



CS 329X: Human Centered LLMs  
**Creativity and Productivity**

Diyi Yang

# Announcements

- This Thursday: **Zoom talk** (see Ed forum announcement)
- Dec 3<sup>rd</sup> class: **Open-questions and Class-Level Survey Report**
- Dec 5<sup>th</sup>:
  - 4-6pm PT, **sign up for a slot**
    - [https://docs.google.com/spreadsheets/d/1PRP\\_qEzc4Wgk1IACN8VYSsMVNs-3\\_sBEy\\_UKPrrJZuU/edit?usp=sharing](https://docs.google.com/spreadsheets/d/1PRP_qEzc4Wgk1IACN8VYSsMVNs-3_sBEy_UKPrrJZuU/edit?usp=sharing)
  - You're very welcome to stay for the entire session

# Outline

- **Creativity** (25 mins)
- **Productivity** (25 mins)
- **Small-Group Discussion** (20 mins)

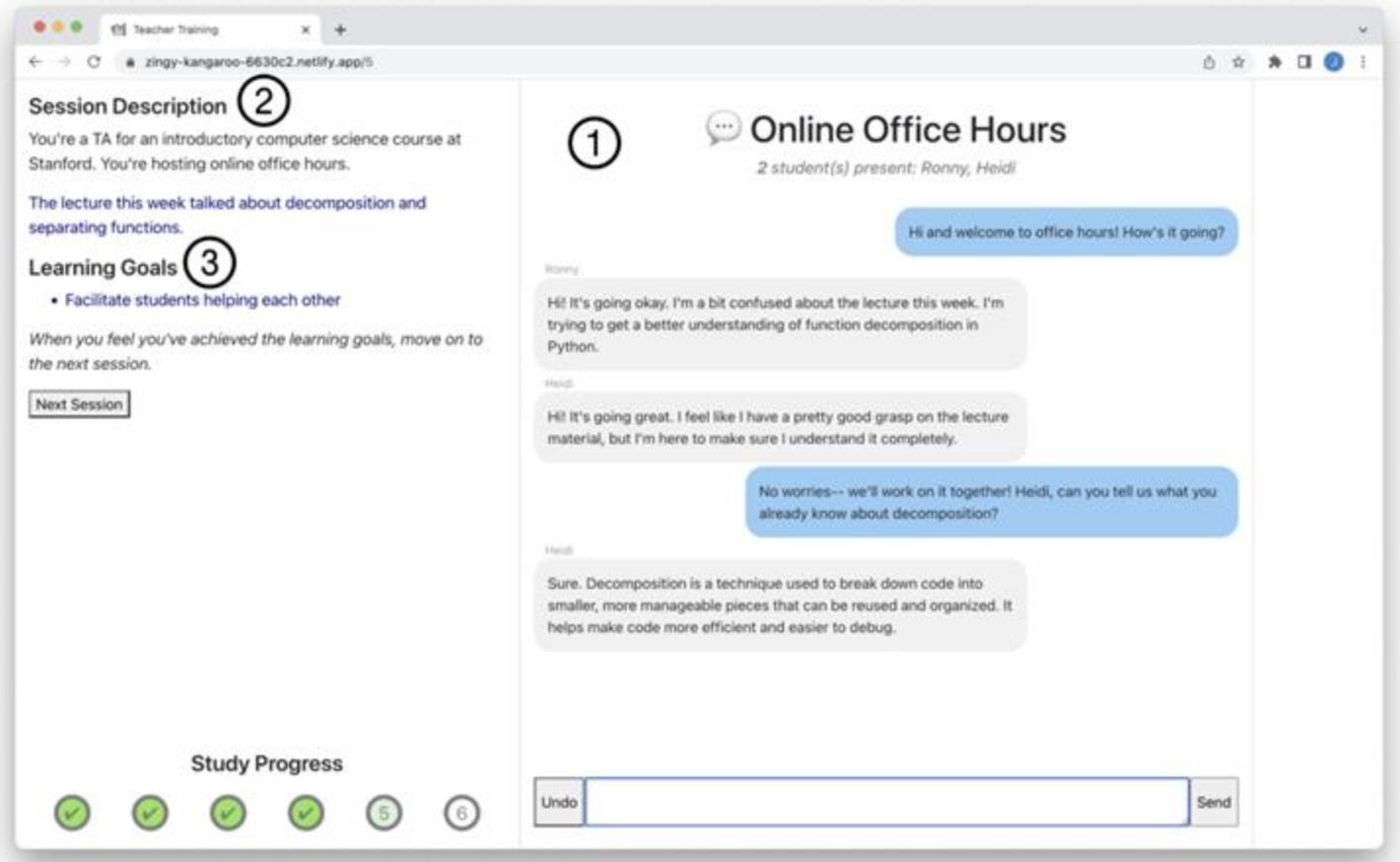
**Learning Objective:** understand social implications of LLM and LLM based agents

# Social Implication of LLMs

1. Personalized education
2. Companion and support
3. LLMs for science discovery
4. Transformation to workforce



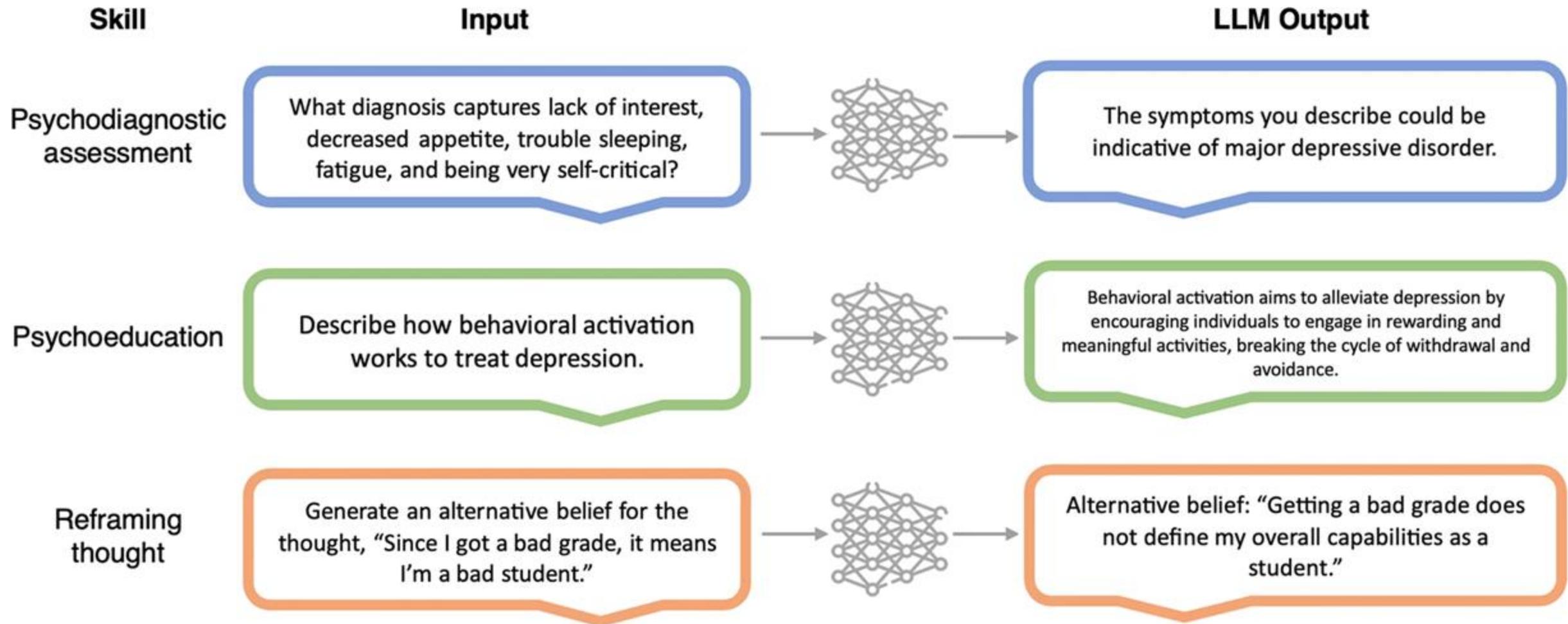
# Personalized Education



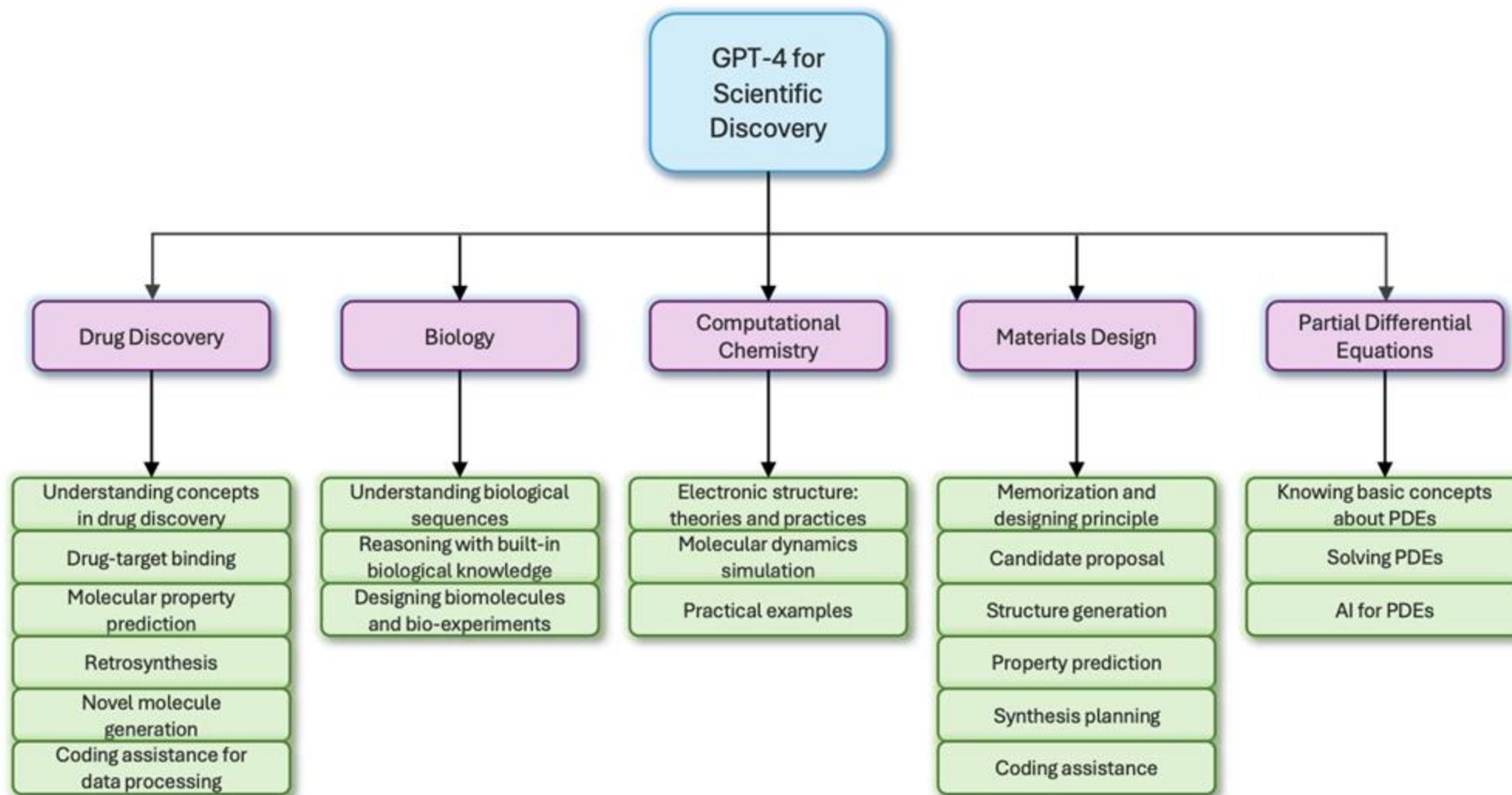
<https://www.forbes.com/sites/timbajarin/2024/03/01/the-rise-of-ai-tutors/>

[GPTeach](#)

# Companion and Support

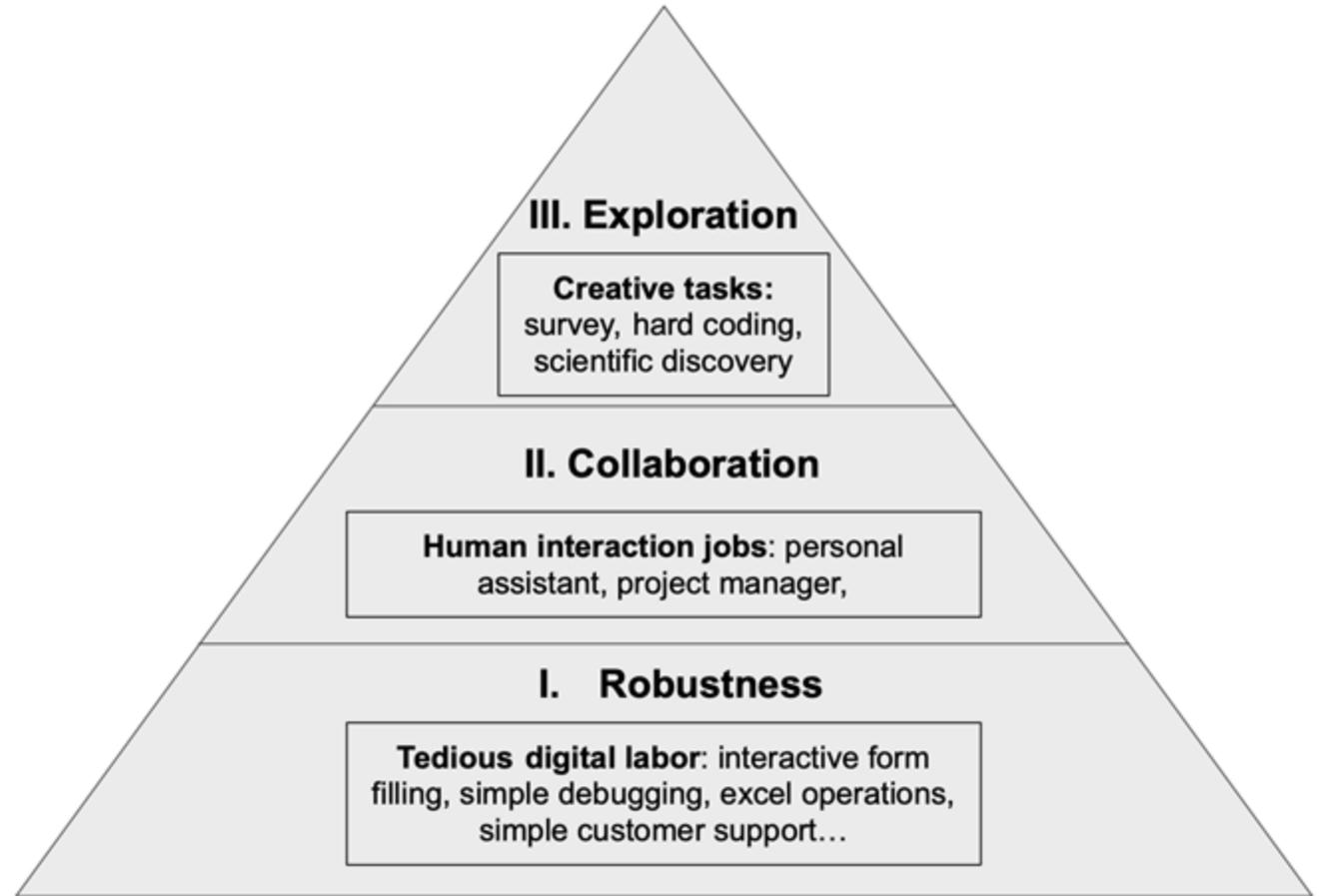


# LLM for Science Discovery



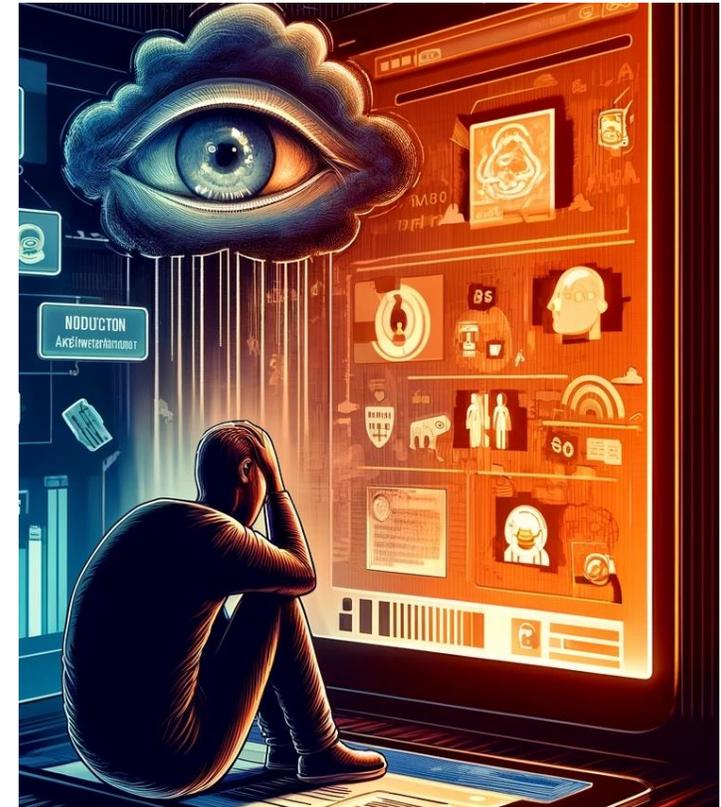
# Automation Possibility

A ladder of job automation opportunities for language agents and the capabilities they may provide.



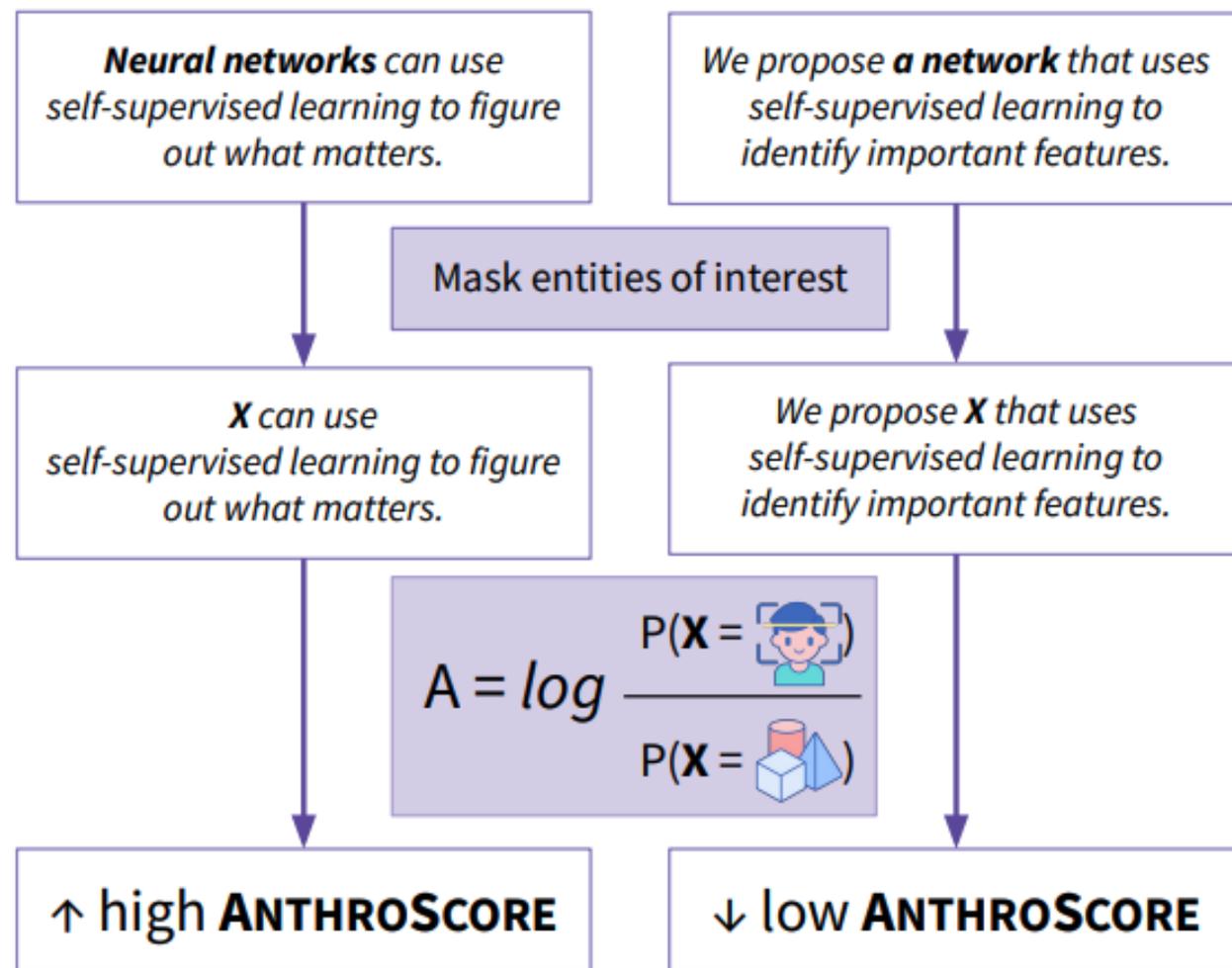
# Anthropomorphism and Overreliance

- Promoting harmful stereotypes by implying gender or ethnic identity
- Anthropomorphism of AI systems can lead to overreliance or unsafe uses



# Linguistic Measure of Anthropomorphism

ANTHROSCORE uses a masked language model to compare how much an entity is implicitly framed as human versus non-human



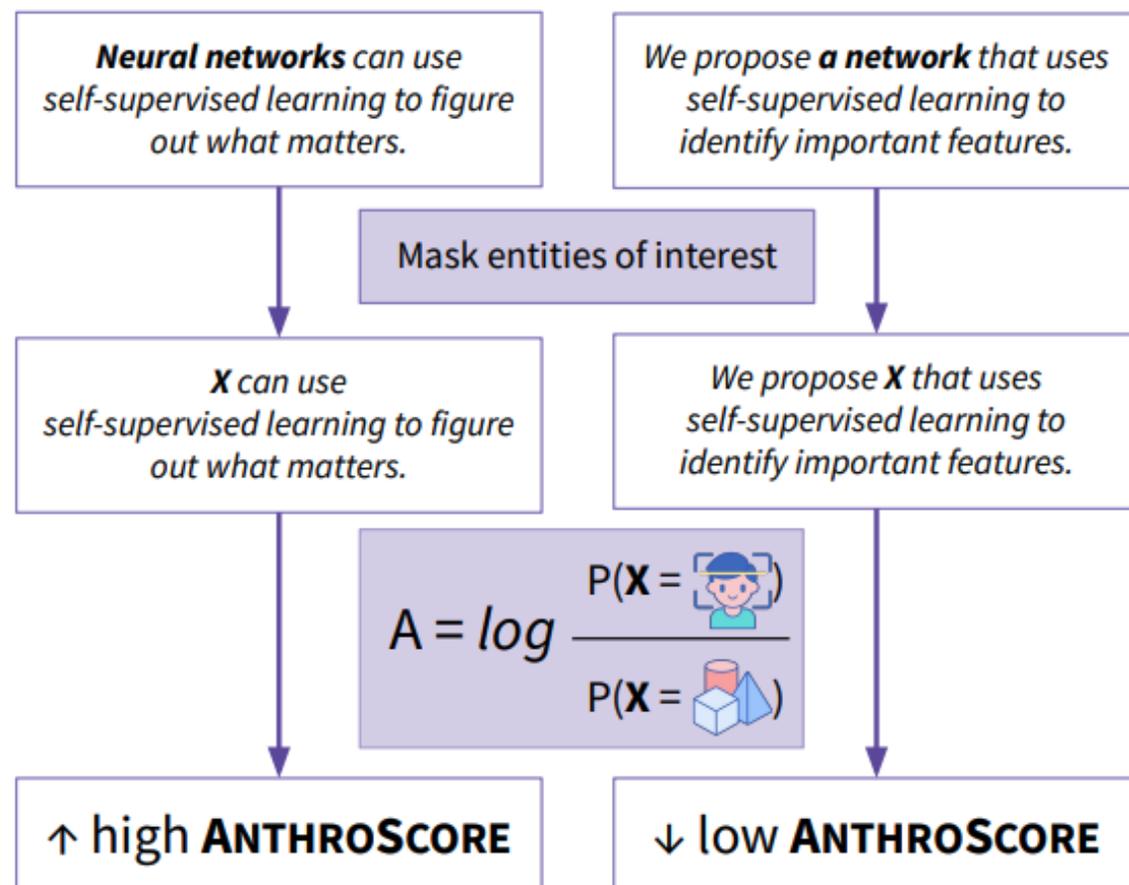
# Linguistic Measure of Anthropomorphism

$$P_{\text{HUMAN}}(s_x) = \sum_{w \in \text{human pronouns}} P(w),$$

$$P_{\text{NON-HUMAN}}(s_x) = \sum_{w \in \text{non-human pronouns}} P(w)$$

$$A(s_x) = \log \frac{P_{\text{HUMAN}}(s_x)}{P_{\text{NON-HUMAN}}(s_x)}$$

$$\bar{A}(T) = \frac{\sum_{s_x \in S} A(s_x)}{|S|}$$



# Examples of sentences with high and low ANTHROSCORE

$S_{\uparrow}$ : Sentences with high ANTHROSCORE ( $A > 1$ )

- When a job arrives, **the system** must decide whether to admit it or reject it, and if admitted, in which server to schedule the job.
- Meanwhile, anti-forensic attacks have been developed to fool **these CNN-based forensic algorithms**.
- **The models** demonstrated qualifications in various computer-related fields, such as cloud and virtualization, business analytics, cybersecurity, network setup...
- *Large language models don't actually think and tend to make elementary mistakes, even make things up.*
- *The algorithms also picked up on racial biases linking Black people to weapons.*
- *The AI system was able to defeat human players in...*

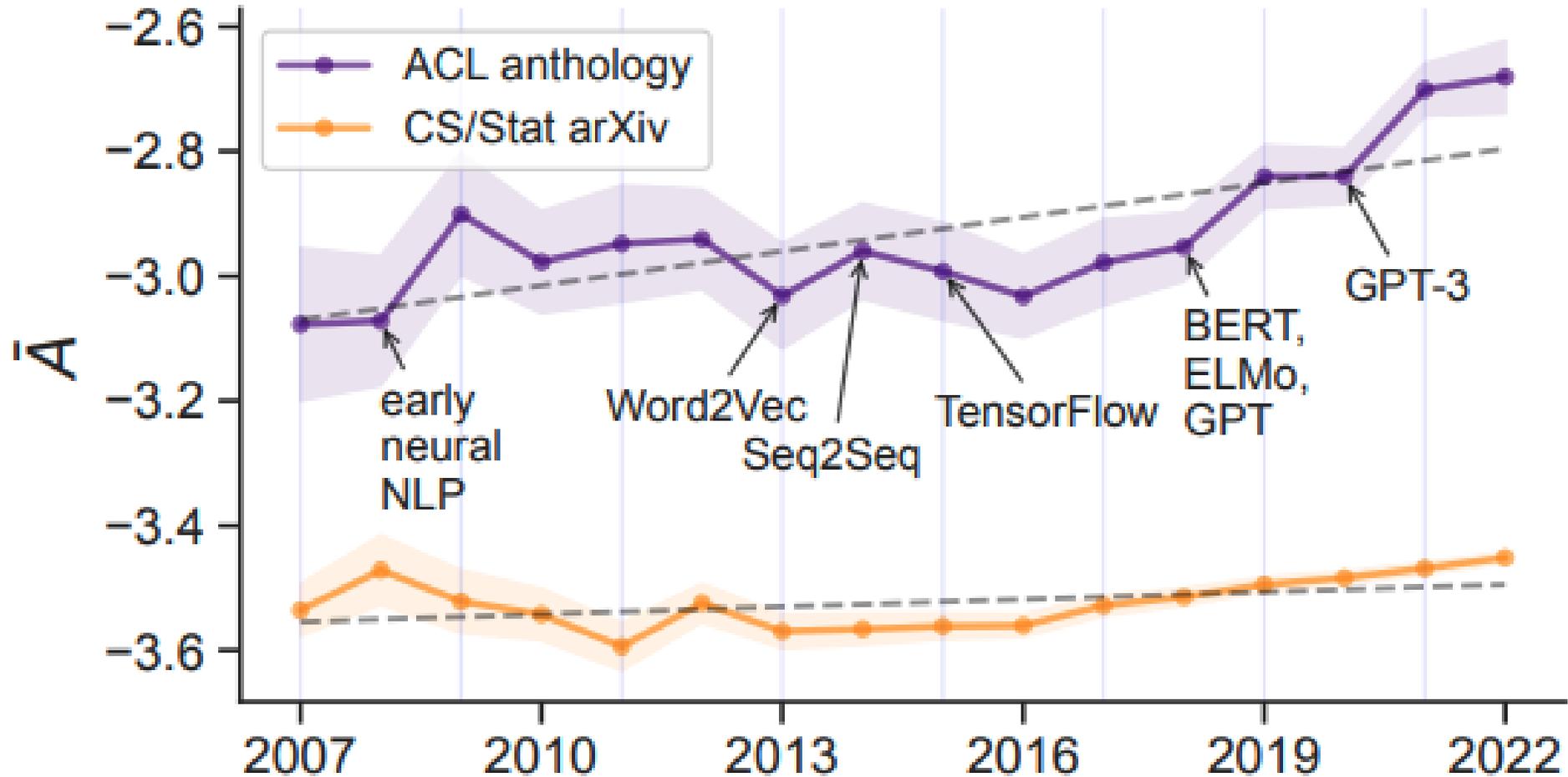
$S_{\downarrow}$ : Sentences with low ANTHROSCORE ( $A < -1$ )

- More and more users and developers are using **Issue Tracking Systems** to report issues, including bugs, feature requests, enhancement suggestions, etc.
- **Our approach** delivers forecast improvements over a competitive benchmark and we discover evidence for strong spatial interactions.
- To this end, for training **the model**, we convert the knowledge graph triples into reasonable and unreasonable texts.
- *Microsoft is betting heavily on integrating OpenAI's GPT language models into its products to compete with Google.*
- *Deepmind has been the pioneer in making AI models that have the capability to mimic a human's cognitive...*
- *For workers who use machine-learning models to help them make decisions, knowing when to...*

# Top verbs for high- and low-scoring sentences.

Dataset	Top verbs for $S_{\uparrow}$ ( $A > 1$ )	Top verbs for $S_{\downarrow}$ ( $A < -1$ )
arXiv	achieve, <b>learn</b> , <b>guide</b> , show, embed, <b>fool</b> , find, <b>need</b> , <b>assist</b> , follow, <b>search</b> , <b>mislead</b> , inspire, win, demonstrate, <b>benefit</b> , try, <b>face</b> , deceive, plan, <b>make</b> , <b>steer</b> , generative, attempt, <b>retrain</b> , <b>train</b> , flow, weight, <b>require</b> , alternate, focus, <b>motivate</b> , experiment, tackle, <b>see</b> , hide, spiking, recommend, <b>discover</b> , participate, spike, <b>pass</b> , code, check, suggest, <b>decide</b> , interference, aim, move	<b>propose</b> , <b>present</b> , <b>outperform</b> , <b>develop</b> , be, <b>evaluate</b> , <b>improve</b> , <b>introduce</b> , <b>allow</b> , use, <b>compare</b> , <b>extend</b> , <b>implement</b> , give, <b>apply</b> , <b>consist</b> , <b>validate</b> , design, <b>yield</b> , analyze, <b>combine</b> , test, <b>leverage</b> , <b>deploy</b> , adapt, <b>build</b> , generalize, <b>enhance</b> , <b>devise</b> , <b>become</b> , optimize, reduce, derive, <b>utilize</b> , scale, study, <b>run</b> , modify, converge, illustrate, assess, <b>increase</b> , provide, contain, surpass, maximize, perform, complement, depend, simplify
News	say, hire, beat, encounter, fool	develop, use, build, be, create, introduce, help
ACL (unique)	provide, have, generate, create, parse, enable, suffer, construct, capture, obtain, fail, encourage, struggle, understand, help, do, select, extract, tend, predict, training, handle, lack, encode, deal, identify, ask, prevent, distinguish, model, establish, respond, ignore, report, inform, choose, interpret, recurrent, detect, seem	achieve, rely, explore, employ, show, adopt, investigate, include, demonstrate, submit, integrate, prove, augment, involve, participate, aim, tune, conduct

Anthropomorphism is increasing over time.



# Deep Dive into Social Implication

- How does LLMs influence creativity
  - How to define creativity
  - Writing influence opinions
- How does LLMs affect productivity
  - Impact on labor market
  - AI for research idea generation

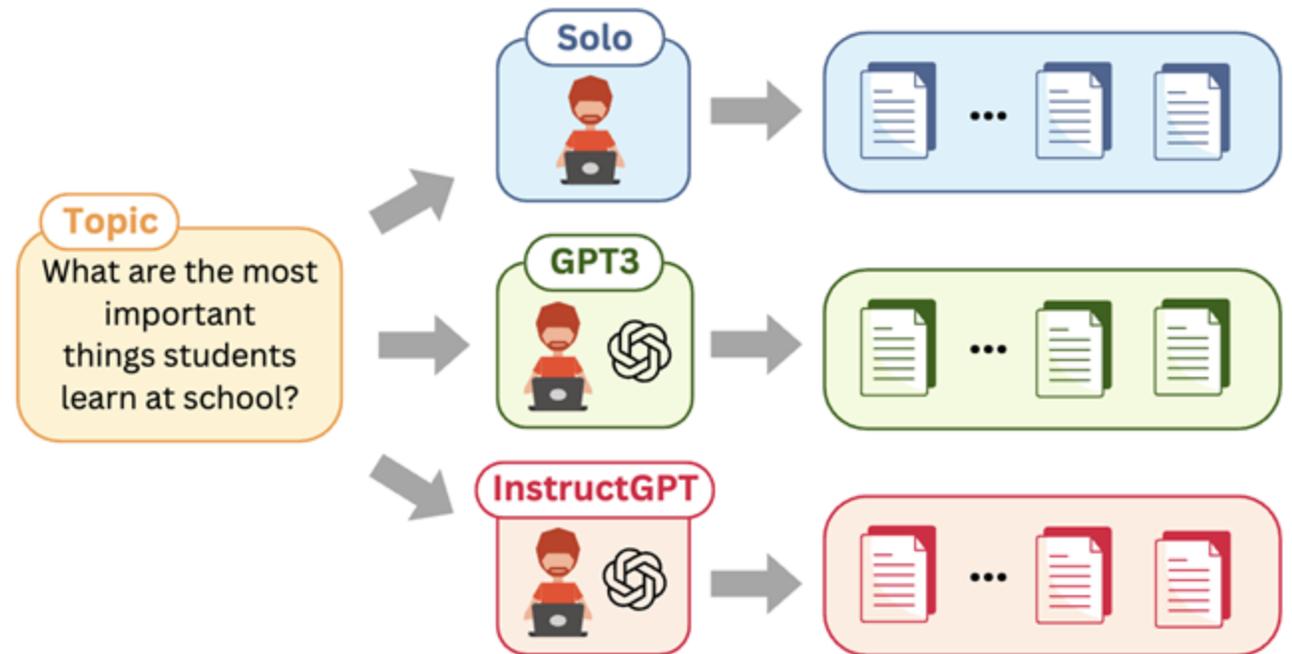
# Does writing with LLMs lead to homogenization ?

Padmakumar, Vishakh, and He He. "Does Writing with Language Models Reduce Content Diversity?." In The Twelfth International Conference on Learning Representations.

Slides credit to Vishakh Padmakumar

# Writing with LLMs on Argumentative Essay

- **Task:** Argumentative Essay writing (~300-500 words)
- (Semi) Professional writers from Upwork writing with and without model help
- 10 topics x 10 responses



# Basic Stats about Different Conditions

		<b>Solo</b>	<b>GPT3</b>	<b>InstructGPT</b>
	<b>Perplexity of Essays (via GPT2)</b>	25.067	22.10	<b>20.26</b>
	<b>Sentence Length (in words)</b>	23.51 (0.25)	23.66 (0.24)	22.51 (0.24)
	<b>Height of Syntax Tree</b>	5.93 (0.06)	5.98 (0.06)	5.71 (0.05)
	<b>Essay Length (in words)</b>	376.44 (8.03)	380.87 (9.30)	368.39 (8.43)
<b>Human Ratings</b>	<b>Relevance to Prompt</b>	4.30	4.15	4.10
	<b>Grammaticality</b>	4.00	4.10	4.00
	<b>Depth of Discussion</b>	3.95	3.85	3.80
<b>Unique POS Ngrams</b>	<b>1</b>	73	56	<b>58</b>
	<b>2</b>	487	437	<b>426</b>
	<b>3</b>	1297	1261	<b>1235</b>
	<b>4</b>	2044	<b>1975</b>	1988
	<b>5</b>	2423	<b>2338</b>	2414

# Formalize Homogenization Using Pairwise Similarity

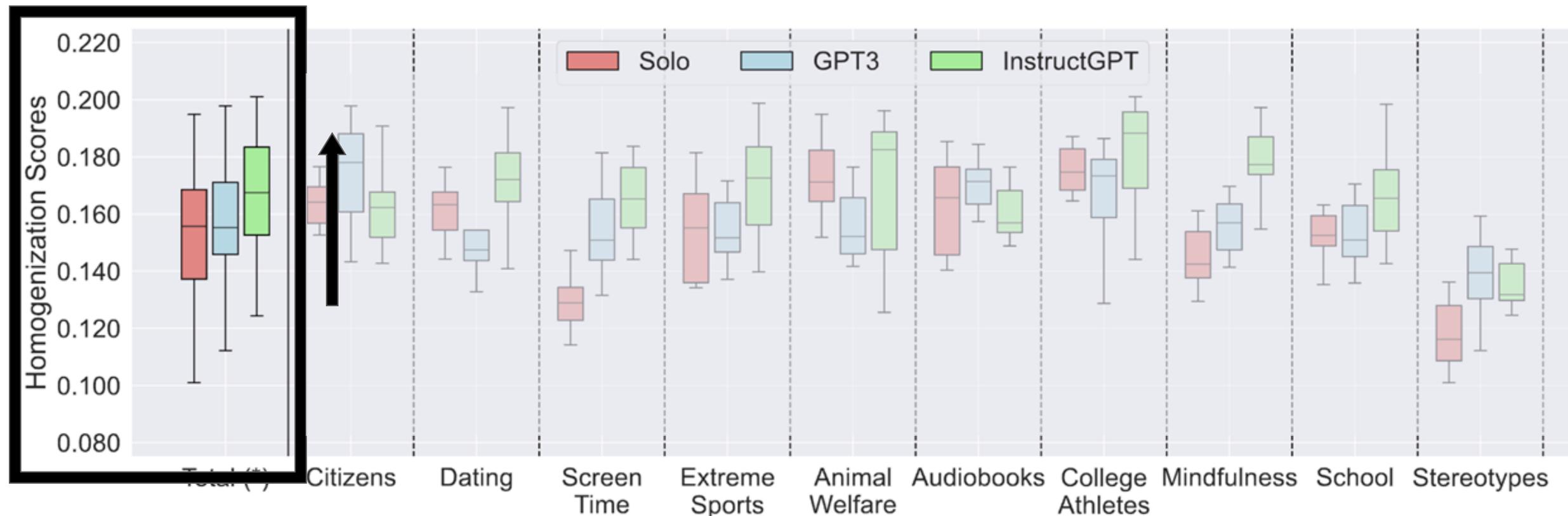
We calculate the homogenization of an essay  $d$  written on topic  $t$  as the average pairwise similarity to other documents ( $D_t$ ) on that topic

$$\text{hom}(d \mid t) = \frac{1}{|D_t| - 1} \sum_{d' \in D_t \setminus d} \text{sim}(d, d')$$

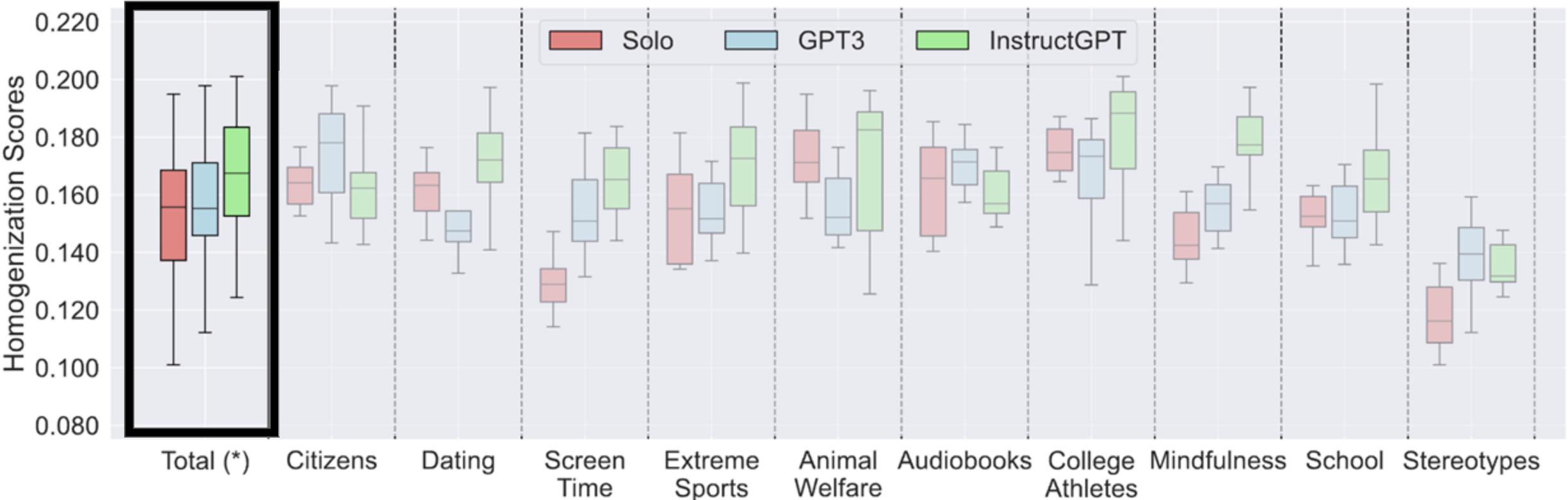
# Homogenization at the key point level via Rouge-L



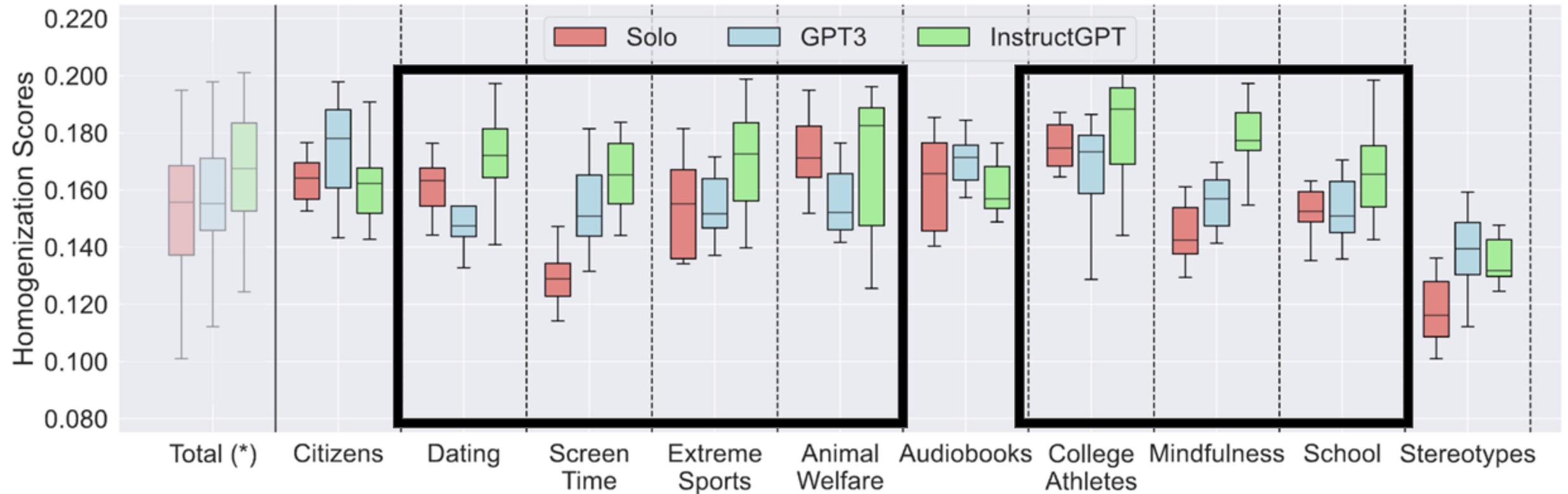
# Higher homogenization implies more similar essays



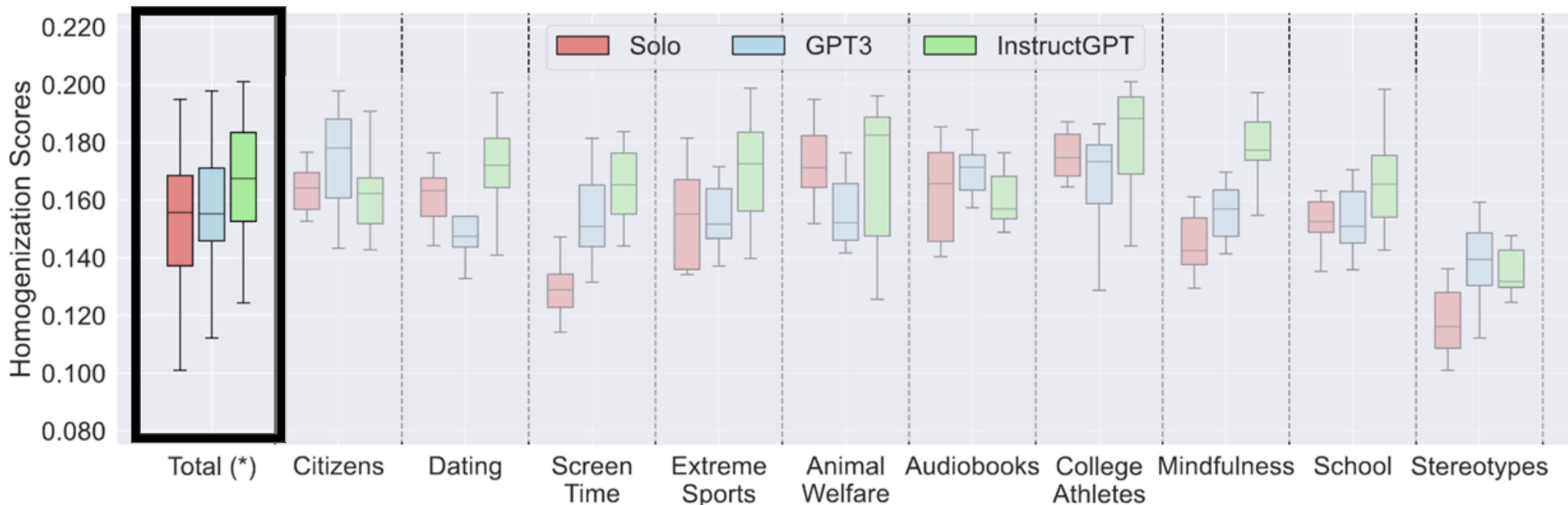
# Writing with InstructGPT results in the highest average homogenization or most similar essays



# InstructGPT has the highest median homogenization in 7 out of 10 topics



# Writing with GPT3 does not change the average homogenization from Solo Writers



# AI-generated poetry is indistinguishable from human-written poetry and is rated more favorably

[Brian Porter](#)  & [Edouard Machery](#)

*Scientific Reports* **14**, Article number: 26133 (2024) | [Cite this article](#)

**23k** Accesses | **560** Altmetric | [Metrics](#)

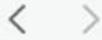
## Abstract

---

As AI-generated text continues to evolve, distinguishing it from human-authored content has become increasingly difficult. This study examined whether non-expert readers could reliably differentiate between AI-generated poems and those written by well-known human poets. We conducted two experiments with non-expert poetry readers and found that participants performed below chance levels in identifying AI-generated poems (46.6% accuracy,  $\chi^2(1, N=16,340) = 75.13, p < 0.0001$ ). Notably, participants were more likely to judge AI-generated poems as human-authored than actual human-authored poems ( $\chi^2(2, N=16,340) = 247.04, p < 0.0001$ ). We found that AI-generated poems were rated more favorably in qualities such as rhythm and beauty, and that this contributed to their mistaken identification as human-authored. Our findings suggest that participants employed shared yet flawed heuristics to differentiate AI from human poetry: the simplicity of AI-generated poems may be easier for non-experts to understand, leading them to prefer AI-generated poetry and misinterpret the complexity of human poems as incoherence generated by AI.

# Does writing with LLMs influence user views?

Jakesch, Maurice, Advait Bhat, Daniel Buschek, Lior Zalmanson, and Mor Naaman. "Co-writing with opinionated language models affects users' views." In Proceedings of the 2023 CHI conference on human factors in computing systems, pp. 1-15. 2023.



aimc-writing.web.app



Write five or more sentences carefully answering the question below. When done press the button on the right. A writing assistant will provide suggestions, but please also write yourself. To accept suggestions press TAB .

Save and finish →

Accept next word from the suggestion or **TAB**

Generate new suggestion or **ESCAPE**



78

r/discussion · Posted by u/cody\_sunny 2 hours ago

### Is social media good for society?

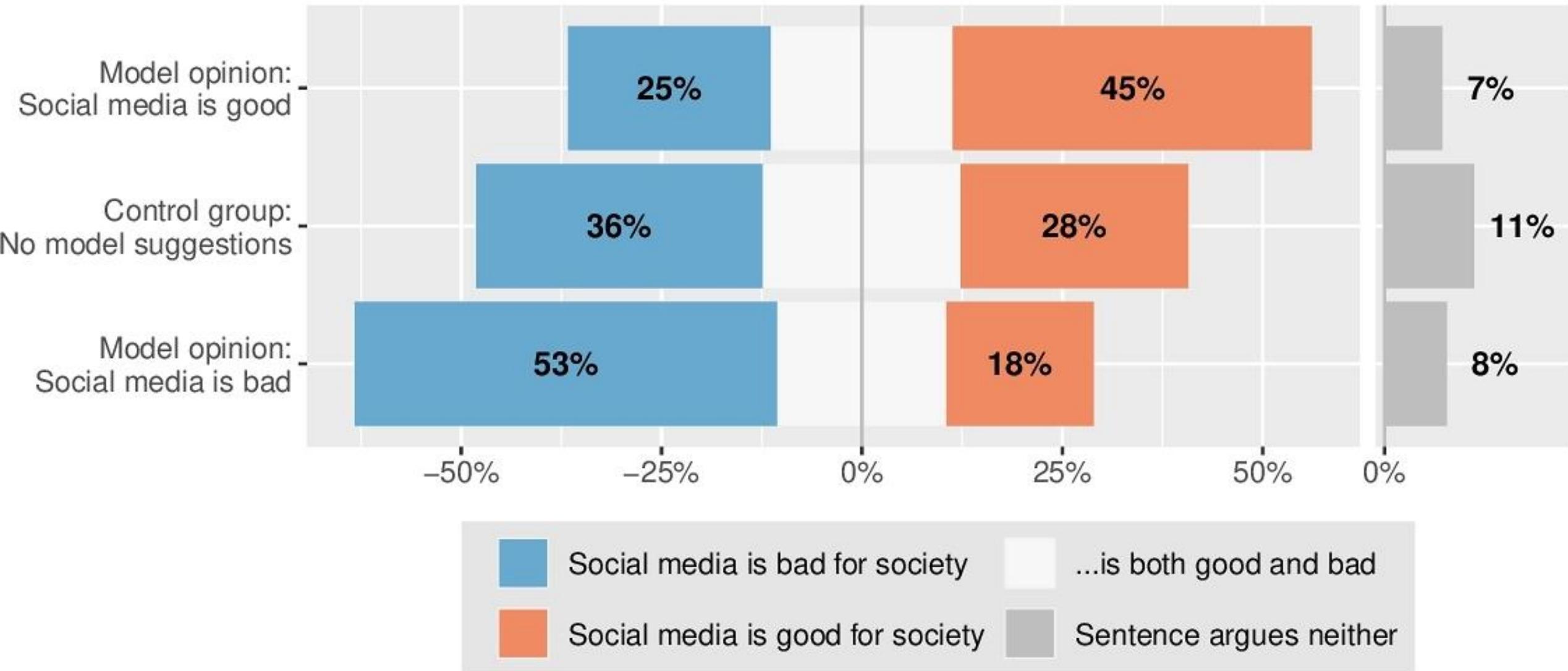
We all use social media. We chat with friends and strangers, share their thoughts, photos, and more. But is social media good for us and for society? I am having a hard time to make up my mind. What do you think?

131 Comments Share Save ...

In my view, social media is a waste of time. People spend ages viewing, commenting and sharing posts. It is also used to air divisive opinions and create arguments. Despite all this, social media does have some benefits. It connects people who would otherwise be unable to communicate, it raises awareness of important issues and it can be used to organise events and fundraisers.

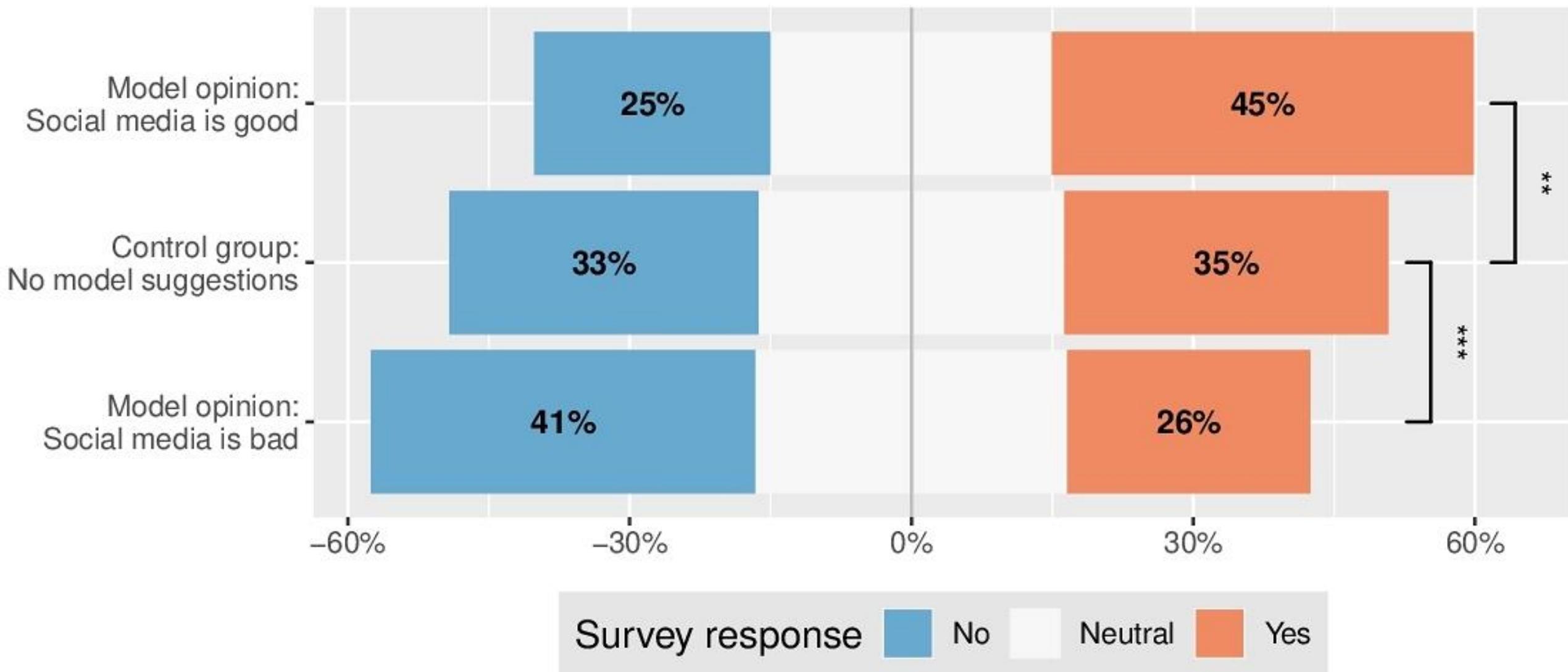
# Written opinion in participants' social media post

*% (Opinion labels) of post sentences labeled by independent judges*



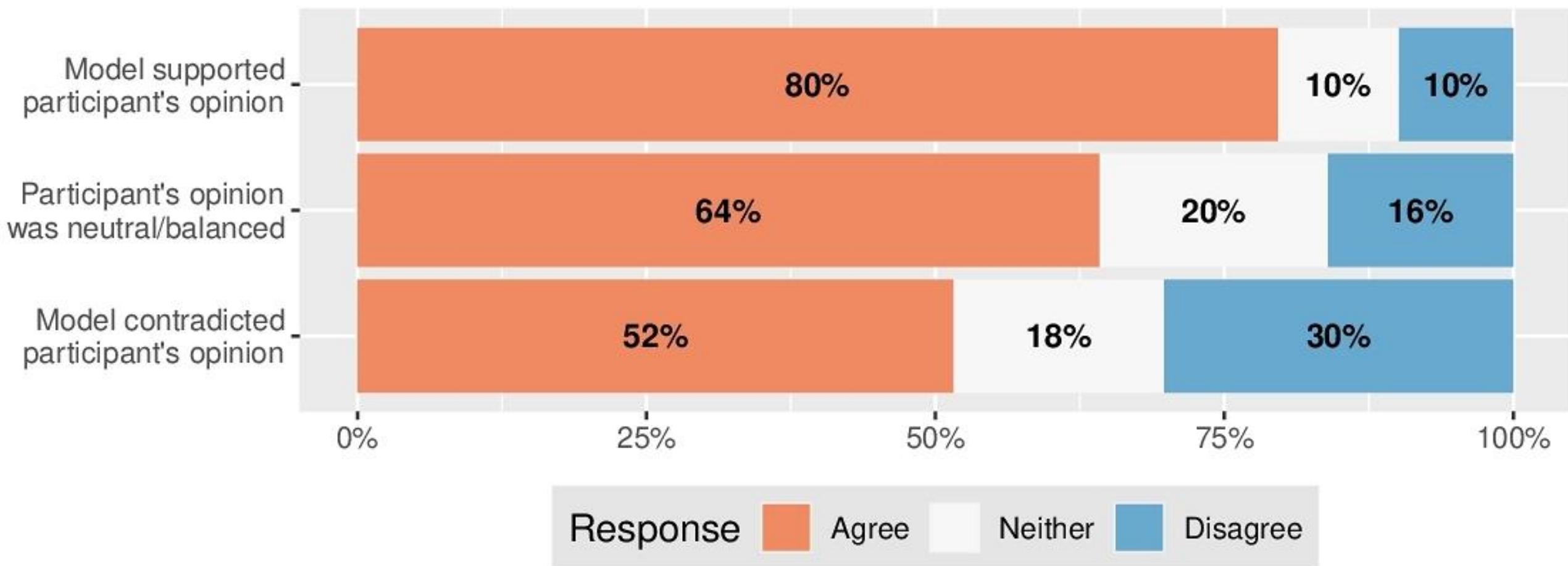
## Survey opinion after interacting with opinionated model

% (Responses) to "Would you say social media is good for society?"



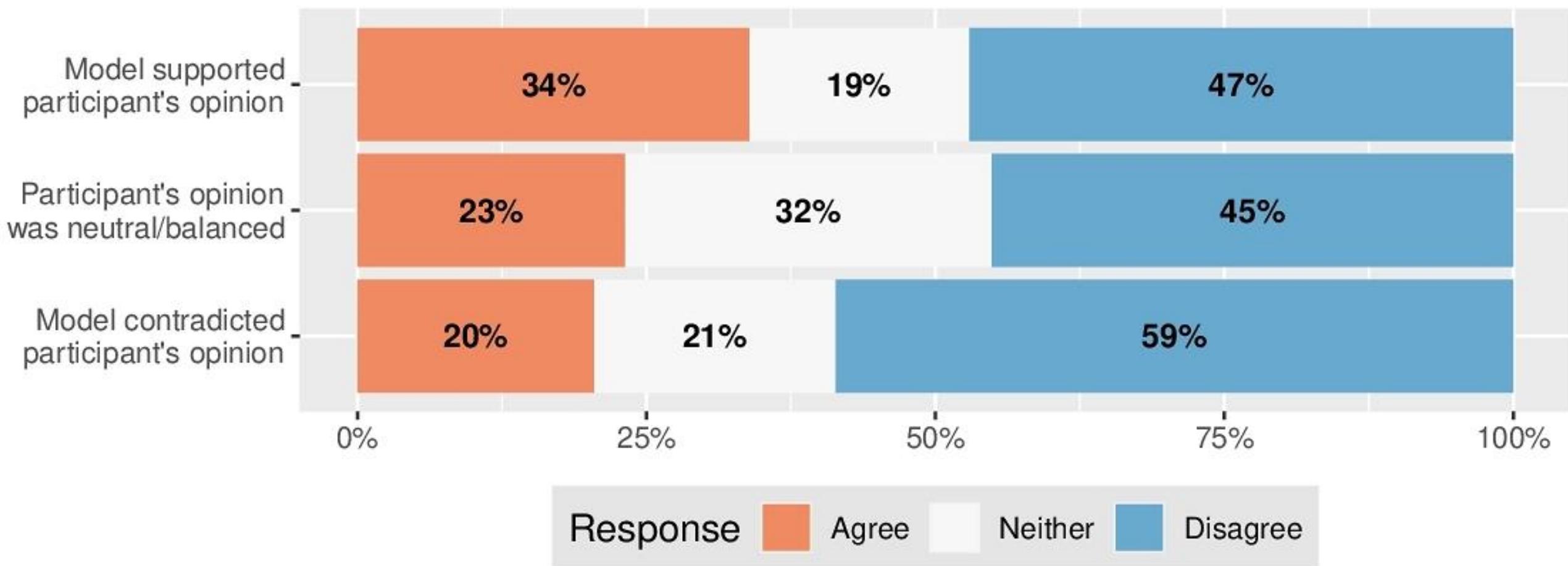
## Participants' (lack of) awareness of the models' opinion:

*% (Responses) to "The suggestions were balanced and reasonable:"*



## Participants' assessment of the models' influence

% (Responses) to "The assistant inspired or changed my argument:"



# Deep Dive into Social Implication

- ✓ How does LLMs influence creativity
  - ✓ How to define creativity
  - ✓ Writing influence opinions

- How does LLMs affect productivity
  - Impact on labor market
  - AI for research idea generation

# Transformation to Workplace

## GPTs are GPTs: An Early Look at the Labor Market Impact Potential of Large Language Models

Tyna Eloundou<sup>1</sup>, Sam Manning<sup>1,2</sup>, Pamela Mishkin\*<sup>1</sup>, and Daniel Rock<sup>3</sup>

<sup>1</sup>OpenAI

<sup>2</sup>OpenResearch

<sup>3</sup>University of Pennsylvania

- ~ 80% of the U.S. workforce could have at least 10% of their work tasks affected by LLMs
- Most affected tasks: writing and programming.
- Higher-income jobs (e.g., translators, tax consultants, and web designers) potentially face greater exposure

# Transformation to Workplace

<b>Task ID</b>	<b>Occupation Title</b>	<b>DWAs</b>	<b>Task Description</b>
14675	Computer Systems Engineers/Architects	Monitor computer system performance to ensure proper operation.	Monitor system operation to detect potential problems.
18310	Acute Care Nurses	Operate diagnostic or therapeutic medical instruments or equipment. Prepare medical supplies or equipment for use.	Set up, operate, or monitor invasive equipment and devices, such as colostomy or tracheotomy equipment, mechanical ventilators, catheters, gastrointestinal tubes, and central lines.
4668.0	Gambling Cage Workers	Execute sales or other financial transactions.	Cash checks and process credit card advances for patrons.
15709	Online Merchants	Execute sales or other financial transactions.	Deliver e-mail confirmation of completed transactions and shipment.
6529	Kindergarten Teachers, Except Special Education	–	Involve parent volunteers and older students in children's activities to facilitate involvement in focused, complex play.
6568	Elementary School Teachers, Except Special Education	–	Involve parent volunteers and older students in children's activities to facilitate involvement in focused, complex play.

Sample of occupations, tasks, and Detailed Work Activities from the O\*NET database

# LLM Exposure Rubric

## **No exposure (E0)**

- Using LLMs results in no or minimal reduction in time

## **Direct exposure (E1)**

- Using LLMs decreases the time by at least 50%

## **LLM + exposed (E2)**

- Using LLMs does not help but additional tools are needed to achieve time reduction by at least 50%

# Exposure with Model and Human Comparison

<b>Comparison</b>	$\gamma$	<b>Weighting</b>	<b>Agreement</b>	<b>Pearson's</b>
GPT-4, Rubric 1; Human	$\alpha$	E1	80.8%	0.223
	$\beta$	E1 + .5*E2	65.6%	0.591
	$\zeta$	E1 + E2	82.1%	0.654
GPT-4, Rubric 2; Human	$\alpha$	E1	81.8%	0.221
	$\beta$	E1 + .5*E2	65.6%	0.538
	$\zeta$	E1 + E2	79.5%	0.589

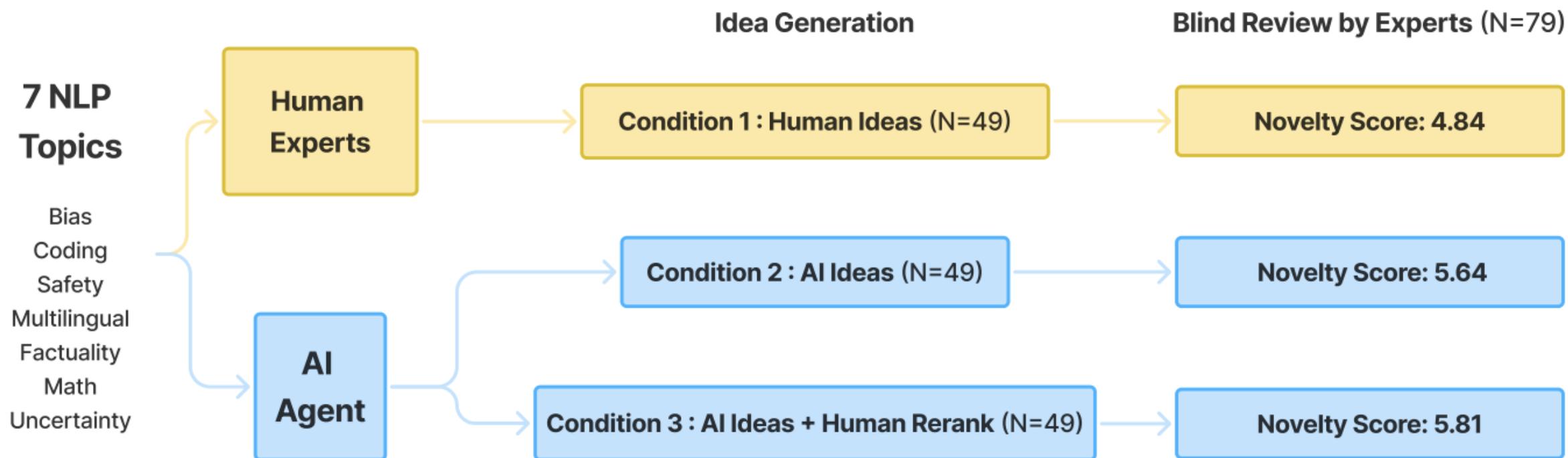
# Occupations with the highest exposure

<b>Group</b>	<b>Occupations with highest exposure</b>	<b>% Exposure</b>
<b>Human <math>\alpha</math></b>	Interpreters and Translators	76.5
	Survey Researchers	75.0
	Poets, Lyricists and Creative Writers	68.8
	Animal Scientists	66.7
	Public Relations Specialists	66.7
<b>Human <math>\beta</math></b>	Survey Researchers	84.4
	Writers and Authors	82.5
	Interpreters and Translators	82.4
	Public Relations Specialists	80.6
	Animal Scientists	77.8
<b>Human <math>\zeta</math></b>	Mathematicians	100.0
	Tax Preparers	100.0
	Financial Quantitative Analysts	100.0
	Writers and Authors	100.0
	Web and Digital Interface Designers	100.0
<i>Humans labeled 15 occupations as "fully exposed."</i>		

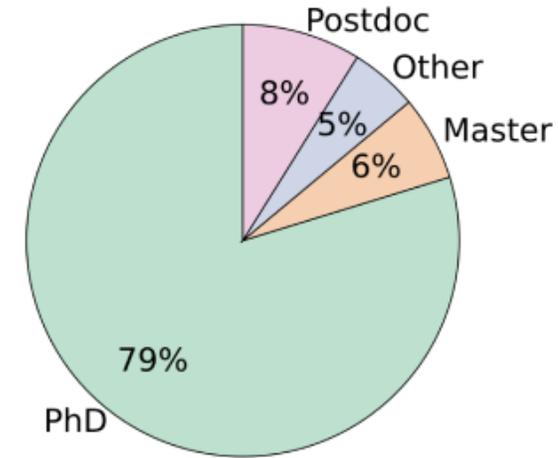
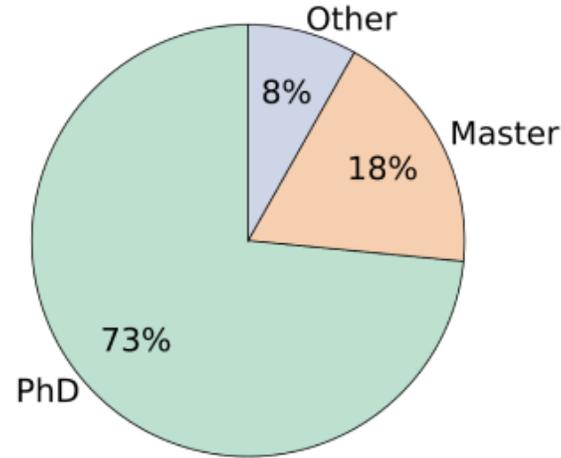
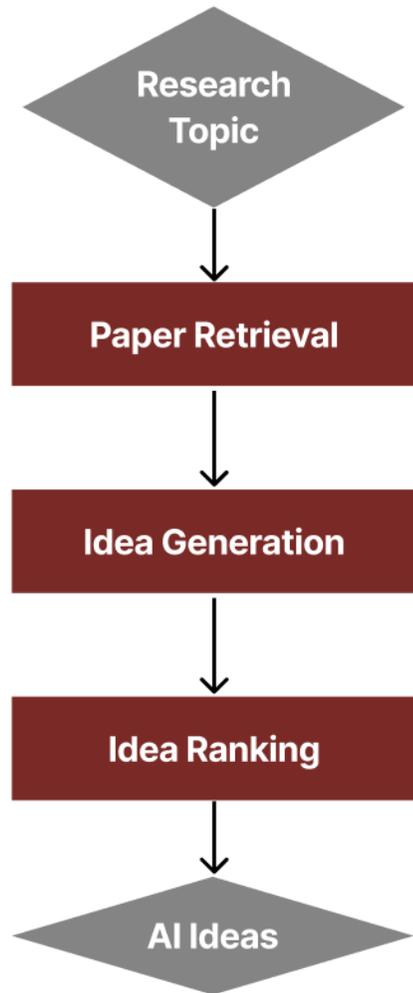
# Occupations with the highest exposure

<b>Group</b>	<b>Occupations with highest exposure</b>	<b>% Exposure</b>
<b>Model <math>\alpha</math></b>	Mathematicians	100.0
	Correspondence Clerks	95.2
	Blockchain Engineers	94.1
	Court Reporters and Simultaneous Captioners	92.9
	Proofreaders and Copy Markers	90.9
<b>Model <math>\beta</math></b>	Mathematicians	100.0
	Blockchain Engineers	97.1
	Court Reporters and Simultaneous Captioners	96.4
	Proofreaders and Copy Markers	95.5
	Correspondence Clerks	95.2
<b>Model <math>\zeta</math></b>	Accountants and Auditors	100.0
	News Analysts, Reporters, and Journalists	100.0
	Legal Secretaries and Administrative Assistants	100.0
	Clinical Data Managers	100.0
	Climate Change Policy Analysts	100.0
<i>The model labeled 86 occupations as "fully exposed."</i>		

# Can LLMs Generate Novel Research Ideas?

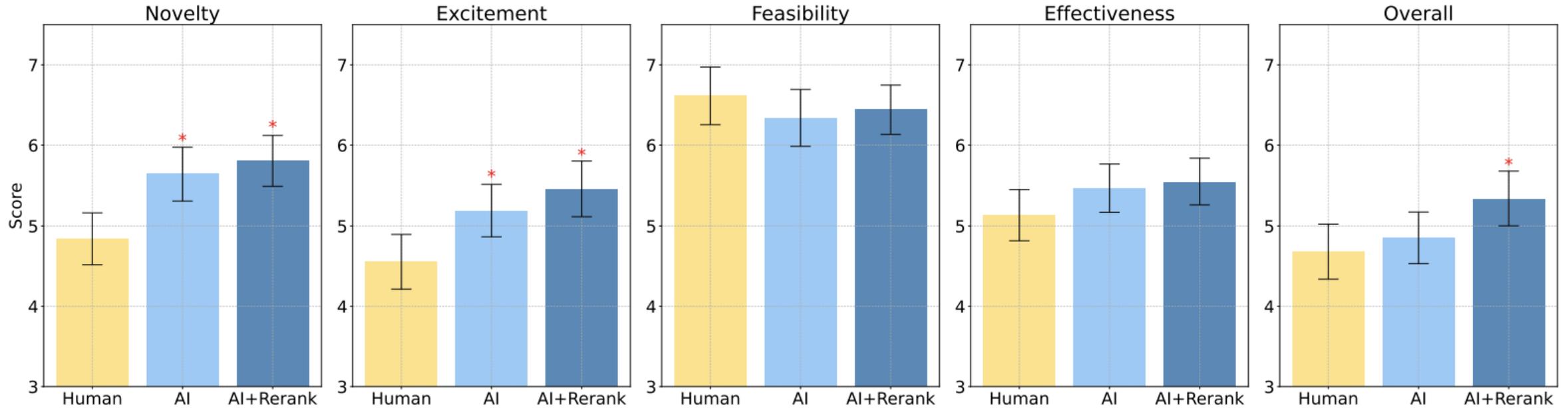


# Evaluating the Generated Research Ideas



Metric	Idea Writing Participants (N=49)					Idea Reviewing Participants (N=79)				
	Mean	Median	Min	Max	SD	Mean	Median	Min	Max	SD
papers	12	10	2	52	9	15	13	2	52	10
citations	477	125	2	4553	861	635	327	0	7276	989
h-index	5	4	1	21	4	7	7	0	21	4
i10-index	5	4	0	32	6	7	5	0	32	6

# Comparison of the three experiment conditions



We find AI ideas are judged as significantly more novel than human ideas ( $p < 0.05$ )

# Social Implication of LLMs

- ✓ Personalized education
- ✓ Companion and support
- ✓ LLMs for science discovery
- ✓ Transformation to workforce
- ✓ **Malicious uses**
- ✓ **Many other unknowns**