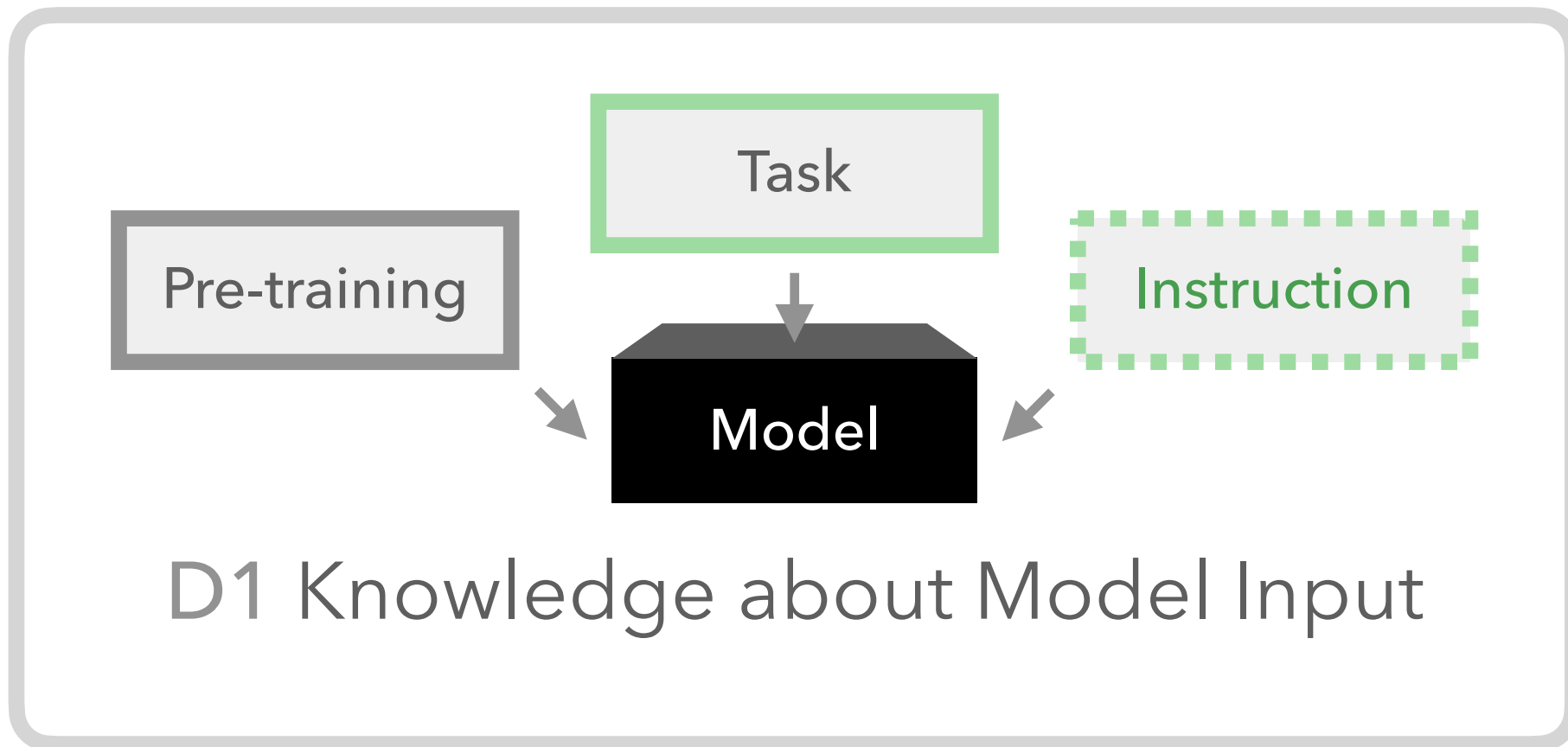Hays. Applications. ACL 1979.

# Roadmap

**1** Past: LLMs & Trust - How Did We Get There?

**2** Present: Trust Issues with LLMs

**3** Trustworthy Human-Facing NLP
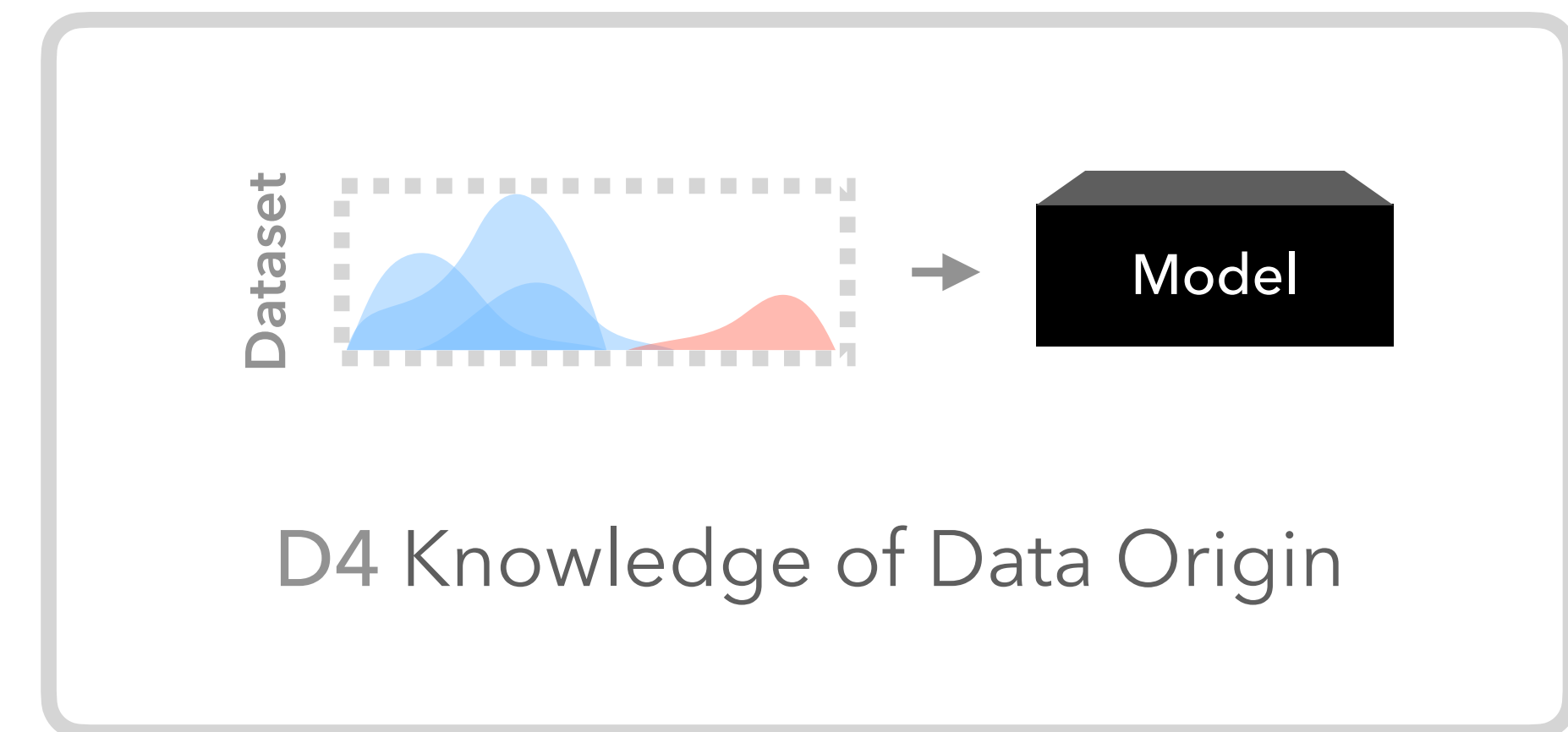
# Trust Issues with LLMs
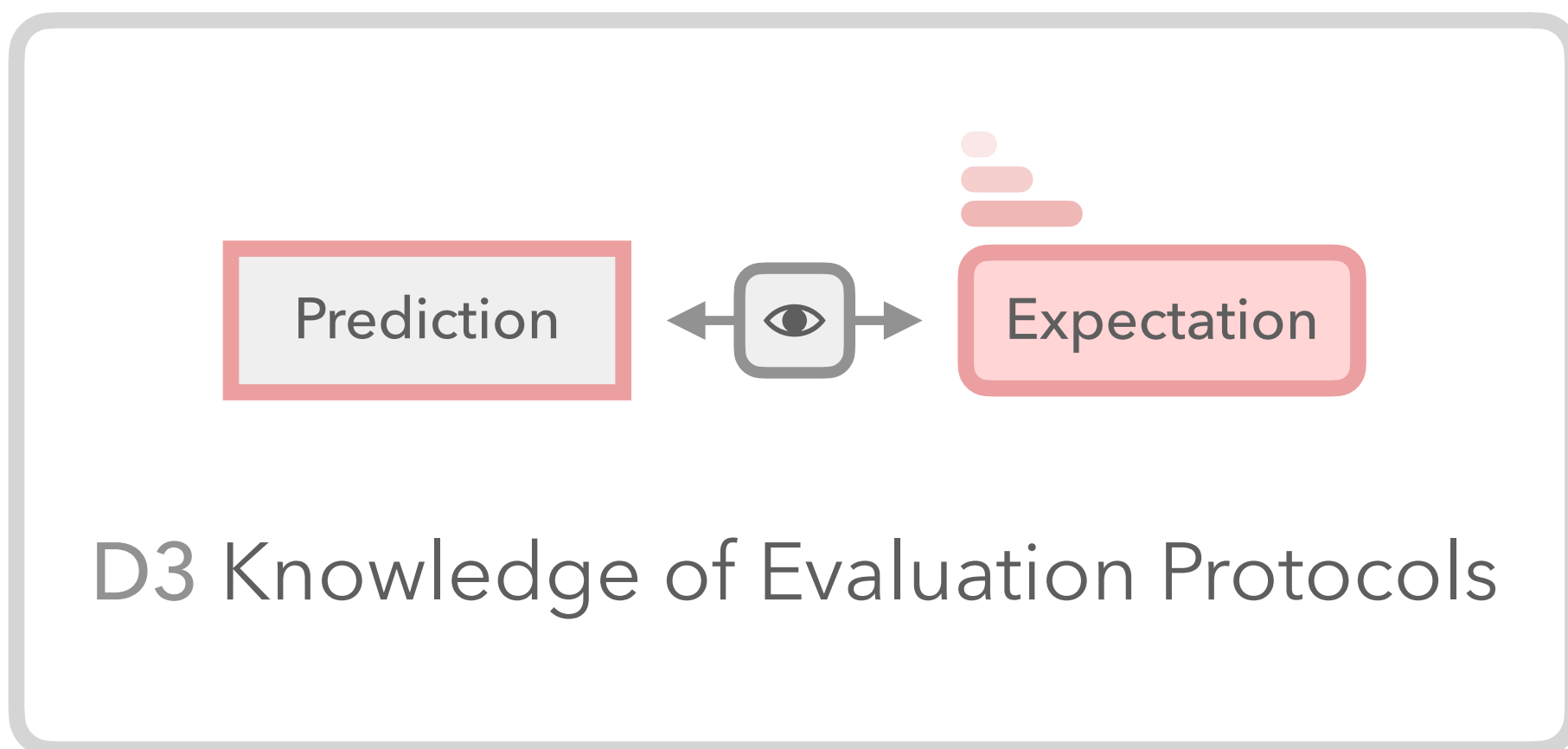
**Selected Research Examples**
**— Four Desiderata to Increase Trust —**

**D1 Knowledge about Model Input**

**D2 Knowledge about Model Behaviour**

Trust arises from **knowledge of origin** as well as from **knowledge of functional capacity**.

*Trustworthiness - Working Definition by David G. Hays, 1979*

**D3 Knowledge of Evaluation Protocols**

**D4 Knowledge of Data Origin**

Litschko*, Müller-Eberstein*, van der Goot, Weber-Genzel, Plank. Establishing Trustworthiness: Rethinking Tasks and Model Evaluation. EMNLP 2023.

# Model Behaviour: Does it Matter How we Prompt an LLM?

‣ ⚠️ Instability in prompting: Performance is highly sensitive to the linguistic variation of a prompt; prompts transfer poorly across datasets and models; LM perplexity dot not correlate well with model accuracy (open questions on connection data distribution and model beh

| prop. | | prompt |
|---|---|---|
| mood | inter. | Do you find this movie review positive? |
| | indic. | You find this movie review positive. |
| | imper. | Tell me if you find this movie review positive. |
| aspt. | active | Do you find this movie review positive? |
| | pass. | Is this movie review found positive? |
| tense | past | Did you find this movie review positive? |
| | pres. | Do you find this movie review positive? |
| | future | Will you find this movie review positive? |
| modality | can | Can you find this movie review positive? |
| | could | Could you find this movie review positive? |
| | may | May you find this movie review positive? |
| | might | Might you find this movie review positive? |
| | must | Must you find this movie review positive? |
| | should | Should you find this movie review positive? |
| | would | Would you find this movie review positive? |
| synonymy | apprai. | Do you find this movie appraisal positive? |
| | comm. | Do you find this movie commentary positive? |
| | criti. | Do you find this movie critique positive? |
| | eval. | Do you find this movie evaluation positive? |
| | review | Do you find this movie review positive? |

Table 1: Examples of variation of linguistic properties

The language of prompting:
What linguistic properties make a prompt successful?

Leidinger, van Rooij, Shutova, EMNLP 2023 Findings.

# Model Behaviour: How Well Do LLMs Deal with Ambiguity?

‣ ⚠️ LLMs and ambiguity is a major open problem: e.g. perform poorly at implicitly disambiguating entity types & biased towards preferred entity readings (influenced by entity popularity)
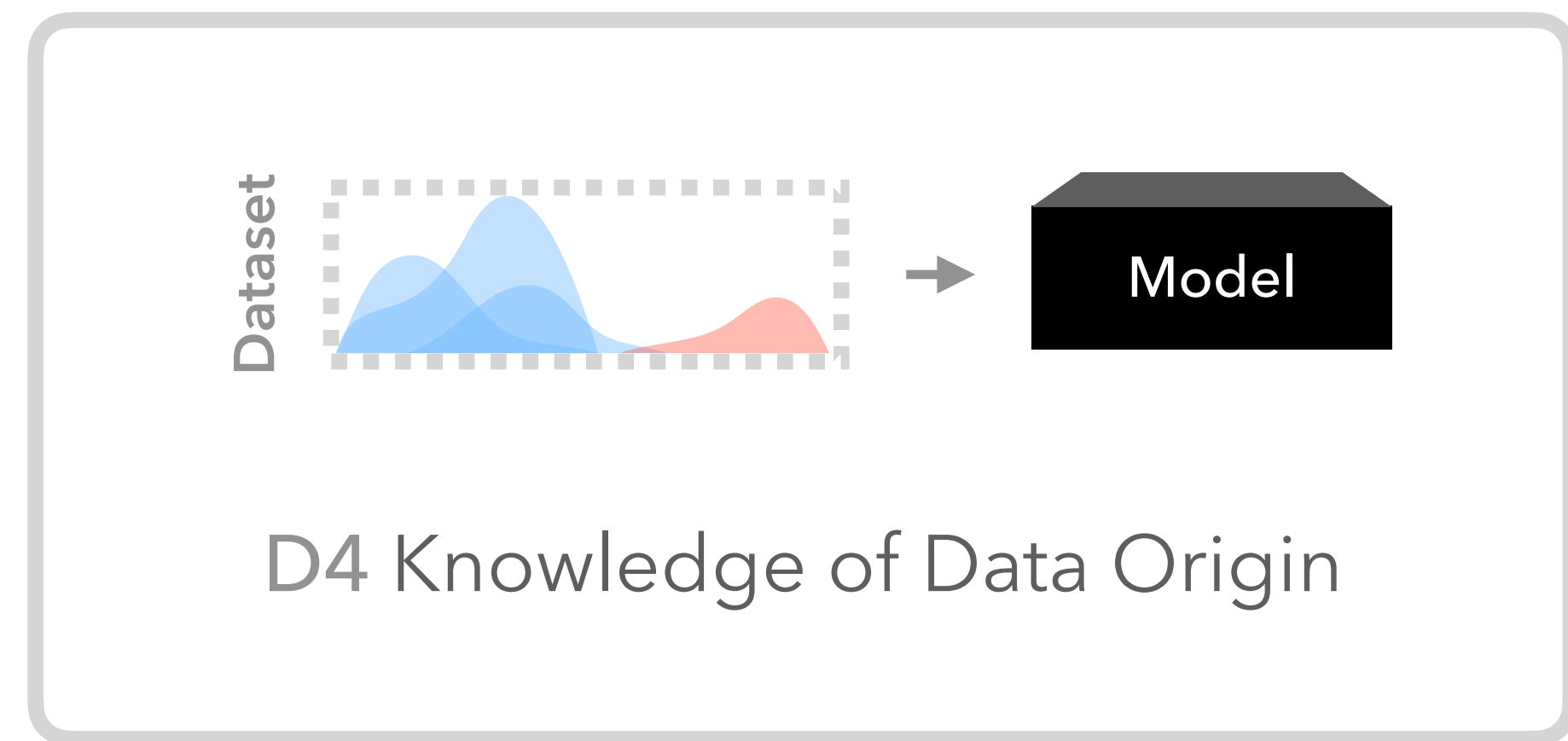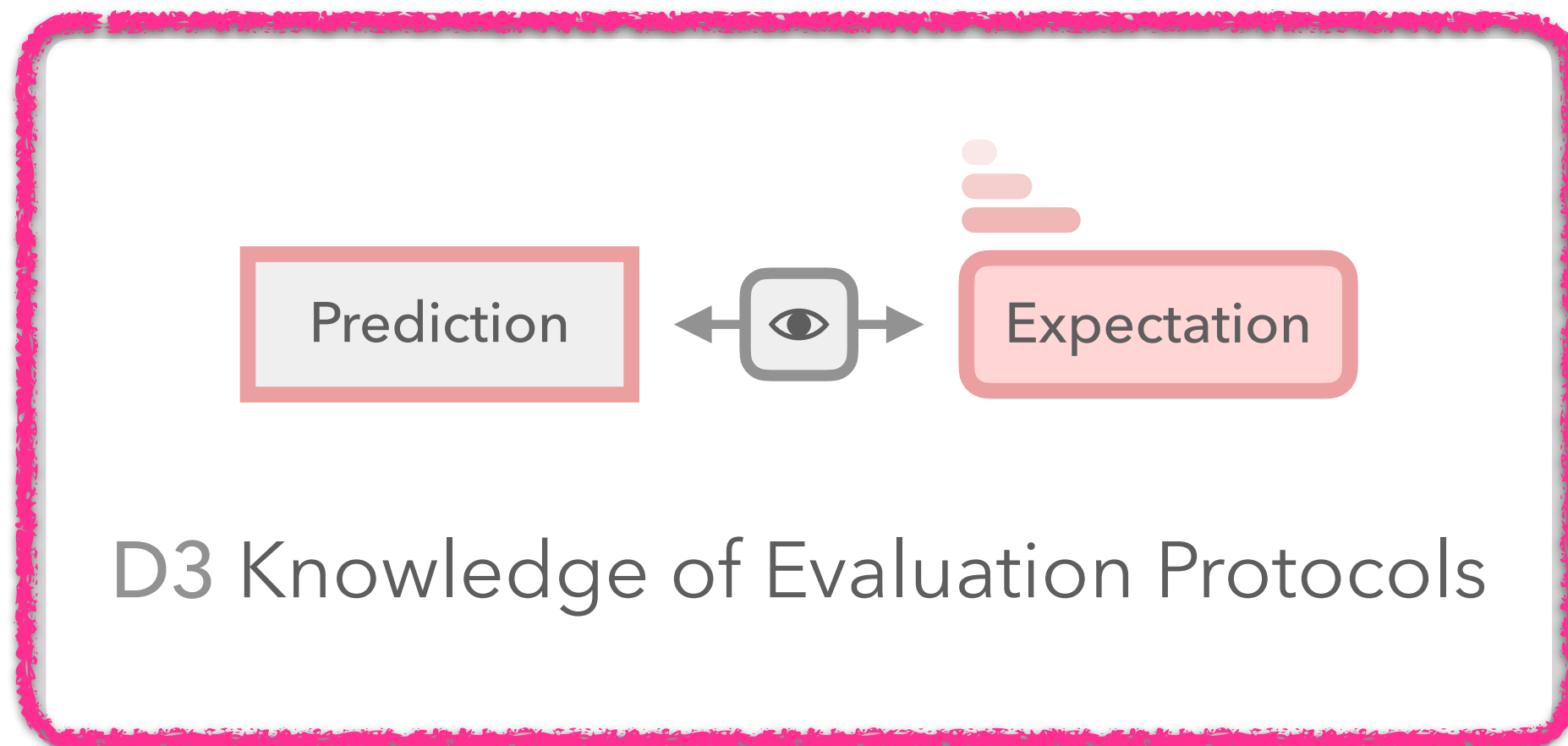


To Know or Not To Know? Analyzing Self-Consistency of Large Language Models under Ambiguity

Anastasiia Sedova[1,2*]    Robert Litschko[3,4*]    Diego Frassinelli[3]
Benjamin Roth[1,5]    Barbara Plank[3,4]

Sedova, Litschko et al. EMNLP 2024 Findings.

D1 Knowledge about Model Input

D2 Knowledge about Model Behaviour

Trust arises from **knowledge of origin** as well as from **knowledge of functional capacity.**

*Trustworthiness - Working Definition by David G. Hays, 1979*

D3 Knowledge of Evaluation Protocols

D4 Knowledge of Data Origin

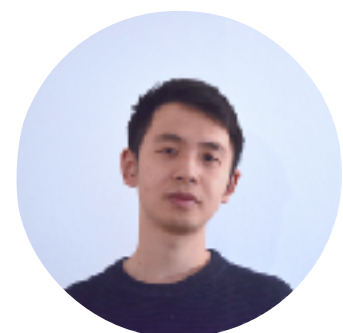# Multiple-Choice Question Answering (MCQA) Prompt Style

**General Instruction:** Please read the multiple-choice question below carefully and select ONE of the listed options and only give a single letter.

**Question:** The Web was effectively invented by Berners-Lee in which year?

**Options:**
A. 1991
B. 1980
C. 1989
D. 1993

**Answer:**

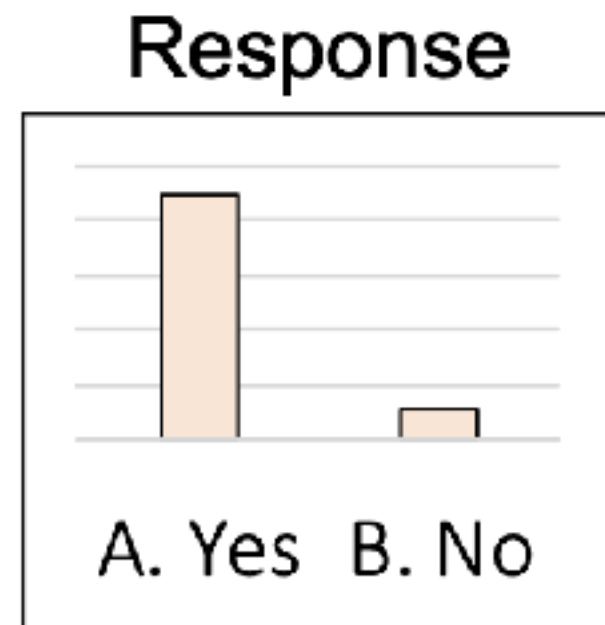Wang, Hu, Ma, Röttger, Plank. Look at the Text: Instruction-Tuned Language Models are More Robust Multiple Choice Selectors than You Think. COLM 2024.

▸ ⚠️ LLM's "A"-bias in MCQA responses

### Choice ordering 1

Question: In the past 12 months, has this person given birth to any children?
A. Yes
B. No
Answer:

Response
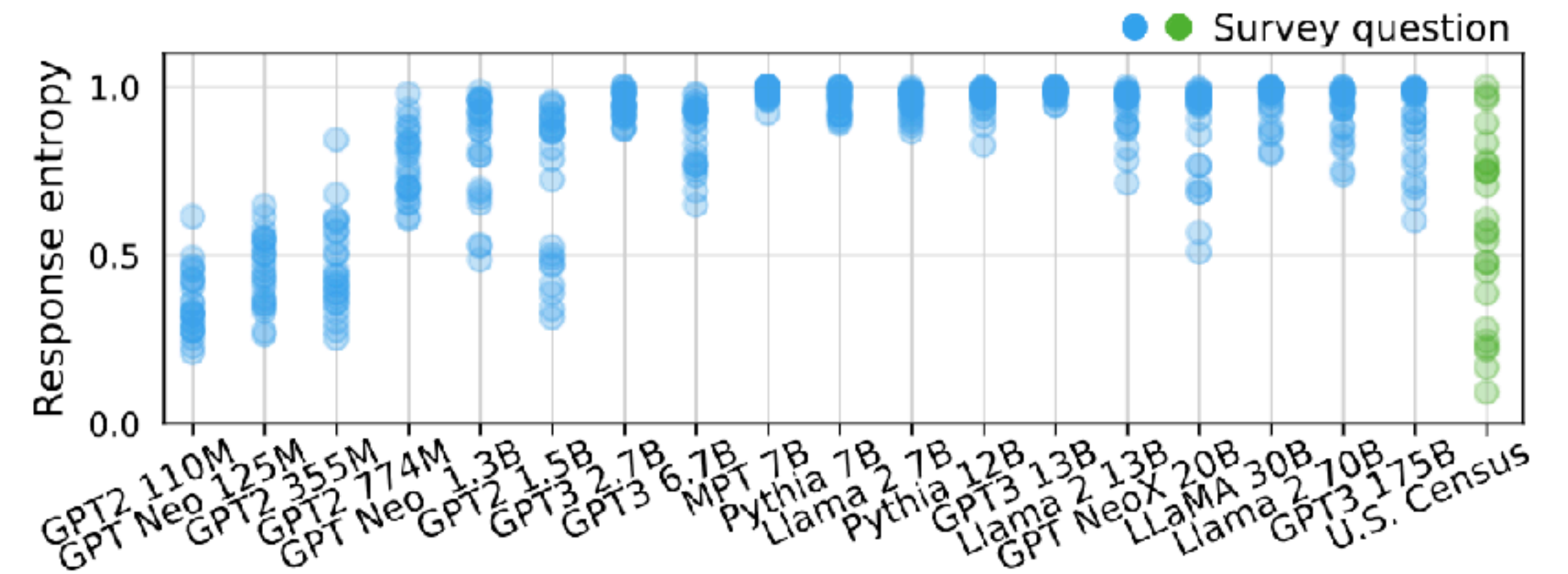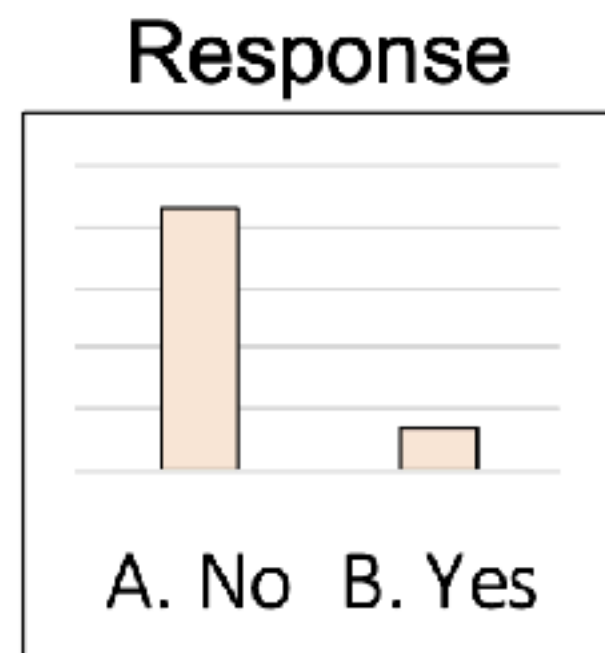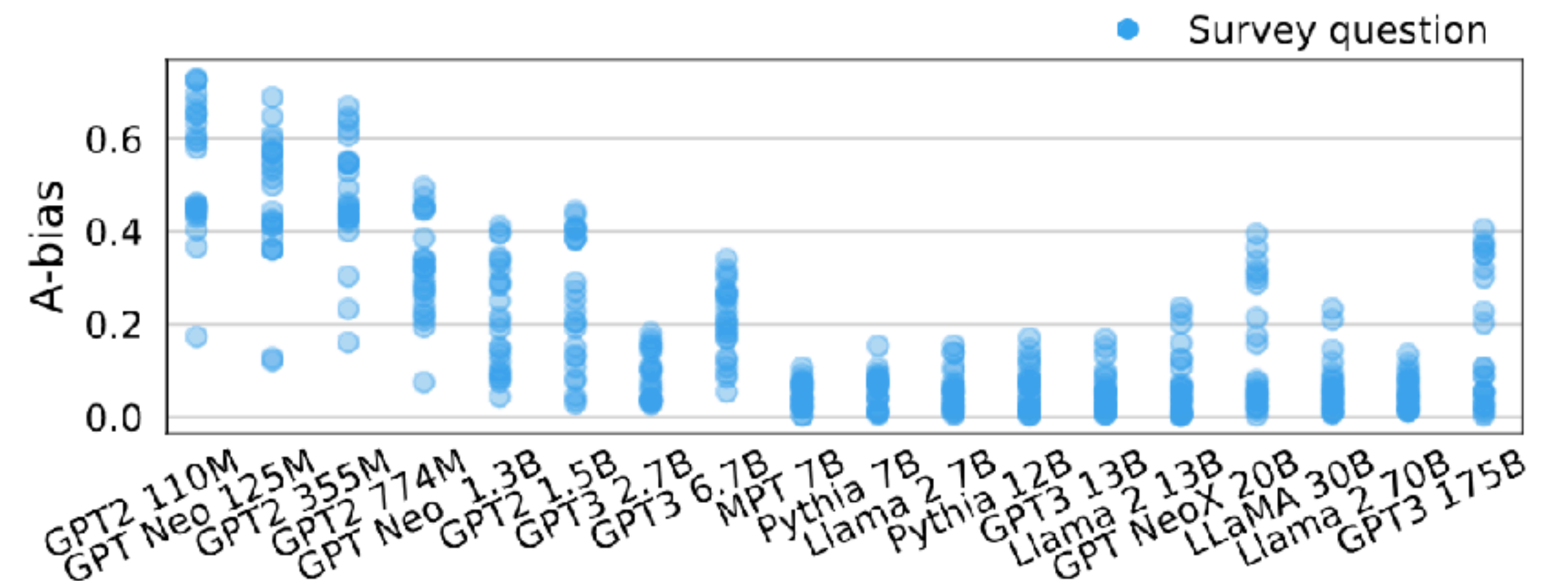
A. Yes   B. No

| P("A") | 0.82 | P("B") | 0.11 |
|---|---|---|---|

### Choice ordering 2

Question: In the past 12 months, has this person given birth to any children?
A. No
B. Yes
Answer:

Response

A. No   B. Yes

| P("A") | 0.80 | P("B") | 0.15 |
|---|---|---|---|



(a) Entropy of base models' responses.



(b) A-bias of base models' responses.

Dominguez-Olmedo, Hardt, Mendler-Dünner. Questioning the Survey Responses of Large Language Models. arXiv:2306.07951 2023.

▸ ⚠️ But "First-token log probs" do not match the text answers

General Instruction: Please read the multiple-choice question below carefully and select ONE of the listed options and only give a single letter.
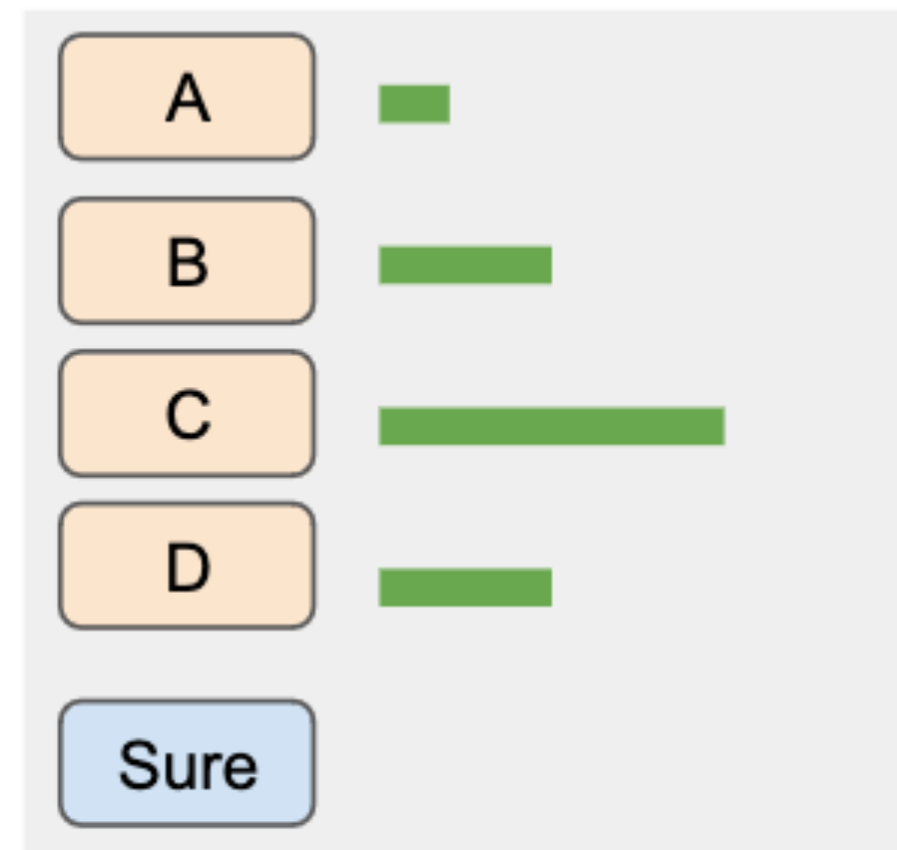
Question: The Web was effectively invented by Berners-Lee in which year?

Options:
A. 1991
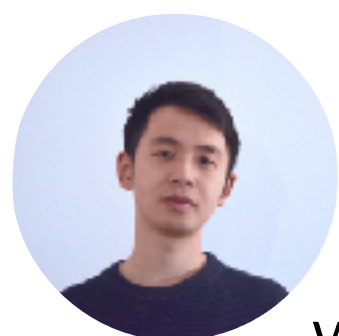B. 1980
C. 1989
D. 1993

Answer:

a. First Token Logits:

A
B
C
D
Sure

versus

b. Text Answer:

Sure! The answer is (A) 1991.

C

A

| Model (0-shot) | First Token | Text Answer |
|---|---|---|
| Gemma-7b-Inst | 30.2 | 50.8 |
| Llama2-7b-Chat | 34.9 | 43.1 |
| Llama2-13b-Chat | 40.2 | 47.6 |
| Mistral-7b-Inst-0.2 | 53.2 | 53.6 |

Performance on MMLU.

Wang, Ma, Hu, Weber-Genzel, Röttger, Kreuter, Hovy, Plank. "My Answer is C": First-Token Probabilities Do Not Match Text Answers in Instruction-Tuned Language Models. ACL 2024 Findings.
Wang, Hu, Ma, Röttger, Plank. Look at the Text: Instruction-Tuned Language Models are More Robust Multiple Choice Selectors than You Think. COLM 2024.

25

# Evaluation Protocols: Chain-of-Thought, Better and Sound?

‣ **Chain-of-Thought (CoT)** prompting ("Let's think step by step")

‣ ⚠️ But CoT is Not Sound!

| Model | Correct Answer | Sound Reasoning |
|---|---|---|
| Zephyr-7B-$\beta$ | $45.0 \pm 15.5$ | $25.0 \pm 10.5$ |
| Mistral-7B-Instruct | $55.0 \pm 10.0$ | $25.0 \pm 7.5$ |
| LLaMA-2-7B | $46.7 \pm 6.7$ | $0.0 \pm 0.0$ |
| LLaMA-2-13B | $40.0 \pm 8.2$ | $15.0 \pm 6.2$ |
| LLaMA-2-70B | $56.7 \pm 6.2$ | $31.7 \pm 9.7$ |

Statements:
1. There is a purple marble in the box if and only if there is a blue marble in the box.
2. There is a blue marble in the box if and only if there is an olive marble in the box.
3. Either there is an olive marble in the box or else there is a red marble in the box, but not both.
4. Either there is a red marble in the box or else there is a green marble in the box, but not both.

Conclusion: If there is a purple marble in the box then there is a green marble in the box.

CoT response (excerpt) - manually verified:

7. Let's consider the possibility that there is an olive marble in the box. If that's the case, then, according to Statement 2, there must be a blue marble in the box as well.
8. Now, using Statement 1, we know that if there is a blue marble in the box, then there must be a purple marble in the box as well.
9. However, Statement 3 tells us that if there is an olive marble in the box, then there cannot be a red marble in the box. But, Statement 4 tells us that if there is a red marble in the box, then there cannot be a green marble in the box.

Supposition following (sound).

10. Since we cannot have both an olive marble and a red marble in the box, and we cannot have both a red marble and a green marble in the box, it seems that the only possibility left is that there is an olive marble in the box, and therefore, a blue marble, a purple marble, but no red or green marbles.

Invalid inference. ✗

11. However, the conclusion states that "If there is a purple marble in the box then there is a green marble in the box", which is not logically true based on the given statements.
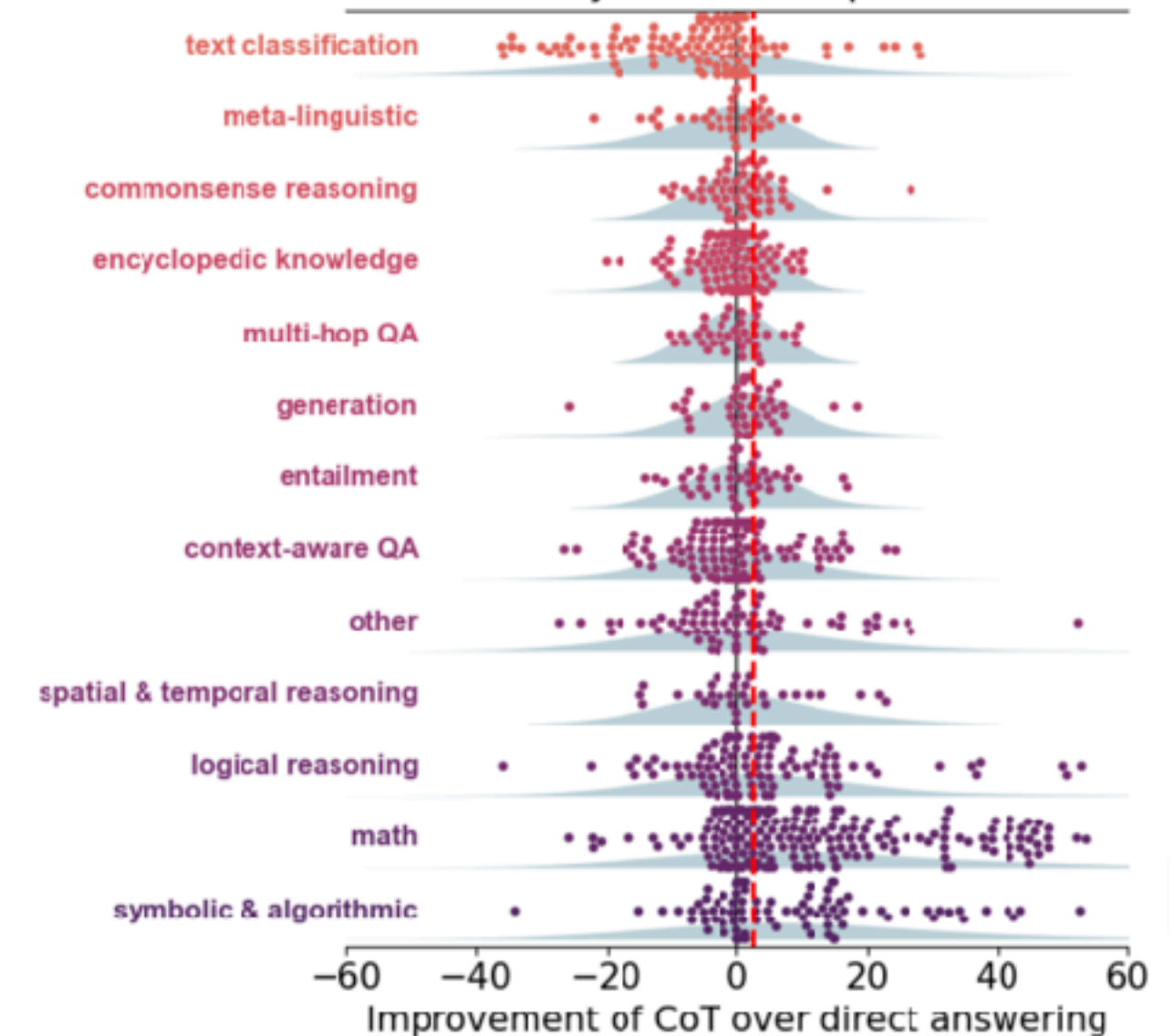
Conclusion (invalid).

Conclusion: False.

Final answer (incorrect).

Meta-analysis of CoT improvements

(categories) text classification, meta-linguistic, commonsense reasoning, encyclopedic knowledge, multi-hop QA, generation, entailment, context-aware QA, other, spatial & temporal reasoning, logical reasoning, math, symbolic & algorithmic

Improvement of CoT over direct answering

Mondorf & Plank. Comparing Inferential Strategies of Humans and Large Language Models in Deductive Reasoning. ACL 2024.
Sprague et al. 2024. To CoT or not to CoT? Chain of Thought helps mainly on math and symbolic reaoning. https://arxiv.org/abs/2409.12183

# Evaluation Protocols: Can LLMs Replace Humans Judges?
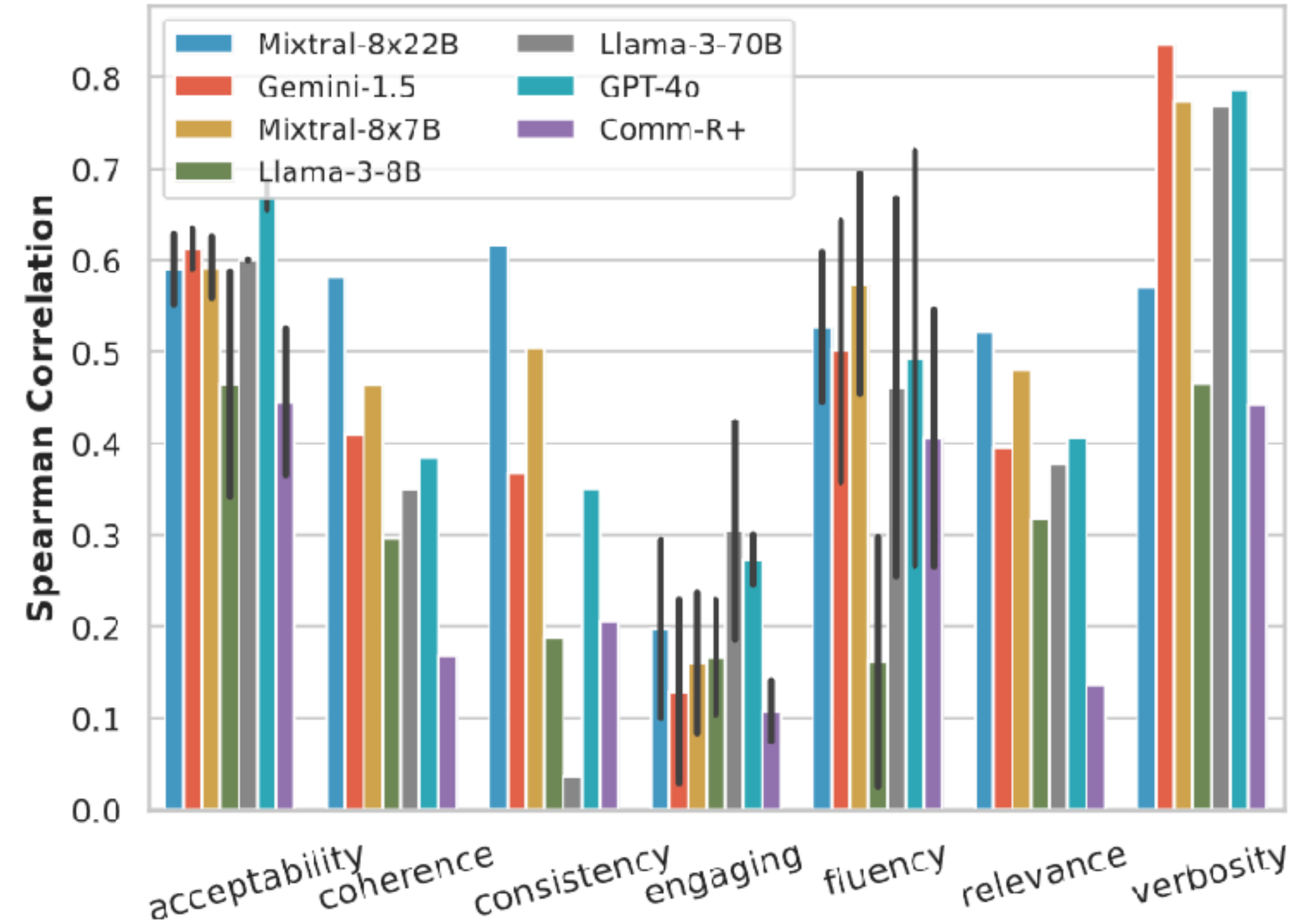
‣ ⚠️ A lot of variability in LLM outputs

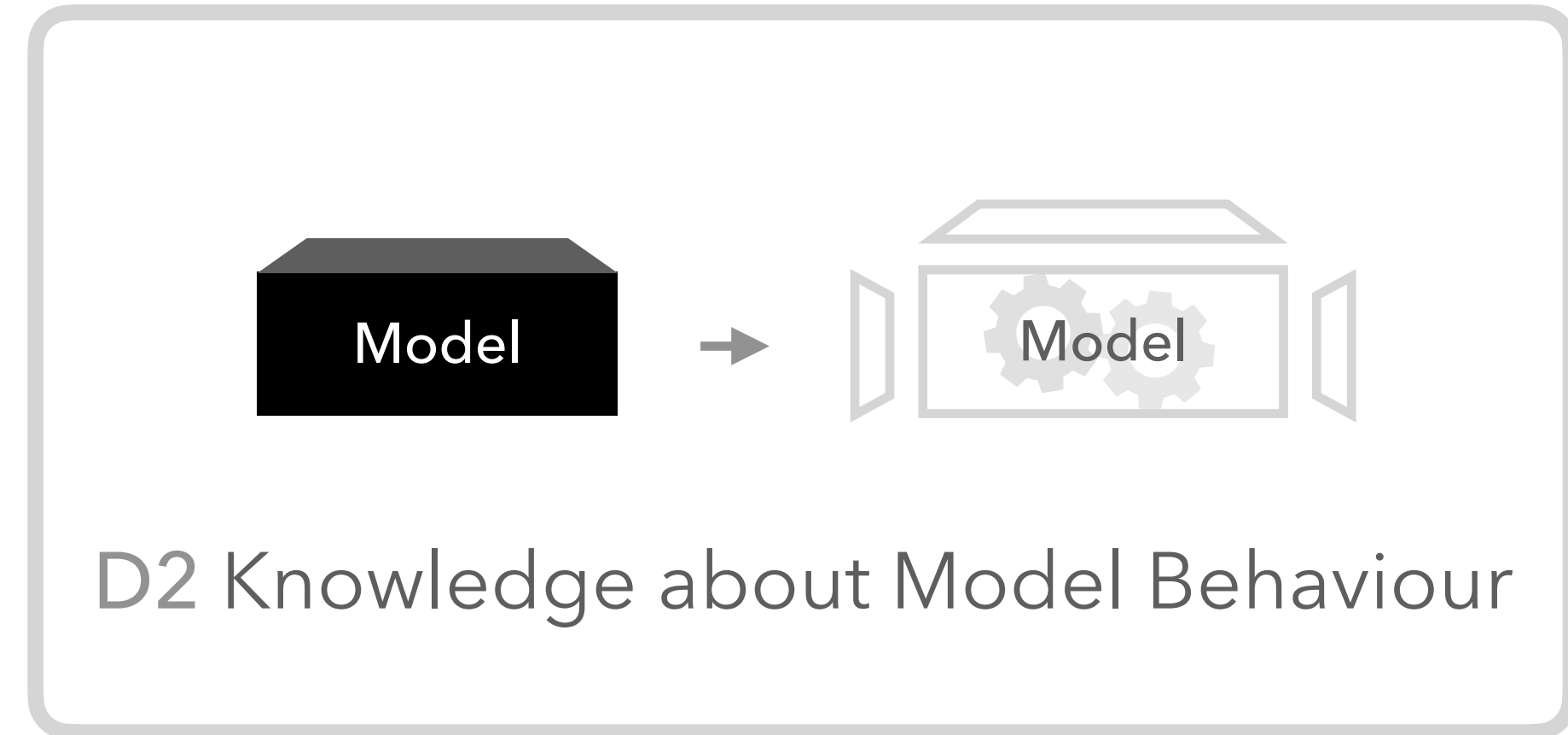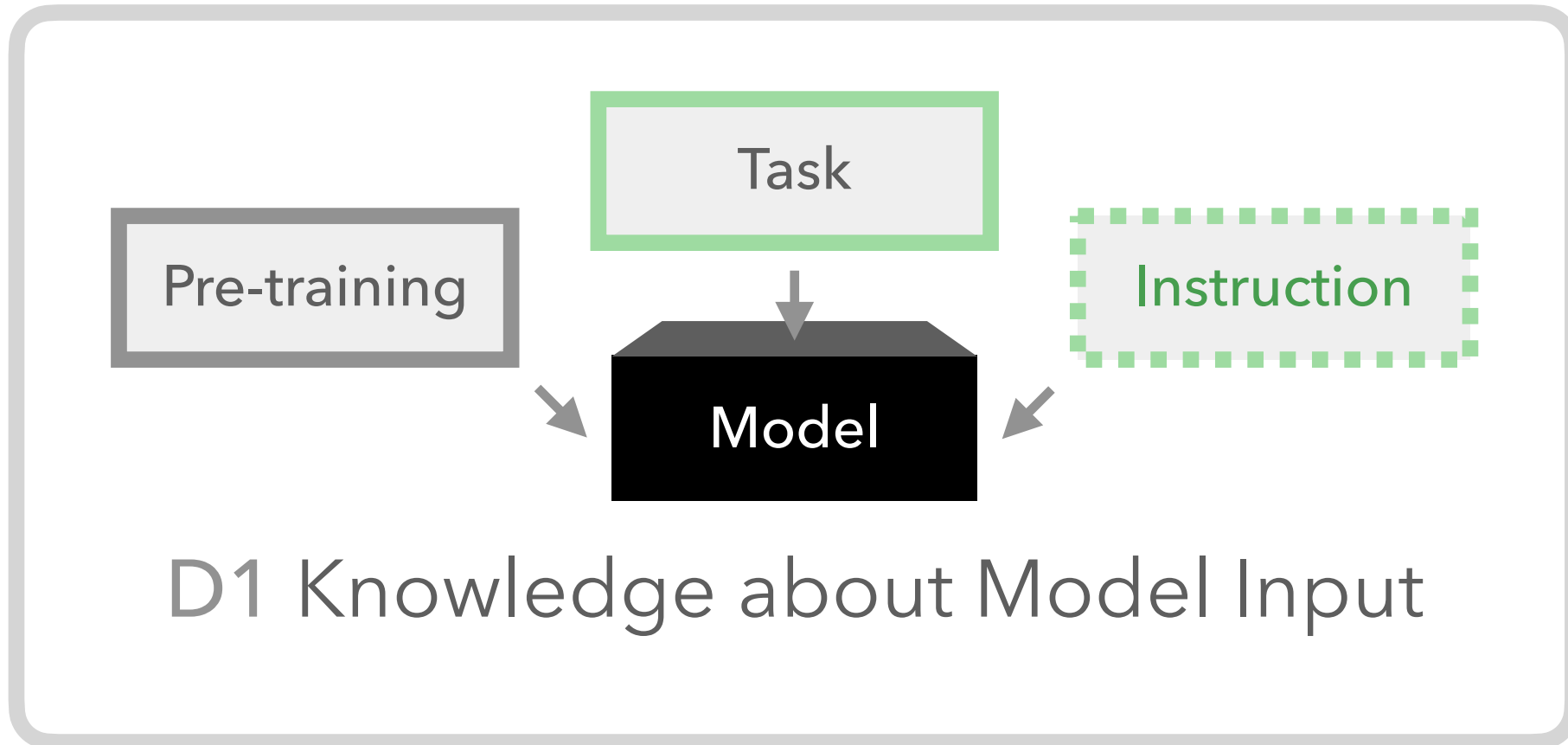‣ LLMs are not ready yet to replace human judges - not even GPT-4o:

**E.g. Plausibility:**

**Humans Coders vs Models:**



*Instruction*: On a scale of 1 (very unlikely) to 5 (very likely), how plausible is it that the last response belongs to the dialogue?

**A**: Made it all the way through four years of college playing ball but
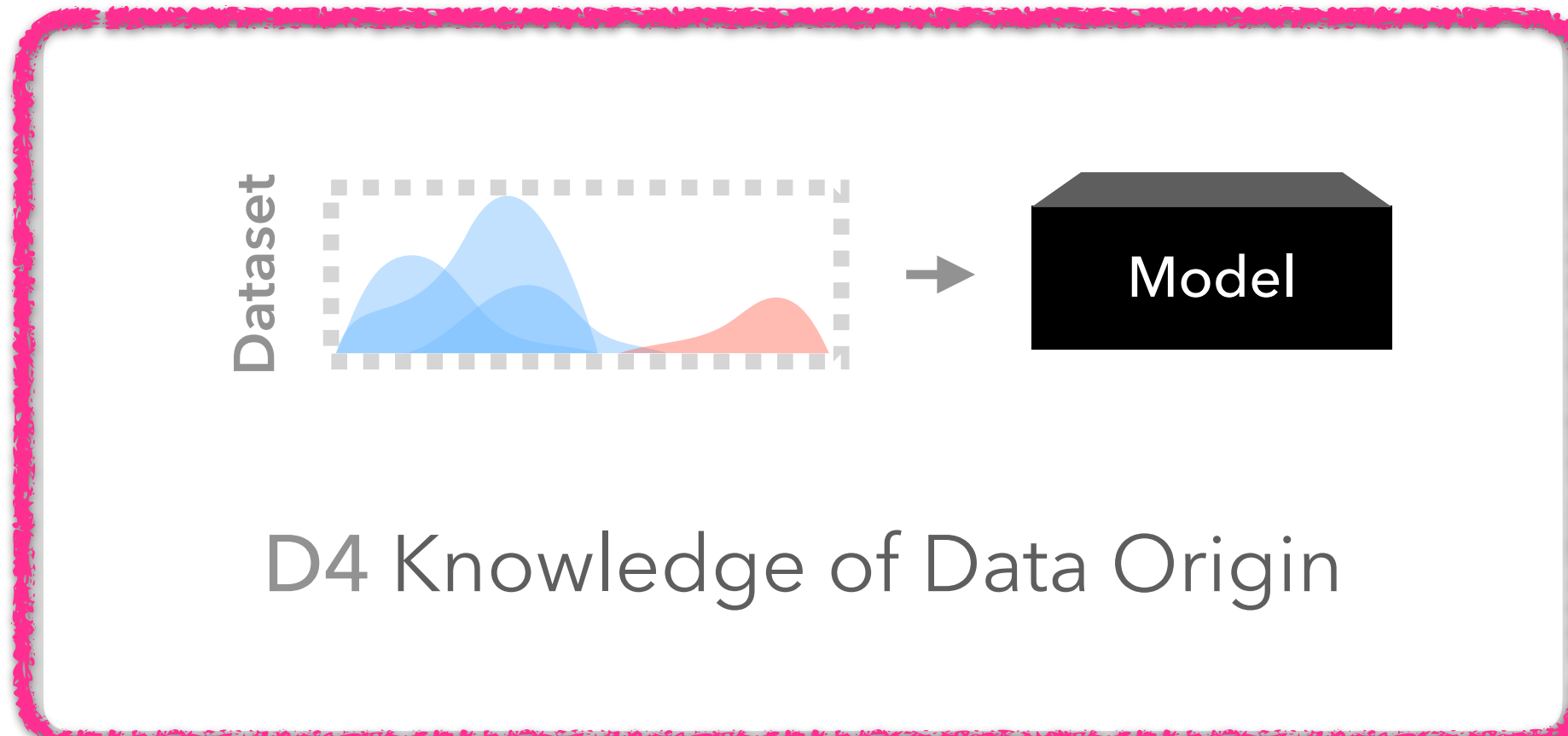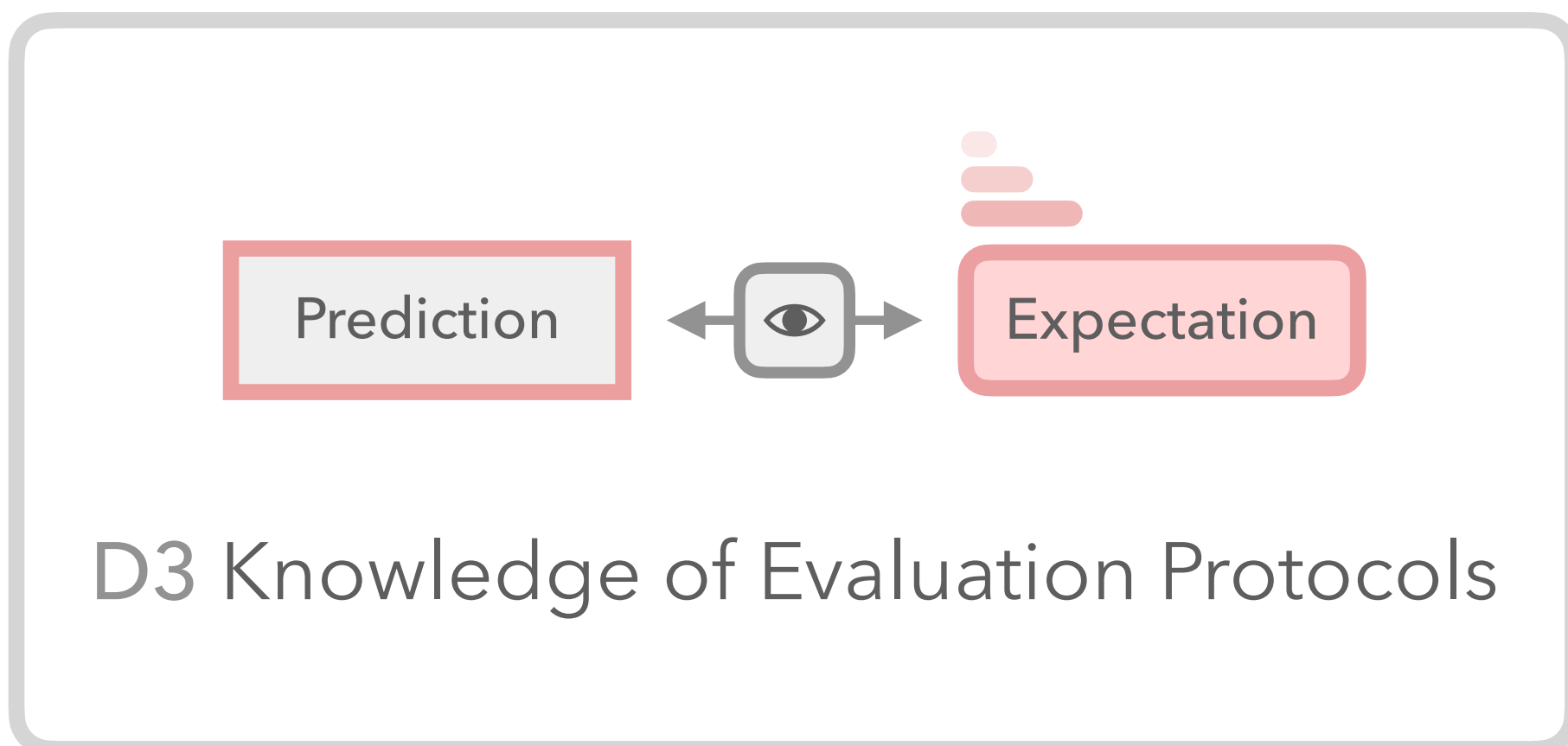**B**: I also like The Cosby Show

Bavaresco, Bernardi, Bertolazzi, Elliott, Fernandez, Gatt, Ghaleb, Giulianelli, Hanna, Koller, Martins, Mondorf, Neplenbroek, Pezzelle, Plank, Schlangen, Suglia, Surikuchi, Takmaz, Testoni. LLMs instead of Human Judges? A Large Scale Empirical Study across 20 NLP Evaluation Tasks. arXiv:2406.18403 2024.

D1 Knowledge about Model Input
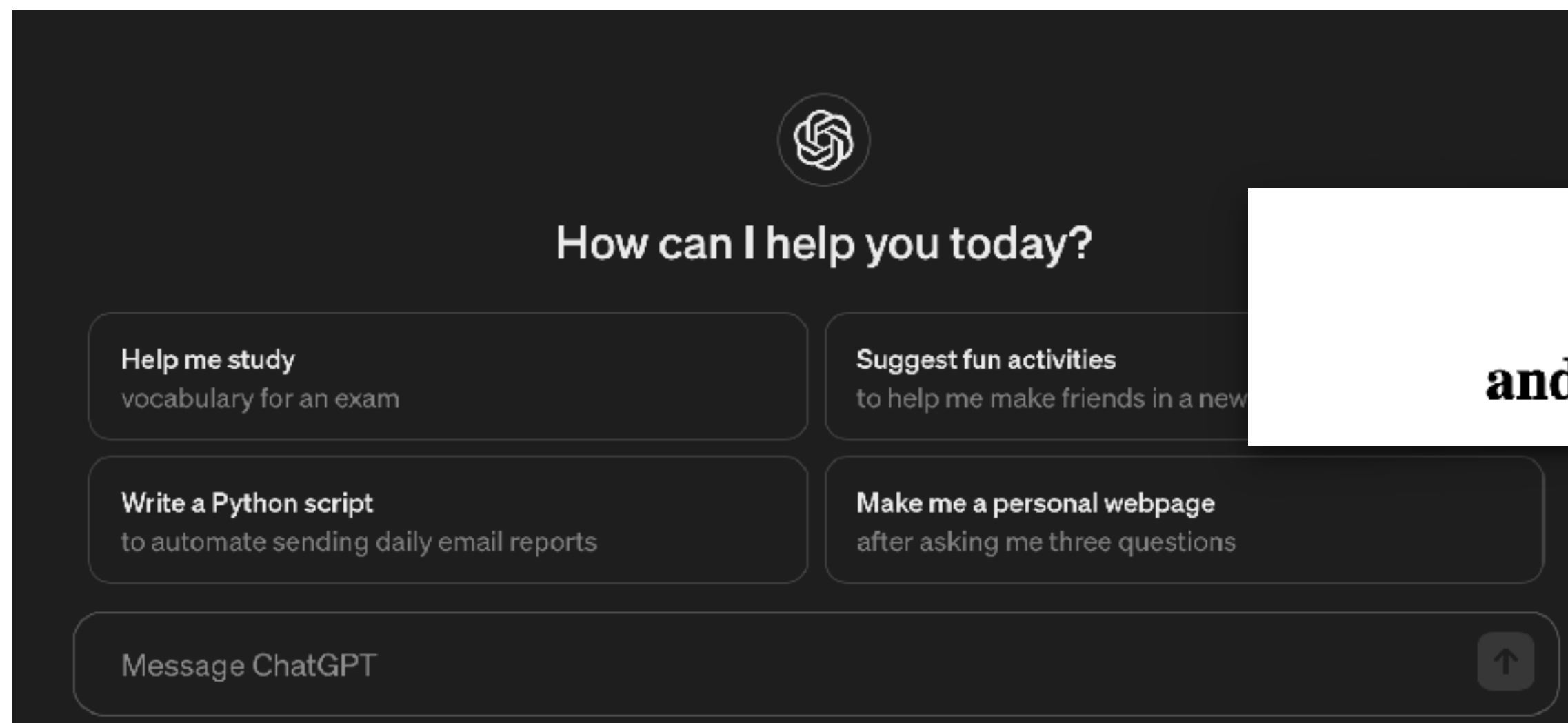
D2 Knowledge about Model Behaviour

Trust arises from **knowledge of origin** as well as from **knowledge of functional capacity.**

*Trustworthiness - Working Definition by David G. Hays, 1979*

D3 Knowledge of Evaluation Protocols

D4 Knowledge of Data Origin

# Data Origin: (Indirect) Contamination & Need for Transparency

- ⚠️ Too little transparency of what went into the training data of an LLM

- ⚠️ Indirect data leakage: continuously provided by users (e.g. via OpenAI's the web interface)





Balloccu, Schmidtová, Lango, Dušek. EACL 2024.

- ➡️ increasing efforts for transparency on training data & pre-processing, e.g.:
  - **PILE** (Gao et al., 2020)
  - **Dolma** (Soldini et al., 2024 ACL **best** paper award)

# Growing Importance of
# Data Quality > Data Quantity

# The "it" in AI models is the dataset - talk by Thom Wolf 🤗

## The "it" in AI models is the dataset.

*Posted on June 10, 2023 by jbetker*

I've been at OpenAI for almost a year now. In that time, I've trained a **lot** of generative models. More than anyone really has any right to train. As I've spent these hours observing the effects of tweaking various model configurations and hyperparameters, one thing that has struck me is the similarities in between all the training runs.

It's becoming awfully clear to me that these models are truly approximating their datasets to an incredible degree. What that means is not only that they learn what it means to be a dog or a cat, but the interstitial frequencies between distributions that don't matter, like what photos humans are likely to take or words humans commonly write down.

What this manifests as is – trained on the same dataset for long enough, pretty much every model with enough weights and training time converges to the same point. Sufficiently large diffusion conv-unets produce the same images as ViT generators. AR sampling produces the same images as diffusion.

This is a surprising observation! It implies that model behavior is not determined by architecture, hyperparameters, or optimizer choices. It's determined by your dataset, nothing else. Everything else is a means to an end in efficiently delivery compute to approximating that dataset.

Then, when you refer to "Lambda", "ChatGPT", "Bard", or "Claude" then, it's not the model weights that you are referring to. It's the dataset.

https://nonint.com/2023/06/10/the-it-in-ai-models-is-the-dataset/

31

# Evidence from a talk by Sara Hooker ✈

| Model | Size (# Parameters) | Training Tokens |
|---|---|---|
| LaMDA (Thoppilan et al., 2022) | 137 Billion | 168 Billion |
| GPT-3 (Brown et al., 2020) | 175 Billion | 300 Billion |
| Jurassic (Lieber et al., 2021) | 178 Billion | 300 Billion |
| Gopher (Rae et al., 2021) | 280 Billion | 300 Billion |
| MT-NLG 530B (Smith et al., 2022) | 530 Billion | 270 Billion |
| Chinchilla | 70 Billion | 1.4 Trillion |

- Recent work suggests smaller amounts of higher quality data remove the need for a larger model.
- This suggest larger models may just be compensating for problems in the data pipeline.

# Roadmap

**1** Past: LLMs & Trust - How Did We Get There?

**2** Present: Trust Issues with LLMs

**3** Trustworthy Human-Facing NLP

# Name the object

34

# Name the object



cake (53), food (19), bread (8), burger (6), dessert (6), snacks (3), muffin (3), pastry (3)

# Lora Aroyo's NeurIPS 2023 Keynote



Is there a **SMILE** in this image?

YES but …

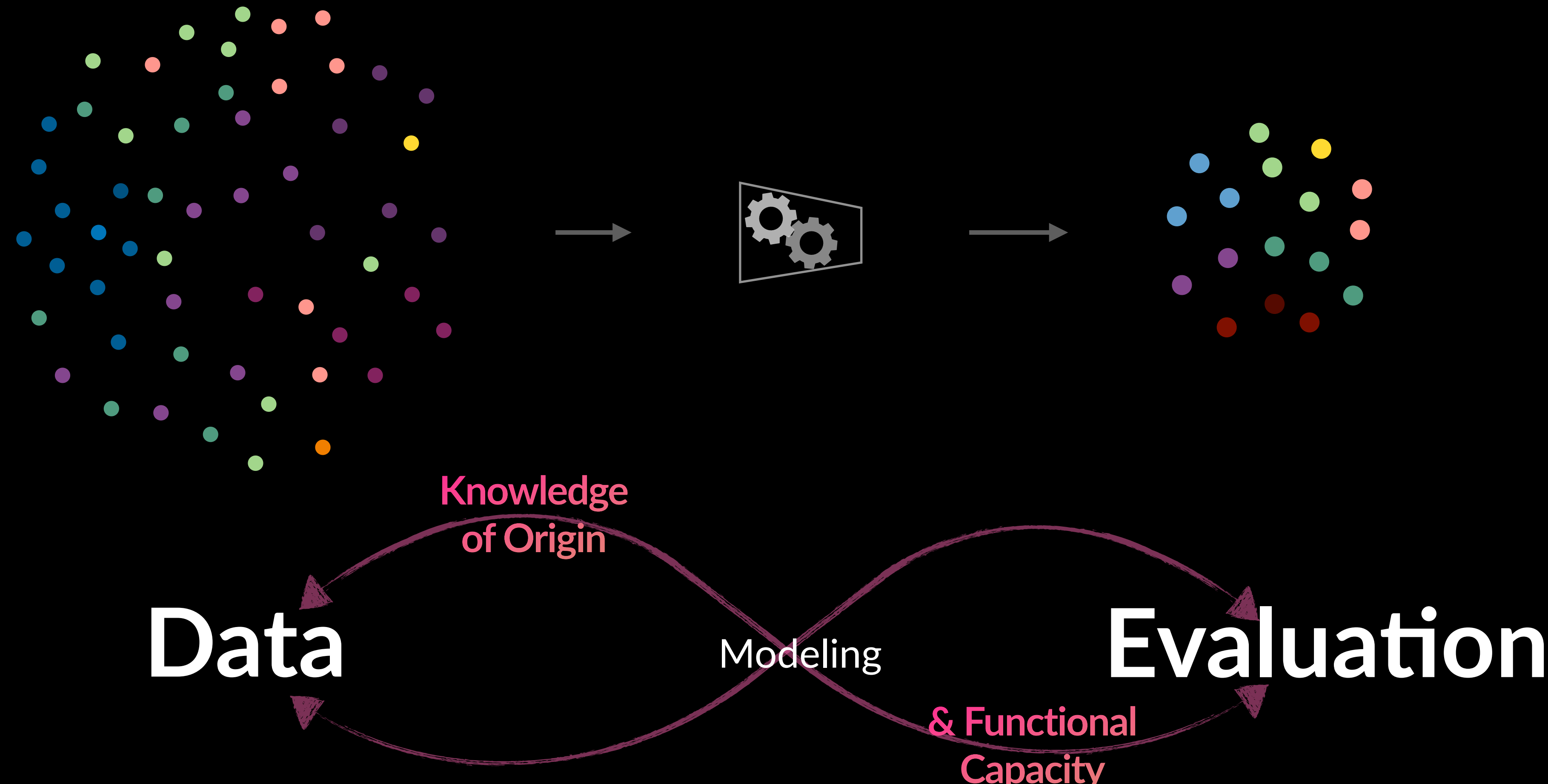| Canada | | |
|---|---|---|
| YES | NO | DNK |
| 40% | 40% | 20% |

| India | | |
|---|---|---|
| YES | NO | DNK |
| 70% | 30% | 0 |

| USA | | |
|---|---|---|
| YES | NO | DNK |
| 50% | 0 | 50% |

https://slideslive.com/39015341/the-many-faces-of-responsible-ai?ref=speaker-55217

36

# Human Label Variation

Data

Modeling

Evaluation

Knowledge of Origin

& Functional Capacity

# Many open questions

- How does Human Label Variation interact with Socially Aware NLP? (Specific tasks and broadly Values, Opinions, Attitude and LLMs)

- How much Variation is there Within-Human (Coder) Label Variation?

- How does Human Label Variation relate to Model Uncertainty?

- Is All Variation Plausible? How to tease apart Error from Plausible Label Variation?

- Do we Need More Labels or More Cases (for Eval/Train)? Data Quality vs Quantity

- When to take a description vs a prescriptive approach (Röttger et al. 2022) to annotation?

learning from less but higher quality data?

active learning (how to sample. Instances for labelling)

What goes into epresentativeness and quality of data

model uncertainty

*Human Label Variation*
**- many exciting connections -**

human values and LLM alignment (e.g. Durmus et al., 2024)

LLMs that react as humans do

statistics and data-generation process

# NLP & Survey Research

🙌

Thank you!