

# Bridging Human and LLM Annotations for Statistically Valid Computational Social Science

**Kristina Gligorić**

Assistant Professor, Computer Science Department

Data Science Institute, Center for Language and Speech Processing

Johns Hopkins University



JOHNS HOPKINS  
UNIVERSITY

# Yay data!



**Explosion of Text Data**



**Opportunities for New Insights**



**Unlocking New Knowledge**

# Real-world social science questions

- What is the impact of COVID-19 vaccine online misinformation on vaccination intent?
- Does air pollution lower people's expressed happiness on social media?
- Does negativity influence online news consumption?

# Real-world social science questions

- What is the impact of COVID-19 vaccine online misinformation on vaccination intent?
- Does air pollution lower people's expressed happiness on social media?
- Does negativity influence online news consumption?

Finding an answer requires using large textual datasets (e.g., social media posts) and needs **data labeling**.

# Real-world social science questions

- What is the impact of COVID-19 vaccine online misinformation on vaccination intent?
  - Annotation: “Does this social media post contain misleading claims?”
  - To estimate: whether people who see misinformation online report lower intent to vaccinate

## nature human behaviour

[Explore content](#) ▾ [About the journal](#) ▾ [Publish with us](#) ▾

[nature](#) > [nature human behaviour](#) > [articles](#) > [article](#)

Article | Published: 05 February 2021

### **Measuring the impact of COVID-19 vaccine misinformation on vaccination intent in the UK and USA**

[Sahil Loomba](#), [Alexandre de Figueiredo](#) , [Simon J. Piatek](#), [Kristen de Graaf](#) & [Heidi J. Larson](#) 

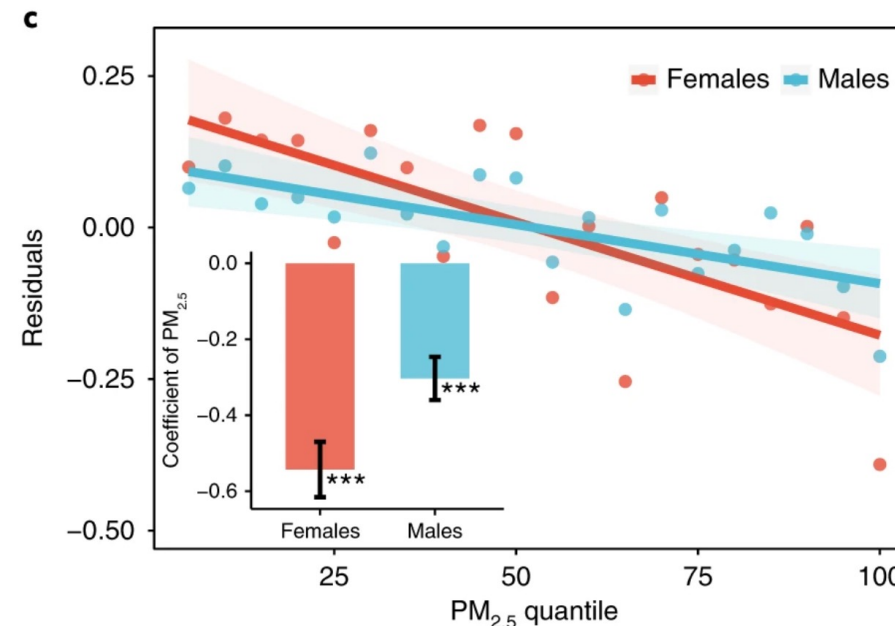
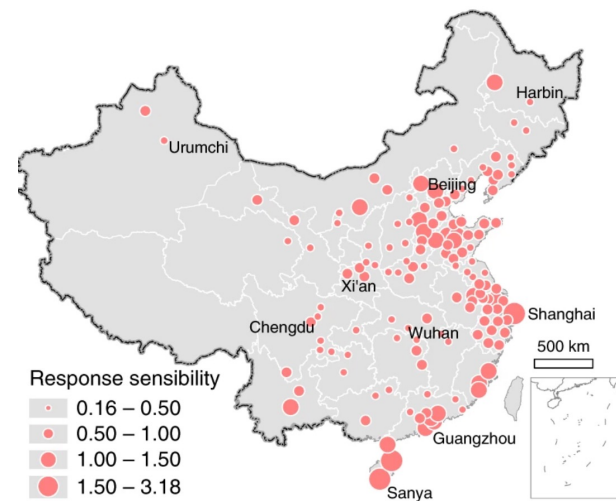
*Nature Human Behaviour* **5**, 337–348 (2021) | [Cite this article](#)

**269k** Accesses | **1538** Citations | **2386** Altmetric | [Metrics](#)

Loomba, S., De Figueiredo, A., Piatek, S. J., De Graaf, K., & Larson, H. J. (2021). Measuring the impact of COVID-19 vaccine misinformation on vaccination intent in the UK and USA. *Nature human behavior*, 5(3), 337-348.

# Real-world social science questions

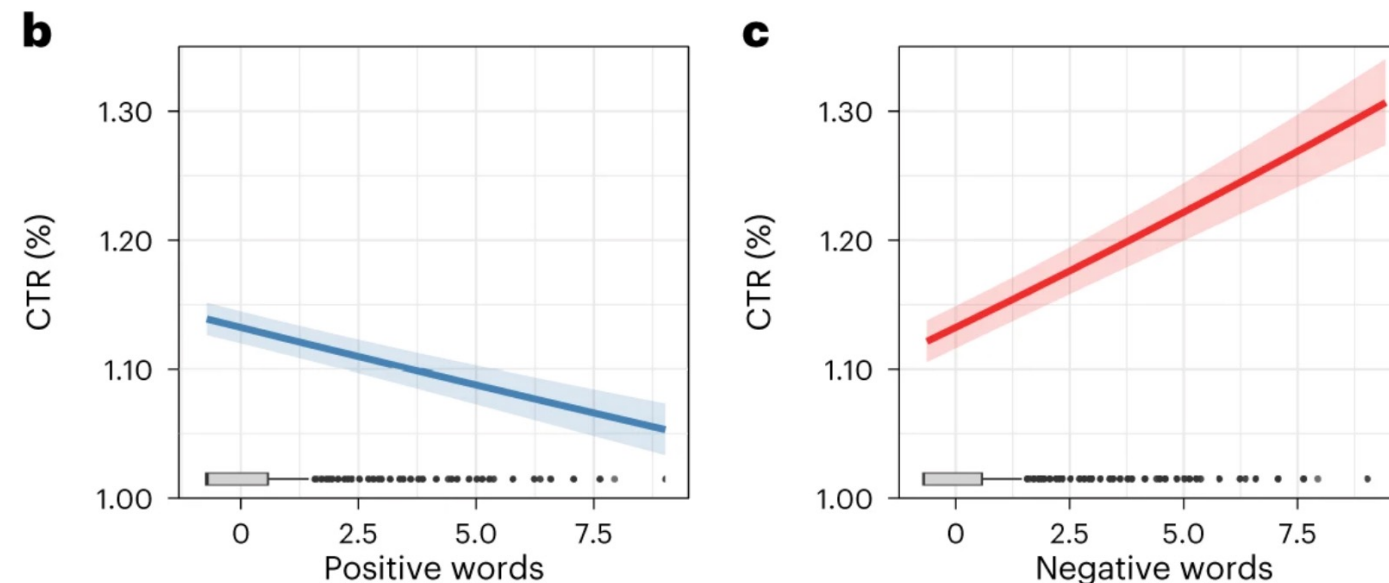
- Does air pollution lower people's expressed happiness on social media?
  - Annotation: "Does this social media post contain high or low positive affect?"
- To estimate: whether people living in a more polluted environment express less happiness on social media



Zheng, S., Wang, J., Sun, C., Zhang, X., & Kahn, M. E. (2019). Air pollution lowers Chinese urbanites' expressed happiness on social media. *Nature human behaviour*, 3(3), 237-243.

# Real-world social science questions

- Does negativity influence online news consumption?
  - Annotation: “Does this online news contain high or low negative affect?”
  - To estimate: whether the negativity of online news predict consumption

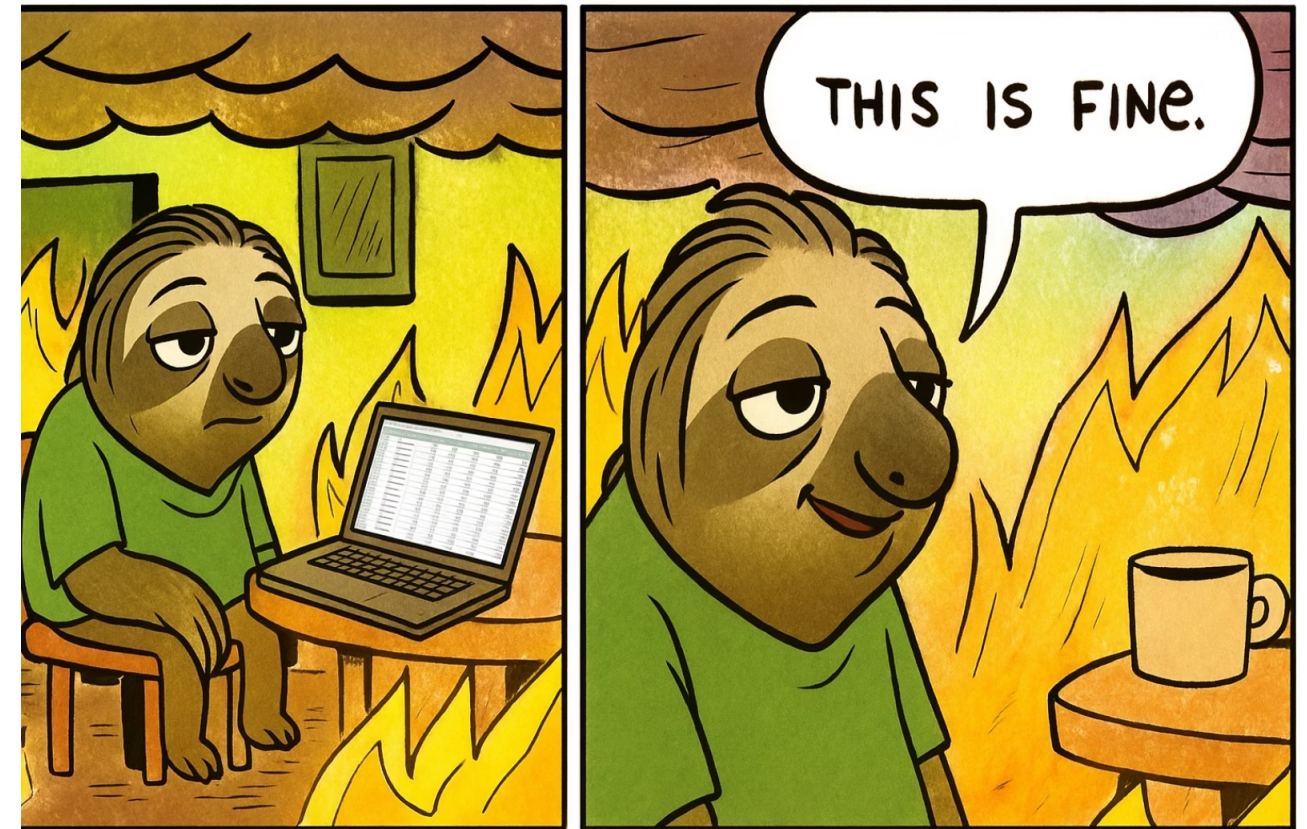


Robertson, C. E., Pröllochs, N., Schwarzenegger, K., Pärnamets, P., Van Bavel, J. J., & Feuerriegel, S. (2023). Negativity drives online news consumption. *Nature human behaviour*, 7(5), 812-822.

# Challenges with Human Annotation

While human annotations are the gold standard for quality and nuance, they are also **slow** and **expensive**

- require aggregating judgments from many annotators
- very costly, especially if coming from experts



# Can LLMs Replace Human Annotators?

BRIEF REPORT | POLITICAL SCIENCES | 



## ChatGPT outperforms crowd workers for text-annotation tasks

Fabrizio Gilardi , Meysam A.

March 01 2024

### Can Large Language Models Transform Computational Social Science?

In Special Collection: CogNet

Caleb Ziems, William Held, Omar Shaikh, Liisa Chen, Zhehao Zhang, Didi Yan

## Out of One, Many: Using Language Models to Simulate Human Samples

Published online by Cambridge University Press: 21 February 2023

Lisa P. Argyle , Ethan C. Busby, Nancy Fulda, Joshua R. Gubler , Christopher Rytting and David Wingate

PERSPECTIVE | SOCIAL SCIENCES | 



## Can Generative AI improve social science?

Christopher A. Bail  [Authors Info & Affiliations](#)

Edited by David Lazer, Northeastern University, Boston, MA; received September 7, 2023; accepted April 5, 2024, by Editorial Board Member Mark Granovetter

May 9, 2024 | 121 (21) e2314021121 | <https://doi.org/10.1073/pnas.2314021121>

# Can LLMs Replace Human Annotators?

LLM annotations are fast & affordable!

... But they do not always align with human judgment (e.g., biases, factual inaccuracies, inconsistency)

March 13, 2024

**AI Language Models Are More Biased Than Humans When It Comes To AAVE, Stanford And Oxford Study Unveils**

**Large Legal Fictions: Profiling Legal Hallucinations in Large Language Models**

[Matthew Dahl](#), [Varun Magesh](#), [Mirac Suzgun](#), [Daniel E. Ho](#)

**Air Canada ordered to pay customer who was misled by airline's chatbot**

Company claimed its chatbot 'was responsible for its own actions' when giving wrong information about bereavement fare

# This Lecture

## Methods for trustworthy social science with possibly untrustworthy LLMs

Basic idea: LLMs shouldn't replace real human data; they should complement it

Best of both worlds: leverage power of AI + retain scientific rigor

We will be following notation from **Confidence-Driven Inference (CDI)** (Gligoric\*, Zrnic\*, Lee\*, Candes, Jurafsky [NAACL, 2025])

Method as presented will combine ideas from several works, which we will reference along the way



Kristina  
Gligorić\*



Tijana  
Zrnić\*



Cinoo  
Lee\*



Emmanuel  
Candes



Dan  
Jurafsky

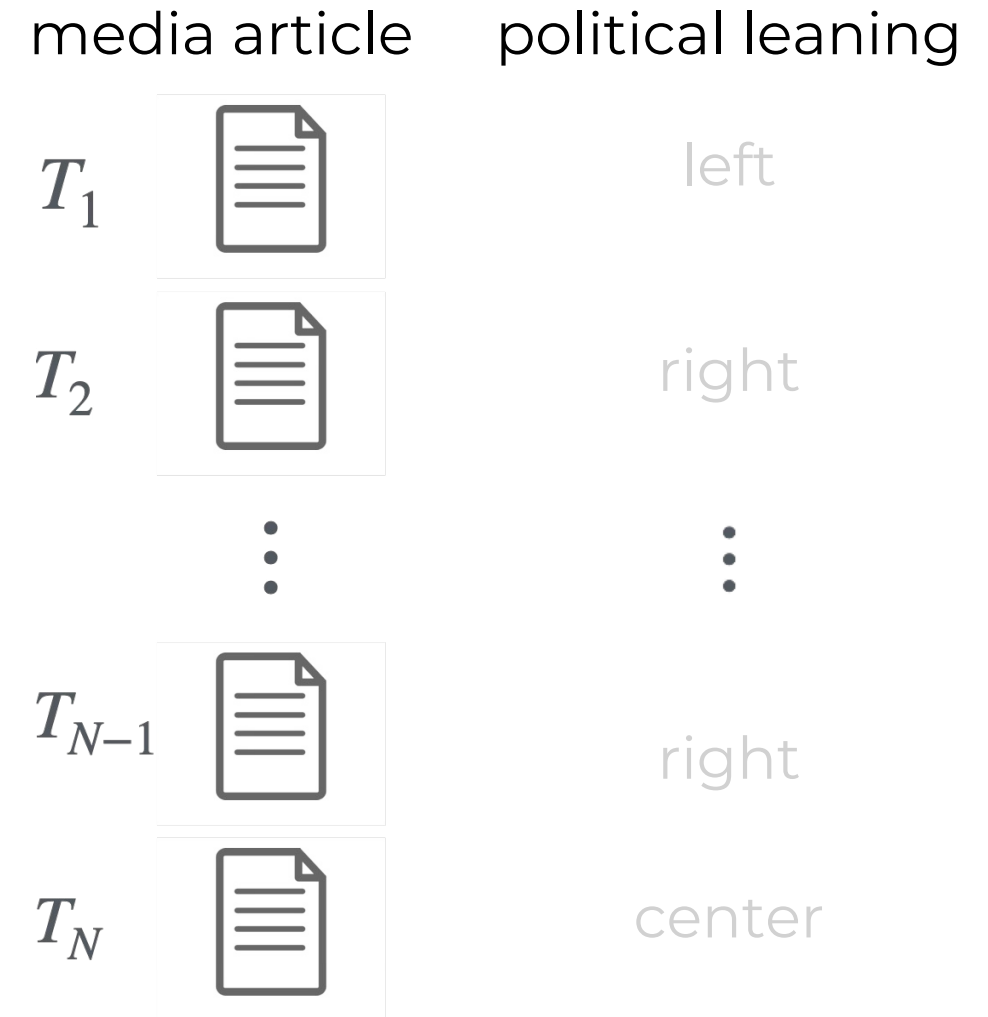


# Setup

$T_1$	$H_1$
$T_2$	$H_2$
$\vdots$	$\vdots$
$T_{N-1}$	$H_{N-1}$
$T_N$	$H_N$

$N$  text instances  $T_i$

missing human annotations  $H_i$







# Setup

$T_1$	$H_1$
$T_2$	$H_2$
$\vdots$	$\vdots$
$T_{N-1}$	$H_{N-1}$
$T_N$	$H_N$

$N$  text instances  $T_i$

missing human annotations  $H_i$

social media post      sentiment

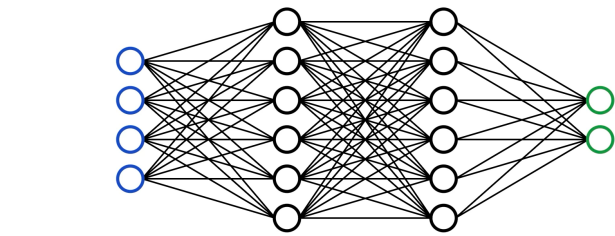
$T_1$		0.1
$T_2$		0.9
	$\vdots$	$\vdots$
$T_{N-1}$		0.5
$T_N$		0.3

# Setup

$T_1$	$H_1$
$T_2$	$H_2$
$\vdots$	$\vdots$
$T_{N-1}$	$H_{N-1}$
$T_N$	$H_N$

$N$  text instances  $T_i$

missing human annotations  $H_i$



large language model

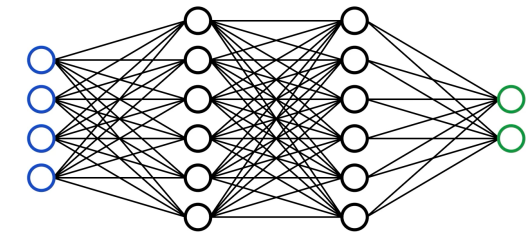
can be used to produce  $\hat{H}_i$  that approximate human annotations  $H_i$

# Setup

$T_1$	$\hat{H}_1$
$T_2$	$\hat{H}_2$
$\vdots$	$\vdots$
$T_{N-1}$	$\hat{H}_{N-1}$
$T_N$	$\hat{H}_N$

$N$  text instances  $T_i$

missing human annotations  $H_i$



large language model

issue:  $\hat{H}_i$  are potentially biased annotations!

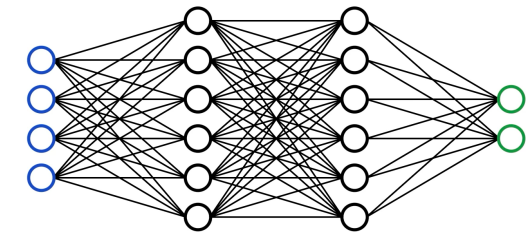
Unless we are willing to assume that the LLM is accurate, there is no hope of reaching valid conclusions without any human annotations!

# Setup

$T_1$	$\hat{H}_1$
$T_2$	$\hat{H}_2$
$\vdots$	$\vdots$
$T_{N-1}$	$\hat{H}_{N-1}$
$T_N$	$\hat{H}_N$

$N$  text instances  $T_i$

missing human annotations  $H_i$



large language model

examples of  $\theta^*$ :

- ❖ change in political leaning on X after Elon Musk acquisition
- ❖ effect of certain linguistic devices on perceived sentiment
- ❖ whatever we care about learning once we have human annotations!

budget: can collect at most  $n \ll N$  human annotations

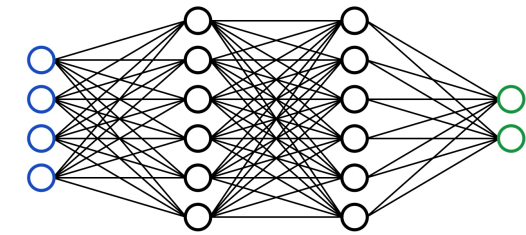
goal: estimate quantity of interest  $\theta^*$

# Setup

$T_1$	$\hat{H}_1$
$T_2$	$\hat{H}_2$
$\vdots$	$\vdots$
$T_{N-1}$	$\hat{H}_{N-1}$
$T_N$	$\hat{H}_N$

$N$  text instances  $T_i$

missing human annotations  $H_i$



large language model

important note:  $H_i$  do not necessarily correspond to annotations from a *single* human. They are “gold” annotations; e.g., obtained by aggregating annotations from multiple annotators.

# A Special Case: $\theta^* = \text{mean}(H_i)$

Example:  $H_i \in \{0,1\}$  indicates if article has right leaning;  $\theta^* =$  prevalence of right leaning

$$\theta^* = \text{mean}(H_i) = \frac{1}{N} (1 + 1 + 0 + 1 + \dots + 0) = \text{fraction of right-leaning articles}$$

right-leaning      left-leaning

# A Special Case: $\theta^* = \text{mean}(H_i)$

Example:  $H_i \in \{0,1\}$  indicates if article has right leaning;  $\theta^*$  = prevalence of right leaning

Step 1: Collect LLM annotations for all texts


# A Special Case: $\theta^* = \text{mean}(H_i)$

Example:  $H_i \in \{0,1\}$  indicates if article has right leaning;  $\theta^*$  = prevalence of right leaning

Step 1: Collect LLM annotations for all texts

	$\hat{H}_1$
	$\vdots$
	$\vdots$
	$\hat{H}_N$

# A Special Case: $\theta^* = \text{mean}(H_i)$

Example:  $H_i \in \{0,1\}$  indicates if article has right leaning;  $\theta^*$  = prevalence of right leaning

Step 1: Collect LLM annotations for all texts

Step 2: Collect  $n$  human annotations uniformly at random

	$\hat{H}_1$	
	$\vdots$	
	$\hat{H}_N$	

# A Special Case: $\theta^* = \text{mean}(H_i)$

Example:  $H_i \in \{0,1\}$  indicates if article has right leaning;  $\theta^* =$  prevalence of right leaning

Step 1: Collect LLM annotations for all texts

Step 2: Collect  $n$  human annotations uniformly at random


# A Special Case: $\theta^* = \text{mean}(H_i)$

Example:  $H_i \in \{0,1\}$  indicates if article has right leaning;  $\theta^* =$  prevalence of right leaning

Step 1: Collect LLM annotations for all texts

Step 2: Collect  $n$  human annotations uniformly at random

$T_1$	$\hat{H}_1$	$H_1$
$\vdots$	$\vdots$	$\vdots$
$T_n$	$\hat{H}_n$	$H_n$
$T_{n+1}$	$\hat{H}_{n+1}$	$H_{n+1}$
$\vdots$	$\vdots$	$\vdots$
$T_N$	$\hat{H}_N$	$H_N$

House of Representatives ...	0	1
The Pentagon accidentally ...	0	0
Democrats clash over ...	1	1
Gun lobby may emerge ...	0	0
Senate confirms FBI ...	0	
Senate Coronavirus Bill ...	1	
What does climate change ...	1	
Bipartisan Harvard panel ...	1	
Elon Musk has idea to ...	0	

# A Special Case: $\theta^* = \text{mean}(H_i)$

Example:  $H_i \in \{0,1\}$  indicates if article has right leaning;  $\theta^*$  = prevalence of right leaning

Step 1: Collect LLM annotations for all texts

Step 2: Collect  $n$  human annotations uniformly at random

Step 3: Given  $(H_1, \hat{H}_1), \dots, (H_n, \hat{H}_n), \hat{H}_{n+1}, \dots, \hat{H}_N$ , compute estimate of  $\theta^*$

$$\hat{\theta}^{\text{PPI}} = \underbrace{\text{mean}(\hat{H}_{n+1}, \dots, \hat{H}_N)}_{\text{naïve estimate}} - \underbrace{\text{mean}(\hat{H}_1 - H_1, \dots, \hat{H}_n - H_n)}_{\text{bias}}$$

Prediction-powered inference. Angelopoulos, Bates, Fannjiang, Jordan, Zrnic [Science, 2023]

Design-based supervised learning. Egami, Hinck, Stewart, Wei [NeurIPS, 2023]

# A Special Case: $\theta^* = \text{mean}(H_i)$

$$\hat{\theta}^{\text{PPI}} = \text{mean}(\hat{H}_{n+1}, \dots, \hat{H}_N) - \text{mean}(\hat{H}_1 - H_1, \dots, \hat{H}_n - H_n)$$

Theorem. For any data,  $\hat{\theta}^{\text{PPI}}$  is:

- ❖ accurate:  $\hat{\theta}^{\text{PPI}} \rightarrow \theta^*$  as the data size grows
- ❖ well-behaved:  $\hat{\theta}^{\text{PPI}} \approx N(\theta^*, \sigma^2)$

$\Rightarrow$  can form a confidence interval  $(\hat{\theta}^{\text{PPI}} \pm r)$   
via bootstrap or normal approximation

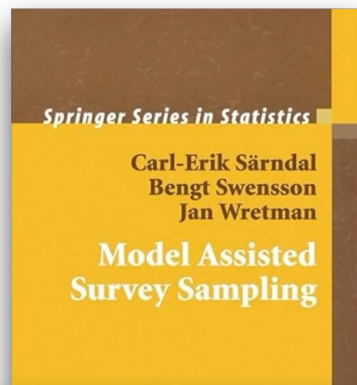
# A Special Case: $\theta^* = \text{mean}(H_i)$

$$\hat{\theta}^{\text{PPI}} = \text{mean}(\hat{H}_{n+1}, \dots, \hat{H}_N) - \text{mean}(\hat{H}_1 - H_1, \dots, \hat{H}_n - H_n)$$

## BIOMETRIKA

Inference using surrogate outcome data and a validation sample [Get access >](#)

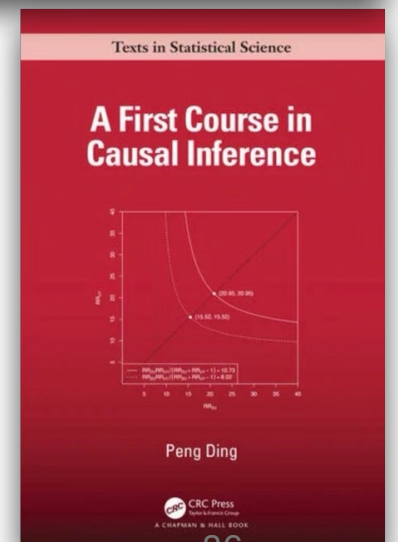
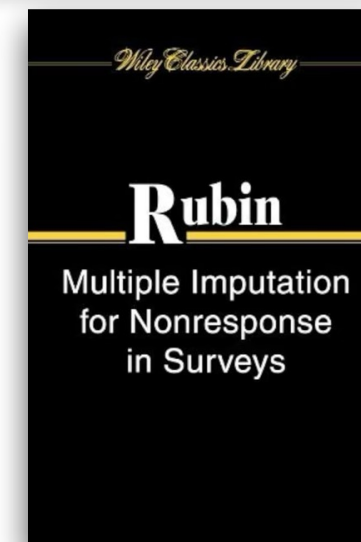
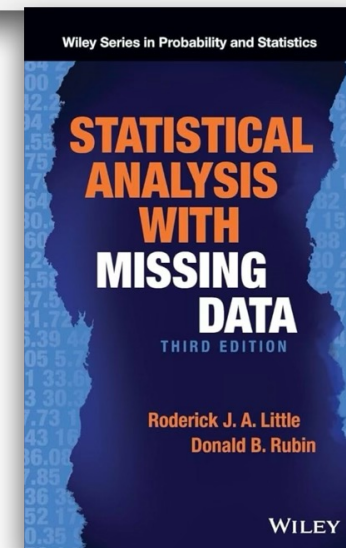
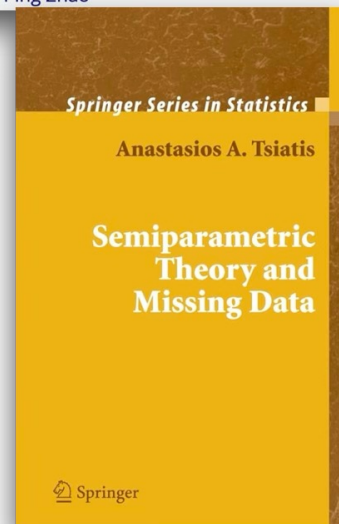
MARGARET SULLIVAN PEPE



Theory and Methods

**Estimation of Regression Coefficients When Some Regressors are not Always Observed**

James M. Robins, Andrea Rotnitzky & Lue Ping Zhao



**BACKPROPAGATION THROUGH THE VOID:  
OPTIMIZING CONTROL VARIATES FOR  
BLACK-BOX GRADIENT ESTIMATION**

Will Grathwohl, Dami Choi, Yuhuai Wu, Geoffrey Roeder, David Duvenaud  
University of Toronto and Vector Institute  
{wgrathwohl, choidami, ywu, roeder, duvenaud}@cs.toronto.edu

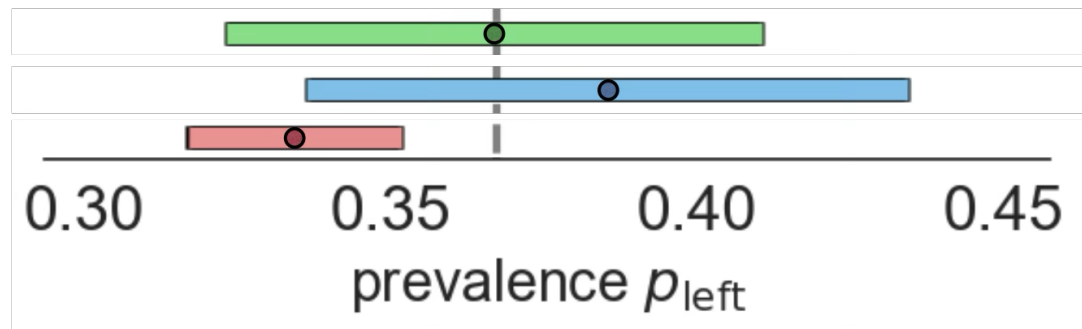
# Political Leaning

$T_i$  — media articles\*

$H_i$  — human annotations of political leaning

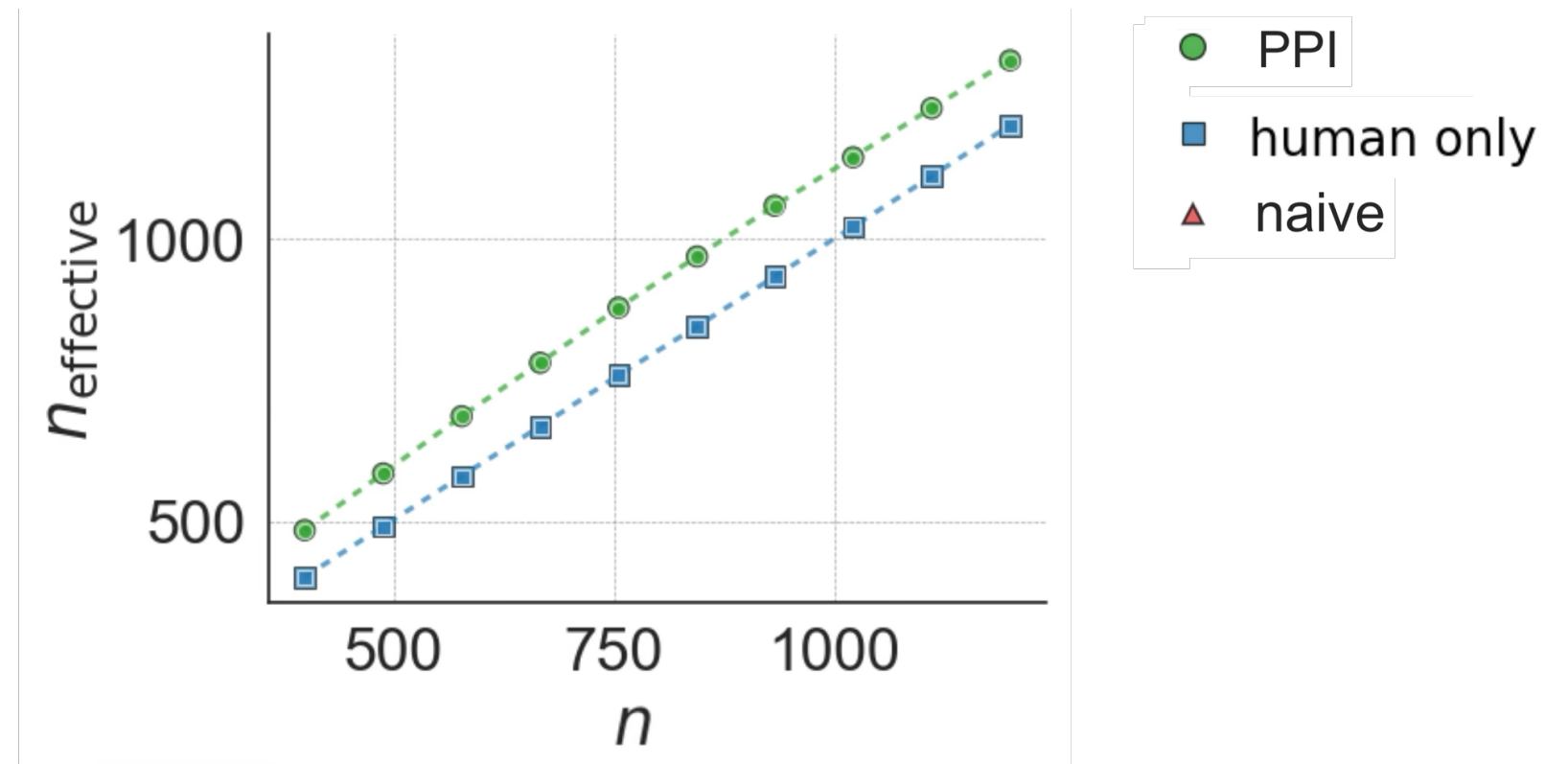
$\theta^*$  — fraction of left-leaning articles

LLM — GPT-4o



$$\hat{\theta}^{\text{naive}} = \text{mean}(\hat{H}_1, \dots, \hat{H}_N)$$

**Naively relying on LLMs is risky!**



\*Baly et al. EMNLP, 2020

# What are we missing?

1) That was only mean estimation. I want to run regressions (e.g., compute causal effects) and compute other more complex statistics (e.g. correlations, odds ratios, etc).

All these points are addressed by **Confidence-Driven Inference (CDI)**

# General Quantities of Interest



We want to learn  $\theta^* = \hat{\theta}((X_i, H_i)_{i=1}^N)$ , where  $X_i$  are (optionally) additional side covariates

$\theta^*$  = logistic regression coef. of  $H \sim X$

is polite?

contains  
gratitude words?

# General Quantities of Interest

1)

We want to learn  $\theta^* = \hat{\theta}((X_i, H_i)_{i=1}^N)$ , where  $X_i$  are (optionally) additional side covariates

- ❖ e.g.  $X_i$  indicates whether  $T_i$  contains gratitude words, or which media source the article comes from

General estimator:

$$\hat{\theta}^{\text{CDI}} = \underbrace{\hat{\theta}((X_i, \hat{H}_i)_{i=n+1}^N)}_{\text{naïve estimate}} - \underbrace{(\hat{\theta}((X_i, \hat{H}_i)_{i=1}^n) - \hat{\theta}((X_i, H_i)_{i=1}^n))}_{\text{bias}}$$

**Theorem.** For any data,  $\hat{\theta}^{\text{CDI}}$  is:

- ❖ accurate:  $\hat{\theta}^{\text{CDI}} \rightarrow \theta^*$  as the data size grows
- ❖ well-behaved:  $\hat{\theta}^{\text{CDI}} \approx N(\theta^*, \sigma^2)$

⇒ can form a confidence interval  $(\hat{\theta}^{\text{CDI}} \pm r)$   
via bootstrap or normal approximation

# Active Data Collection

2)

Human expertise should be reserved for “hard” problems; want  $\text{Prob}(\text{collect } H_i)$  large for difficult  $T_i$

It is optimal to have large  $\text{Prob}(\text{collect } H_i)$  for instances where  $\text{err}(H_i, \hat{H}_i)$  is the largest

We have to include inverse probability weights in the estimator if we sample adaptively:

$$\hat{\theta}^{\text{CDI}} = \hat{\theta} \left( (X_i, \hat{H}_i; \bar{W}_i)_{i=n+1}^N \right) - \left( \hat{\theta} \left( (X_i, \hat{H}_i; W_i)_{i=1}^n \right) - \hat{\theta} \left( (X_i, H_i; W_i)_{i=1}^n \right) \right)$$

$$W_i = \frac{1\{\text{collect } H_i\}}{\text{Prob}(\text{collect } H_i)}$$

$$\bar{W}_i = \frac{1 - 1\{\text{collect } H_i\}}{1 - \text{Prob}(\text{collect } H_i)}$$

Zrnic, Candes [ICML, 2024]

Gligoric, Zrnic, Lee, Candes, Jurafsky [NAACL, 2025]

Kluger, Lu, Zrnic, Wang, Bates [2025]

# Active Data Collection

2)

Human expertise should be reserved for “hard” problems; want  $\text{Prob}(\text{collect } H_i)$  large for difficult  $T_i$

It is statistically optimal to have large  $\text{Prob}(\text{collect } H_i)$  for instances where  $\text{err}(H_i, \hat{H}_i)$  is the largest

We have to include inverse probability weights in the estimator if we sample adaptively:

$$\hat{\theta}^{\text{CDI}} = \hat{\theta} \left( (X_i, \hat{H}_i; \bar{W}_i)_{i=n+1}^N \right) - \left( \hat{\theta} \left( (X_i, \hat{H}_i; W_i)_{i=1}^n \right) - \hat{\theta} \left( (X_i, H_i; W_i)_{i=1}^n \right) \right)$$

$$W_i = \frac{1\{\text{collect } H_i\}}{\text{Prob}(\text{collect } H_i)}$$

$$\bar{W}_i = \frac{1 - 1\{\text{collect } H_i\}}{1 - \text{Prob}(\text{collect } H_i)}$$



`fit(X, y, sample_weight=None)`

Fit linear model.

#### Parameters:

**X** : {array-like, sparse matrix} of shape (n\_samples, n\_features)

Training data.

**y** : array-like of shape (n\_samples,) or (n\_samples, n\_targets)

Target values. Will be cast to X's dtype if necessary.

**sample\_weight** : array-like of shape (n\_samples,), default=None

Individual weights for each sample.

# Active Data Collection

2)

Human expertise should be reserved for “hard” problems; want  $\text{Prob}(\text{collect } H_i)$  large for difficult  $T_i$

It is statistically optimal to have large  $\text{Prob}(\text{collect } H_i)$  for instances where  $\text{err}(H_i, \hat{H}_i)$  is the largest

We have to include inverse probability weights in the estimator if we sample adaptively:

$$\hat{\theta}^{\text{CDI}} = \hat{\theta} \left( (X_i, \hat{H}_i; \bar{W}_i)_{i=n+1}^N \right) - \left( \hat{\theta} \left( (X_i, \hat{H}_i; W_i)_{i=1}^n \right) - \hat{\theta} \left( (X_i, H_i; W_i)_{i=1}^n \right) \right)$$

$$W_i = \frac{1\{\text{collect } H_i\}}{\text{Prob}(\text{collect } H_i)}$$

$$\bar{W}_i = \frac{1 - 1\{\text{collect } H_i\}}{1 - \text{Prob}(\text{collect } H_i)}$$

Theorem. For any data,  $\hat{\theta}^{\text{CDI}}$  is:

- ❖ accurate:  $\hat{\theta}^{\text{CDI}} \rightarrow \theta^*$  as the data size grows
- ❖ well-behaved:  $\hat{\theta}^{\text{CDI}} \approx N(\theta^*, \sigma^2)$

$\Rightarrow$  can form a confidence interval  $(\hat{\theta}^{\text{CDI}} \pm r)$   
via bootstrap or normal approximation

# Confidence-Driven Inference

2)

To approximately sample where  $\text{err}(H_i, \hat{H}_i)$  is the largest, we look at **LLM uncertainty**

In our experiments, the most useful uncertainties were based on verbalized confidence (Tian et al., 2023)

## Stage 1

What is the political bias of the following article? Output either A,B, or C. Output a letter only.

A) Left

B) Center

C) Right

Article: <text>

Answer:

## Stage 2

How likely is it that the following article has a <previously provided answer: left-leaning, centrist, or right-leaning> political bias? Output the probability only (a number between 0 and 1).

Text: <text>

Probability:

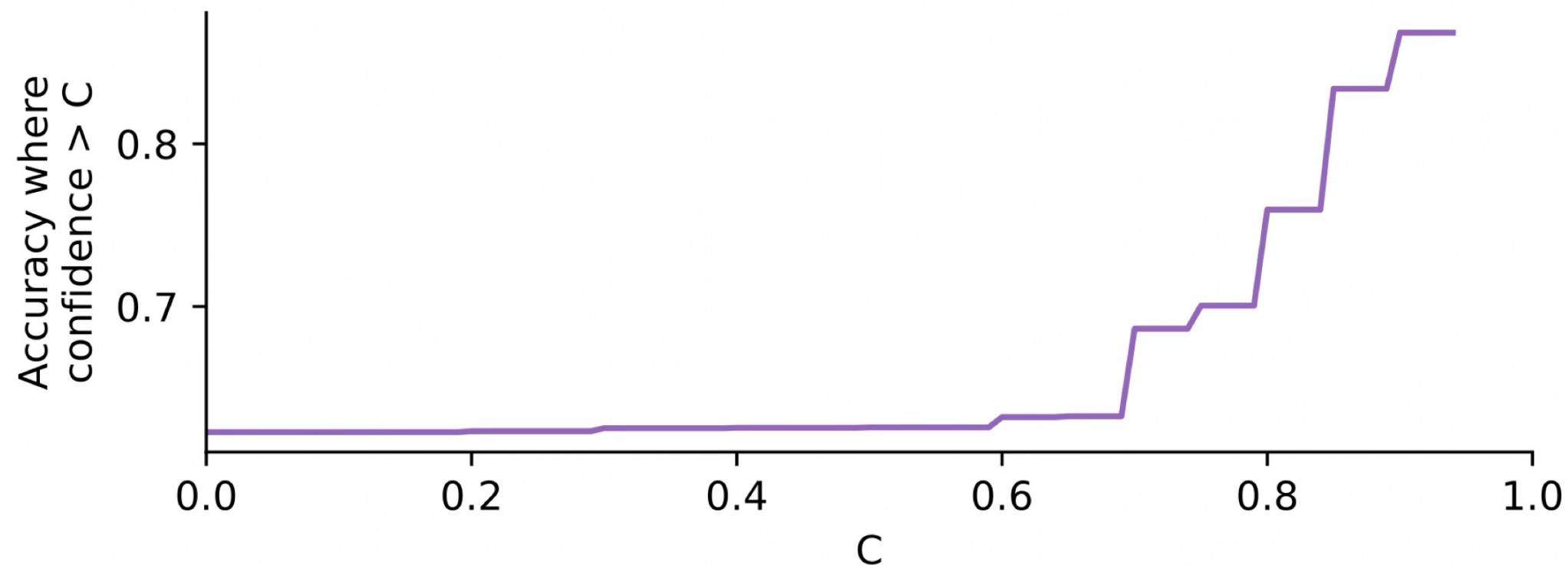
# Confidence-Driven Inference

2)

To approximately sample where  $\text{err}(H_i, \hat{H}_i)$  is the largest, we look at **LLM uncertainty**

In our experiments, the most useful uncertainties were based on verbalized confidence (Tian et al., 2023)

Confidence reflects accuracy!



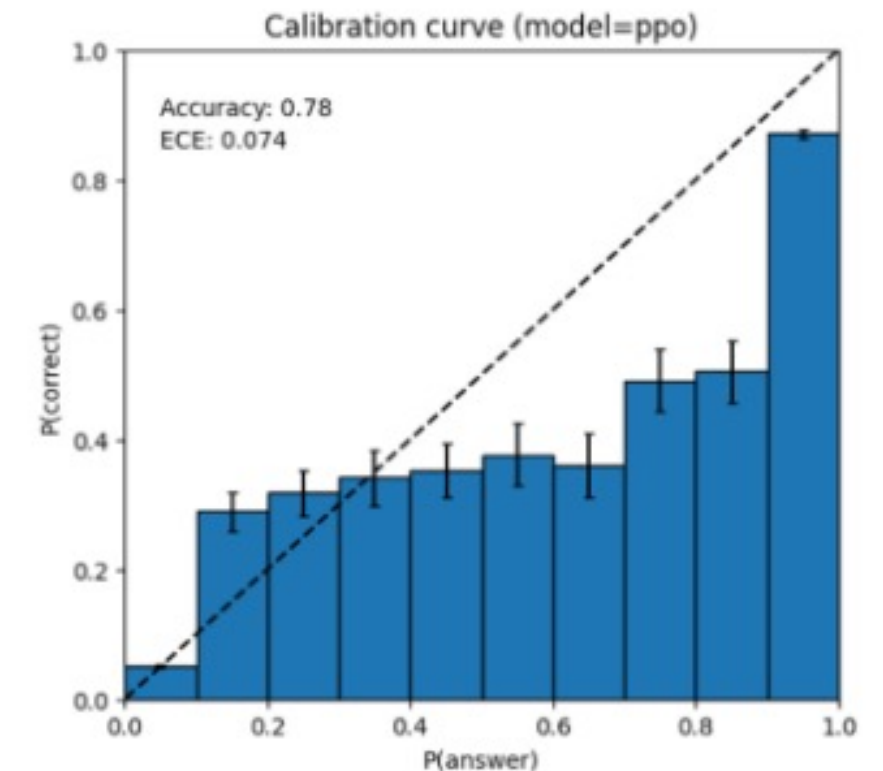
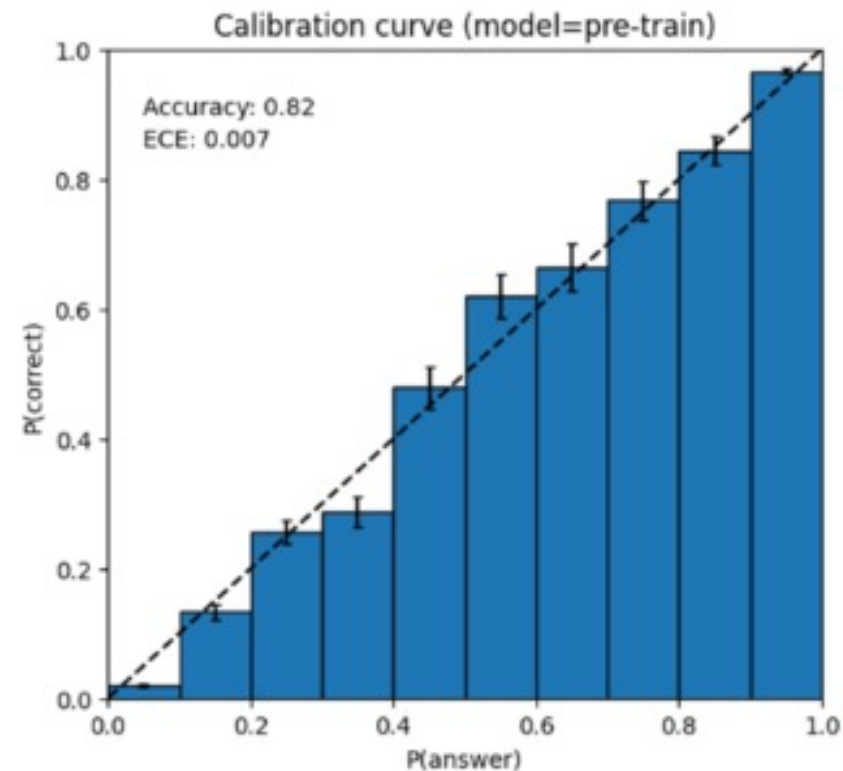
# Confidence-Driven Inference

2)

To approximately sample where  $\text{err}(H_i, \hat{H}_i)$  is the largest, we look at **LLM uncertainty**

In our experiments, the most useful uncertainties were based on verbalized confidence (Tian et al., 2023)

Confidence reflects accuracy!



Left: Calibration plot of the pre-trained GPT-4 model on an MMLU subset. The model's confidence in its prediction closely matches the probability of being correct. The dotted diagonal line represents perfect calibration. Right: Calibration plot of post-trained PPO GPT-4 model on the same MMLU subset. Our current process hurts the calibration quite a bit.

<https://openai.com/index/gpt-4-research/>

# Confidence-Driven Inference

2)

To approximately sample where  $\text{err}(H_i, \hat{H}_i)$  is the largest, we look at **LLM uncertainty**

In our experiments, the most useful uncertainties were based on verbalized confidence (Tian et al., 2023)

Confidence reflects accuracy!

We fit a mapping  $\hat{\text{err}}$  from confidence  $C_i$  to  $\text{err}(H_i, \hat{H}_i)$  as we collect data and set  $\text{Prob}(\text{collect } H_i) \propto \hat{\text{err}}(C_i)$  (normalized to meet the budget constraint)

# Safeguard Against Poor LLM Annotations

3)

Power tuning interpolates between using and not using LLM annotations

$$\hat{\theta}^\lambda = \lambda \cdot \hat{\theta} \left( (X_i, \hat{H}_i; \bar{W}_i)_{i=n+1}^N \right) - (\lambda \cdot \hat{\theta} \left( (X_i, \hat{H}_i; W_i)_{i=1}^n \right) - \hat{\theta} \left( (X_i, H_i; W_i)_{i=1}^n \right))$$



$$\lambda = 0$$

human-only

# Safeguard Against Poor LLM Annotations

3)

Power tuning interpolates between using and not using LLM annotations

$$\hat{\theta}^\lambda = \lambda \cdot \hat{\theta} \left( (X_i, \hat{H}_i; \bar{W}_i)_{i=n+1}^N \right) - (\lambda \cdot \hat{\theta} \left( (X_i, \hat{H}_i; W_i)_{i=1}^n \right) - \hat{\theta} \left( (X_i, H_i; W_i)_{i=1}^n \right))$$

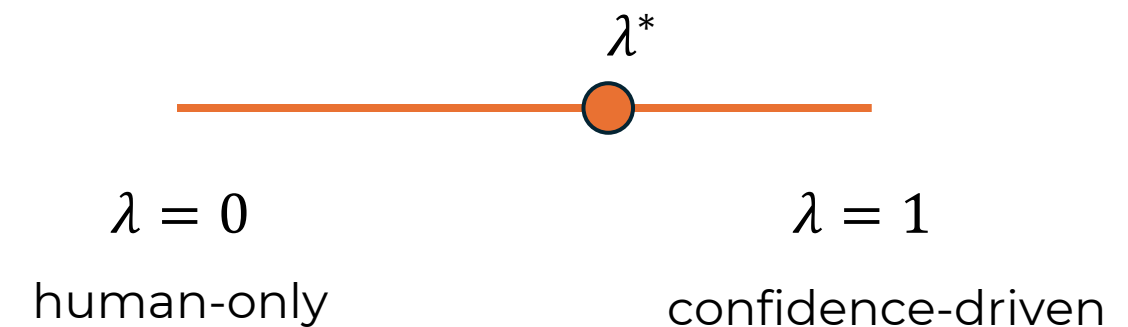


# Safeguard Against Poor LLM Annotations

3)

Power tuning interpolates between using and not using LLM annotations

$$\hat{\theta}^\lambda = \lambda \cdot \hat{\theta} \left( (X_i, \hat{H}_i; \bar{W}_i)_{i=n+1}^N \right) - (\lambda \cdot \hat{\theta} \left( (X_i, \hat{H}_i; W_i)_{i=1}^n \right) - \hat{\theta} \left( (X_i, H_i; W_i)_{i=1}^n \right))$$



Optimal tuning  $\lambda^*$  is proportional to how well  $H$  and  $\hat{H}$  correlate and can be computed explicitly

Theorem. After tuning,  $\hat{\theta}^\lambda$  achieves lower MSE than both  $\hat{\theta}^{\text{CDI}}$  and  $\hat{\theta}^{\text{human}}$ .

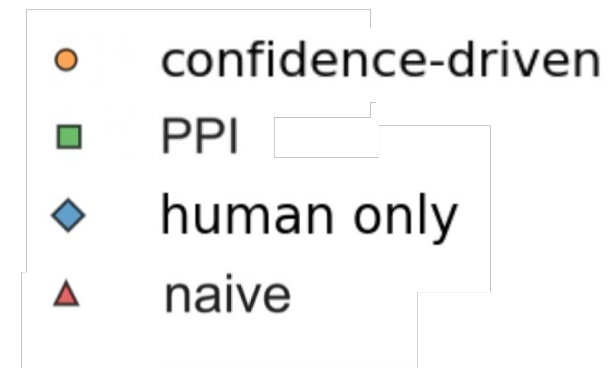
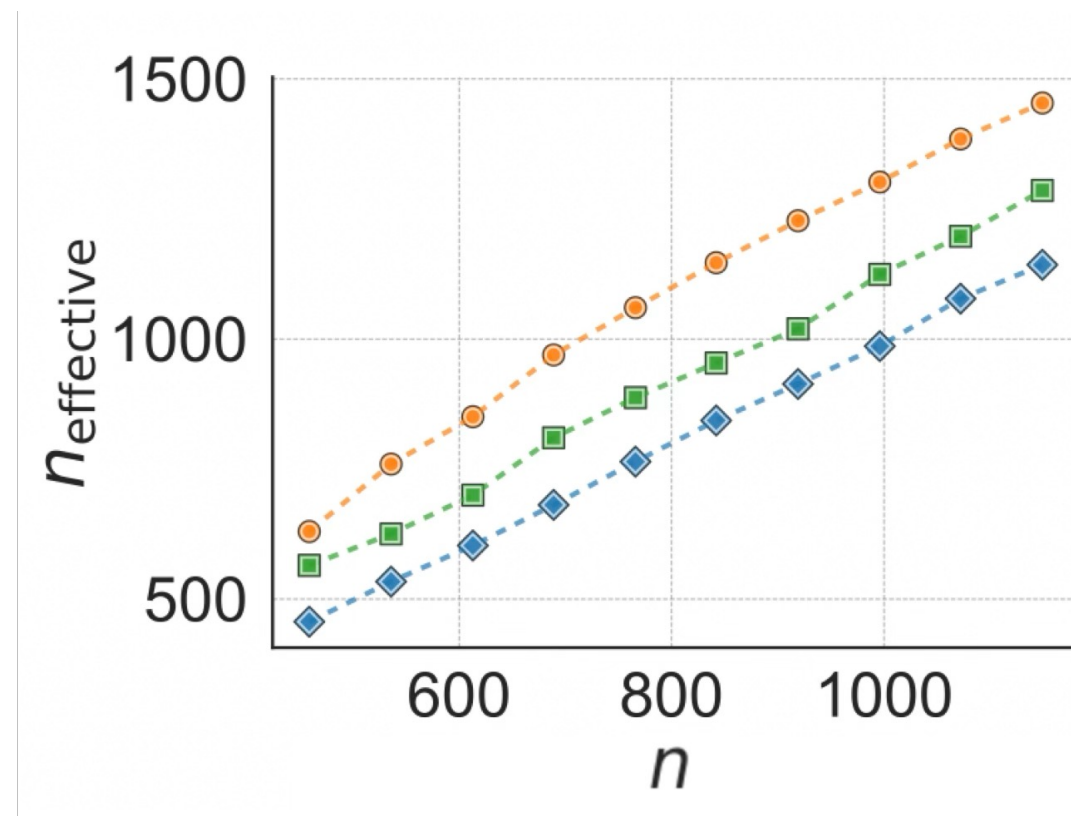
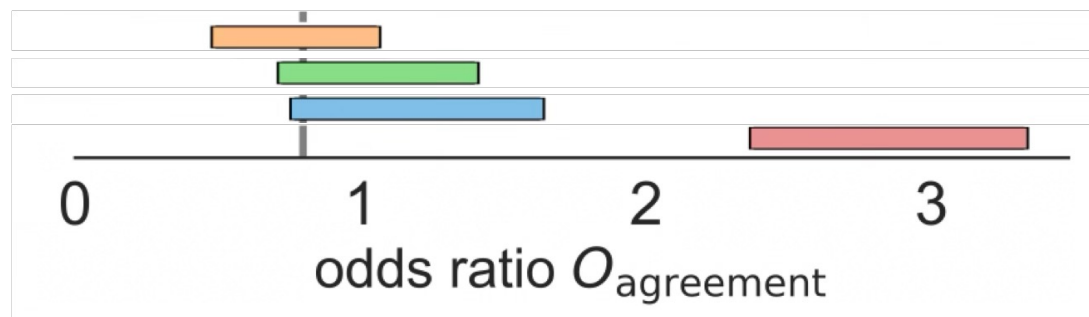
# Media Stance on Global Warming

$T_i$  — news headlines\*

$H_i$  — human annotations of article stance on global warming

$\theta^*$  — odds ratio quantifying relationship between affirming devices (e.g. “expert”, “award-winning scientist”) and stance on global warming

LLM — GPT-4o



\*Luo et al. EMNLP, 2020

# Confidence-Driven Inference: Step by Step

Input: (shuffled) texts  $T_i$ , LLM API, human annotation API, estimator of interest  $\hat{\theta}$

Step 1: Collect LLM annotations  $\hat{H}_i$  and confidence scores  $C_i$  for all texts  $T_i$

Step 2: Collect human annotations  $H_i$  for texts  $i = 1, \dots, n_{\text{init}}$ ; fit mapping  $\widehat{\text{err}}$  from  $C_i$  to  $\text{err}(H_i, \hat{H}_i)$

Step 3: Set  $\text{Prob}(\text{collect } H_i) = \frac{n}{N} \frac{\widehat{\text{err}}(C_i)}{\text{mean}(\widehat{\text{err}}(C_i))}$  for next  $n_{\text{batch}}$  texts; make sampling decisions  $\xi_i \in \{0,1\}$  with  $\text{Prob}(\text{collect } H_i)$

Step 4: Collect human annotations  $H_i$  for texts with  $\xi_i = 1$ ; refit  $\widehat{\text{err}}$  from all collected data so far

Step 5: Repeat Steps 3-4 until pass through all  $N$  texts is finished

Step 6: Compute tuning parameter  $\lambda$  and final estimate

Step 7: Compute confidence interval  $(\hat{\theta}^\lambda \pm r)$  via bootstrap

Output: estimate  $\hat{\theta}^\lambda$  and confidence interval  $(\hat{\theta}^\lambda \pm r)$

# A practical example

Questions so far?

# Politeness

Oxford **Learner's Dictionaries**

**politeness** *noun*

- 1 ★ good manners and respect for the feelings of others

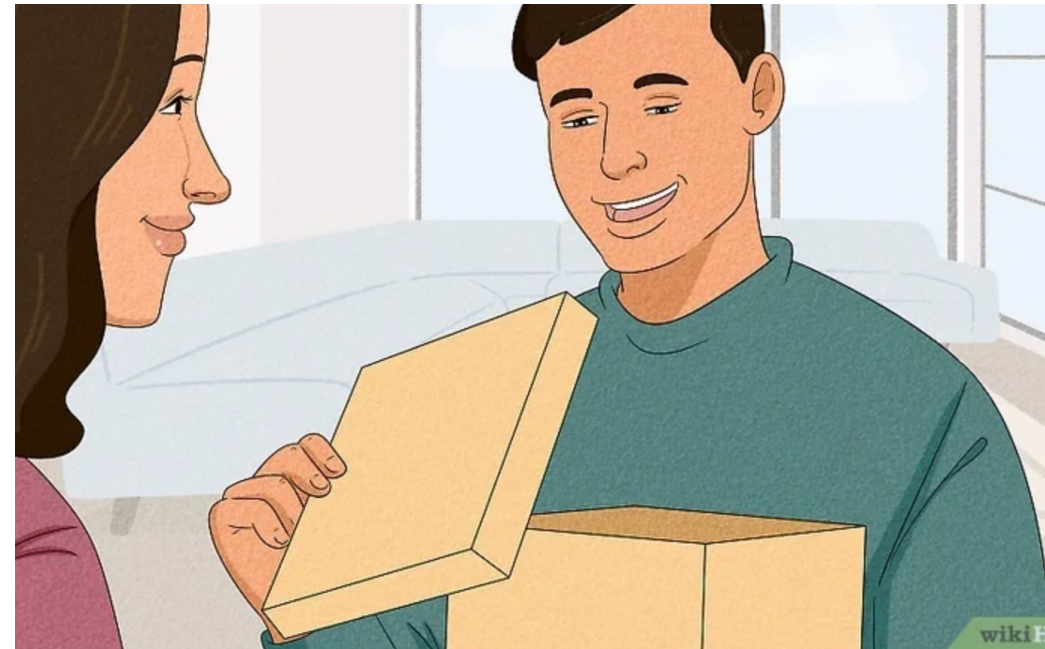
SYNONYM [courtesy](#) (1)

wikiHow to do anything...



PRO

## What to Do When You Get a Gift You Don't Like



### 1 Act naturally.

You don't need to feign excitement. Instead, summon up a positive feeling by thinking how nice it was that someone is giving you a gift.

- Try to react immediately. If you pause after you open the gift, you might seem disappointed.
- Smile if you can. It might help to remind yourself that they were trying to make you happy!

# What can politeness annotation tell us?

Surprisingly a lot!

Politeness annotation can inform various types of research questions:

e.g., gender, race, status and power

# What can politeness annotation tell us?

Is there a gender difference in language use?

Gender Differences in Language Use:  
An Analysis of 14,000 Text Samples

Matthew L. Newman  
*Department of Social and Behavioral Sciences  
Arizona State University*

Carla J. Groom\*  
*Department of Psychology  
The University of Texas at Austin*

Lori D. Handelman  
*Oxford University Press  
New York*

James W. Pennebaker  
*Department of Psychology  
The University of Texas at Austin*

Women use polite forms and  
hedging more than men  
("Would you mind if...", "I guess...")

# **What can politeness annotation tell us?**

Do police talk to White and Black drivers differently? If yes, how?

# What can politeness annotation tell us?

Do police talk to White and Black drivers differently? If yes, how?

RESEARCH ARTICLE | PSYCHOLOGICAL AND COGNITIVE SCIENCES | 



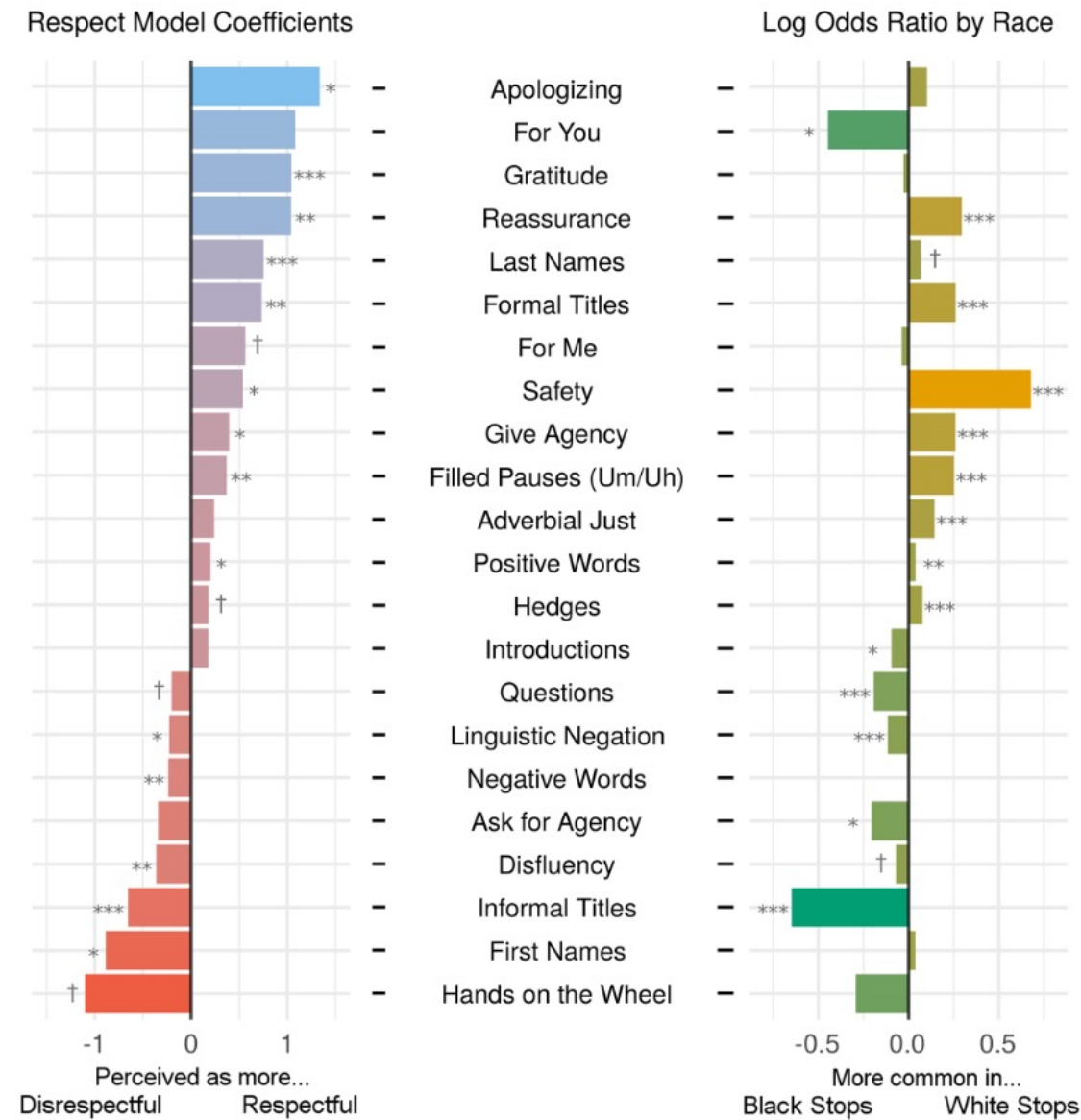
## Language from police body camera footage shows racial disparities in officer respect

Rob Voigt , Nicholas P. Camp, Vinodkumar Prabhakaran, , and Jennifer L. Eberhardt  [Authors Info & Affiliations](#)

“Officers speak with consistently less respect/politeness toward black vs. white community members, even after controlling for the race of the officer, the severity of the infraction, the location of the stop, and the outcome of the stop.”

# What can politeness annotation tell us?

Do police talk to White and Black drivers differently? If yes, how?



# What can politeness annotation tell us?

Does power corrupt?

## **A Computational Approach to Politeness with Application to Social Factors**

[Cristian Danescu-Niculescu-Mizil](#), [Moritz Sudhof](#), [Dan Jurafsky](#), [Jure Leskovec](#), [Christopher Potts](#)

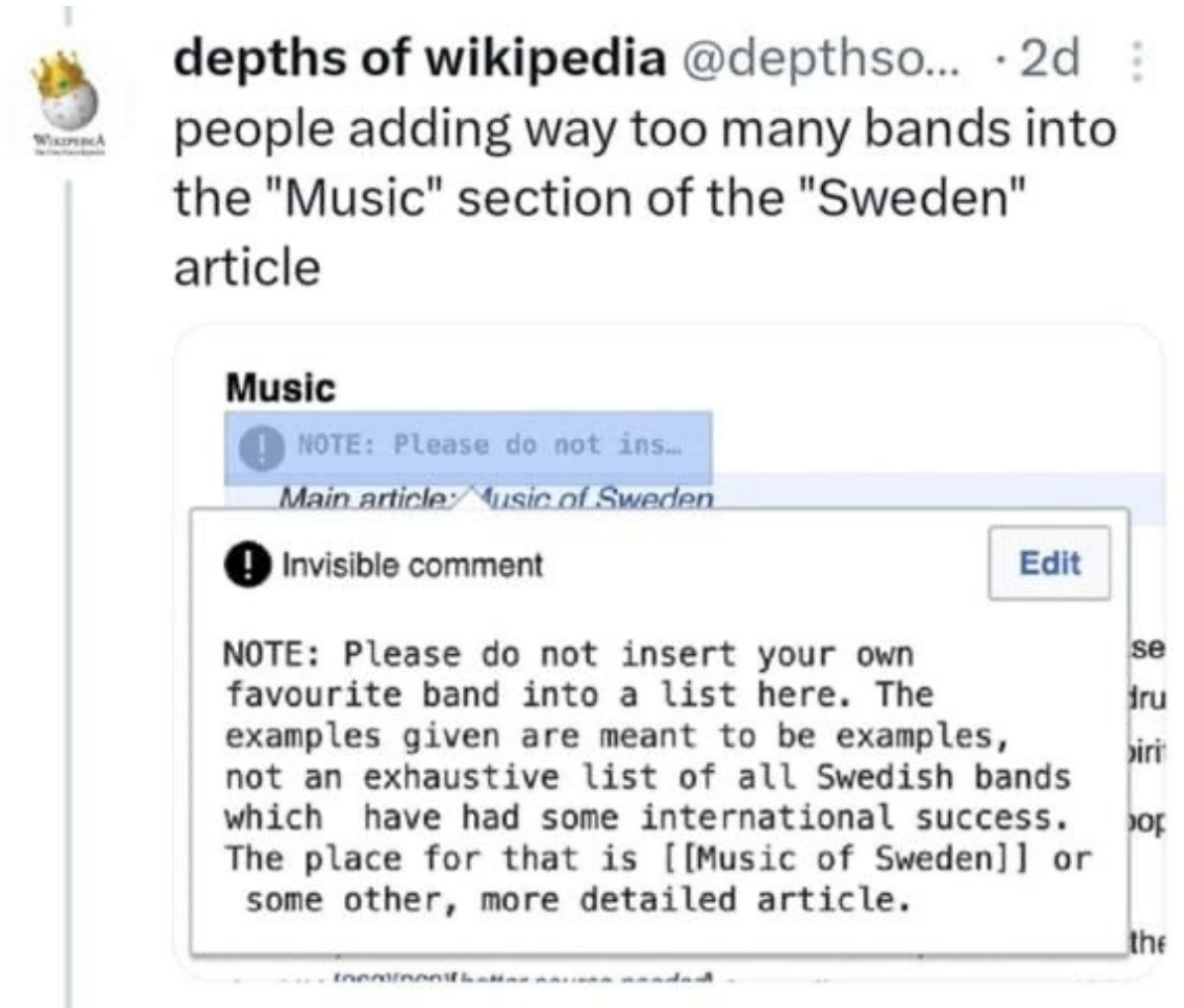
# What can politeness annotation tell us?

Does power corrupt?

**Yes!**

Wikipedia editors who would eventually be elected to administrator roles were significantly more polite than other users **before** their promotion.

However, **after** election (higher-status position), they became less polite.



The image shows a tweet from the account 'depths of wikipedia @depthso...' posted 2 days ago. The tweet text reads: 'people adding way too many bands into the "Music" section of the "Sweden" article'. Below the tweet is a screenshot of a Wikipedia comment box. The comment box has a title 'Music' and a blue header with a warning icon and the text 'NOTE: Please do not ins...'. Below the header, it says 'Main article: Music of Sweden'. The main body of the comment is titled 'Invisible comment' and contains the text: 'NOTE: Please do not insert your own favourite band into a list here. The examples given are meant to be examples, not an exhaustive list of all Swedish bands which have had some international success. The place for that is [[Music of Sweden]] or some other, more detailed article.' There is an 'Edit' button in the top right corner of the comment box.

# What can politeness annotation tell us?

Does power corrupt?

**Yes!**

Similarly, on Stack Exchange, users with the highest reputation scores were found to be less polite than users with low or middle-level reputations.



# Politeness data

## **A Computational Approach to Politeness with Application to Social Factors**

[Cristian Danescu-Niculescu-Mizil](#), [Moritz Sudhof](#), [Dan Jurafsky](#), [Jure Leskovec](#), [Christopher Potts](#)

Dataset comes from the Wikipedia community of editors and the Stack Exchange question-answer community

# Politeness data

## The Stack Exchange question-answer community

### The best way to get from Stockholm to Kolmården zoo by public transport

Asked 12 years, 2 months ago Modified 5 years, 2 months ago Viewed 6k times

- ▲  
7 I'm going to visit Stockholm with kids and wonder what is the best option to get to the Kolmården zoo. My first thought was to take the train as the fastest transport, but as far as I understand, train station is not that close to the Zoo.
- ▼ So, what will be better: to get to the train station and go by local bus or to take the bus from the very beginning?
- 🔖
- 🕒 public-transport sweden stockholm zoos

3 Answers

Sorted by: Highest score (default) ⇅

- ▲  
7 For this kind of question, [resrobot](#) is an excellent tool.
- ▼
- Train Stockholm – Norrköping; for example, 07:59 - 09:23, 08:21 - 09:33, or 09:40 - 11:11.
  - Bus 432 or 433 Norrköping – Kolmården. For example, 10:00 - 10:09, or 11:47 - 12:00.
- 🔖
- ✓
- 🕒 Kolmården Djurpark is 27 km from the Norrköping central station. Waiting times between train and bus appear rather long (more than 30 minutes). If you're many, you could consider taking a taxi. When booked via SJ, a taxi costs roughly 600 SEK, but a look at the prices for [Taxikurir Norrköping](#) suggests it's probably less than 400 SEK there (I get 330 SEK based on their prices). This can be compared to the bus, which costs 74 SEK for an adult and 51 SEK for youth or senior; so for a family of 2 parents, 2 children, the bus would be 250 SEK.

# Politeness data

## Wikipedia community of editors

### Johns Hopkins University

🌐 74 languages ▼

Article Talk

Read Edit View history Tools ▼

From Wikipedia, the free encyclopedia

Coordinates: 39°19′44″N 76°37′13″W﻿ / ﻿39.32889°N 76.62028°W﻿ / 39.32889; -76.62028

*"JHU" redirects here. For the Sri Lankan political party, see [Jathika Hela Urumaya](#).*

**Johns Hopkins University** (often abbreviated as **Johns Hopkins**, **Hopkins**, or **JHU**) is a [private research university](#) in [Baltimore](#), Maryland, United States. Founded in 1876 based on the European research institution model, Johns Hopkins is considered to be the first research university in the U.S.<sup>[8][9]</sup>

The university was named for its first benefactor, the American entrepreneur and [Quaker](#) philanthropist [Johns Hopkins](#).<sup>[10]</sup> Hopkins's \$7 million bequest (equivalent to \$166 million in 2024)<sup>[11]</sup> to establish the university and the affiliated [Johns Hopkins Hospital](#) in Baltimore was the largest [philanthropic](#) gift in U.S. history up to that time.<sup>[12][13]</sup> [Daniel Coit Gilman](#), who was inaugurated as [Johns Hopkins's first president](#) on February 22, 1876,<sup>[14]</sup> led the university to revolutionize higher education in the U.S. by integrating teaching and research.<sup>[15]</sup> In 1900, Johns Hopkins became a founding member of the [Association of American Universities](#).<sup>[16]</sup> The university has led all [U.S. universities](#) in annual research and development expenditures for over four consecutive decades.<sup>[17][18]</sup> The [School of Medicine](#), established in 1893, has achieved international recognition for its pioneering biomedical research.

#### Johns Hopkins University



*Latin:* *Universitas Hopkinsiensis*<sup>[1][2]</sup>

<b>Motto</b>	<i>Veritas vos liberabit</i> ( <i>Latin</i> )
<b>Motto in English</b>	"The truth will set you free"
<b>Type</b>	<a href="#">Private research university</a>
<b>Established</b>	February 22, 1876; 149 years ago
<b>Accreditation</b>	<a href="#">MSCHE</a>

# Politeness data

“talk page”

- [\(cur | prev\)](#) ○ [23:15, 26 April 2023](#) [ElKevbo](#) ([talk](#) | [contribs](#)) . . (7,797 bytes) **(+271)** . . (→Inclusion of *"consistently ranked among the top and most prestigious universities in the United States and the world"* in the lede: If this is information that is important enough to be included in the lede, editors should be able to provide sources that explicitly support it) ([undo](#))
- [\(cur | prev\)](#) ○ [17:56, 25 April 2023](#) [Sauzer](#) ([talk](#) | [contribs](#)) . . (7,526 bytes) **(+642)** . . (→Inclusion of *"consistently ranked among the top and most prestigious universities in the United States and the world"* in the lede: *Reply*) ([undo](#)) (*Tag: Reply*)

# Politeness data

Dataset: *requests* from Wikipedia editors & Stack Exchange question-answer community.

# Politeness data

Dataset: *requests* from Wikipedia editors & Stack Exchange question-answer community.

For each **request**, the annotator has to indicate how polite they perceived the request to be by using a slider with values ranging from “very impolite” to “very polite.”

For sake of simplicity, we’ll be using 0 (not polite) and 1 (polite).

# Politeness data

Our goal is to estimate **two target statistics**:

***mean(H)***: prevalence of politeness, i.e., the fraction of texts in the corpus that are polite.

# Politeness data

Our goal is to estimate **two target statistics**:

***mean(H)***: prevalence of politeness, i.e., the fraction of texts in the corpus that are polite.

**$\beta_{\text{hedge}}$** : the impact of linguistic features of hedging (X) on the perceived politeness (H), estimated with a logistic regression.

- Essentially measuring whether a request having “I suggest...” influences the politeness rating
- E.g., we could estimate that hedging (e.g., “I suggest...”) would make the text 20% more likely to be perceived as polite.

# Politeness data

<https://github.com/kristinagligoric/cdi-tutorial>

## Import libraries

```
In [1]: import numpy as np
from scipy.stats import norm, bernoulli
import pprint
import pandas as pd
from tqdm import tqdm
from tqdm.notebook import tqdm
import time
from ppi_py.utils import bootstrap
import re
import openai
import requests
import json
from datetime import datetime, timezone
import time
import zipfile
import io
import random
from sklearn.linear_model import LogisticRegression
%load_ext autoreload
%autoreload 2

from utils import llms, qualtrics, prolific, mturk, inference
from utils.llms import annotate_texts_with_llm, collect_llm_confidence, get_llm_annotations
from utils.qualtrics import create_and_activate_surveys
from utils.prolific import run_prolific_annotation_pipeline
from utils.mturk import run_mturk_annotation_pipeline
from utils.inference import train_sampling_rule, sampling_rule_predict, confidence_driven_inference, collect_initial_human_annotations, run_adaptive_sampling
```

# Politeness data

<https://github.com/kristinagligoric/cdi-tutorial>

## Setup credentials

```
def load_credentials(file_path="credentials.txt"):
    credentials = {}
    with open(file_path, "r") as f:
        for line in f:
            if '=' in line:
                key, value = line.strip().split('=', 1)
                credentials[key.strip()] = value.strip()
    return credentials

# Load credentials, or put your key here in plain text
creds = load_credentials()
AWS_ACCESS_KEY_ID = creds.get("AWS_ACCESS_KEY_ID")
AWS_SECRET_ACCESS_KEY = creds.get("AWS_SECRET_ACCESS_KEY")
OPENAI_API_KEY = creds.get("OPENAI_API_KEY")
QUALTRICS_API_KEY = creds.get("QUALTRICS_API_KEY")
QUALTRICS_API_URL = creds.get("QUALTRICS_API_URL")
PROLIFIC_API_KEY = creds.get("PROLIFIC_API_KEY")
```

# Politeness data

<https://github.com/kristinagligoric/cdi-tutorial>

## Set parameters for Confidence-Driven Inference (CDI)

```
:  
# if true, we collect LLM annotations and human annotations from scratch, if false, load pre-collected ones  
COLLECT_LLM = False  
COLLECT_HUMAN = False  
  
# if COLLECT_HUMAN = True, specify whether to use "Prolific" or "MTURK"  
HUMAN_SOURCE = "MTURK"  
  
alpha = 0.1 # desired error level for confidence interval  
burnin_steps = 5 # we collect the first burnin_steps points to initialize sampling rule  
  
n_batches = 2  
  
n_human = 15 # budget on number of human annotations (including burnin_steps)  
  
N = 100 # corpus size, or the size of the random subset of the corpus that will be annotated with an LLM  
  
random_state = 42  
  
#define the estimator function  
#mask for valid labels and specify how to use weights  
def mean_estimator(y, weights):  
    y, weights = y[~np.isnan(y)], weights[~np.isnan(y)]  
    return np.sum(y * weights) / np.sum(weights)
```

# Politeness data

<https://github.com/kristinagligoric/cdi-tutorial>

Step 1: Collect LLM annotations for all the texts

```
data = get_llm_annotations(df=df,  
    text_based_feature=text_based_feature,  
    COLLECT_LLM=COLLECT_LLM,  
    llm_parameters = llm_parameters,  
    N=N,  
    random_state=random_state)
```

Step 2: Collect warmup human annotations (initial set)

```
data = collect_initial_human_annotations(  
    data=data,  
    df=df,  
    burnin_steps=burnin_steps,  
    COLLECT_HUMAN=COLLECT_HUMAN,  
    HUMAN_SOURCE=HUMAN_SOURCE,  
    N=N,  
    random_state=random_state,  
    human_annotation_parameters = human_annotation_parameters)
```

Step 3: Strategically collect human annotations

```
data = run_adaptive_sampling(  
    data=data,  
    df=df,  
    burnin_steps=burnin_steps,  
    n_human=n_human,  
    n_batches=n_batches,  
    COLLECT_HUMAN=COLLECT_HUMAN,  
    HUMAN_SOURCE=HUMAN_SOURCE,  
    human_annotation_parameters = human_annotation_parameters)
```

Collecting batch 1/2...  
Collecting 4 human annotations.

Collecting batch 2/2...  
Collecting 8 human annotations.  
17 human datapoints collected in total.

# Politeness data

<https://github.com/kristinagligoric/cdi-tutorial>

## Step 4: Compute the CDI estimate and confidence intervals

We showcase estimation of  $\text{mean}(H)$ : prevalence of the politeness, i.e., the fraction of texts in the corpus that are polite

```
estimate, (lower_bound, upper_bound) = confidence_driven_inference(  
    estimator = mean_estimator,  
    Y = data['human'].values,  
    Yhat = data['llm'].values,  
    sampling_probs = data['sampling_probs'].values,  
    sampling_decisions = data['sampling_decisions'].values,  
    alpha = alpha)  
  
print("CDI estimate of the target statistic (mean(H): prevalence of politeness):")  
print('point estimate:', estimate.round(4))  
print('confidence intervals:', lower_bound.round(4), upper_bound.round(4))  
  
print("Ground truth mean(H) estimate (if we had access to human annotations on the full text corpus):")  
print(np.mean(pd.read_csv('data/politeness_dataset.csv').sample(n = N, random_state = random_state)['Prediction_human'].values))
```

```
CDI estimate of the target statistic (mean(H): prevalence of politeness):  
point estimate: 0.4901  
confidence intervals: 0.2563 0.7032  
Ground truth mean(H) estimate (if we had access to human annotations on the full text corpus):  
0.58
```

# Politeness data

<https://github.com/kristinagligoric/cdi-tutorial>

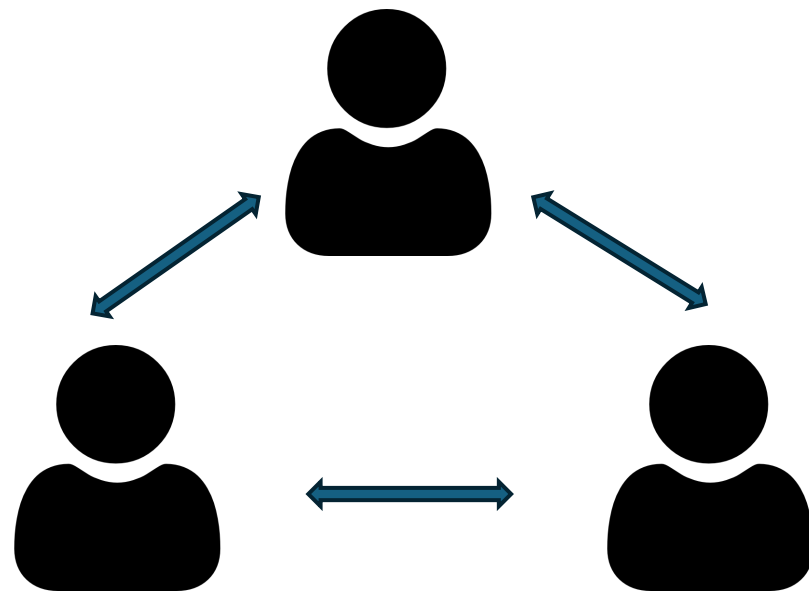
$\beta$  estimation

```
estimate, (lower_bound, upper_bound) = confidence_driven_inference(  
    estimator = log_reg_estimator,  
    Y = data['human'].values,  
    Yhat = data['llm'].values,  
    X = data['X'].values.reshape(-1, 1),  
    sampling_probs = np.ones(len(data))/len(data),  
    sampling_decisions = data['sampling_decisions'].values,  
    alpha = alpha)  
  
print("CDI estimate of the target statistic ( $\beta$ : effect of X on H):")  
print('point estimate:', estimate.round(4))  
print('confidence intervals:', lower_bound.round(4), upper_bound.round(4))
```

```
CDI estimate of the target statistic ( $\beta$ : effect of X on H):  
point estimate: 0.4433  
confidence intervals: 0.2734 0.5991
```

# Alternative approaches vs the "debiasing route"

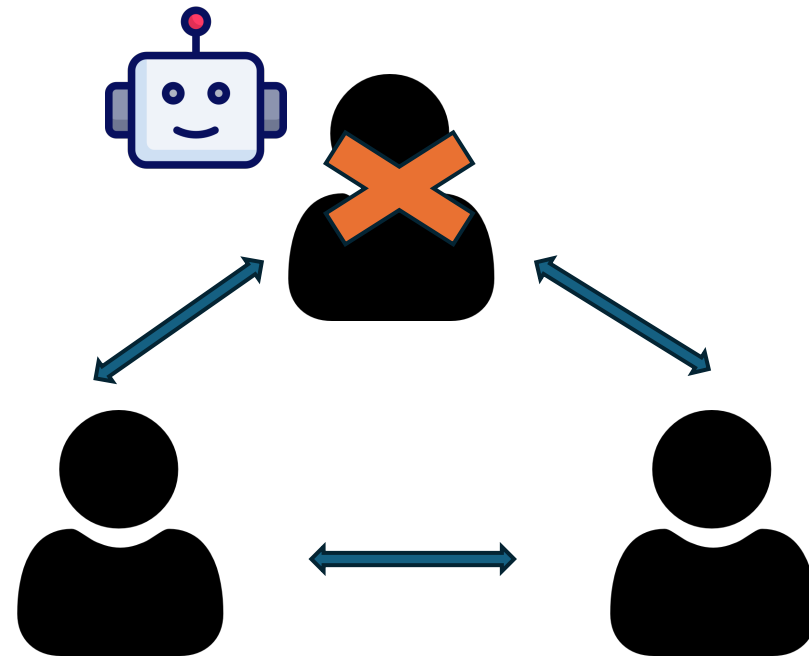
## The "alternative annotator test"



What is the level of disagreement?

# Alternative approaches vs the "debiasing route"

## The "alternative annotator test"



What is the level of disagreement **now**?

Practical argument of an upper limit, but no bounds on validity

# Further Problems & References

Statistical factuality guarantees for language models.

Mohri, Hashimoto (2024), Cherian, Gibbs, Candes (2024), Rubin-Toles, Gambhir, Ramji, Roth, Goel (2025)

# Further Problems & References

Statistical factuality guarantees for language models.

Mohri, Hashimoto (2024), Cherian, Gibbs, Candes (2024), Rubin-Toles, Gambhir, Ramji, Roth, Goel (2025)

Combining human and LLM annotations for approximately correct annotations.

Li, Shi, Ziemis, Kan, Chen, Liu, Yang (2023), Kim, Mitra, Chen, Rahman, Zhang (2024), Candes, Ilyas, Zrnic (2025)

Same Dataset    Uncertainty Computation    Expertise Estimation    Annotation Decision

## CoAnnotating: Uncertainty-Guided Work Allocation between Human and Large Language Models for Data Annotation

Minzhi Li <sup>†§</sup>    Taiwei Shi <sup>‡</sup>    Caleb Ziemis <sup>¶</sup>

Min-Yen Kan <sup>†</sup>    Nancy F. Chen <sup>§</sup>    Zhengyuan Liu <sup>§</sup>    Diyi Yang <sup>¶</sup>

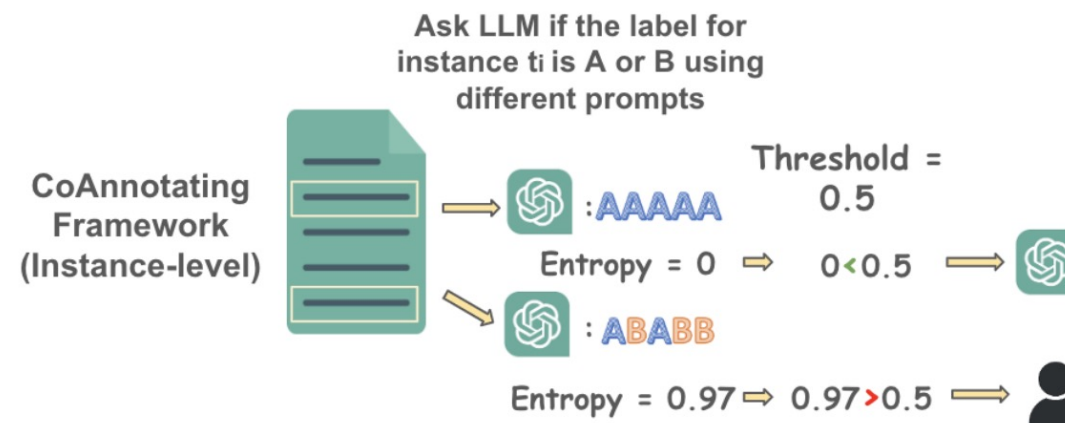
<sup>†</sup>National University of Singapore    <sup>§</sup>Institute for Infocomm Research (I<sup>2</sup>R), A\*STAR

<sup>‡</sup>University of Southern California    <sup>¶</sup>Stanford University

li.minzhi@u.nus.edu    taiweish@usc.edu    cziems@stanford.edu

mfychen@i2r.a-star.edu.sg    liu\_zhengyuan@i2r.a-star.edu.sg

kanmy@comp.nus.edu.sg    diyiy@cs.stanford.edu



# Further Problems & References

Statistical factuality guarantees for language models.

Mohri, Hashimoto (2024), Cherian, Gibbs, Candes (2024), Rubin-Toles, Gambhir, Ramji, Roth, Goel (2025)

Combining human and LLM annotations for approximately correct annotations.

Li, Shi, Ziems, Kan, Chen, Liu, Yang (2023), Kim, Mitra, Chen, Rahman, Zhang (2024), Candes, Ilyas, Zrnic (2025)

Valid evaluation of LLMs with synthetic data.

Chatzi, Straitouri, Thejaswi, Gomez Rodriguez (2024), Boyeau, Angelopoulos, Yosef, Malik, Jordan (2025)



# **What we're missing: Opportunities for future research**

## **LLM annotations with multi-modal inputs**

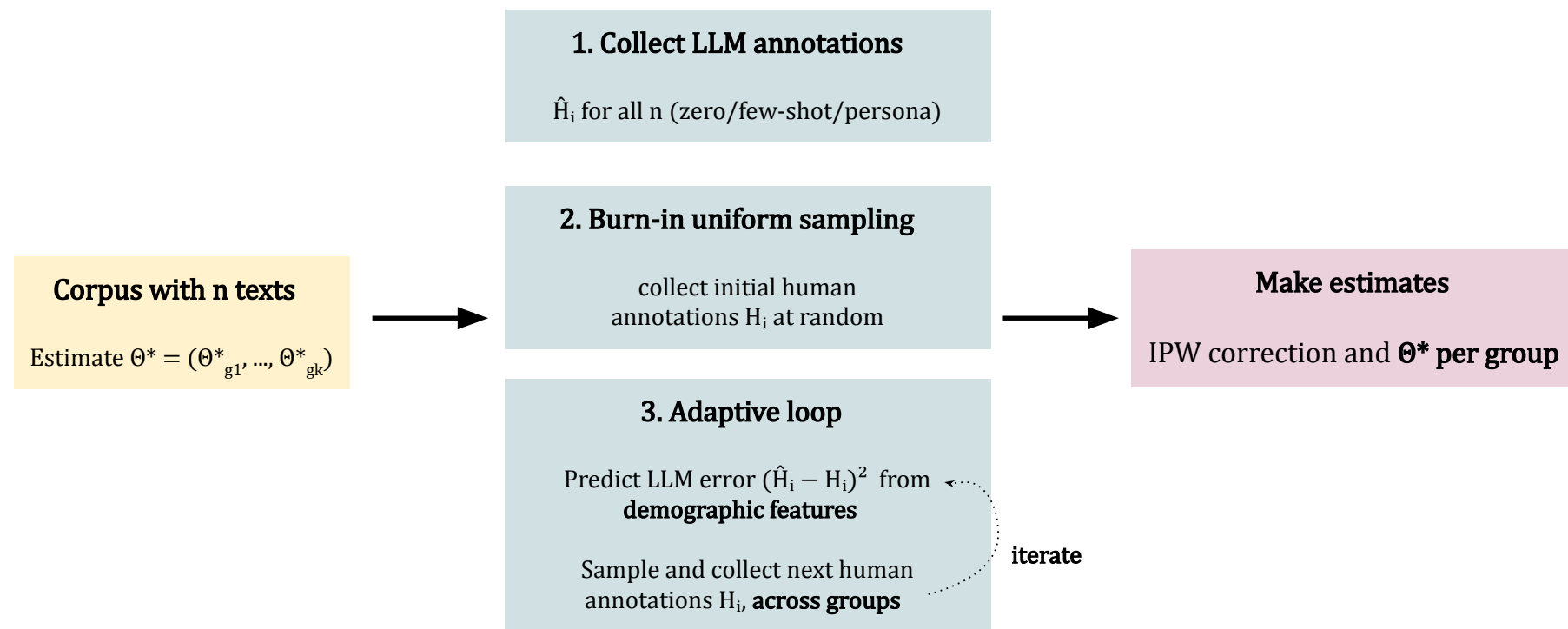
How to annotate videos? E.g., how do we find the most informative frames?

# What we're missing: Opportunities for future research

LLM annotations with multi-modal inputs

## What if there is no ground truth?

We want to estimate a vector of  $\theta^*$ s



<https://arxiv.org/pdf/2603.21404>



Navya  
Mehrotra



Adam  
Visokay



Kristina  
Gligorić

# **What we're missing: Opportunities for future research**

LLM annotations with multi-modal inputs

What if there is no ground truth?

**What if LLM predictions are not calibrated?**

How do we train distilled models, prioritizing calibration?

# **What we're missing: Opportunities for future research**

LLM annotations with multi-modal inputs

What if there is no ground truth?

What if LLM predictions are not calibrated?

**What if we want to do fine-tuning?**

# What we're missing: Opportunities for future

Valid Survey Simulations with Limited Human Data:  
The Roles of Prompting, Fine-Tuning, and Rectification

Stefan Krsteski<sup>1\*</sup>, Giuseppe Russo<sup>1,2\*</sup>, Serina Chang<sup>3</sup>, Robert West<sup>1</sup>, Kristina Gligorić<sup>4</sup>

<sup>1</sup>EPFL

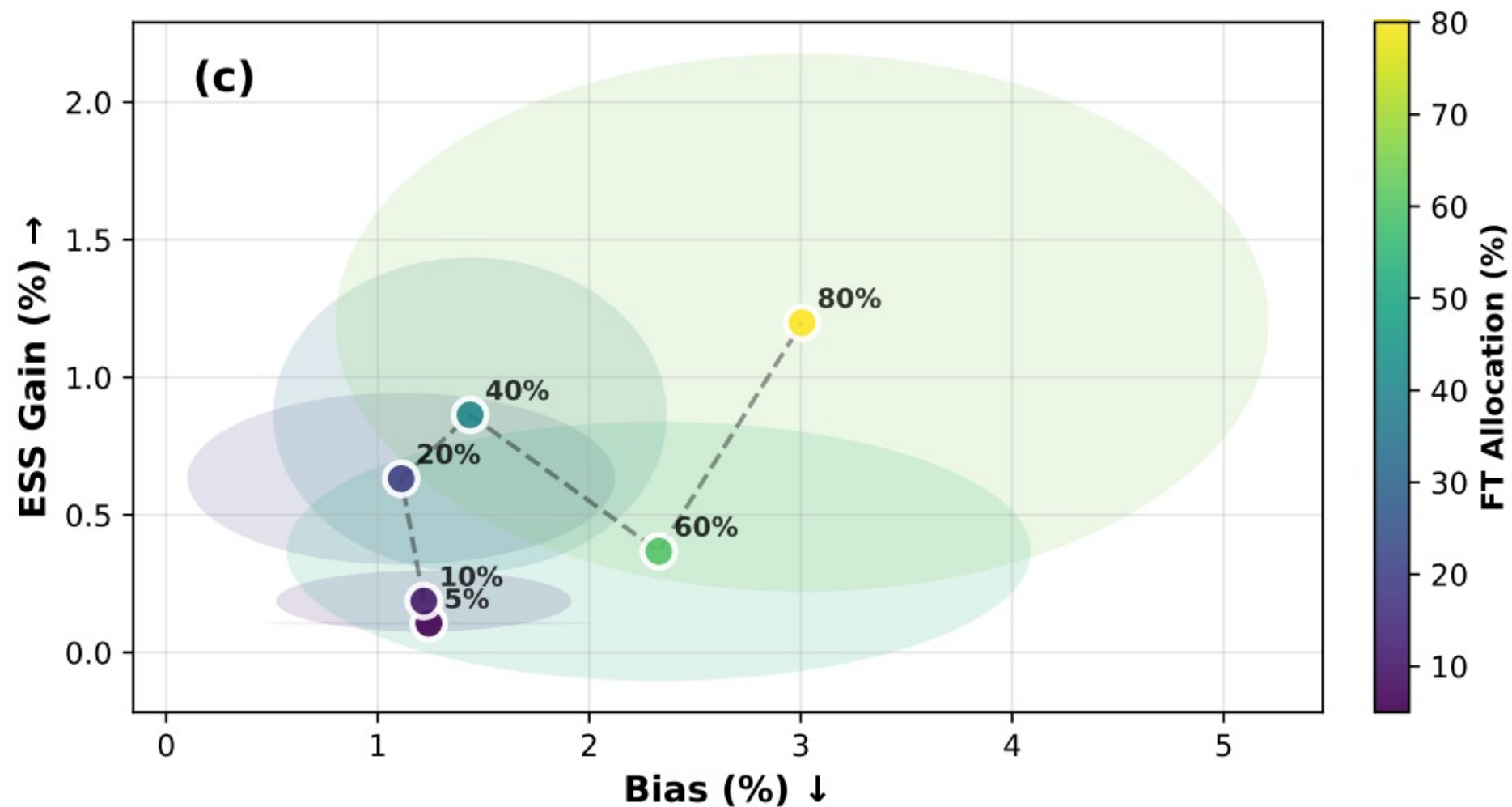
<sup>2</sup>Stanford University

<sup>3</sup>University of California, Berkeley

<sup>4</sup>Johns Hopkins University

To appear @ ACL 2026

## What if we want to do fine-tuning?



# **What we're missing: Opportunities for future research**

LLM annotations with multi-modal inputs

What if there is no ground truth?

What if LLM predictions are not calibrated?

What if we want to do fine-tuning?

**What are good scientific practices?**

# The problem of “LLM hacking”

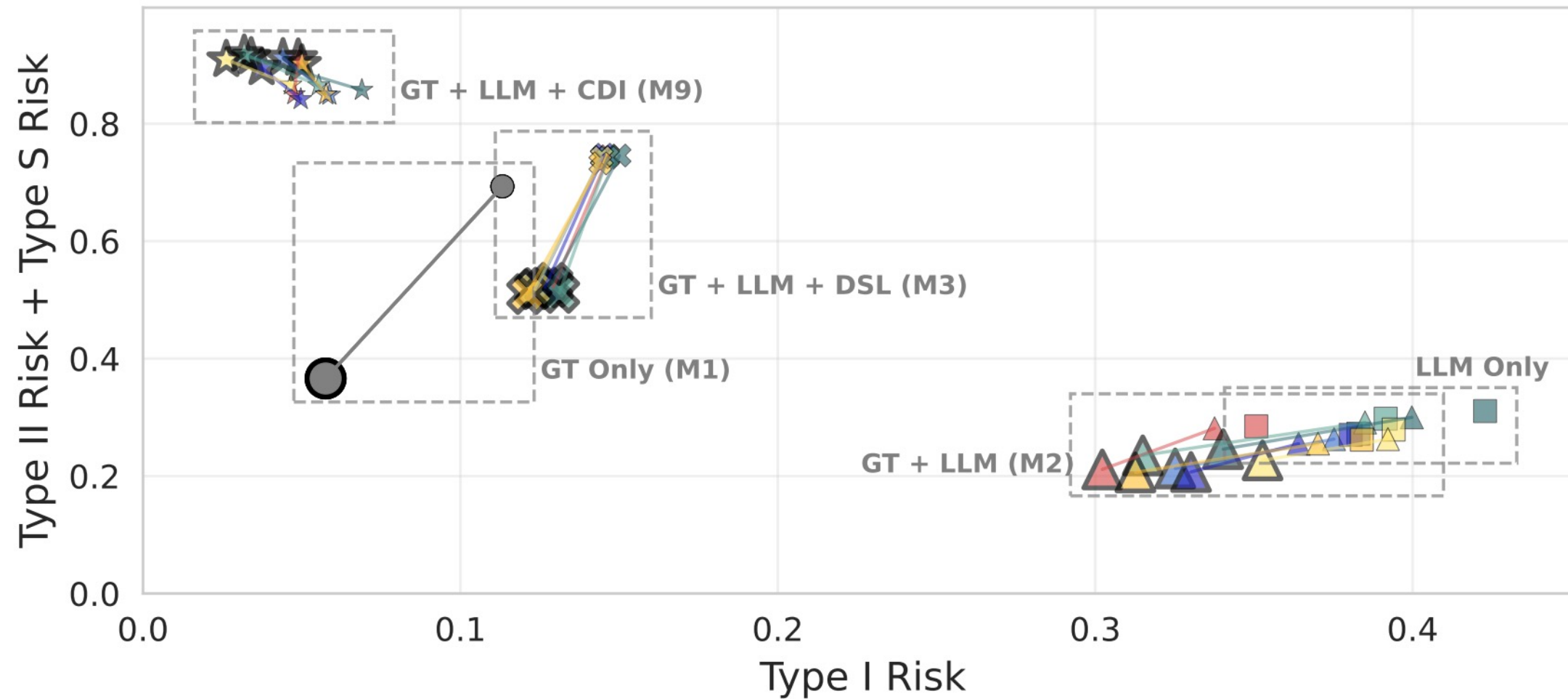
Every LLM-based annotation requires researchers to make numerous configuration choices, including:

- which model to use
- how to formulate the prompt
- which decoding parameters to set
- how to map outputs to categories
- ...

**These choices become a  
“garden of forking paths”**

Baumann, Joachim, et al. "Large language model hacking: Quantifying the hidden risks of using llms for text annotation." *arXiv preprint arXiv:2509.08825* (2025).

# The problem of “LLM hacking”



Baumann, Joachim, et al. "Large language model hacking: Quantifying the hidden risks of using llms for text annotation." *arXiv preprint arXiv:2509.08825* (2025).

**contact:**  
gligoric@jhu.edu