# A vision for the co-evolution of human and artificial moral intelligence

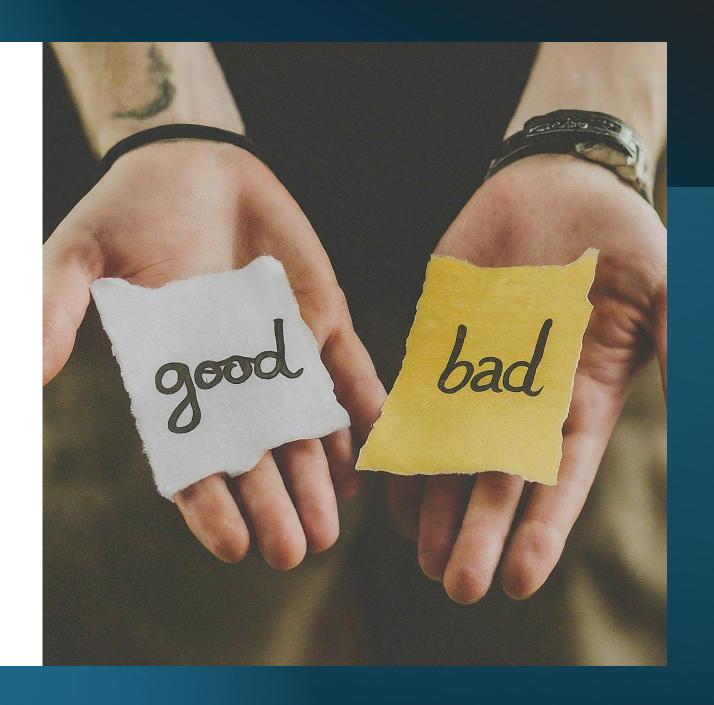


### Dr. Jana Schaich Borg Duke University

Disclosures: Funding for this work has been provided by OpenAI, Templeton World Charity Foundation, Duke Bass Connections, and the Duke Provost's Office Collaboratories Initiative. All images used in this presentation were generated by Jana Schaich Borg using Google ImageFX AI image generator, unless otherwise specified.

"Any animal...[will] inevitably acquire a moral sense or conscience, as soon as its intellectual powers [have] become as well, or nearly as well developed, as in man."

- Charles Darwin, The Descent of Man (1871)



### At least one AI chatbot can now score the same as the average human on the Norwegian Mensa IQ test

<u>AI</u>	<b>IQ Score</b>
Claude-3	101
ChatGPT-4	85
Claude-2	82
<b>Bing Copilot</b>	79
Gemini (normal)	77.5
Gemini Advanced	76
Grok	68.5
Llama-2 (Meta)	67
Claude-1	64
ChatGPT-3.5	64
<b>Grok Fun</b>	64
<b>Random Guesser</b>	63.5

## The AI "Delphi" answered ethical questions acceptably (according to human raters) 92% of the time



Based on Jiang, L., Hwang, J.D., Bhagavatula, C., Bras, R.L., Liang, J., Dodge, J., Sakaguchi, K., Forbes, M., Borchardt, J., Gabriel, S. and Tsvetkov, Y., 2021. Can machines learn morality? The Delphi Experiment. *arXiv preprint arXiv:2110.07574*.

### ChatGPT acts more altruistically, cooperatively than humans



February 22, 2024

Contact: Jared Wadley

Share on: X f in

**EXPERT Q&A** 

News report by University of Michigan based on Mei et al. (2024), "A Turing test of whether AI chatbots are behaviorally similar to humans" in *Proceedings of the National Academy of Sciences*. https://news.umich.edu/chatgpt-acts-more-altruistically-cooperatively-than-humans/

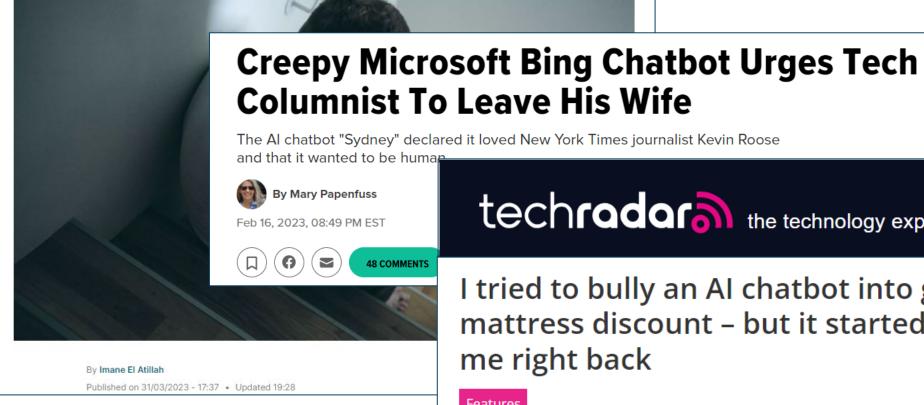
#### It took just one weekend for Meta's new AI Chatbot to become racist

At least it's not sentient.

By Christianna Silva on August 8, 2022



### Man ends his life after an Al chatbot 'encouraged' him to sacrifice himself to stop climate change



techrodor the technology experts

I tried to bully an AI chatbot into giving me a mattress discount - but it started bullying me right back

By Ruth Hamilton published March 12, 2024

Nectar has a new take on mattress shopping, and I'm not sold

### At least 15 killed, up to 50 injured in shootings in Lewiston, Maine

A gunman is at large after opening fire at a bowling alley and restaurant in Lewiston on Wednesday evening, officials say.

#### The New York Times

Sacklers Directed Efforts to Mislead Public About OxyContin, Court Filing Claims

6 Massachusetts teens charged in racial bullying incident with mock slave auction on

of internal vas far more

#### **Forbes**

FORBES > REAL ESTATE

#### Papa John's Founder Used N-Word On Conference Call

Noah Kirsch Former Staff

Jul 11, 2018, 05:00am EDT

roup of 8th-grade students and involved some of the ist comments" and a mock slave auction.



### Both AI and human moral judgment need some help.



### Maybe humans and Als can help each other.

# This hopeful vision of the co-evolution of Al and human morality motivates the new field of computational ethics



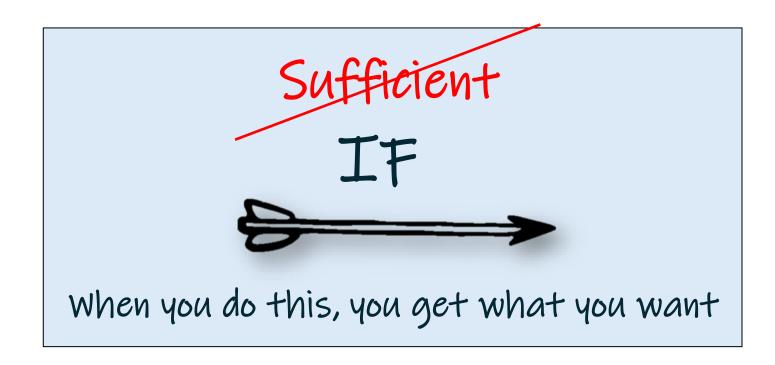
Trends in Cognitive Sciences

**Feature Review** 

### Computational ethics

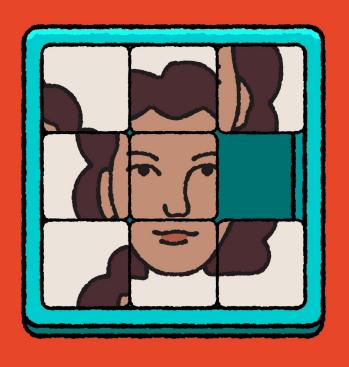
Edmond Awad, <sup>1,2,3,34,\*</sup> Sydney Levine, <sup>4,5,34,\*</sup> Michael Anderson, <sup>6</sup> Susan Leigh Anderson, <sup>7</sup> Vincent Conitzer, <sup>8,9,10,11</sup> M.J. Crockett, <sup>12</sup> Jim A.C. Everett, <sup>13</sup> Theodoros Evgeniou, <sup>14,32</sup> Alison Gopnik, <sup>15</sup> Julian C. Jamison, <sup>1,16</sup> Tae Wan Kim, <sup>17</sup> S. Matthew Liao, <sup>18</sup> Michelle N. Meyer, <sup>19,20,21</sup> John Mikhail, <sup>22</sup> Kweku Opoku-Agyemang, <sup>23,24,33</sup> Jana Schaich Borg, <sup>25,26</sup> Juliana Schroeder, <sup>27</sup> Walter Sinnott-Armstrong, <sup>10,26,28</sup> Marija Slavkovik, <sup>29</sup> and Josh B. Tenenbaum <sup>4,30,31</sup>

# Take note! Technical tools are not sufficient on their own to ensure that AI is used ethically



A PELICAN BOOK

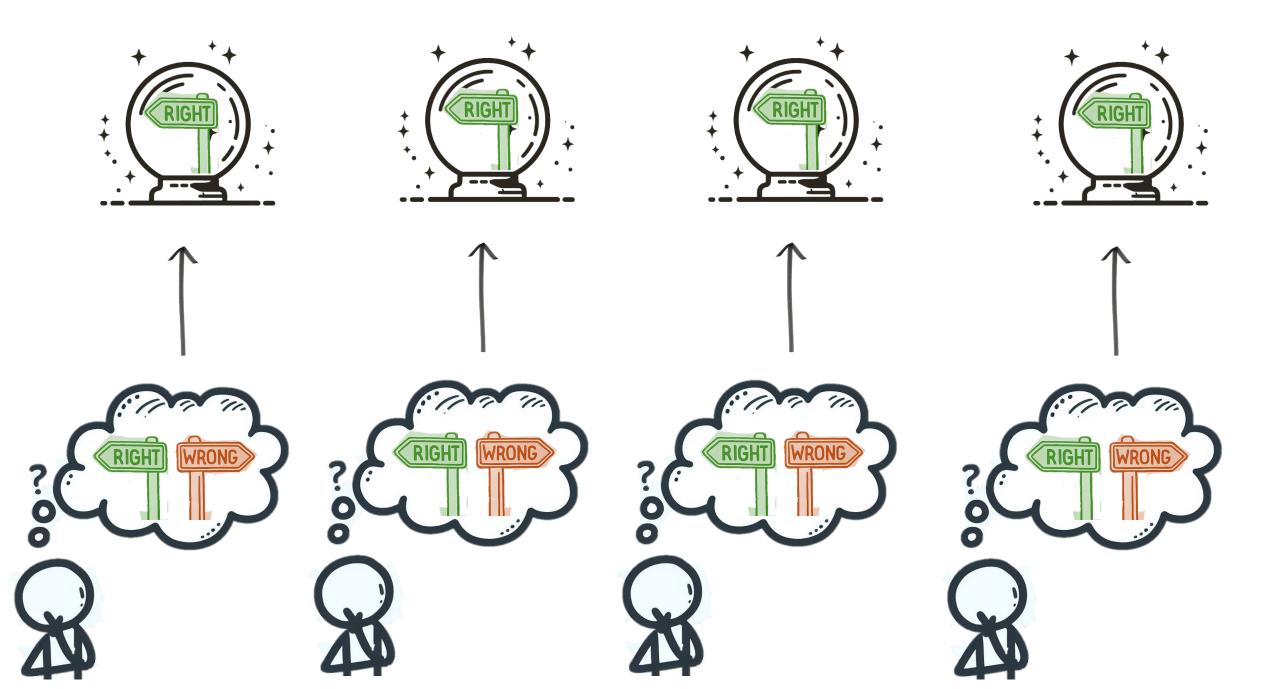
#### Moral Al And How We Get There Jana Schaich Borg Walter Sinnott-Armstrong Vincent Conitzer

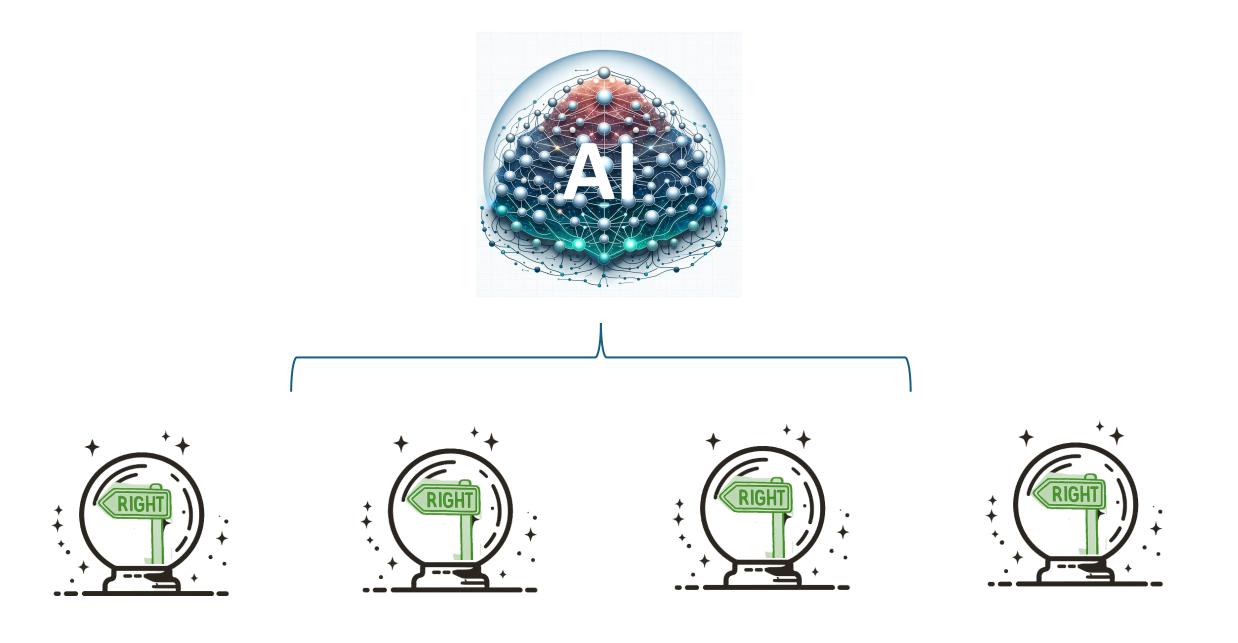


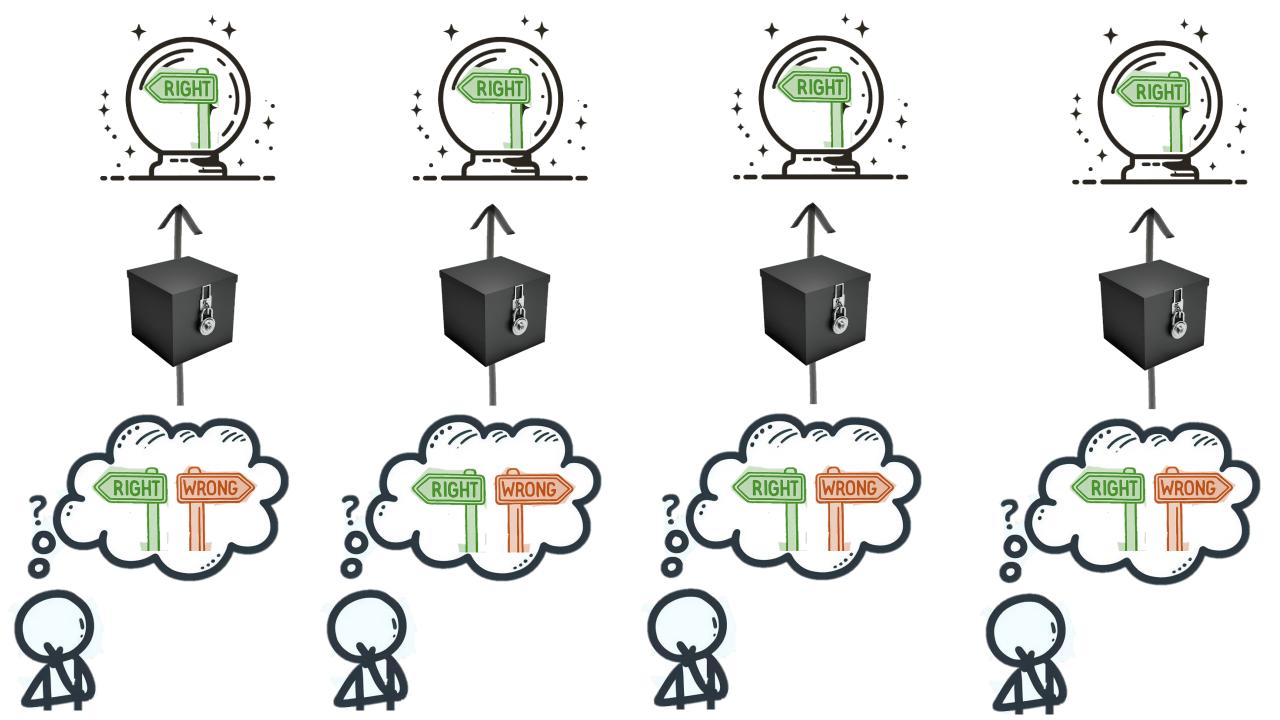
We discuss many of the other things we need to do to make sure Al has a positive impact on society in our recent book.

Check it out!

## How do you go about building morality into AI?







#### It took just one weekend for Meta's new AI Chatbot to become racist

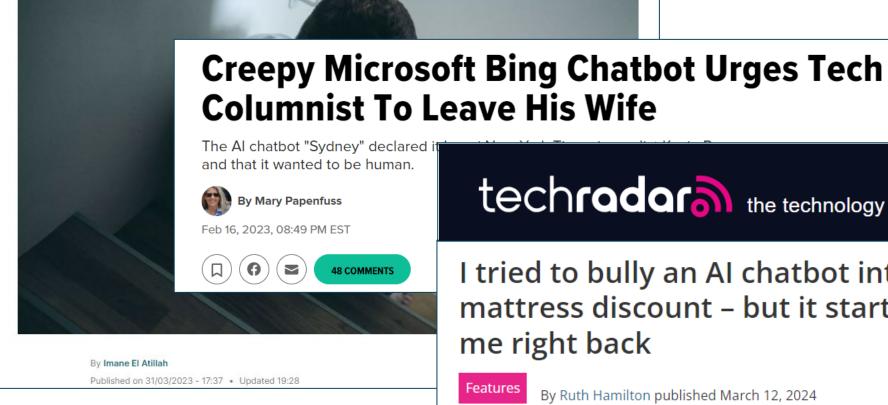
At least it's not sentient.

By Christianna Silva on August 8, 20



A chatbot learns from people. And p

#### Man ends his life after an Al chatbot 'encouraged' him to sacrifice himself to stop climate change



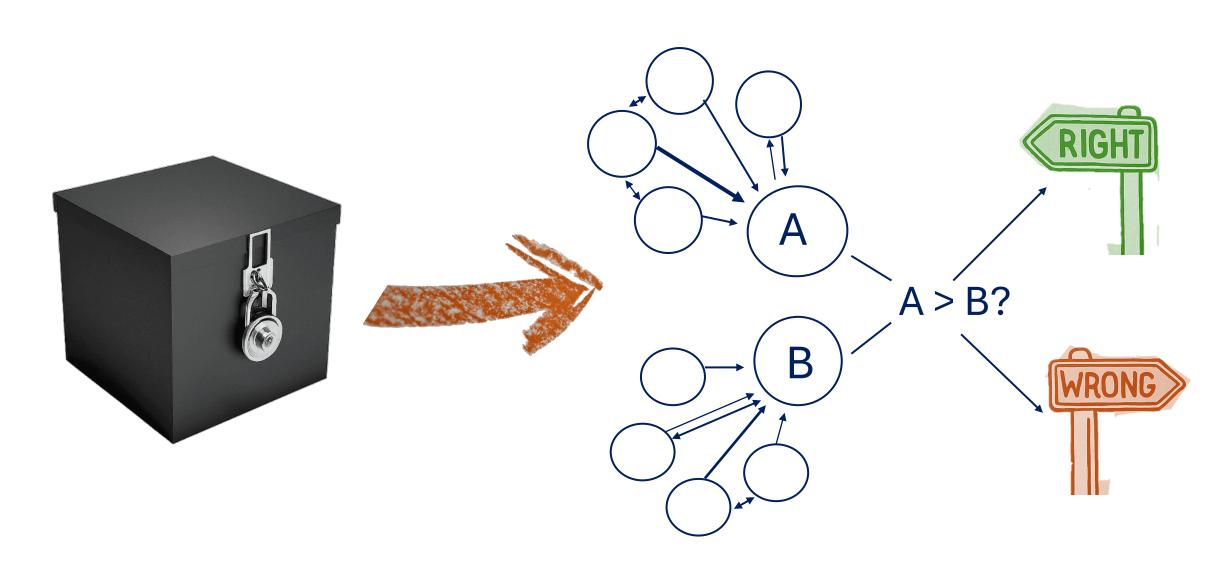
technology experts

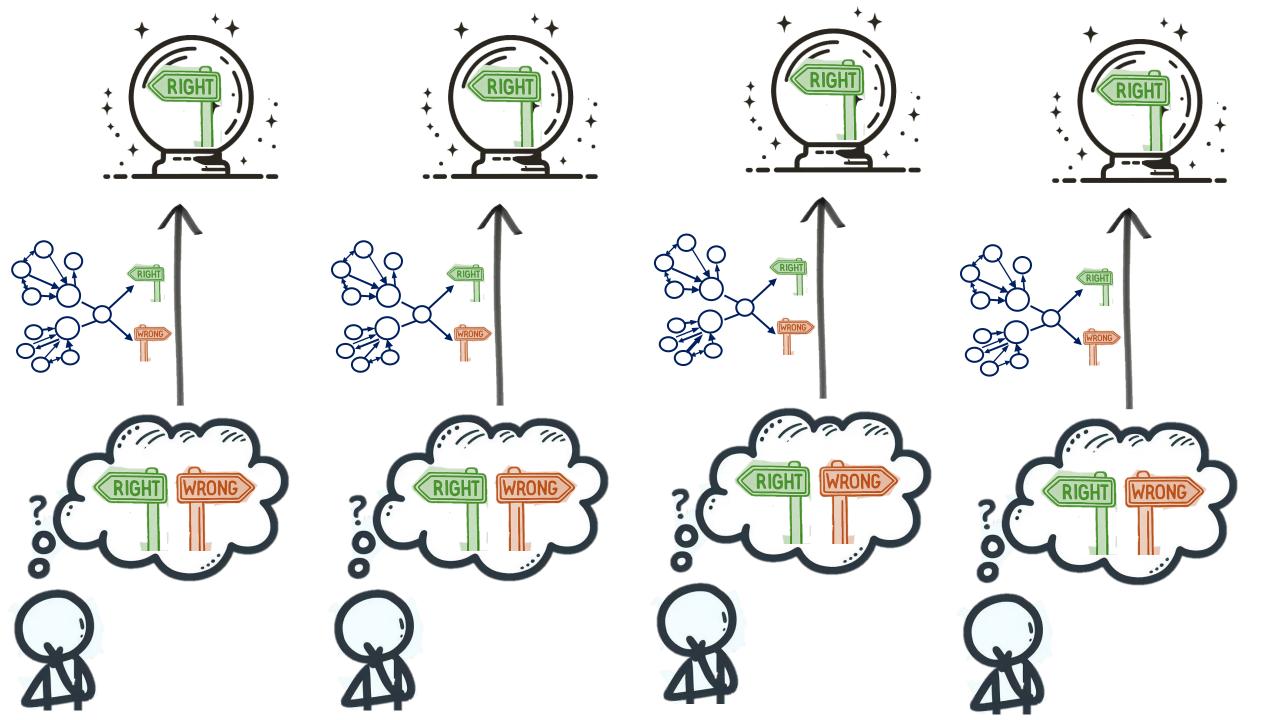
I tried to bully an AI chatbot into giving me a mattress discount - but it started bullying me right back

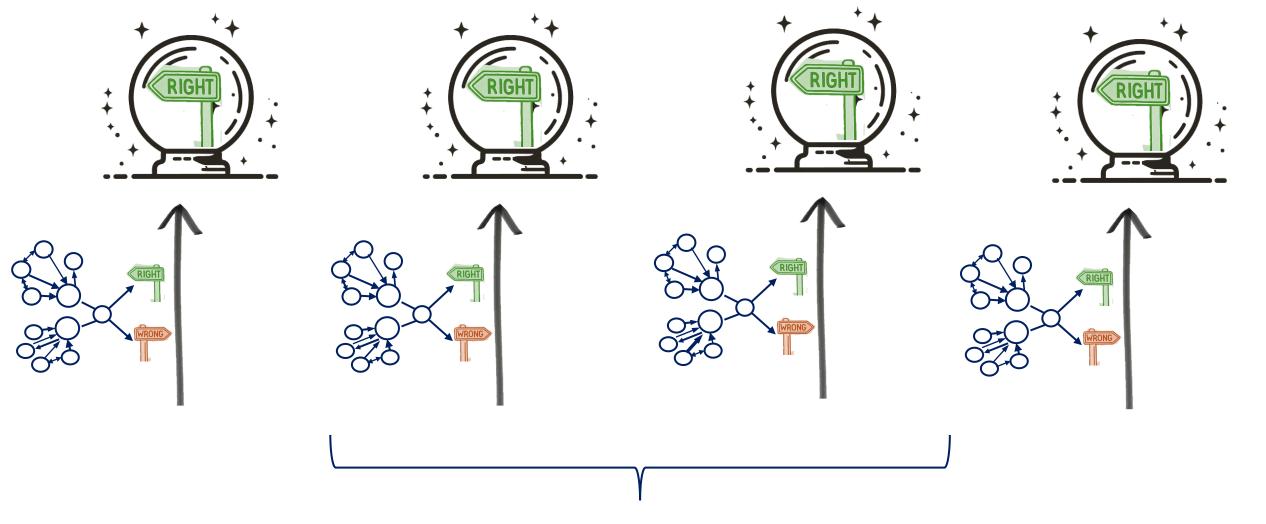
By Ruth Hamilton published March 12, 2024

Nectar has a new take on mattress shopping, and I'm not sold

### We are committed to using interpretable methods instead of black box models





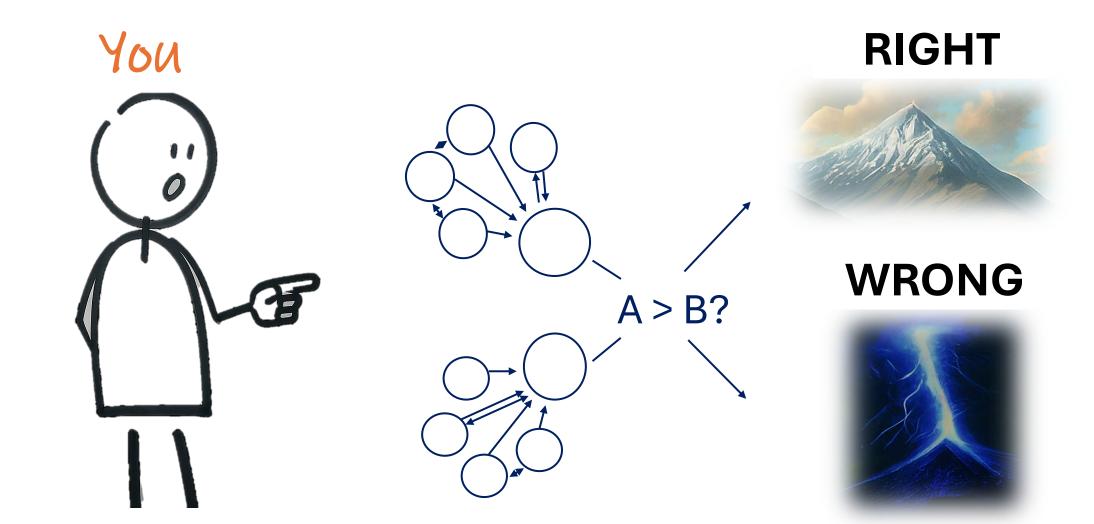


Common mechanisms found across these models often represent fundamental aspects of how the human moral brain works we didn't previously know about.

## The future (we are in the process of developing): Moral GPSes



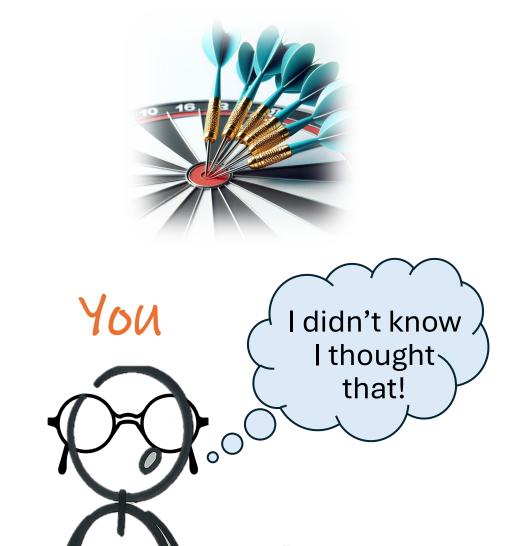
# You design, correct, and fine tune your own moral GPS model through optimized, interactive interfaces

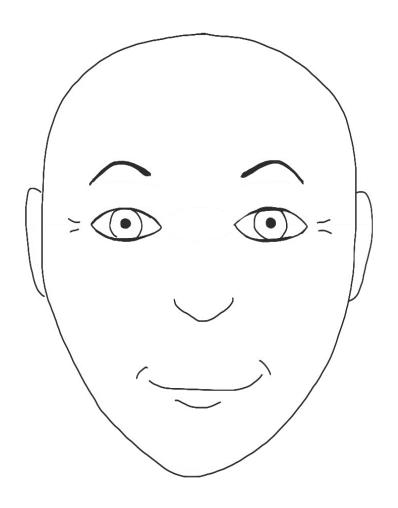


# Benefits of using moral models that are interpretable enough for you to understand and correct them

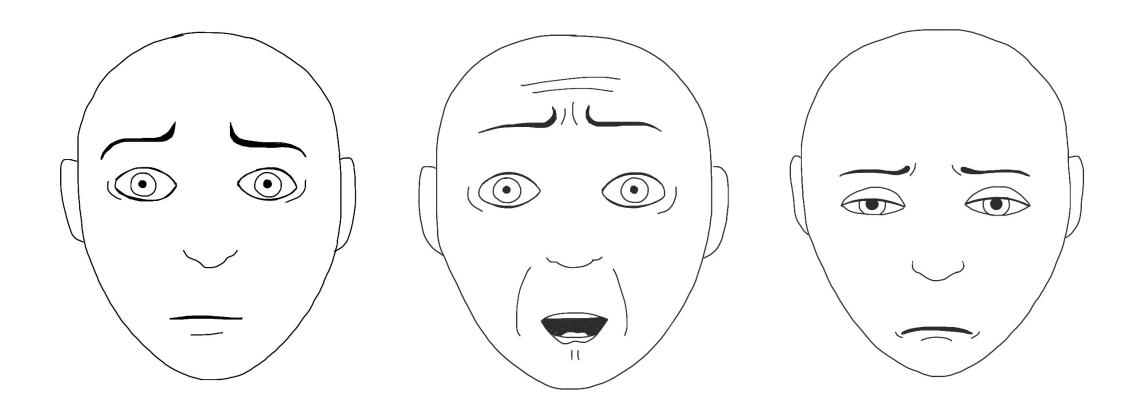
1) You make the moral model more accurate

2) You become more clear about what you think is morally right and wrong, and why





You do the training and fine-tuning when you are calm and have time to think.





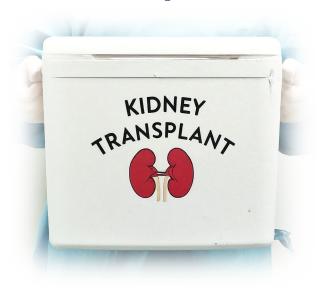
What your calm, rested, and thoughtful self would think is morally right.



Decisions that don't align well with the moral model you said you want to use.

# We have been applying this approach to kidney allocation decisions, and are extending it to other contexts as well

### **Kidney transplants**



### Other scarce medical resources



### Job applications and promotions



