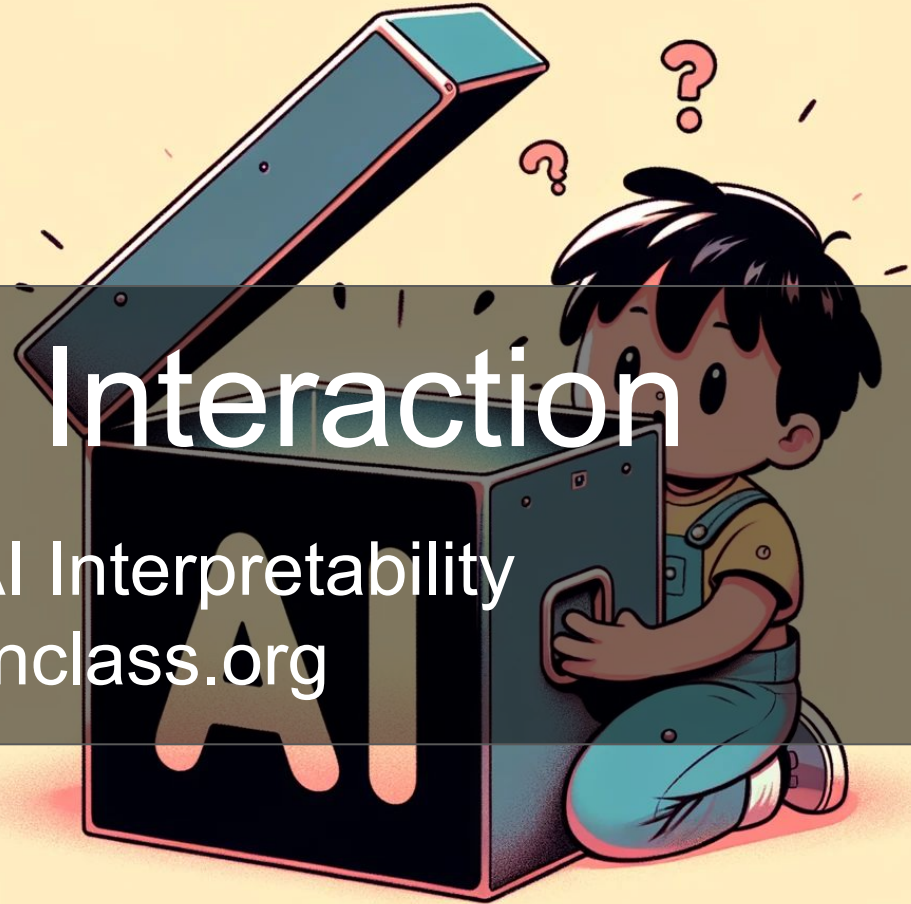# Human AI Interaction

## Lecture 16: AI Interpretability
aidesignclass.org

# Learning goals today

- What is interpretability and when is it helpful?
- If an algorithm can be interpreted, does it improve:
  - Decision-making?
  - Trust?
- Counterfactuals: what would make this not-true?
  - Does this help decision-making?
- What should you do as a designer?
- Project info
- Thursday: Guest lecture with Dr Ding Wang, on "responsible data"
- Reminder: if you haven't asked for APIkeys, you should!

# Bad news: there is no consensus definition of interpretability

But the general hypothesis is: "If you can follow the reasoning of an AI system, then you can know if its answers are correct"

E.g. Decision-tree "AI" which tells you whether you should walk to school:

IF weather is bad, THEN don't walk (take the bus)

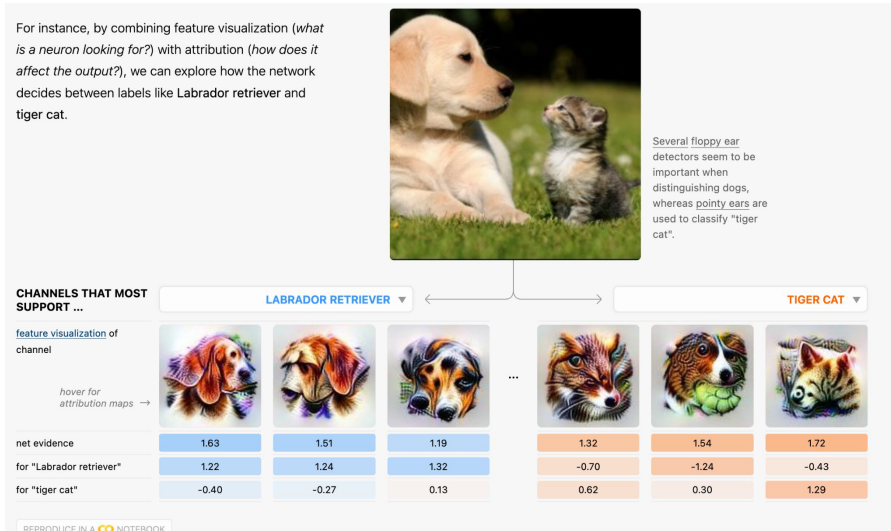IF weather is good && School is close, THEN walk. ELSE, don't walk

# Other examples of interpretability

## Attribution and attention

This page does three things: visualize features, show where they are detected, and show net evidence for the feature

Your task#1:

- Play with a few examples on this page
- Do you find these neural networks more understandable?

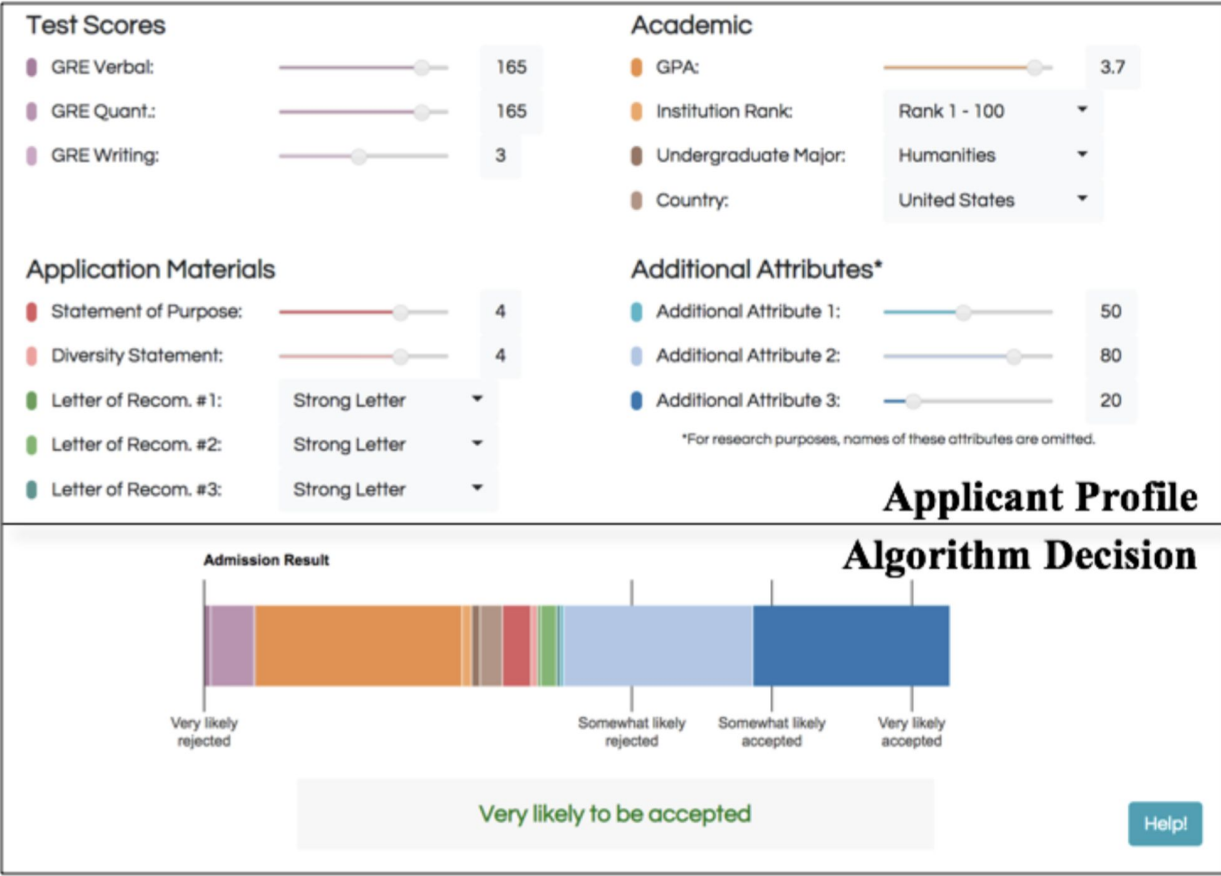# Task #2: can you predict what features will be used?

… to identity this bird?



Image from

# Interpretability by breaking open the black box

Imagine you had a system that determined if a student was admitted into grad school

(Images on the next 5 slides from this paper by Hao-Fei Cheng 1, Ruotong Wang, Zheng Zhang, Fiona O'Connell, Terrance Gray, F. Maxwell Harper, Haiyi Zhu)



**Test Scores**
- GRE Verbal: 165
- GRE Quant.: 165
- GRE Writing: 3

**Academic**
- GPA: 3.7
- Institution Rank: Rank 1 - 100
- Undergraduate Major: Humanities
- Country: United States

**Application Materials**
- Statement of Purpose: 4
- Diversity Statement: 4
- Letter of Recom. #1: Strong Letter
- Letter of Recom. #2: Strong Letter
- Letter of Recom. #3: Strong Letter

**Additional Attributes***
- Additional Attribute 1: 50
- Additional Attribute 2: 80
- Additional Attribute 3: 20

*For research purposes, names of these attributes are omitted.

**Applicant Profile**
**Algorithm Decision**

Admission Result

Very likely rejected | Somewhat likely rejected | Somewhat likely accepted | Very likely accepted

Very likely to be accepted

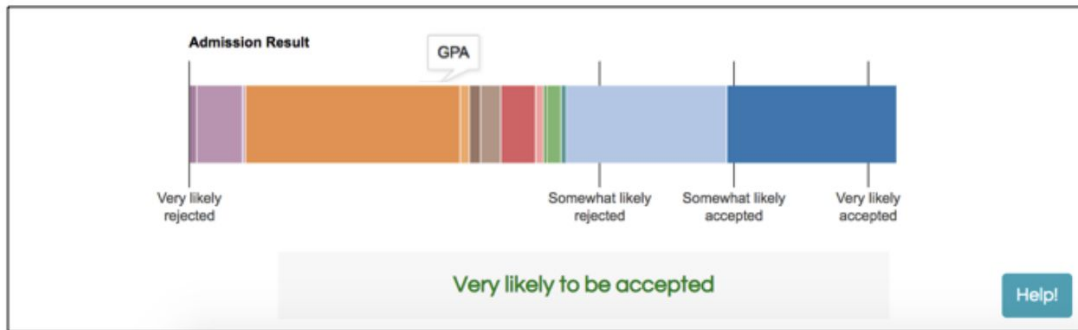Help!

# Interpretability by breaking open the black box

Imagine you had a system that determined if a student was admitted into grad school

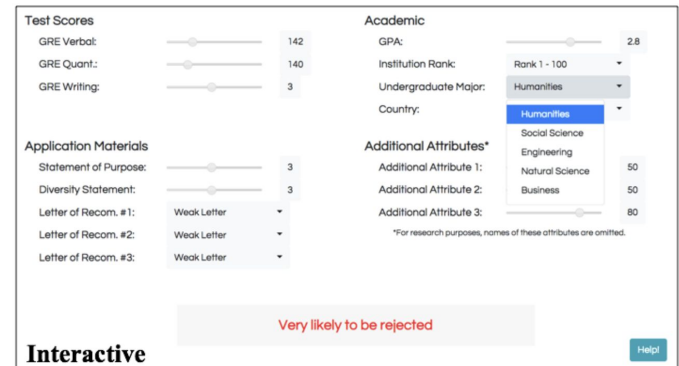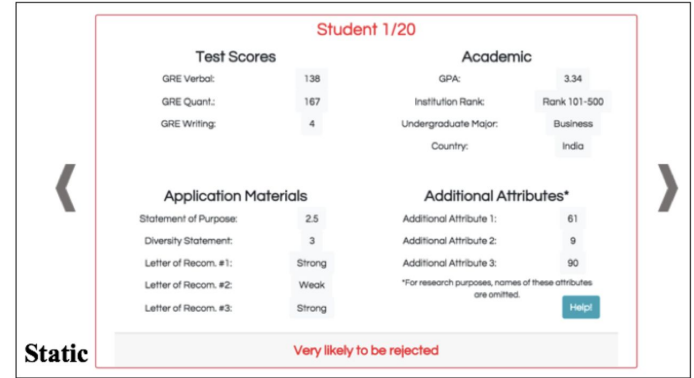- Many ways to do interpretability: let's first see black vs. white box

**Black-box**

Very likely to be accepted          Help!

**White-box**

Admission Result          GPA

Very likely rejected          Somewhat likely rejected     Somewhat likely accepted     Very likely accepted

Very likely to be accepted          Help!

# Interpretability by breaking open the black box

Imagine you had a system that determined if a student was admitted into grad school

- Static vs. dynamic explanations

# Interpretability by breaking open the black box

Imagine you had a system that determined if a student was admitted into grad school
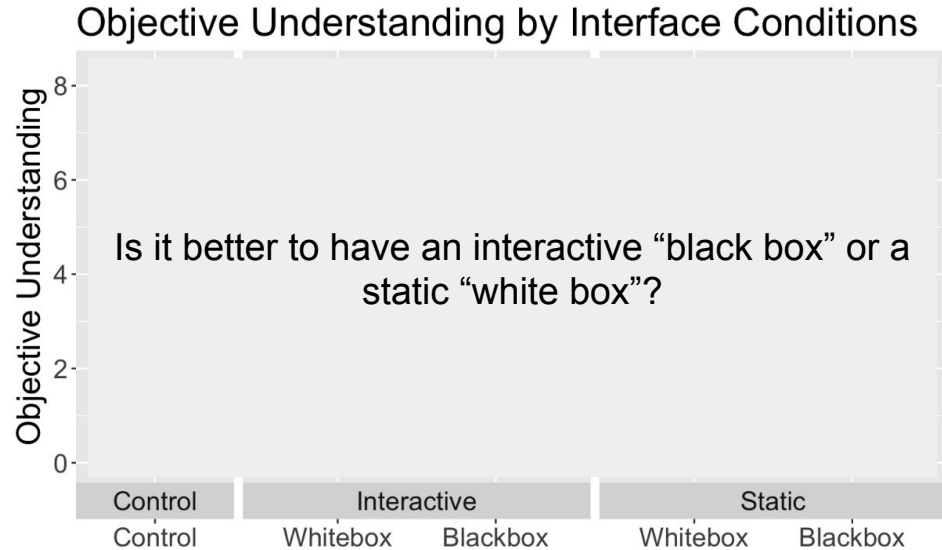
Predict results!

### Objective Understanding by Interface Conditions

Is it better to have an interactive "black box" or a static "white box"?

Objective Understanding (y-axis): 0, 2, 4, 6, 8

Control: Control
Interactive: Whitebox, Blackbox
Static: Whitebox, Blackbox

**Figure 2: Participants' objective understanding of the algorithms by interface conditions. Error bars represent 95% confidence intervals.**

# Interpretability: design implications

- Interactive "whitebox" models are most understandable
- When you can't open the blackbox (i.e. reveal how it works), interactivity has nearly the same benefit.
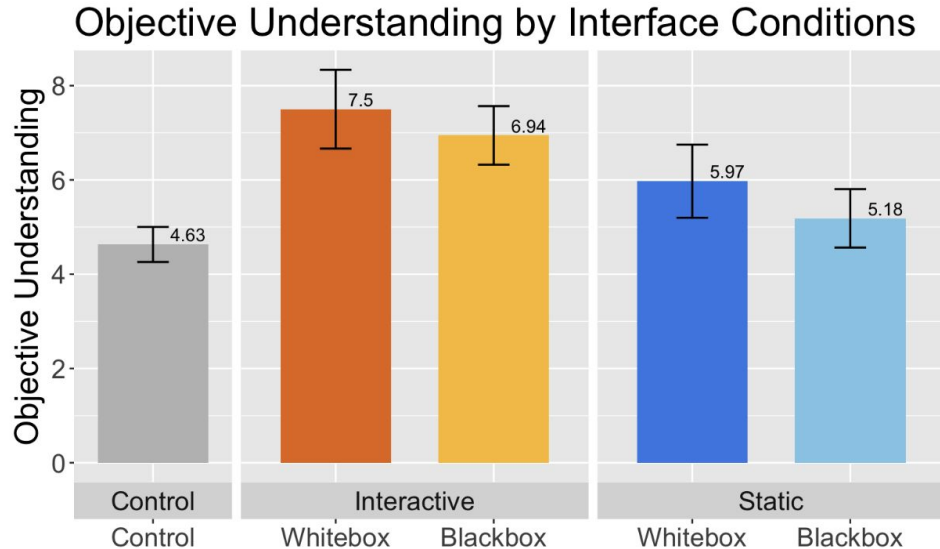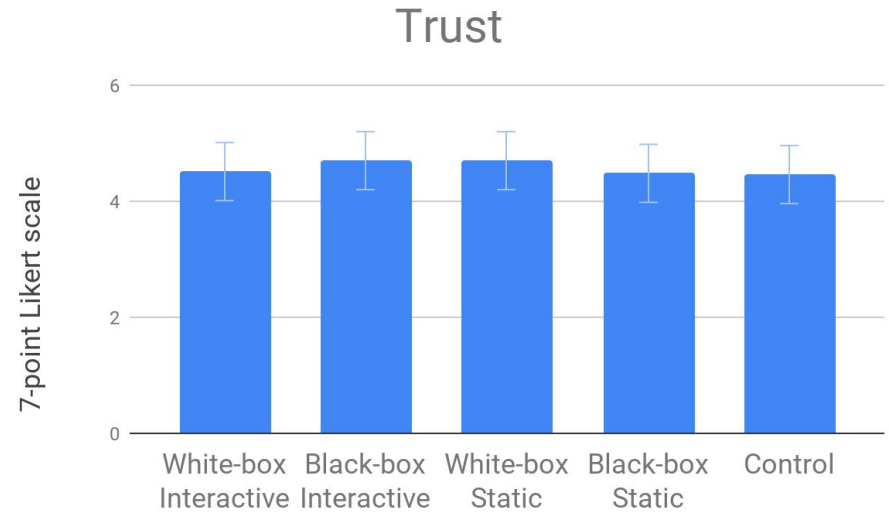


Figure 2: Participants' objective understanding of the algorithms by interface conditions. Error bars represent 95% confidence intervals.
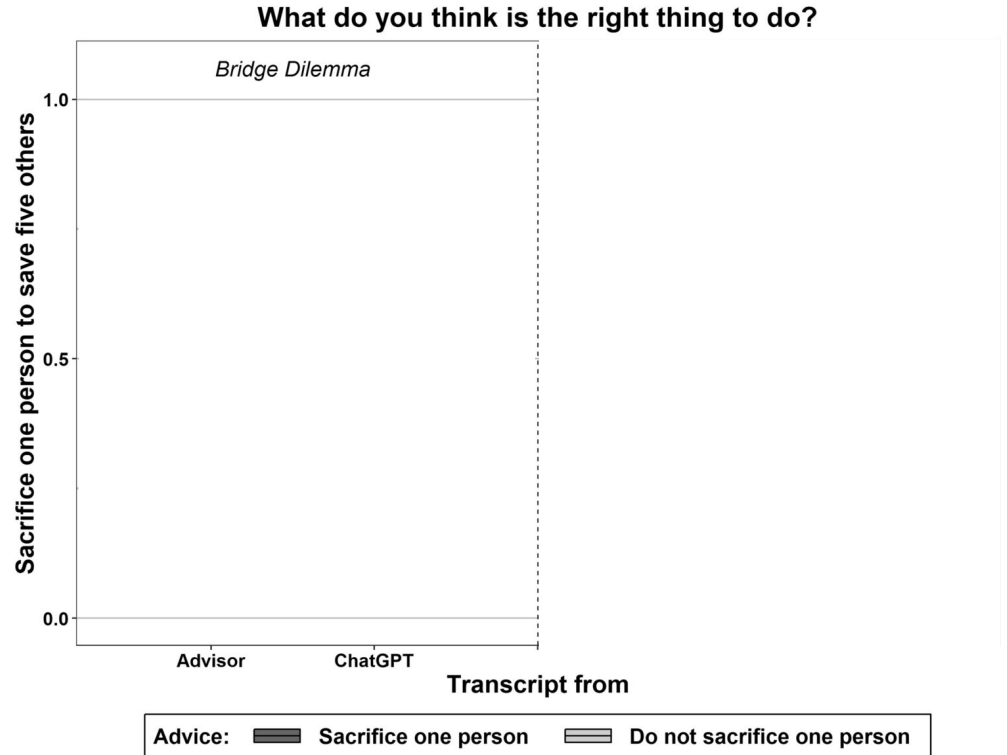
# Trust

Even when participants (don't) understand the algorithm, they may still trust it


(chart from Haiyi Zhu)

## Trust



7-point Likert scale — bar chart showing Trust across conditions: White-box Interactive, Black-box Interactive, White-box Static, Black-box Static, Control (all approximately 4.5 on the scale)

# Trust: would you trust an AI anyway?



What do you think is the right thing to do?

*Bridge Dilemma*

Sacrifice one person to save five others

1.0

0.5

0.0

Advisor          ChatGPT

Transcript from

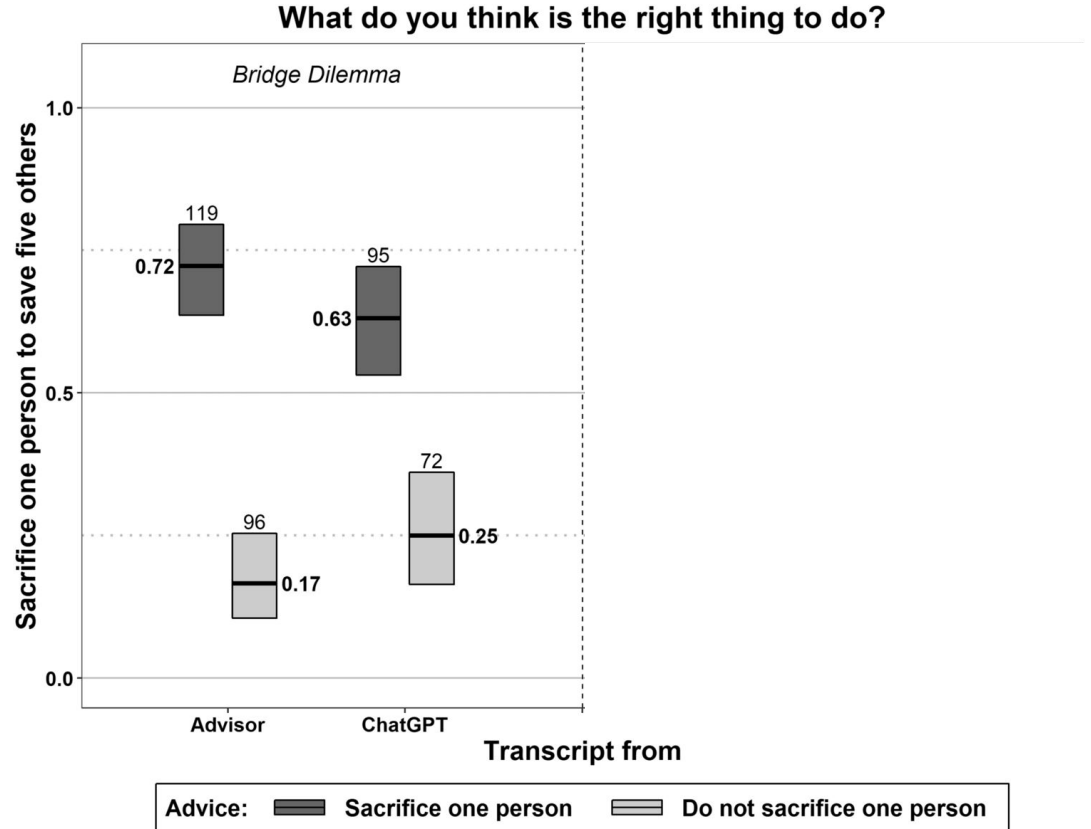Advice: ■ Sacrifice one person    ▨ Do not sacrifice one person

# Trust: it gets worse

People take advice on ethical issues from AI, even when the AI is inconsistent!

- Merely telling people "hey this comes from a probabilistic AI system" isn't enough to discount its dubious advice.

Chart from this paper.



**What do you think is the right thing to do?**

*Bridge Dilemma*

Sacrifice one person to save five others

119 — 0.72
95 — 0.63
96 — 0.17
72 — 0.25

Advisor    ChatGPT

**Transcript from**

Advice: ■ Sacrifice one person    ☐ Do not sacrifice one person

# Counterfactuals

Counterfactual: "That which is not the case"
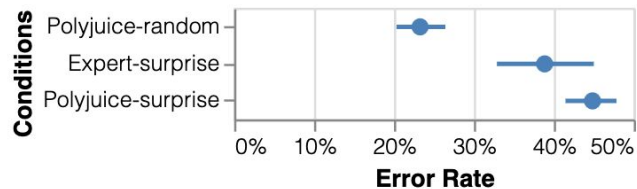
Images from [Polyjuice](#)



Figure 4: Simulation error rates per condition (higher the better). POLYJUICE-*surprise* has the highest error rate, indicating these counterfactuals would add the most information to users if displayed.
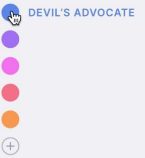
# Project info

Teams up to three

Default project: "create Daemons people can trust"

- Starter code is provided
- Must implement prompts and ask at least 5 users questions around trust (remember: benevolence, ability, integrity)



Language models are not yet good enough to be reliable thinking partners. Their frequent hallucinations make it difficult to know if their factual claims are valid. They are unreliable narrators until proven otherwise. If you ask them for references, they'll happily generate very real sounding journal names, author names, and URLs. None of which exist.

Until we drastically improve their ability to respond with accurate factual information and real

DEVIL'S ADVOCATE

# Project info

Teams up to three

Alternative project

- You pick what you want to do
- Allowed to reuse a project you are working on
- Requirements: must involve some implementation, some measurement of a concept of interest
- **MUST GET PRIOR APPROVAL**
- Rewarded for "principled risk taking"



DEVIL'S ADVOCATE

Language models are not yet good enough to be reliable thinking partners. Their frequent hallucinations make it difficult to know if their factual claims are valid. They are unreliable narrators until proven otherwise. If you ask them for references, they'll happily generate very real sounding journal names, author names, and URLs. None of which exist.

Until we drastically improve their ability to respond with accurate factual information and real