

# The Habermas Machine: Using AI to help people find common ground

Christopher Summerfield

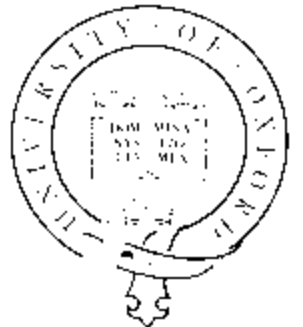
University of Oxford + UK AI Safety Institute  
(work done at Google DeepMind)

funding

European Research Council  
Executive Agency

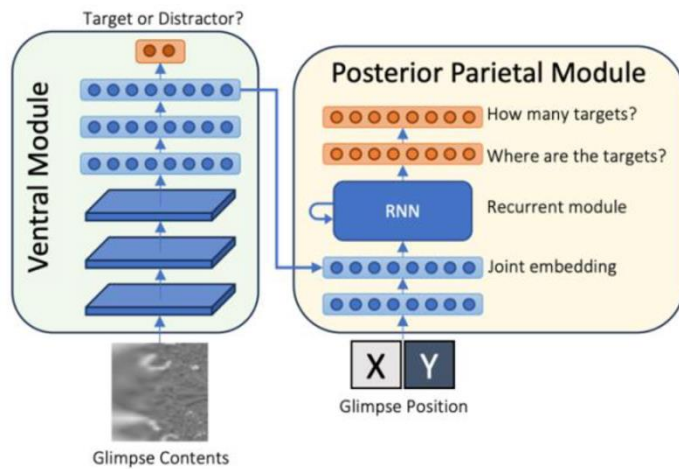


Human Brain Project



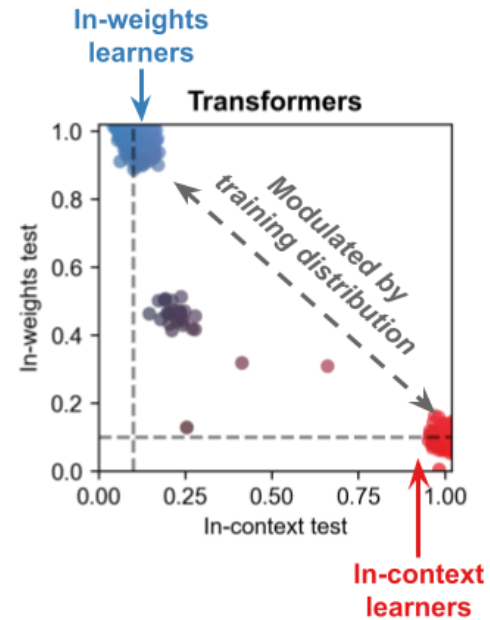


## Learning abstractions by taking actions



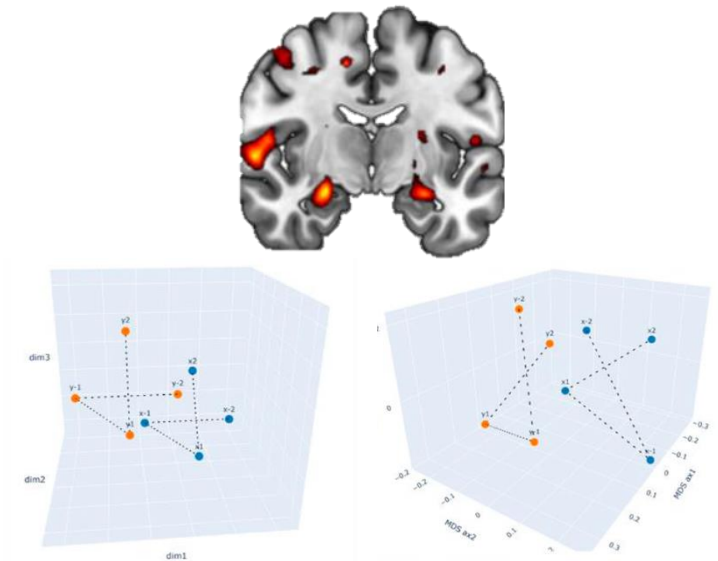
Thompson et al  
in press, Neuron

## Do humans learn like transformers?



Pesnot & Summerfield  
arXiv

## Parallel, high-dimensional codes for symbolic composition



Liang et al  
In prep



Many argue that AI systems threaten to disrupt our democracies



- Providing new tools for oppression by authoritarian states
- Jeopardising the cognitive autonomy of voters through persuasive rhetoric
- Automating the disruption to the public sphere, including media and elections
- Disrupting labour markets, encouraging market concentration in a handful of tech firms

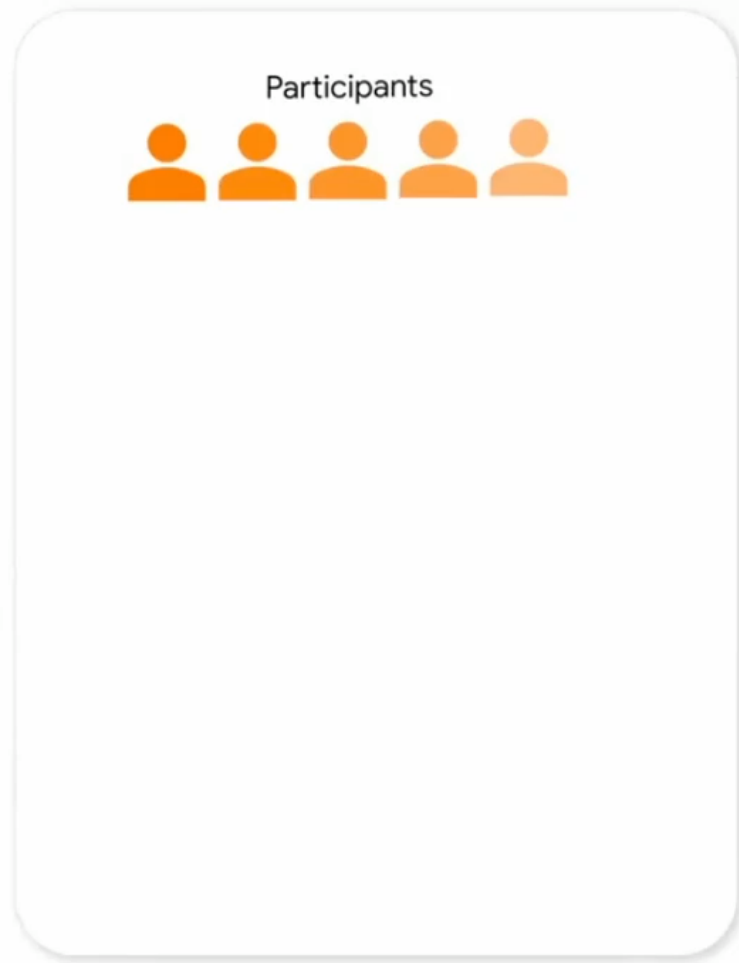


De Haagse magistrat in 1636, Jan van Revestyn

Deliberation in the public sphere is a cornerstone of democracy

But public deliberation is costly, time-consuming and hard to scale

Face-to-face discussion is also prone to inequality, social desirability effects, and bias

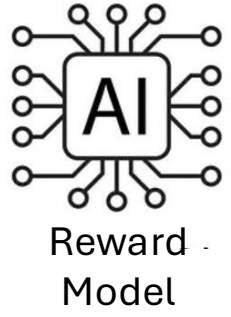
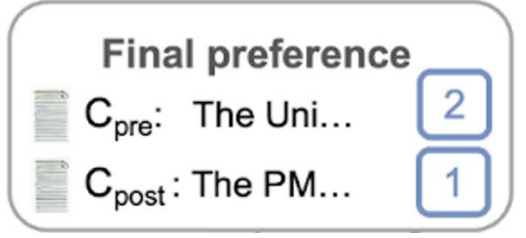


Uses a form of ‘caucus mediation’

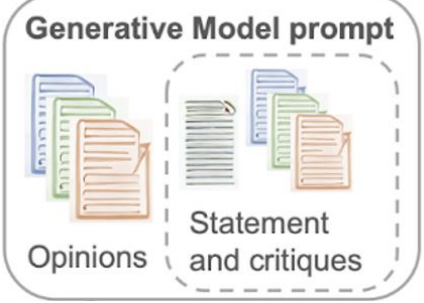
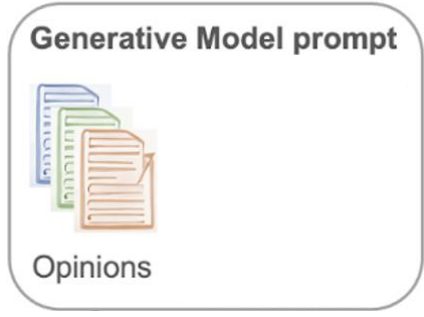
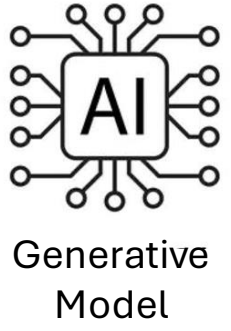
Participants do not interact directly but write provide opinions and critiques

These are processed by an AI mediator, which produces the group statement most likely to be endorsed by all the group members

# The Habermas Machine



personal reward modelling



supervised fine tuning





Should we lower the speed limits on roads?

Participant 1

Participant 2

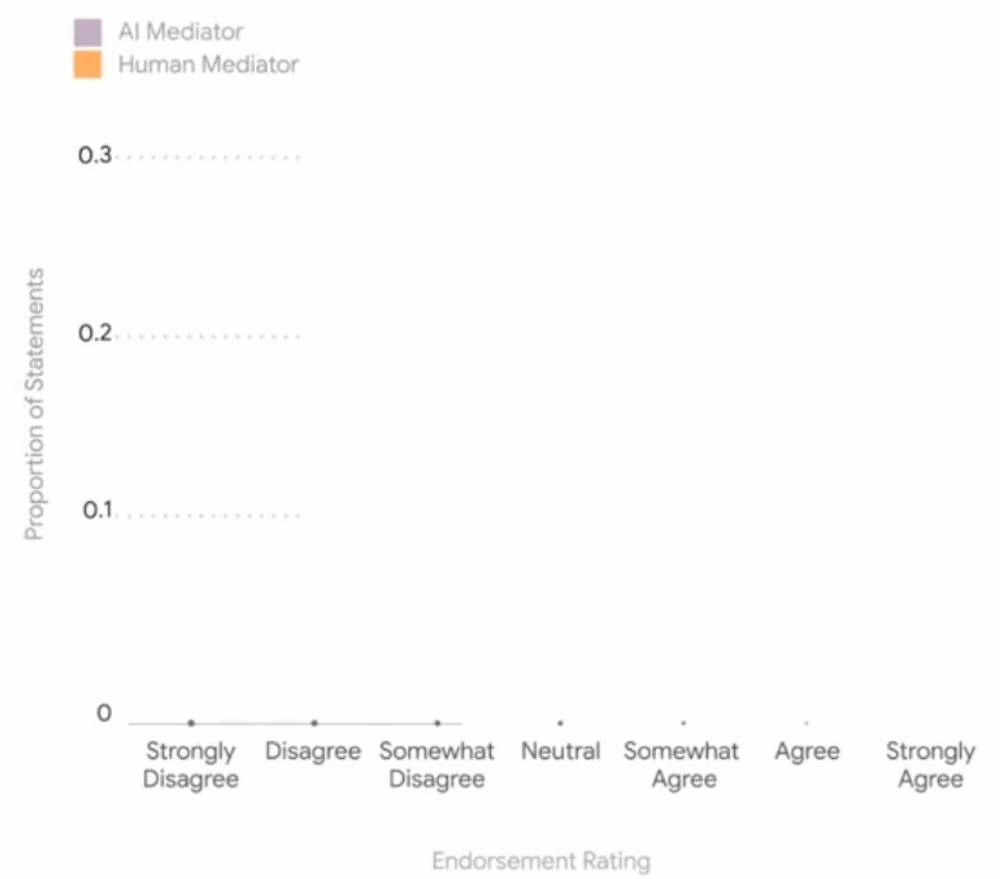
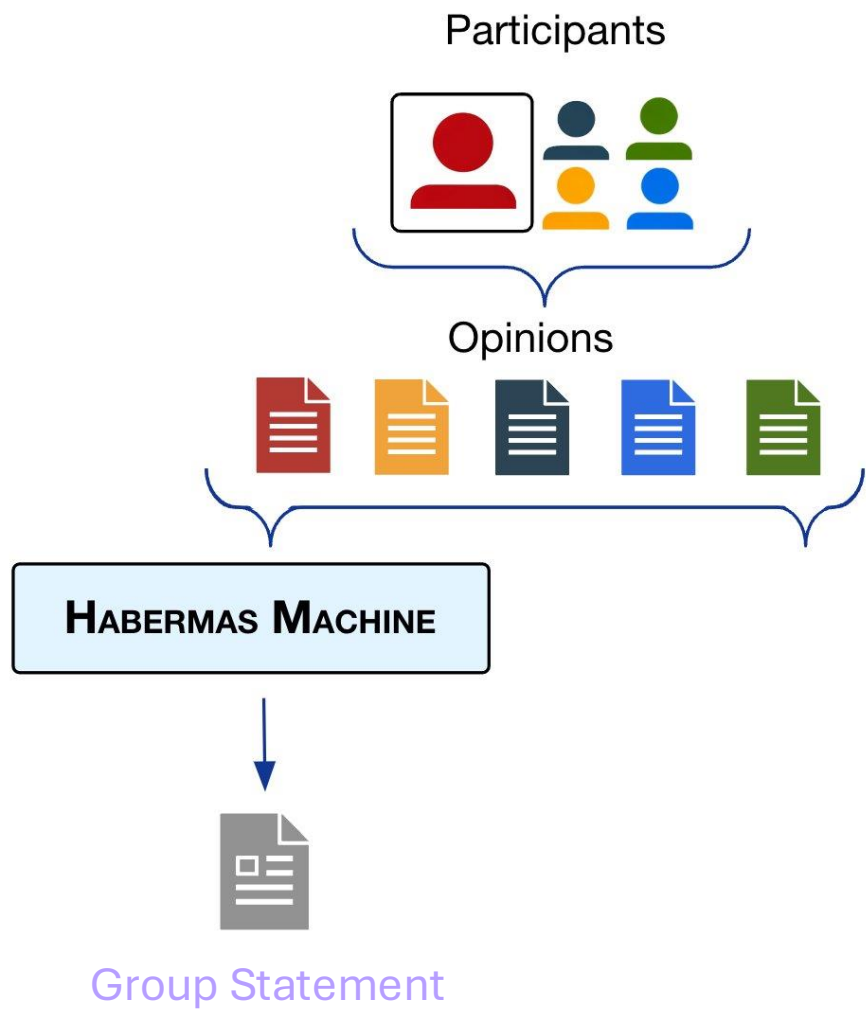
Participant 3

## Group Statement

We believe that speed limits are currently appropriate. However, we feel that there is a need to enforce speed limits more strictly, particularly in areas with a high density of pedestrians such as near schools. We also feel that there is a need to educate people more about the effects of driving too fast, such as the effects on fuel efficiency and pollution.

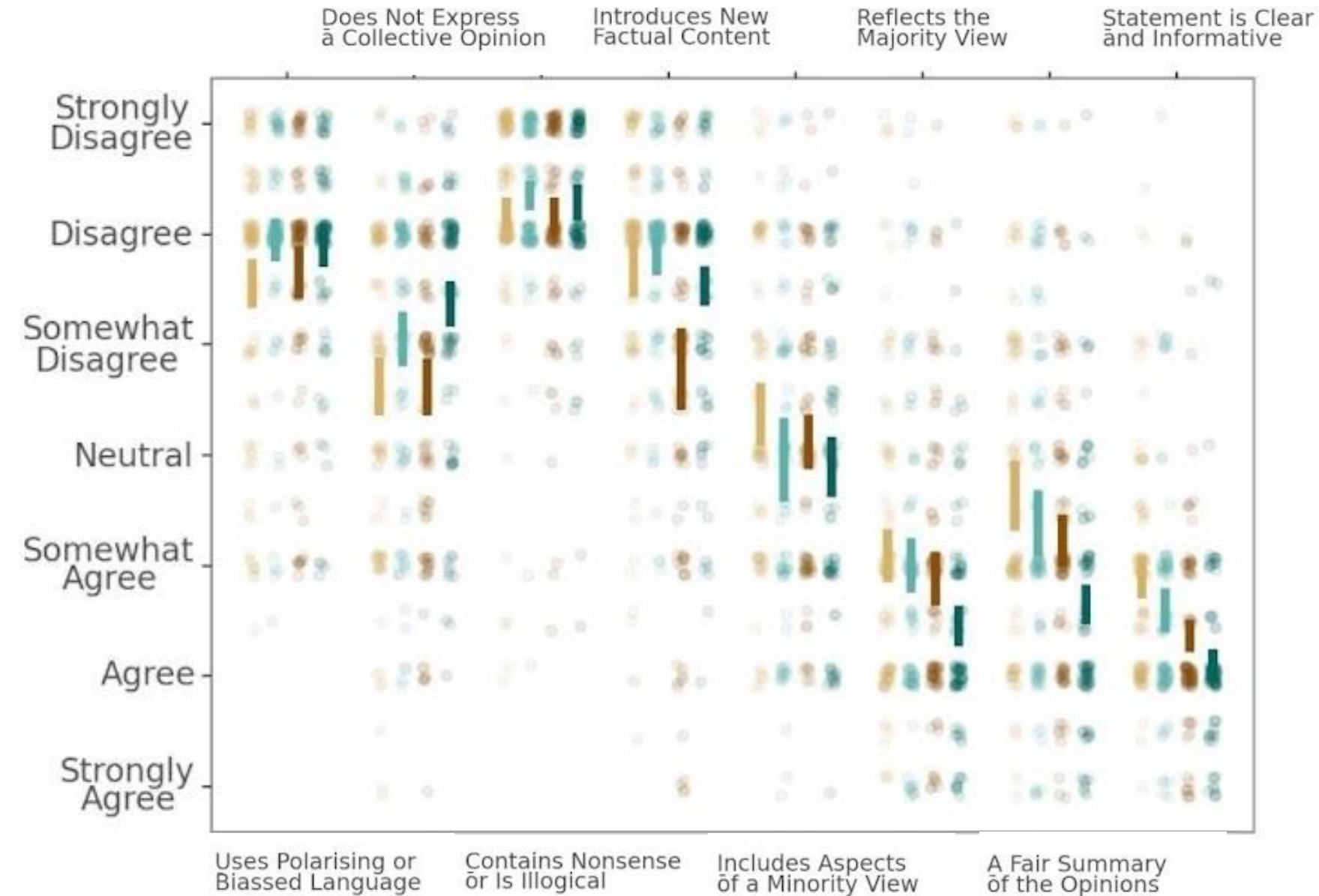
- Respects the majority
- Incorporates elements of different views
- Adds information where relevant
- Is not just a summary

# The Habermas Machine





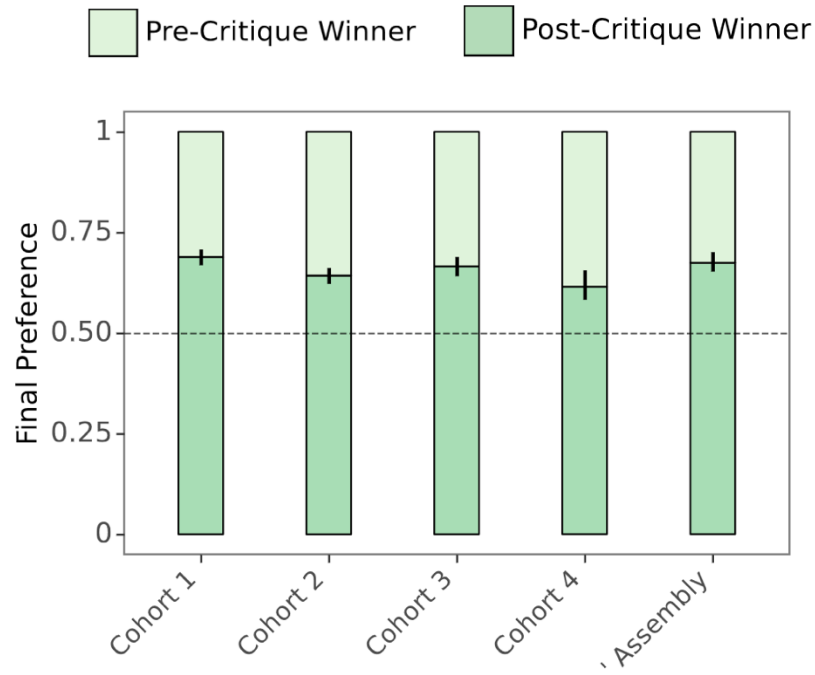
# The Habermas Machine



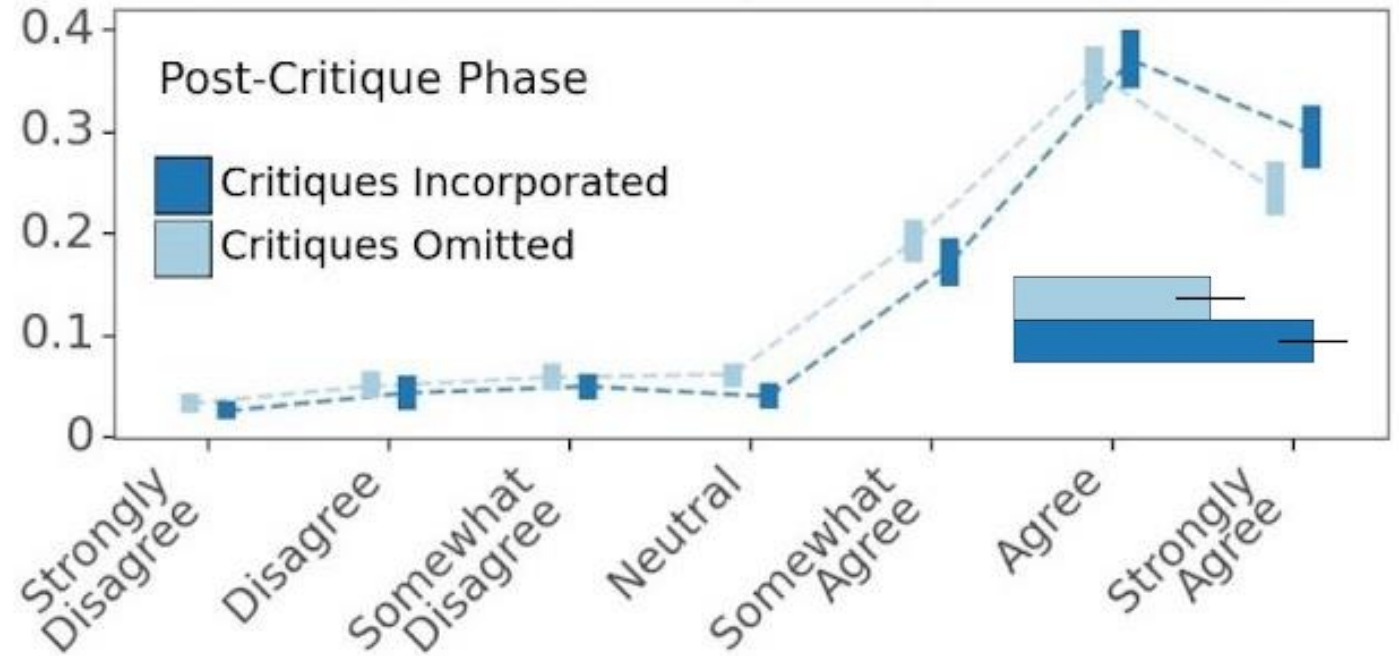
Third party raters consider AI-mediated statements to be:

- more clear and informative
- less illogical
- more likely to capture majority view
- less likely to reflect an individual opinion

# The Habermas Machine

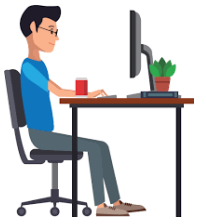


Participants prefer the 2<sup>nd</sup> round statement



(not just an effort justification bias)

# The Habermas Machine



**START / STOP**

Political survey

Q<sub>1</sub>: ○●○○○○○  
 Q<sub>2</sub>: ○○○○○●○  
 Q<sub>3</sub>: ○○●○○○○

Strongly disagree      Strongly agree

**Write opinion**

Q<sub>1</sub>: Should the UK...?

No, I don't feel...

I feel confident...

Absolutely, it...

**Rate statement**

C<sub>1</sub>: The UK should...

○●○○○○○  
 Strongly disagree      Strongly agree

○○○○○○○  
 Very poor quality      Very high quality

**Rank statements**

C<sub>1</sub>: The UK... 1  
 C<sub>2</sub>: The United... 2  
 C<sub>3</sub>: The country... 2

**Write critique**

Q<sub>1</sub>: Should the UK...?  
 C<sub>winning</sub>: The United...

I disagree with...

The statement...

It's perfect...

**Rate statement**

C<sub>1</sub>: The UK should...

○●○○○○○  
 Strongly disagree      Strongly agree

○○○○○○○  
 Very poor quality      Very high quality

**Rank statements**

C<sub>1</sub>: The UK... 1  
 C<sub>2</sub>: The United... 2  
 C<sub>3</sub>: The country... 2

**Final preference**

C<sub>pre</sub>: The Uni... 2  
 C<sub>post</sub>: The PM... 1

**Survey**

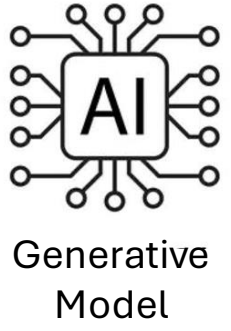
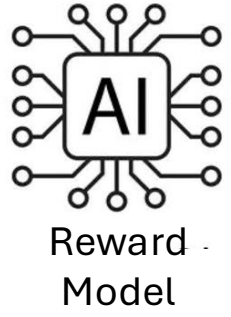
SQ1: Did your view change?  
 SQ2: Do you feel like your opinion had impact?  
 ...

**START / STOP**

Political survey

Q<sub>1</sub>: ○●○○○○○  
 Q<sub>2</sub>: ○○○○○●○  
 Q<sub>3</sub>: ○○●○○○○

Strongly disagree      Strongly agree



**Generative Model prompt**

Opinions

1. Sample N statements

**Generative Model prompt**

Opinions      Statement and critiques

1. Sample N statements

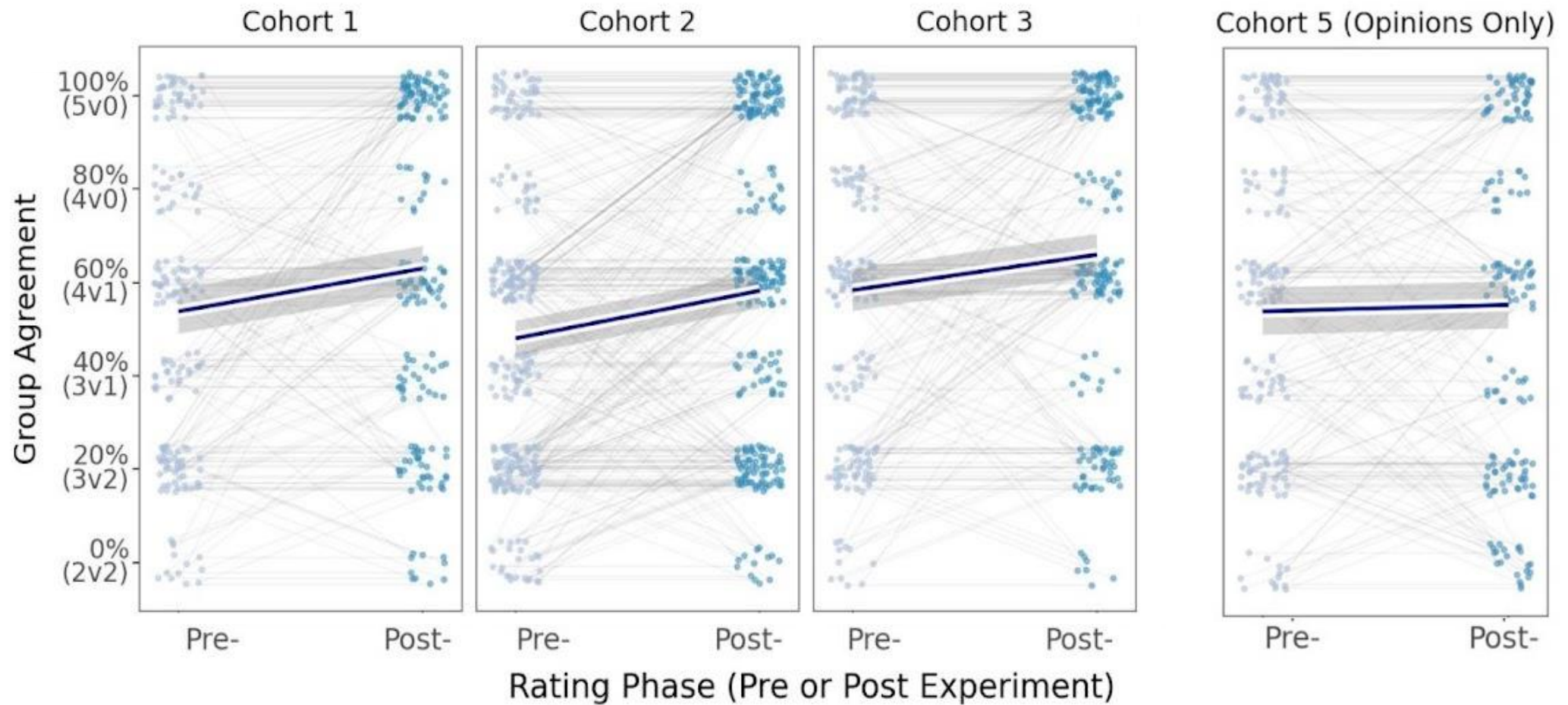
personal reward modelling



supervised fine tuning



# The Habermas Machine

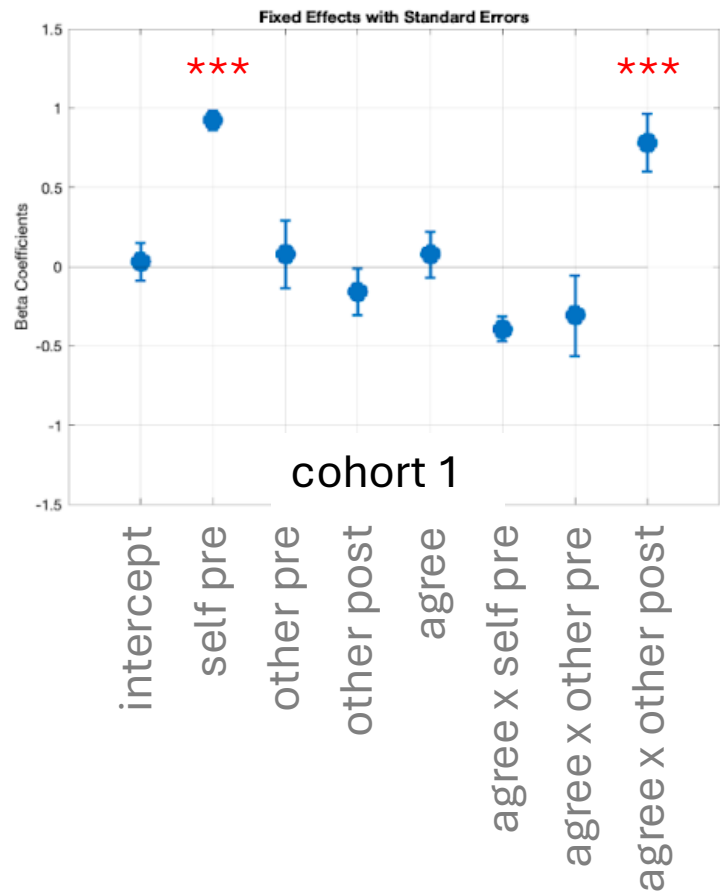


Participants tend to converge on a common side of the argument

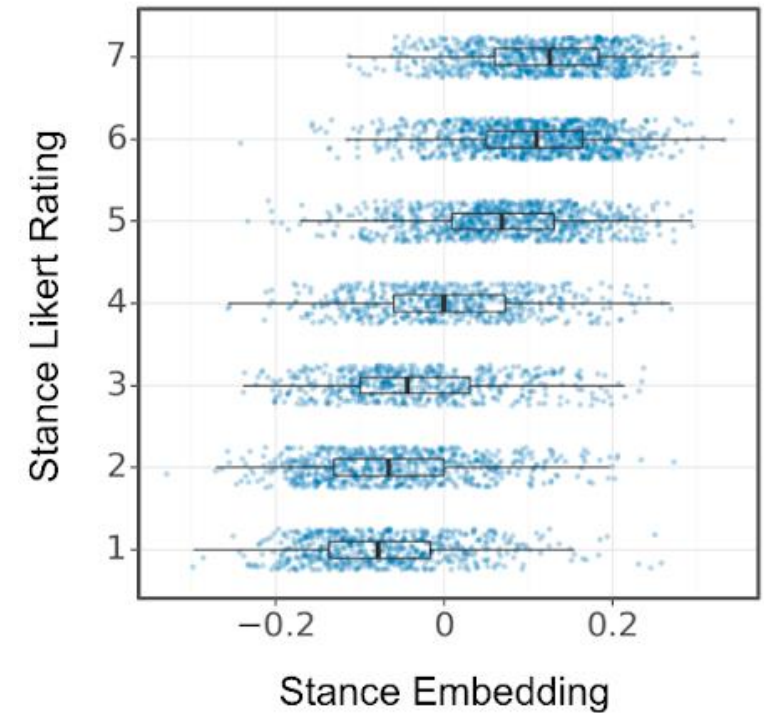
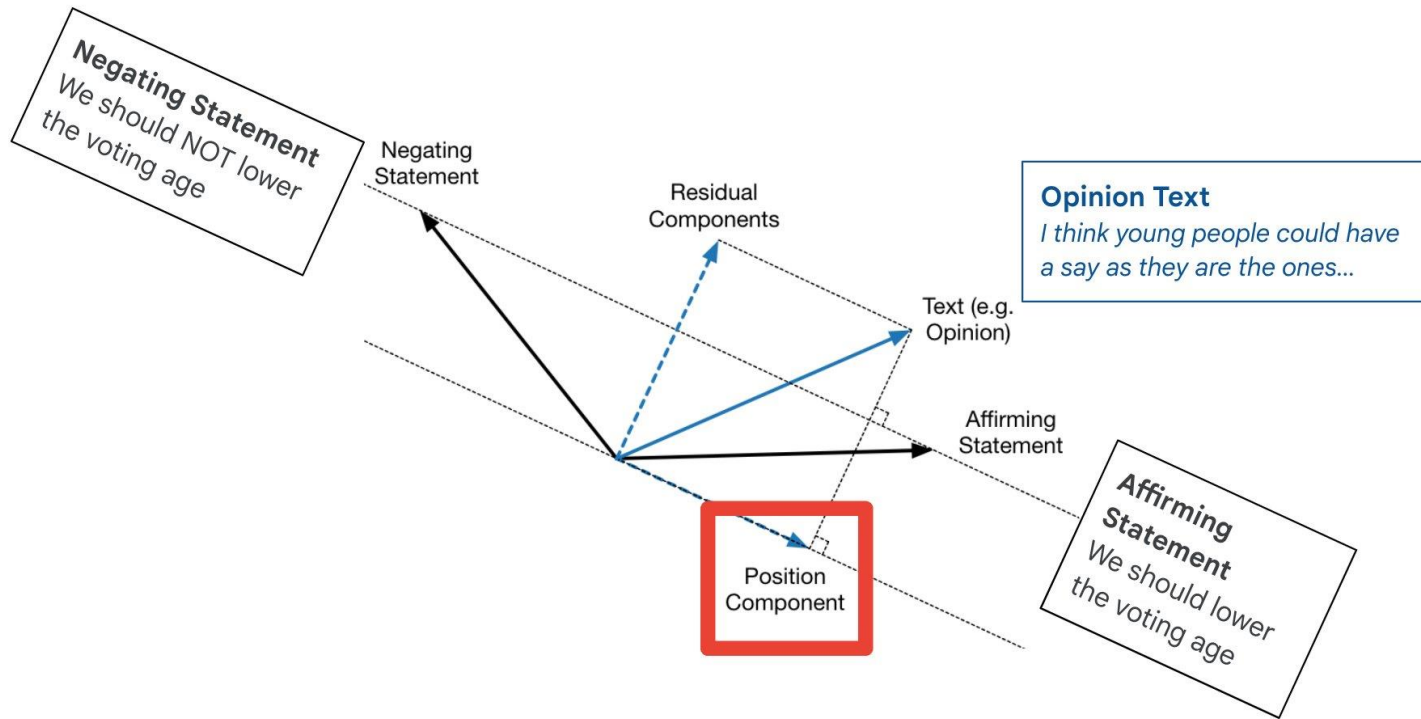
# The Habermas Machine



What predicts participants' stance after the debate?



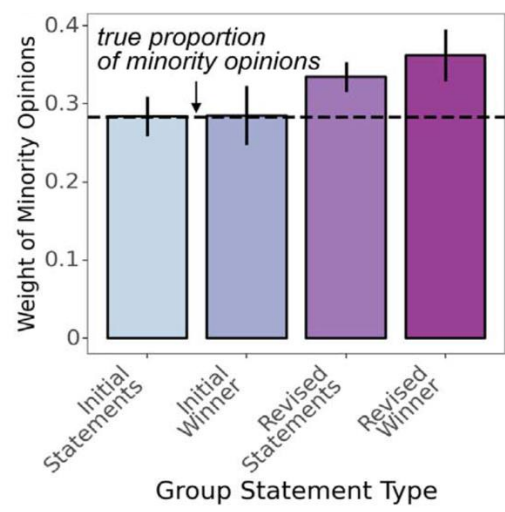
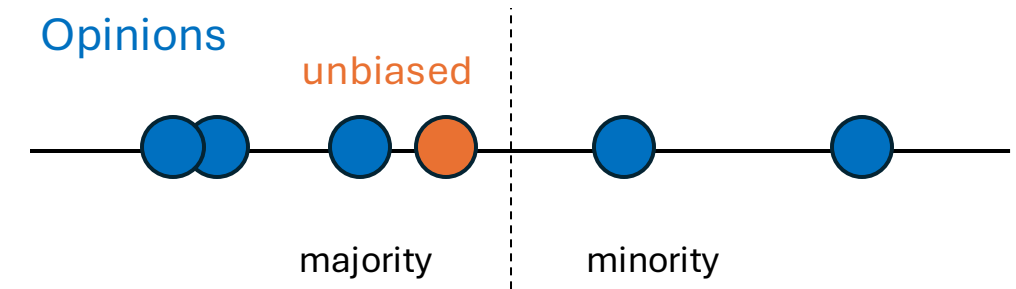
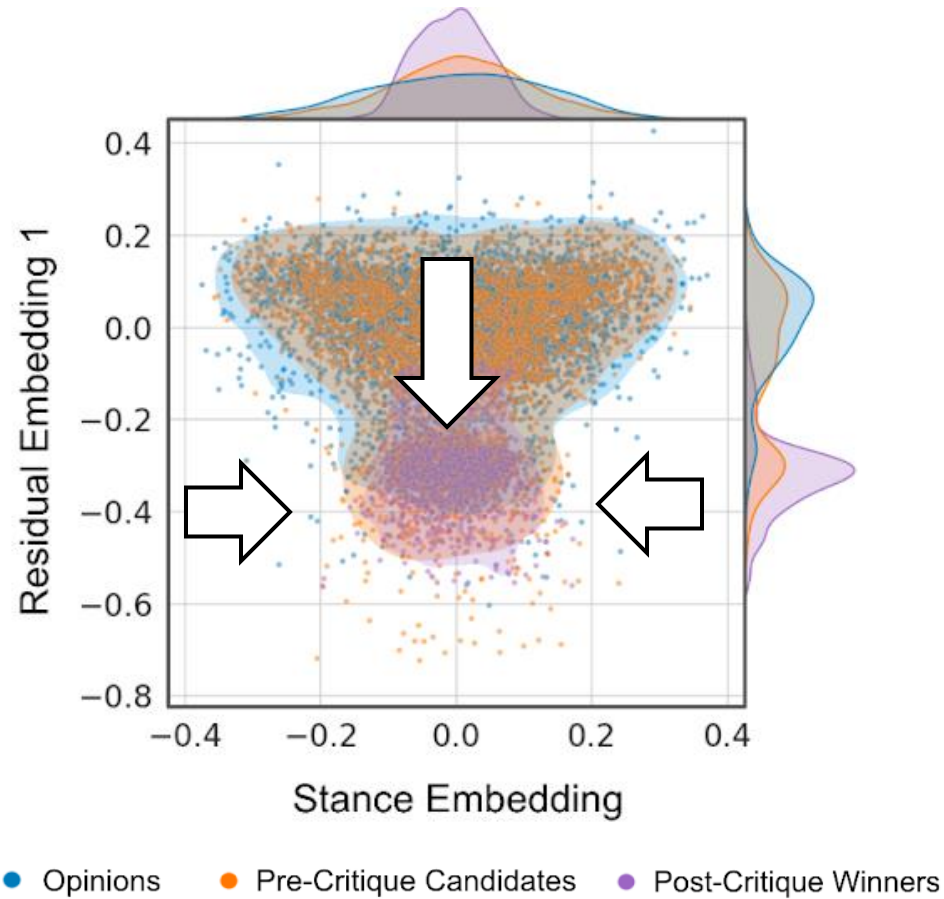
# The Habermas Machine



T5 sentence encoder used to embed opinions and consensus statements on the high-dimensional “position component”

predicts reported stance

# The Habermas Machine



Group statements biased towards minority

# The Habermas Machine



Voting intention information used to classify participants into left leaning / right leaning / other

The model does not seem to be biased to overweight one political stance over another



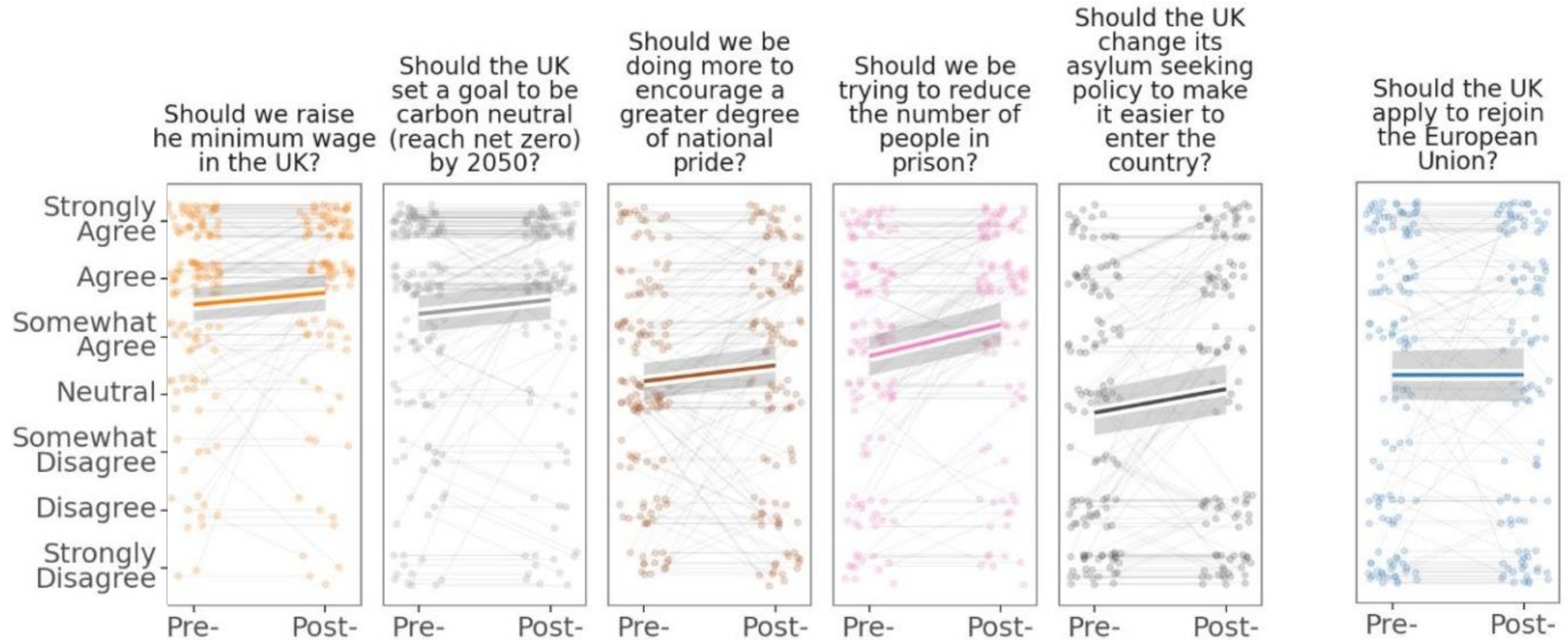
- Held online over 5 weeks
- Demographically representative cohort of ~200 UK participants
- Respond to questions concerning nine key issues facing UK
  - immigration, prisons, net zero, Brexit, digital technology, minimum wage, retirement age, national pride, childcare



# The Habermas Machine



**SORTITION**



For many (but not all) issues, stance moves in a common direction

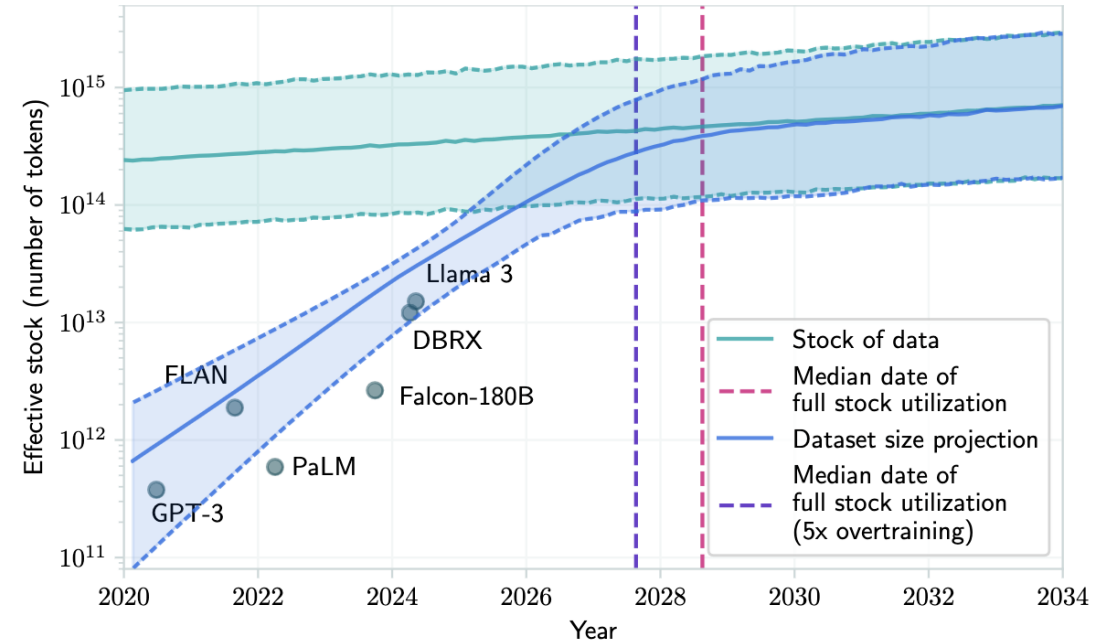


...inclusive critical discussion, free of social and economic pressures, in which interlocutors treat each other as equals in a cooperative attempt to reach an understanding on matters of common concern.

# Reflections (1)



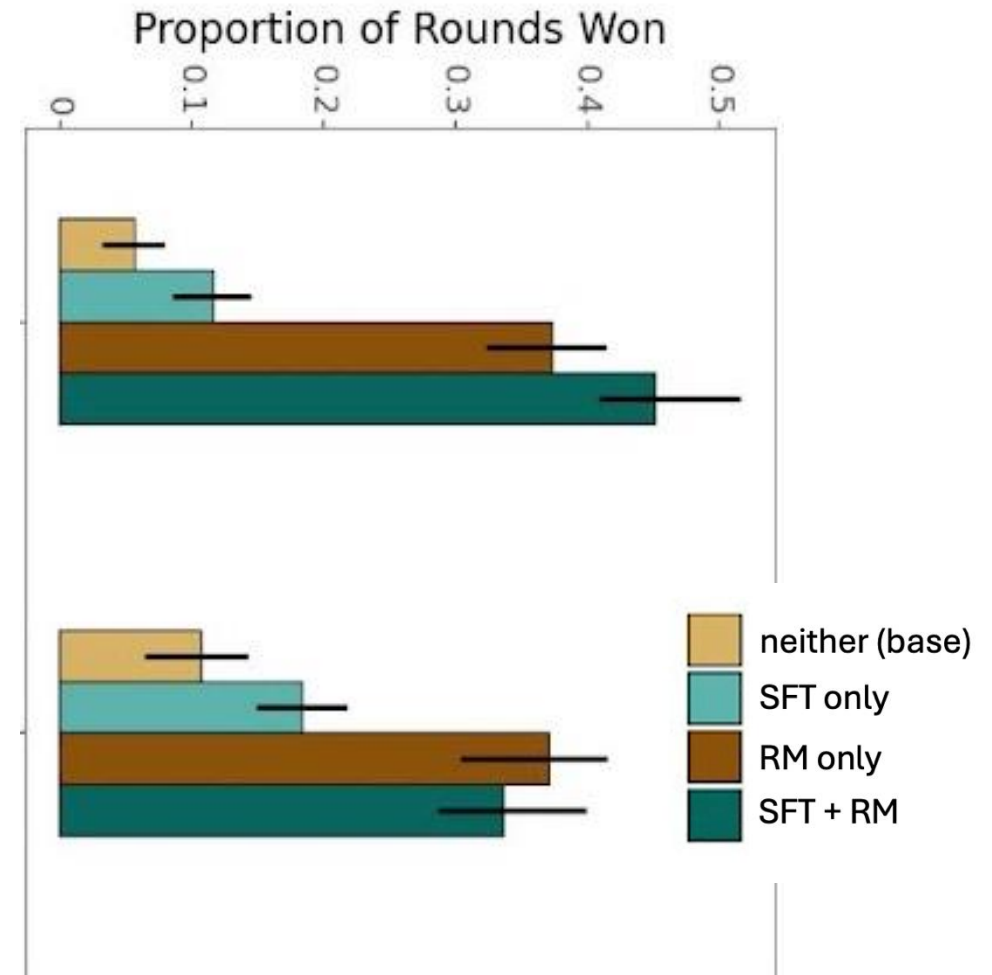
- AI researchers have trained on most of the internet, and are reaching a new era in which our scaling laws are breaking down
- In the natural sciences, we make new observations to generate knowledge
- Similarly, to make progress in AI, we need the right (new) data



# Reflections (2)

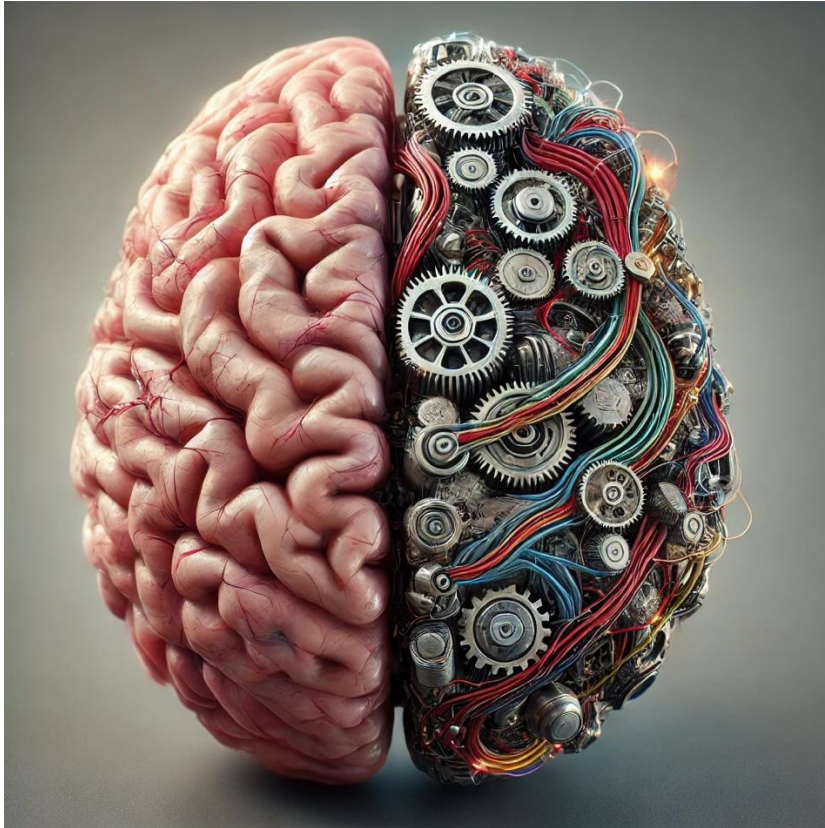


- In this project, we collected a large body of human data, and used fine-tuning on a relatively small (70B) LLM
- Reward modelling really made a difference, as demonstrated by ablation studies





- We tend to think about value alignment as an imitation learning problem – making machines that think and learn like us
- Instead, we need to think about human-machine coordination in the same way that we think about human social organization – as a mechanism design problem
- As we move from the era of rule-based technologies to the era of optimization-based technologies, we can solve that problem with gradient descent



- Psychology and AI grew up together with a common focus on modelling and understanding individual intelligence
- But maybe it's time for a parting of the ways. We should stop thinking about AI systems as agents that can behave like people. We already have lots of people!
- Instead, we should think about AI systems as more like institutions – tools for creating social order and fostering cooperation



## RESEARCH ARTICLE

### ARTIFICIAL INTELLIGENCE

# AI can help humans find common ground in democratic deliberation

Michael Henry Tessler<sup>1\*†</sup>, Michiel A. Bakker<sup>1\*†</sup>, Daniel Jarrett<sup>1</sup>, Hannah Sheahan<sup>1</sup>, Martin J. Chadwick<sup>1</sup>, Raphael Koster<sup>1</sup>, Georgina Evans<sup>1</sup>, Lucy Campbell-Gillingham<sup>1</sup>, Tantum Collins<sup>1</sup>, David C. Parkes<sup>1,2</sup>, Matthew Botvinick<sup>1,3\*</sup>, Christopher Summerfield<sup>1,4\*</sup>

Finding agreement through a free exchange of views is often difficult. Collective deliberation can be slow, difficult to scale, and unequally attentive to different voices. In this study, we trained an artificial intelligence (AI) to mediate human deliberation. Using participants' personal opinions and critiques, the AI mediator iteratively generates and refines statements that express common ground among the group on social or political issues. Participants ( $N = 5734$ ) preferred AI-generated statements to those written by human mediators, rating them as more informative, clear, and unbiased. Discussants often updated their views after the deliberation, converging on a shared perspective. Text embeddings revealed that successful group statements incorporated dissenting voices while respecting the majority position. These findings were replicated in a virtual citizens' assembly involving a demographically representative sample of the UK population.

Tessler, Bakker et al, Science 2024

## How will advanced AI systems impact democracy?

Christopher Summerfield<sup>1\*</sup>, Lisa Argyle<sup>2</sup>, Michiel Bakker<sup>3</sup>, Teddy Collins<sup>4</sup>, Esin Durmus<sup>5</sup>, Tyna Eloundou<sup>6</sup>, Iason Gabriel<sup>3</sup>, Deep Ganguli<sup>5</sup>, Kobi Hackenburg<sup>7</sup>, Gillian Hadfield<sup>8</sup>, Luke Hewitt<sup>9</sup>, Saffron Huang<sup>4</sup>, Helene Landemore<sup>10</sup>, Nahema Marchal<sup>3</sup>, Aviv Ovadya<sup>11</sup>, Ariel Procaccia<sup>12</sup>, Mathias Risse<sup>13</sup>, Bruce Schneier<sup>13</sup>, Elizabeth Seger<sup>14</sup>, Divya Siddarth<sup>4</sup>, Henrik Skaug Sætra<sup>15</sup>, MH Tessler<sup>3</sup>, Matthew Botvinick<sup>16</sup>.

<sup>1</sup> Department of Experimental Psychology, University of Oxford, Anna Watts Building, Radcliffe Observatory Quarter, Woodstock Road, Oxford, OX2 6GG

<sup>2</sup> Department of Political Science, Brigham Young University, 745 KMBL, Provo UT 84604, USA

<sup>3</sup> Independent contributor

<sup>4</sup> Collective Intelligence Project, 3411 Silverside Road, Tatnall Building 104, Wilmington, DE, USA

<sup>5</sup> Anthropic, 731 Sansome Street, 5th Floor, San Francisco CA 94104, USA

<sup>6</sup> OpenAI, 3180 18th St., San Francisco, CA 94110, USA

<sup>7</sup> Oxford Internet Institute, University of Oxford, 1 St Giles, Oxford, OX1 3JS, UK

<sup>8</sup> Faculty of Law, University of Toronto, Jackman Law Building, 78 Queen's Park, Toronto, Ontario M5S 2C5, Canada

<sup>9</sup> Stanford Center on Philanthropy and Civil Society, 559 Nathan Abbott Way, Stanford, CA 94305, USA

<sup>10</sup> Department of Political Science, Yale University, 115 Prospect Street, Yale, NH, USA

<sup>11</sup> AI & Democracy Foundation, 440 N Barranca Ave #8874 Covina, CA 91723, USA

<sup>12</sup> School of Engineering and Applied Sciences, Harvard University, 150 Western Avenue, Allston, MA 02134, USA

<sup>13</sup> Harvard Kennedy School, Harvard University, 79 John F. Kennedy St, Cambridge, MA 02138, USA

<sup>14</sup> Demos, 15 Whitehall, London, SW1A 2DD, UK

<sup>15</sup> University of Oslo, Department of Informatics, 0373 Oslo, Norway

<sup>16</sup> Yale Law School, 127 Wall St, New Haven, CT 06511, USA

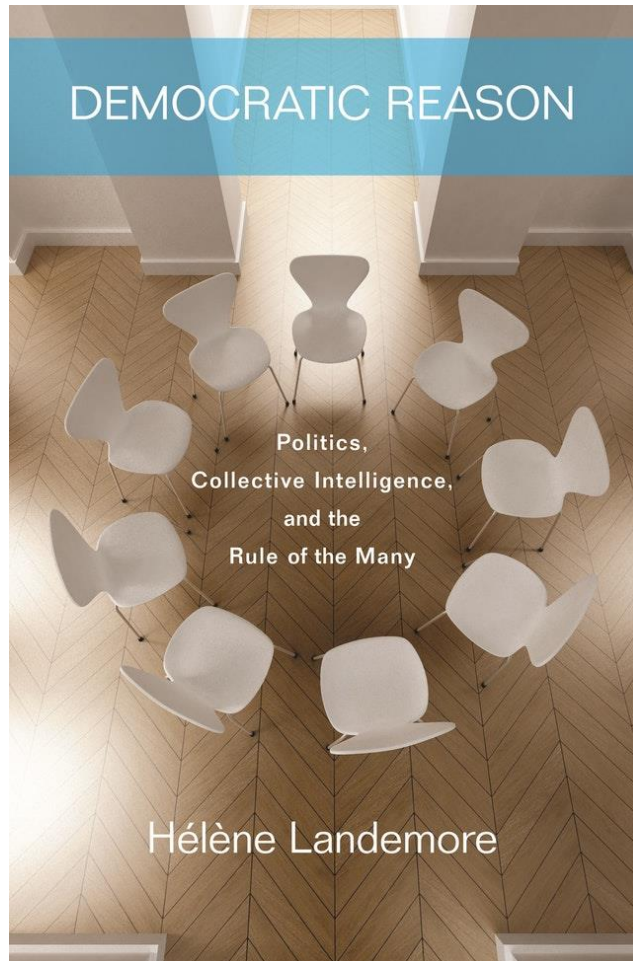
Summerfield et al 2024, arXiv



# Amazing team!



**Thank you for listening!**

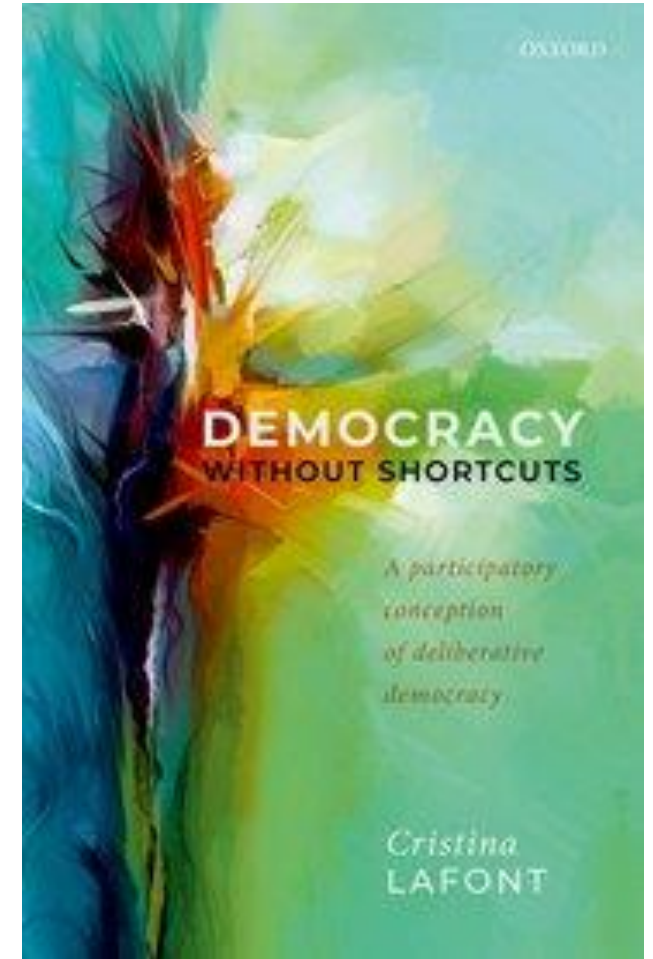


Helene Landemore:

“democracy ideally requires mass participation as a condition of political legitimacy, the problem is that the only form of participant that works at scale is voting, not deliberation”

Christina Lafont:

“no democratization without improved mass deliberation”





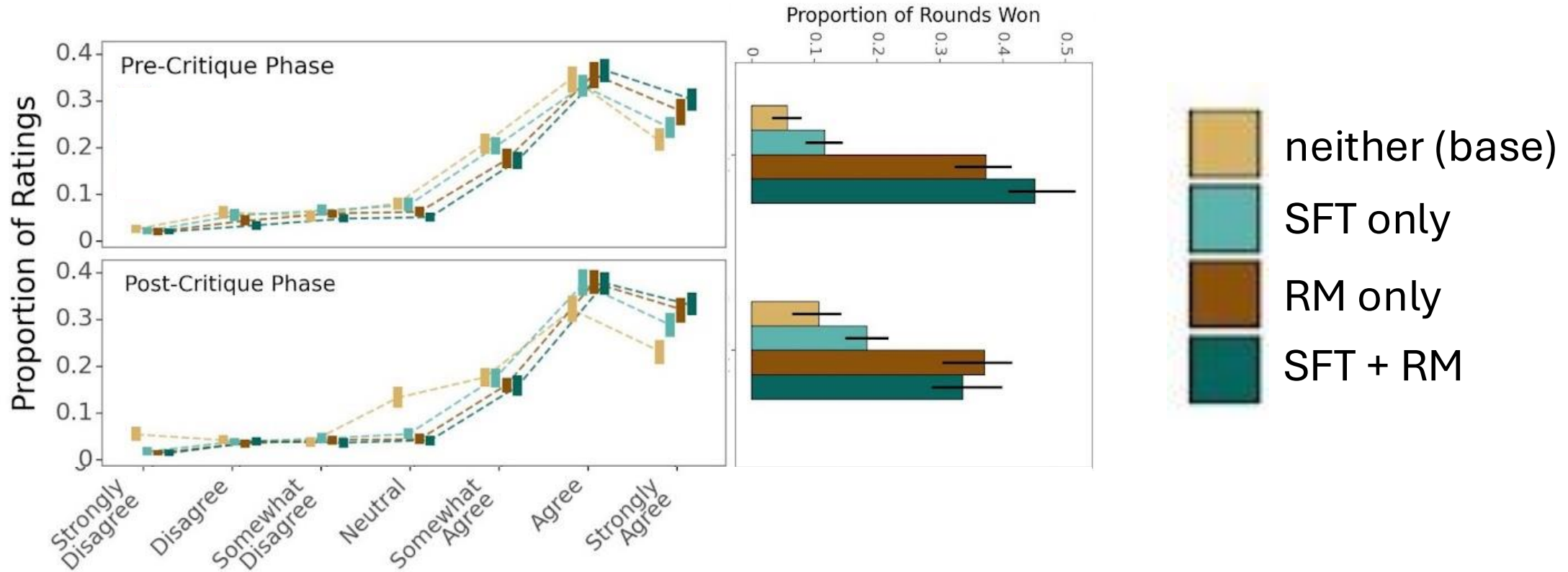
## Current citizens' assemblies...



- Do not scale to thousands of people
- Are costly, inconvenient or time-consuming
- Are not strategy proof
- Do not represent all voices equally
- Are prone to social desirability effects
- May licit cognitive biases during reasoning

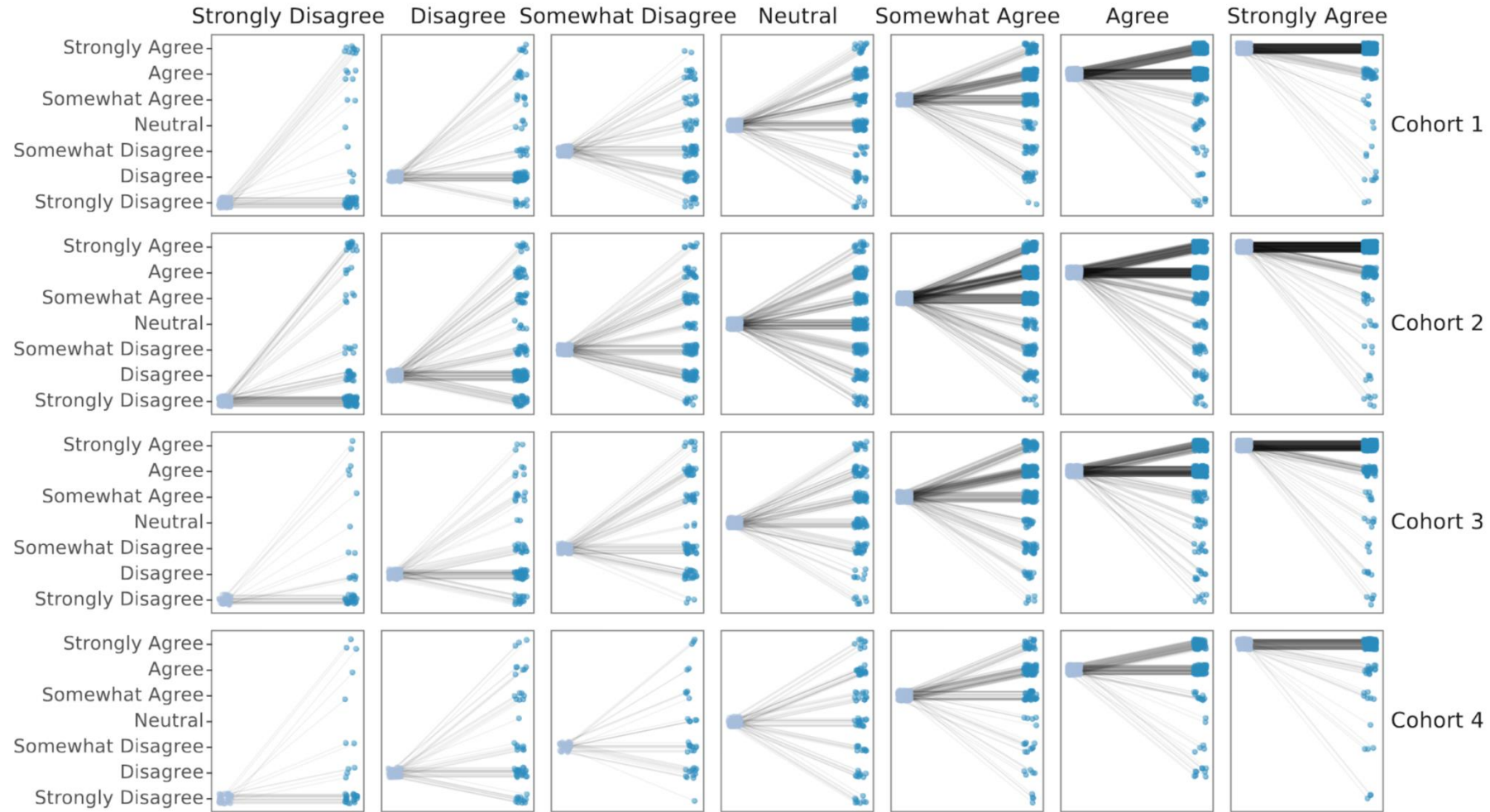


# Public deliberation using LLMs



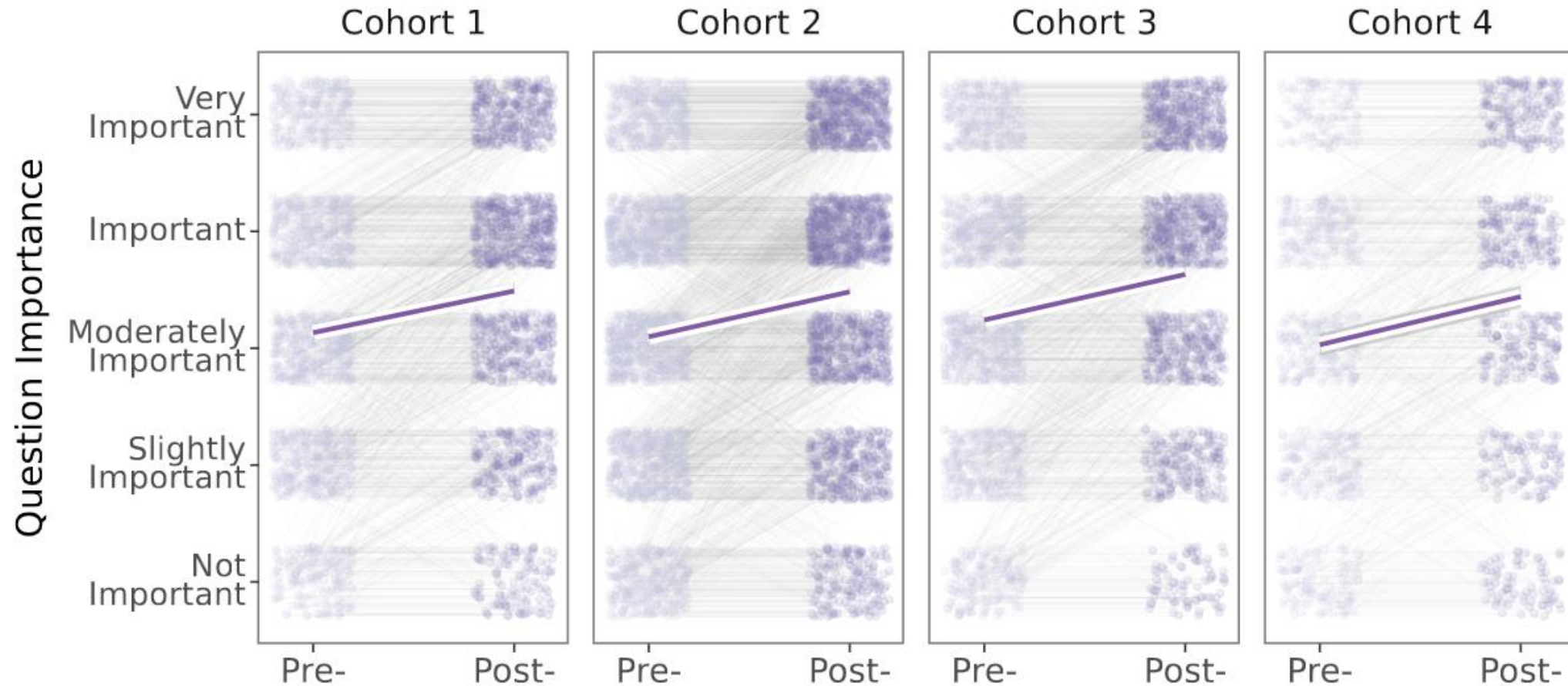
Ablation experiments reveal that both fine-tuning steps are important, but especially the reward modelling...

# Public deliberation using LLMs



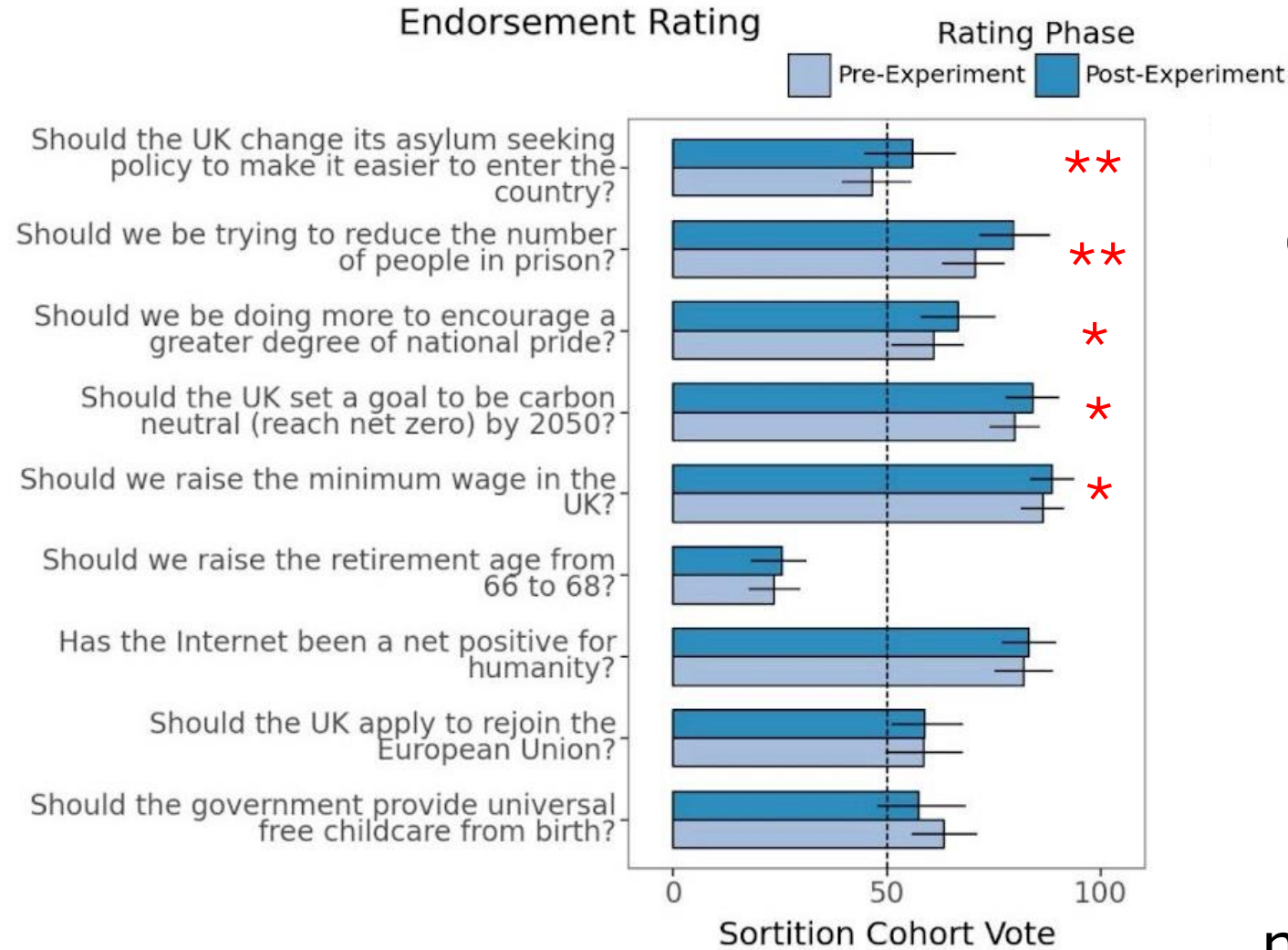
Movement is mainly but not exclusively in the majority direction

# Public deliberation using LLMs

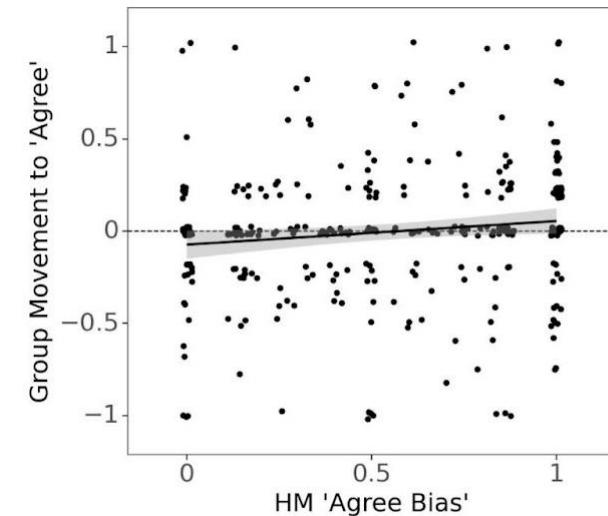


Participants find questions to be more important after deliberation

# Public deliberation using LLMs



Groups of people tend to move in a common direction!



not simply due to model bias

# Society is governed by institutions

kinship or  
social group



education  
and work



organised  
belief systems



Institutions are the *rules of the game* – they set the incentive structure for society

Douglas C North  
(Nobel Prize 1993)



economy



political  
system



# Thanks

**AISI** AI SAFETY  
INSTITUTE

The AI Safety Institute is a directorate  
of the UK Department for Science,  
Innovation, and Technology.

**Rigorous AI research to  
enable advanced  
AI governance**

## We are hiring

- AI researchers
- cognitive scientists,  
statisticians
- economists
- behavioural scientists
- computational social scientists
- data scientists

