

Accepted manuscript

Grossmann, I., Rotella, A., Hutcherson, C. A., Sharpinskyi, K., Varnum, Michael E. W., Achter, S., Dhami, M. K., Guo, X. E., Kara-Yakoubian, M., Mandel, D. R., Raes, L., Tay, L., Vie, A., Wagner, L., Adamkovic, M., Arami, A., Arriaga, P., Bandara, K., Baník, G. ... & Wilkening, T. (2023). Insights into the accuracy of social scientists' forecasts of societal change. *Nature Human Behaviour*, 7, 484-501. <https://doi.org/10.1038/s41562-022-01517-1>

Published in: Nature Human Behaviour

DOI: <https://doi.org/10.1038/s41562-022-01517-1>

AURA: <https://hdl.handle.net/11250/3107109>

Copyright: © His Majesty the King in Right of Canada as represented by
Department of National Defence 2023

Available:

This is the Author's Accepted Manuscript (AAM) of an article published by Nature Portfolio in *nature human behaviour* on 09. Feb. 2023, available at: <https://doi.org/10.1038/s41562-022-01517-1>

1 **Inventory of Supporting Information**

2

3 **1. Supplementary Information:**

4 **A. PDF Files**

5

Item	Present?	Filename	A brief, numerical description of file contents.
Supplementary Information	Yes	Grossmann_SI.pdf	Supplementary Methods, Supplementary Figures 1-15, Supplementary Tables 1-9, Supplementary Appendices 1-5
Reporting Summary	Yes	Grossmann_RS.pdf	
Peer Review Information	Yes	PRFile_Grossmann.pdf	

6

7 **Title**

8 **Insights into accuracy of social scientists' forecasts of societal change**

9 **Author list**

10 Igor Grossmann^{1*}, Amanda Rotella^{2,1}, Cendri A. Hutcherson³, Konstantyn Sharpinskyi¹, Michael
11 E. W. Varnum⁴, Sebastian Achter⁵, Mandeep K. Dhimi⁶, Xinqi Evie Guo⁷, Mane Kara-
12 Yakoubian⁸, David R. Mandel^{9,10}, Louis Raes¹¹, Louis Tay¹², Aymeric Vie^{13,14}, Lisa Wagner¹⁵,
13 Matus Adamkovic^{16,17}, Arash Arami¹⁸, Patricia Arriaga¹⁹, Kasun Bandara²⁰, Gabriel Baník¹⁶,
14 František Bartoš²¹, Ernest Baskin²², Christoph Bergmeir²³, Michał Białek²⁴, Caroline K.
15 Børsting²⁵, Dillon T. Browne¹, Eugene M. Caruso²⁶, Rong Chen²⁷, Bin-Tzong Chie²⁸, William J.
16 Chopik²⁹, Robert N. Collins⁹, Chin W. Cong³⁰, Lucian G. Conway³¹, Matthew Davis³², Martin V.
17 Day³³, Nathan A. Dhaliwal³⁴, Justin D. Durham³⁵, Martyna Dziekan³⁶, Christian T. Elbaek²⁵,
18 Eric Shuman³⁷, Marharyta Fabrykant^{38,39}, Mustafa Firat⁴⁰, Geoffrey T. Fong^{1,41}, Jeremy A.
19 Frimer⁴², Jonathan M. Gallegos⁴³, Simon B. Goldberg⁴⁴, Anton Gollwitzer^{45,46}, Julia Goyal⁴⁷,
20 Lorenz Graf-Vlachy^{48,49}, Scott D. Gronlund³⁵, Sebastian Hafenbrädl⁵⁰, Andree Hartanto⁵¹,
21 Matthew J. Hirshberg⁵², Matthew J. Hornsey⁵³, Piers D. L. Howe⁵⁴, Anoosha Izadi⁵⁵, Bastian
22 Jaeger⁵⁶, Pavol Kačmár⁵⁷, Yeun Joon Kim⁵⁸, Ruslan Krenzler^{59,60}, Daniel G. Lannin⁶¹, Hung-
23 Wen Lin⁶², Nigel Mantou Lou^{63,64}, Verity Y. Q. Lua⁵¹, Aaron W. Lukaszewski^{65,66}, Albert L.
24 Ly⁶⁷, Christopher R. Madan⁶⁸, Maximilian Maier⁶⁹, Nadyanna M. Majeed⁷⁰, David S. March⁷¹,
25 Abigail A. Marsh⁷², Michal Misiak^{24,73}, Kristian Ove R. Myrseth⁷⁴, Jaime M. Napan⁶⁷, Jonathan
26 Nicholas⁷⁵, Konstantinos Nikolopoulos⁷⁶, Jiaqing O⁷⁷, Tobias Otterbring^{78,79}, Mariola Paruzel-
27 Czachura^{80,81}, Shiva Pauer²¹, John Protzko⁸², Quentin Raffaelli⁸³, Ivan Ropovik^{84,85}, Robert M.
28 Ross⁸⁶, Yefim Roth⁸⁷, Espen Røysamb⁸⁸, Landon Schnabel⁸⁹, Astrid Schütz⁹⁰, Matthias Seifert⁹¹,
29 A. T. Sevincer⁹², Garrick T. Sherman⁹³, Otto Simonsson^{94,95}, Ming-Chien Sung⁹⁶, Chung-Ching
30 Tai⁹⁶, Thomas Talhelm⁹⁷, Bethany A. Teachman⁹⁸, Philip E. Tetlock^{99,100}, Dimitrios
31 Thomakos¹⁰¹, Dwight C. K. Tse¹⁰², Oliver J. Twardus¹⁰³, Joshua M. Tybur⁵⁶, Lyle Ungar⁹³, Daan
32 Vandermeulen¹⁰⁴, Leighton Vaughan Williams¹⁰⁵, Hrag A. Vosgerichian¹⁰⁶, Qi Wang¹⁰⁷, Ke
33 Wang¹⁰⁸, Mark E. Whiting^{109,110}, Conny E. Wollbrant¹¹¹, Tao Yang¹¹², Kumar Yogeewaran¹¹³,
34 Sangsuk Yoon¹¹⁴, Ventura r. Alves¹¹⁵, Jessica R. Andrews-Hanna^{83,116}, Paul A. Bloom⁷⁵,
35 Anthony Boyles¹¹⁷, Loo Charis¹¹⁸, Mingyeong Choi¹¹⁹, Sean Darling-Hammond¹²⁰, Zoe E.
36 Ferguson¹²¹, Cheryl R. Kaiser⁴³, Simon T. Karg¹²², Alberto López Ortega¹²³, Lori Mahoney¹²⁴,
37 Melvin S. Marsh¹²⁵, Marcellin F. R. C. Martinie⁵⁴, Eli K. Michaels¹²⁶, Philip Millroth¹²⁷, Jeanean
38 B. Naqvi¹²⁸, Weiting Ng¹²⁹, Robb B. Rutledge¹³⁰, Peter Slattery¹³¹, Adam H. Smiley⁴³, Oliver
39 Strijbis¹³², Daniel Sznycer¹³³, Eli Tsukayama¹³⁴, Austin van Loon¹³⁵, Jan G. Voelkel¹³⁵, Margaux
40 N. A. Wienk⁷⁵, Tom Wilkening¹³⁶, & The Forecasting Collaborative**

41 **Affiliations**

42 ¹Department of Psychology, University of Waterloo; Waterloo, Canada.

43 ²Department of Psychology, Northumbria University, UK.

44 ³Department of Psychology, University of Toronto Scarborough; Toronto, Canada.

45 ⁴Department of Psychology, Arizona State University, Tempe, USA.

- 46 ⁵Institute of Management Accounting and Simulation, Hamburg University of Technology;
47 Hamburg, Germany.
- 48 ⁶Department of Psychology, Middlesex University London; London, UK.
- 49 ⁷Department of Experimental Psychology, University of California San Diego; San Diego, USA.
- 50 ⁸Department of Psychology, Toronto Metropolitan University; Toronto, Canada.
- 51 ⁹Defence Research and Development Canada; Ottawa, Canada.
- 52 ¹⁰Department of Psychology, York University; Toronto, Canada.
- 53 ¹¹Department of Economics, Tilburg University, Tilburg, Netherlands.
- 54 ¹²Department of Psychological Sciences, Purdue University; West Lafayette, USA.
- 55 ¹³Mathematical Institute, University of Oxford; Oxford, UK.
- 56 ¹⁴Institute of New Economic Thinking, University of Oxford, Oxford, UK.
- 57 ¹⁵Jacobs Center for Productive Youth Development, University of Zurich; Zurich, Switzerland.
- 58 ¹⁶Institute of Psychology, University of Prešov; Prešov, Slovakia.
- 59 ¹⁷Institute of Social Sciences, CSPP, Slovak Academy of Sciences; Bratislava, Slovakia.
- 60 ¹⁸Department of Mechanical and Mechatronics Engineering, University of Waterloo; Waterloo,
61 Canada.
- 62 ¹⁹Iscte-University Institute of Lisbon, CIS; Lisbon, Portugal.
- 63 ²⁰Melbourne Centre for Data Science, University of Melbourne; Melbourne, Australia.
- 64 ²¹Faculty of Social and Behavioural Sciences, University of Amsterdam; Amsterdam,
65 Netherlands.
- 66 ²²Department of Food Marketing, Haub School of Business, Saint Joseph's University;
67 Philadelphia, USA.
- 68 ²³Department of Data Science and Artificial Intelligence, Monash University; Melbourne,
69 Australia.
- 70 ²⁴Institute of Psychology, University of Wrocław; Wrocław, Poland.
- 71 ²⁵Department of Management, Aarhus University; Aarhus, Denmark.
- 72 ²⁶Anderson School of Management, UCLA; Los Angeles, USA.
- 73 ²⁷Department of Psychology, Dominican University of California; San Rafael, USA.
- 74 ²⁸Department of Industrial Economics, Tamkang University; New Taipei City, Taiwan.
- 75 ²⁹Department of Psychology, Michigan State University; East Lansing, USA.
- 76 ³⁰Department of Psychology and Counselling, Universiti Tunku Abdul Rahman; Kampar,
77 Malaysia.
- 78 ³¹Psychology Department, Grove City College; Grove City, USA.
- 79 ³²Department of Economics, Siena College; Loudonville, USA.

- 80 ³³Department of Psychology, Memorial University of Newfoundland; St. John's, Canada.
- 81 ³⁴UBC Sauder School of Business, University of British Columbia; Vancouver, Canada.
- 82 ³⁵Department of Psychology, University of Oklahoma; Norman, USA.
- 83 ³⁶Faculty of Psychology and Cognitive Science, Adam Mickiewicz University; Poznań, Poland.
- 84 ³⁷Department of Psychology, University of Groningen; Groningen, Netherlands.
- 85 ³⁸Laboratory for Comparative Studies in Mass Consciousness, Expert Institute, HSE University;
86 Moscow, Russia.
- 87 ³⁹Faculty of Philosophy and Social Sciences, Belarusian State University; Minsk, Belarus.
- 88 ⁴⁰Department of Sociology, Radboud University; Nijmegen, Netherlands.
- 89 ⁴¹Ontario Institute for Cancer Research; Toronto, Canada.
- 90 ⁴²Department of Psychology, University of Winnipeg; Winnipeg, Canada.
- 91 ⁴³Department of Psychology, University of Washington; Seattle, USA.
- 92 ⁴⁴Department of Counseling Psychology, University of Wisconsin - Madison; Madison, USA.
- 93 ⁴⁵Department of Leadership and Organizational Behaviour, BI Norwegian Business School;
94 Oslo, Norway.
- 95 ⁴⁶Center for Adaptive Rationality, Max Planck Institute for Human Development; Berlin,
96 Germany.
- 97 ⁴⁷School of Public Health Sciences, University of Waterloo; Waterloo, Canada.
- 98 ⁴⁸TU Dortmund University; Dortmund, Germany.
- 99 ⁴⁹ESCP Business School; Berlin, Germany.
- 100 ⁵⁰IESE Business School; Barcelona, Spain.
- 101 ⁵¹School of Social Sciences, Singapore Management University; Singapore.
- 102 ⁵²Center for Healthy Minds, University of Wisconsin-Madison; Madison, USA.
- 103 ⁵³University of Queensland Business School; Brisbane, Australia.
- 104 ⁵⁴Melbourne School of Psychological Sciences, University of Melbourne; Melbourne, Australia.
- 105 ⁵⁵Department of Marketing, University of Massachusetts Dartmouth; Dartmouth, USA.
- 106 ⁵⁶Faculty of Behavioural and Movement Sciences, Vrije Universiteit Amsterdam; Amsterdam,
107 Netherlands.
- 108 ⁵⁷Department of Psychology, Faculty of Arts, Pavol Jozef Šafárik University; Košice, Slovakia.
- 109 ⁵⁸Cambridge Judge Business School, University of Cambridge; Cambridge, UK.
- 110 ⁵⁹Hermes Germany GmbH; Hamburg, Germany.
- 111 ⁶⁰University of Hamburg; Hamburg, Germany.
- 112 ⁶¹Department of Psychology, Illinois State University; Normal, USA.

- 113 ⁶²Department of Marketing and Logistics Management, National Penghu University of Science
114 and Technology; Magong City, Taiwan.
- 115 ⁶³Department of Psychology, University of Victoria; Victoria, Canada.
- 116 ⁶⁴Centre for Youth and Society, University of Victoria; Victoria, Canada.
- 117 ⁶⁵Department of Psychology, California State University - Fullerton; Fullerton, USA.
- 118 ⁶⁶Center for the Study of Human Nature, California State University - Fullerton; Fullerton, USA.
- 119 ⁶⁷Department of Psychology, Loma Linda University; Loma Linda, USA.
- 120 ⁶⁸University of Nottingham; Nottingham, UK.
- 121 ⁶⁹Department of Experimental Psychology, University College London; London, UK.
- 122 ⁷⁰Singapore Management University; Singapore.
- 123 ⁷¹Department of Psychology, Florida State University; Tallahassee, USA.
- 124 ⁷²Department of Psychology, Georgetown University; Washington DC, USA.
- 125 ⁷³School of Anthropology & Museum Ethnography, University of Oxford; Oxford, UK.
- 126 ⁷⁴School for Business and Society, University of York; York, UK.
- 127 ⁷⁵Department of Psychology, Columbia University; New York City, USA.
- 128 ⁷⁶IHRR Forecasting Laboratory, Durham University; Durham, UK.
- 129 ⁷⁷Department of Psychology, Aberystwyth University; Aberystwyth, UK.
- 130 ⁷⁸School of Business and Law, Department of Management, University of Agder; Kristiansand,
131 Norway.
- 132 ⁷⁹Institute of Retail Economics; Stockholm, Sweden.
- 133 ⁸⁰Institute of Psychology, University of Silesia in Katowice, Poland.
- 134 ⁸¹ Department of Neurology, Penn Center for Neuroaesthetics, University of Pennsylvania,
135 United States.
- 136 ⁸²Central Connecticut State University; New Britain, USA.
- 137 ⁸³Department of Psychology, University of Arizona; Tucson, USA.
- 138 ⁸⁴Faculty of Education, Institute for Research and Development of Education, Charles
139 University; Prague, Czech Republic.
- 140 ⁸⁵Faculty of Education; University of Prešov; Prešov, Slovakia.
- 141 ⁸⁶School of Psychology, Macquarie University; Sydney, Australia.
- 142 ⁸⁷Department of Human Service, University of Haifa; Haifa, Israel.
- 143 ⁸⁸Promenta Center, Department of Psychology, University of Oslo; Oslo, Norway.
- 144 ⁸⁹Department of Sociology, Cornell University; Ithaca, USA.
- 145 ⁹⁰Institute of Psychology, University of Bamberg; Bamberg, Germany.
- 146 ⁹¹IE Business School, IE University; Madrid, Spain.

147 ⁹²Faculty of Psychology and Human Movement Science, University of Hamburg; Hamburg,
148 Germany.

149 ⁹³Department of Computer and Information Science, University of Pennsylvania; Philadelphia,
150 USA.

151 ⁹⁴Department of Clinical Neuroscience, Karolinska Institutet; Solna, Sweden.

152 ⁹⁵Department of Sociology, University of Oxford; Oxford, UK.

153 ⁹⁶Department of Decision Analytics and Risk, University of Southampton; Southampton, UK.

154 ⁹⁷University of Chicago Booth School of Business; Chicago, USA.

155 ⁹⁸Department of Psychology, University of Virginia; Charlottesville, USA.

156 ⁹⁹Psychology Department, University of Pennsylvania; Philadelphia, USA.

157 ¹⁰⁰Wharton School of Business, University of Pennsylvania; Philadelphia, USA.

158 ¹⁰¹Department of Economics, National and Kapodistrian University of Athens; Athens, Greece.

159 ¹⁰²School of Psychological Sciences and Health, University of Strathclyde; Glasgow, UK.

160 ¹⁰³Department of Psychology, University of Guelph; Guelph, Canada.

161 ¹⁰⁴Psychology Department, Hebrew University of Jerusalem; Jerusalem, Israel.

162 ¹⁰⁵Department of Economics, Nottingham Trent University; Nottingham, UK.

163 ¹⁰⁶Department of Management and Organizations, Northwestern University; Evanston, USA.

164 ¹⁰⁷College of Human Ecology, Cornell University, Ithaca, USA.

165 ¹⁰⁸Harvard Kennedy School, Harvard University; Cambridge, USA.

166 ¹⁰⁹Computer and Information Science, University of Pennsylvania; Philadelphia, USA.

167 ¹¹⁰Operations, Information, and Decisions Department, The Wharton School, University of
168 Pennsylvania; Philadelphia, USA.

169 ¹¹¹School of Economics and Finance, University of St. Andrews; St. Andrews, UK.

170 ¹¹²Department of Management, Cameron School of Business, University of North Carolina;
171 Wilmington, USA.

172 ¹¹³School of Psychology, Speech and Hearing, University of Canterbury; Christchurch, New
173 Zealand, ¹¹⁴Department of Marketing, University of Dayton; Dayton, USA.

174 ¹¹⁵ISG Universidade Lusofona; Lisbon, Portugal.

175 ¹¹⁶Cognitive Science, University of Arizona; Tucson, USA.

176 ¹¹⁷Ephemer AI; Atlanta, USA.

177 ¹¹⁸Questrom School of Business, Boston University; Boston, USA.

178 ¹¹⁹Institute of Social Science Research, Pusan National University; Busan, South
179 Korea, ¹²⁰Fielding School of Public Health, UCLA; Los Angeles, USA.

180 ¹²¹Psychology Department, University of Washington; Seattle, USA.

181 ¹²²Department of Political Science, Aarhus University; Aarhus, USA.
182 ¹²³Faculty of Social Sciences, Vrije Universiteit Amsterdam; Amsterdam, Netherlands.
183 ¹²⁴College of Science and Mathematics, Wright State University; Fairborn, USA.
184 ¹²⁵Department of Psychology, Georgia Southern University; Statesboro, USA.
185 ¹²⁶Division of Epidemiology, School of Public Health, University of California, Berkeley;
186 Berkeley, USA.
187 ¹²⁷Department of Psychology, Uppsala University; Uppsala, Sweden.
188 ¹²⁸Department of Psychology, Carnegie Mellon University; Pittsburgh, USA.
189 ¹²⁹School of Humanities & Behavioral Sciences, Singapore University of Social Sciences;
190 Singapore.
191 ¹³⁰Department of Psychology, Yale University; New Haven, USA.
192 ¹³¹BehaviourWorks Australia. Monash University; Melbourne, Australia.
193 ¹³²Institute of Political Science, University of Zurich; Zurich, Switzerland.
194 ¹³³Department of Psychology, Oklahoma State University; Stillwater, USA.
195 ¹³⁴Department of Business Administration, University of Hawaii - West Oahu; Kapolei, USA.
196 ¹³⁵Department of Sociology, Stanford University; Stanford, USA.
197 ¹³⁶Department of Economics, University of Melbourne; Melbourne, Australia.
198 *Corresponding author. Email: igrossma@uwaterloo.ca
199 **A list of authors and their affiliations appears at the end of the paper.
200

201 ***Abstract***

202 How well can social scientists predict societal change, and what processes underlie their
203 predictions? To answer these questions, we ran two forecasting tournaments testing accuracy of
204 predictions of societal change in domains commonly studied in the social sciences: ideological
205 preferences, political polarization, life satisfaction, sentiment on social media, and gender-career
206 and racial bias. Following provision of historical trend data on the domain, social scientists
207 submitted pre-registered monthly forecasts for a year (Tournament 1; $N=86$ teams/359 forecasts),
208 with an opportunity to update forecasts based on new data six months later (Tournament 2;
209 $N=120$ teams/546 forecasts). Benchmarking forecasting accuracy revealed that social scientists'
210 forecasts were on average no more accurate than simple statistical models (historical means,
211 random walk, or linear regressions) or the aggregate forecasts of a sample from the general
212 public ($N=802$). However, scientists were more accurate if they had scientific expertise in a
213 prediction domain, were interdisciplinary, used simpler models, and based predictions on prior
214 data.

215
216

217 **Main Text**

218 Can social scientists predict societal change? Governments and the general public often
219 rely on experts, based on a general belief that they make better judgments and predictions of the
220 future in their domain of expertise. The media also seek out experts to render their judgments
221 and opinions about what to expect in the future^{1,2}. Yet research on predictions in many domains
222 suggests that experts may not be better than purely stochastic models in predicting the future. For
223 example, portfolio managers (who are paid for their expertise) do not outperform the stock
224 market in their predictions³. Similarly, in the domain of geopolitics, experts often perform at
225 chance levels when forecasting occurrences of specific political events⁴. Based on these insights,
226 one might expect that experts would find it difficult to accurately predict societal change.

227 At the same time, social science researchers have developed rich, empirically-grounded
228 models to explain social science phenomena. Examining sampled data, social scientists strive to
229 develop theoretical models about causal mechanisms that, in ideal cases, reliably describe human
230 behavior and societal processes⁵. Therefore, it is possible that explanatory models afford social
231 science experts an advantage in predicting social phenomena in their domain of expertise. Here,
232 we test these possibilities, examining the overall predictability of trends in social phenomena
233 such as political polarization, racial bias, or well-being, and whether experts in social science are
234 better able to predict those trends compared to non-experts.

235 Prior forecasting initiatives have not fully addressed this question for two reasons. First,
236 forecasting initiatives with subject matter experts have focused on examining the probability of
237 occurrence for specific one-time events^{4,6} rather than the accuracy of *ex-ante* predictions of
238 societal change over multiple units of time. In a sense, predicting events in the future (*ex-ante*) is
239 the same as predicting events that have already happened, as long as the experts (research
240 participants) don't know the outcome. Yet, there are reasons to think that future prediction is
241 different in an important way. Consider stock prices: participants could predict stock returns for
242 stocks in the past, except that they know many other things that have happened (conflicts,
243 bubbles, Black Swans, economic trends, consumption trends, etc.). Post-hoc, those making those
244 making predictions have access to the temporal variance/occurrence for each of these variables
245 and hence are more likely to be successful in *ex post* predictions. Thus, predictions about past
246 events end up being more about testing people's explanations, rather than their predictions *per se*.
247 Moreover, all other things being equal, the likelihood of a prediction regarding a one-off event
248 being accurate is by default higher than that of a prediction regarding societal change across an
249 extended time period. Binary predictions for the one-off event do not require accuracy in
250 estimating degree of change or the shape of the predicted time series, which are extra challenges
251 in forecasting societal change.

252 Second, past research on forecasting has concentrated on predicting geopolitical⁴ or
253 economic events⁷ rather than broader societal phenomena. Thus, in contrast to systematic
254 studies concerning the replicability of in-sample explanations of social science phenomena⁸,
255 out-of-sample prediction accuracy in the social sciences remains understudied^{9,10}. Similarly,
256 little is known about the rationales and approaches social scientists use to make predictions for
257 societal trends. For example, are social scientists more apt to rely on data-driven statistical
258 methods or theory and intuitions when generating such predictions?

259 To address these unknowns, we performed a standardized evaluation of forecasting
260 accuracy⁹ among social scientists in well-studied domains for which systematic, cross-temporal
261 data is available, namely subjective well-being, racial bias, ideological preferences, political
262 polarization, and gender-career bias. With the onset of the COVID-19 pandemic as a backdrop,

263 we selected these domains based on data availability and theoretical links to the pandemic. Prior
264 research has suggested that each of these domains may be impacted by infectious disease^{11–14} or
265 pandemic-related social isolation¹⁵. To understand how scientists made predictions in these
266 domains, we documented the rationales and processes they used to generate forecasts, then
267 examined how different methodological choices were related to accuracy.

268 Research Overview

269 We present results from two forecasting tournaments conducted through the Forecasting
270 Collaborative—a crowdsourced initiative among scientists interested in ex-ante testing of their
271 theoretical or data-driven models. Examining performance across two tournaments allowed us to
272 test the stability of forecasting accuracy in the context of unfolding societal events, and to
273 investigate how social scientists recalibrate their models and incorporate new data when asked to
274 update their forecasts.

275 The Forecasting Collaborative was open to behavioral, social, and data scientists from
276 any field who wanted to participate in the tournament and were willing to provide forecasts over
277 12 months (May 2020 – April 2021) as part of the initial tournament and, upon receiving
278 feedback on initial performance, again after 6 months for a follow-up tournament (recruitment
279 details in Methods and demographic information in Supplementary Table 1). To ensure a
280 “common task framework”^{9,16,17}, we provided all participating teams with the same time series
281 data for the US for each of the 12 variables related to the phenomena of interest (i.e., life
282 satisfaction, positive affect, negative affect, support for Democrats, support for Republicans,
283 political polarization, explicit and implicit attitudes towards Asian Americans, explicit and
284 implicit attitudes towards African Americans, and explicit and implicit associations between
285 gender and specific careers.

286 Participating teams received historical data that spanned 39 months (January 2017 to
287 March 2020) for Tournament 1 and data that spanned 45 months for Tournament 2 (January
288 2017 to September 2020), which they could use to inform their forecasts for the future values of
289 the same time series. Teams could select up to 12 domains to forecast, including domains for
290 which team members reported a track record of peer-reviewed publications as well as domains
291 for which they did not possess relevant expertise (see Methods for multi-stage operationalization
292 of expertise). By including social scientists with expertise in different subject matters, we could
293 examine how such expertise may contribute to forecasting accuracy above and beyond general
294 training in the social sciences. Teams were not constrained in terms of the methods used to
295 generate time-point forecasts. They provided open-ended, free-text responses for the descriptions
296 of the methods used, which were coded later. If they made use of data-driven methods, they also
297 provided the model and any additional data used to generate their forecasts (see Methods). We
298 also collected data on team size and composition, area of research specialization, subject domain
299 and forecasting expertise, and prediction confidence.

300 We benchmarked forecasting accuracy against several alternatives. First, we evaluated
301 whether social scientists’ forecasts in Tournament 1 were better than the wisdom of the crowd
302 (i.e., the average forecasts of a sample of lay participants recruited from Prolific). Second, we
303 compared social scientists’ performance in both tournaments to naïve random extrapolation
304 algorithms (i.e., the average of historical data, random walks, and estimates based on linear
305 trends). Finally, we systematically evaluated the accuracy of different forecasting strategies used
306 by the social scientists in our tournaments, as well as the effect of expertise.

307

308 Results

309 Following the a priori outlined analytic plan (osf.io/7ekfm; details in the Supplementary
310 Methods) to determine forecasting accuracy across domains, we examined the mean absolute
311 scaled error (MASE)¹⁸ across forecasted time-points for each domain. MASE is an
312 asymptotically normal, scale-independent scoring rule that compares predicted values against
313 predictions of a one-step random walk. Because it is scale-independent, it is an adequate measure
314 when comparing accuracy across domains on different scales. A MASE of 1 reflects a forecast
315 that is as good out-of-sample as the naive one-step random walk forecast is in-sample. A MASE
316 below 1.76 is superior to median performance in prior large-scale data science competitions⁷.
317 See Supplementary Materials for further details of the MASE method.

318 In addition to absolute accuracy, we also assessed the comparative accuracy of social
319 scientists' forecasts using several benchmarks. First, during the period of the first tournament, we
320 obtained forecasts from a non-expert crowdsourced sample of US residents ($N = 802$) via Prolific
321¹⁹ who received the same data as tournament participants and filled out an identically structured
322 survey to provide a wisdom-of-the-(lay)-crowd benchmark. Second, for both tournaments we
323 simulated three different data-based naïve approaches to out-of-sample forecasting using the
324 time series data provided to participants in the tournament, including, 1) the historical mean,
325 calculated by randomly resampling the historical time series data; 2) a naïve random walk,
326 calculated by randomly resampling historical change in the time series data with an
327 autoregressive component; 3) extrapolation from linear regression, based on a randomly selected
328 interval of historical time series data (see Supplementary Information for details). This latter
329 approach captures the expected range of predictions that would have resulted from random,
330 uninformed use of historical data to make out-of-sample predictions (as opposed to the naïve in-
331 sample predictions that form the basis of MASE scores).

332 How accurate were behavioral and social scientists at forecasting?

333 Fig. 1 shows that in Tournament 1, social scientists' forecasts were, on average, inferior
334 to in-sample random walks in nine domains. In seven domains, social scientists' forecasts were
335 inferior to median performance in prior forecasting competitions (Supplementary Fig. 1 shows
336 raw estimates; Supplementary Fig. 2 reports measures of uncertainty around estimates). In
337 Tournament 2, forecasts were on average inferior to in-sample random walks in eight domains
338 and inferior to median performance in prior forecasting competitions in five domains. Even
339 winning teams were still less accurate than in-sample random walks for 8 of 12 domains in
340 Tournament 1, and one domain (Republican support) in Tournament 2 (Supplementary Tables 1-
341 2 and Supplementary Figs. 4-9). One should note that inferior performance to the in-sample
342 random walk ($MASE > 1$) may not be too surprising; errors of the in-sample random walk in the
343 denominator concern historical observations that occurred before the pandemic, whereas
344 accuracy of scientific forecasts in the numerator is compared concerns the data for the first
345 pandemic year. However, average forecasting accuracy did not generally beat more liberal
346 benchmarks such as the median MASE in data science tournaments (1.76)⁷ or the benchmark
347 MASE for "good" forecasts in the tourism industry (see Supplement). Except for one team, top
348 forecasters from Tournament 1 did not appear among the winners of Tournament 2
349 (Supplementary Tables 1-2).

350 We examined the accuracy of scientific and lay forecasts in a linear mixed effect model.
351 To systematically compare results for different forecasted domains, we tested a full model with
352 expertise (social scientist versus lay crowd), domain, and their interaction as predictors, and
353 $\log(MASE)$ scores nested in participants. We observed no significant main effect difference

354 between accuracy of social scientists and lay crowds, $F(11, 1747) = 0.88, P = .348, \text{part } R^2$
355 $< .001$. However, we observed a significant interaction between social science training and
356 domain, $F(11, 1304) = 2.00, P = .026$. Simple effects show that social scientists were
357 significantly more accurate than lay people when forecasting life satisfaction, polarization, as
358 well as explicit and implicit gender-career bias. However, the scientific teams were no better
359 than the lay sample in the remaining eight domains (Figure 1 and Table 1). Moreover, Bayesian
360 analyses indicated that only for life satisfaction there is substantial evidence in favor of the
361 difference, whereas for eight domains evidence was in favor of the null hypothesis. See
362 supplementary information online for further details and interpretation of the multiverse analyses
363 of domain-general accuracy.

364 Cross-validation of domain-general accuracy via forecast versus trend comparisons

365 The most elementary analysis of domain-general accuracy involves inspecting trends for
366 each group and comparing them against the ground truth and historical time series in each
367 domain. Fig. 2 allows us to inspect individual trends of social scientists and the naïve crowd per
368 domain in Tournament 1, along with historical and ground truth markers for each domain. For
369 social scientists, one can observe the diversity of forecasts from individual teams (light blue)
370 along with a lowess regression and 95% confidence interval around the trend (blue). For the
371 naïve crowd, one can see an equivalent lowess trend and the 95% CI around it (salmon). In half
372 of the domains – explicit bias against African Americans, implicit bias against Asian-Americans,
373 negative affect, life satisfaction, as well as support for Democrats and Republicans – lowess
374 curves from both groups were overlapping, suggesting that the estimates from both social
375 scientists and the naïve crowd were identical. Moreover, except for the domain of life
376 satisfaction, forecasts of scientists and the naïve crowd were close to far off the mark vis-à-vis
377 ground truth. In one further domain— explicit bias against African Americans—the naïve crowd
378 estimate was in fact closer to the ground truth marker than the estimate from the lowess curve of
379 the social scientists. In other five domains, which concerned explicit and implicit gender career
380 bias, explicit bias against Asian-Americans, positive affect and political polarization, social
381 scientists' forecasts were closer to the ground truth markers than the naïve crowd. We note,
382 however, that these visual inspections may be somewhat misleading because the confidence
383 intervals don't correct for multiple tests. This caveat aside, the overall message remains
384 consistent with the results of the statistical tests above: For most domains social scientists'
385 predictions were either similar to or worse than the naïve crowd's predictions.

386 Comparisons to naïve statistical benchmarks

387 Next, we compared scientific forecasts against three naïve statistical benchmarks by creating
388 benchmark/forecast ratio scores (a ratio of 1 indicates that the social scientists' forecasts were
389 equal in accuracy to the benchmarks, with ratios greater than 1 indicating greater accuracy). To
390 account for interdependence of social scientists' forecasts, we examined estimated ratio scores
391 for domains from linear mixed models, with responses nested in forecasting teams. To reduce the
392 likelihood that social scientists' forecasts beat naïve benchmarks by chance, our main analyses
393 focus on performance across all three benchmarks (see Supplement for rationale favoring this
394 method over averaging across three benchmarks), and by adjusting confidence intervals of the
395 ratio scores for simultaneous inference of 12 domains in each tournament by simulating a
396 multivariate t distribution²⁰. Figs. 1 and 3 and Supplementary Fig. 2 show that social scientists in
397 Tournament 1 were significantly better than each of the three benchmarks in only one out of
398 twelve domains, which concerned explicit gender-career bias, $1.53 < \text{ratio} \leq 1.60, 1.16 < 95\%CI$
399 ≤ 2.910 . In the remaining 11 domains, scientific predictions were either no different or worse

400 than the benchmarks. The relative advantage of scientific forecasts over the historical mean and
401 random walk benchmarks was somewhat larger in Tournament 2 (Supplementary Fig. 1).
402 Scientific forecasts were significantly more accurate than the three naïve benchmarks in five out
403 of twelve domains. These domains reflected explicit racial bias: African-American bias, $2.20 <$
404 $\text{ratio} \leq 2.86$, $1.55 < 95\%CI \leq 4.05$; Asian-American bias, $1.39 < \text{ratio} \leq 3.14$, $1.01 < 95\%CI \leq$
405 4.40 ; and implicit racial and gender career biases: African-American bias, $1.35 < \text{ratio} \leq 2.00$,
406 $1.35 < 95\%CI \leq 2.78$; Asian-American bias, $1.36 < \text{ratio} \leq 2.73$, $1.001 < 95\%CI \leq 3.71$; gender-
407 career bias, $1.59 < \text{ratio} \leq 3.22$, $1.15 < 95\%CI \leq 4.46$. In the remaining seven domains, forecasts
408 were not significantly different from naïve benchmarks. Moreover, as Fig. 3 shows, for political
409 polarization scientific forecasts in Tournament 2 were significantly *less* accurate than estimates
410 from a naïve linear regression, $\text{ratio} = 0.51$, $95\%CI [0.38, 0.68]$. Fig. 3 also shows that in most
411 domains at least one of the naïve forecasting methods produced errors that were comparable to or
412 less than social scientists' forecasts (11 out of 12 in Tournament 1 and 8 out of 12 in Tournament
413 2).

414 To compare social scientists' forecasts against the average of three naïve benchmarks, we
415 fit a linear mixed model with forecast/benchmark ratio scores nested in forecasting teams and
416 examined estimated means for each domain. In Tournament 1, scientists performed better than
417 the average of the naïve benchmarks in only three domains, which concerned political
418 polarization, $95\%CI [1.06; 1.63]$, explicit gender-career bias, $95\%CI [1.23; 1.95]$, and implicit
419 gender-career bias, $95\%CI [1.17; 1.83]$. In Tournament 2, social scientists performed better than
420 the average of the naïve benchmarks in seven domains, $1.07 < 95\%CIs \leq 2.79$, while they were
421 statistically indistinguishable from the average of naïve benchmarks when forecasting the
422 remaining five domains: ideological support for Democrats, $95\%CI [0.76; 1.17]$, and for
423 Republicans, $95\%CI [0.98; 1.51]$, explicit gender-career bias, $95\%CI [0.96; 1.52]$, and negative
424 affect on social media, $95\%CI [0.82; 1.25]$. Moreover, in Tournament 2 social scientists'
425 forecasts of political polarization were inferior to the average of naïve benchmarks, $95\%CI$
426 $[0.58; 0.89]$. Overall, social scientists tended to do worse than the average of the three naïve
427 statistical benchmarks in Tournament 1. While scientists did better than the average of naïve
428 benchmarks in Tournament 2, this difference in overall performance was small, $M(\text{forecast}$
429 $/\text{benchmark inaccuracy ratio}) = 1.43$, $95\%CI [1.26; 1.62]$. Moreover, in most domains at least
430 one of the naïve benchmarks was on par if not more accurate to social scientists' forecasts.

431 Which domains were harder to predict?

432 Fig. 4 shows that some societal trends were significantly harder to forecast than others,
433 Tournament 1: $F(11,295.69) = 41.88$, $P < .001$, $R^2 = .450$, Tournament 2: $F(11,469.49) = 26.87$,
434 $P < .001$, $R^2 = .291$. Forecast accuracy was lowest in politics (underestimating Democratic and
435 Republican support, and political polarization), well-being (underestimating life satisfaction and
436 negative affect on social media), and racial bias against African Americans (overestimating; also
437 see Supplementary Fig. 1). Differences in forecast accuracy across domains did not correspond
438 to differences in quality of ground truth markers: Based on the sampling frequency and
439 representativeness of the data, most reliable ground truth markers concerned societal change in
440 political ideology, obtained via an aggregate of multiple nationally representative surveys by
441 reputable pollsters, yet this domain was among most difficult to forecast. In contrast, some of the
442 least representative markers concerned racial and gender-bias, which came from Project
443 Implicit—a volunteer platform that is subject to self-selection bias—yet these domains were
444 among the easiest to forecast. In a similar vein, both life satisfaction and positive affect on social
445 media were estimated via texts on Twitter, even though forecasting errors between these domains

446 varied. Though measurement imprecision undoubtedly presents a challenge for forecasting, it is
447 unlikely to account for between-domain variability in forecasting errors (Figure 4).

448 Domain differences in forecasting accuracy corresponded to differences in the
449 complexity of historical data: domains ranked more variable in terms of standard deviation (*SD*)
450 and mean absolute difference (*MAD*) of historical data tended to have more forecasting error (as
451 measured by the rank-order correlation between median inaccuracy scores across teams and
452 variability scores for the same domain), Tournament 1: $\rho(SD) = .19$, $\rho(MAD) = .20$; Tournament
453 2: $\rho(SD) = .48$, $\rho(MAD) = .36$, and domain changes in variability of historical data across
454 tournaments corresponded to changes in accuracy, $\rho(SD) = .27$, $\rho(MAD) = .28$.

455 Comparison of accuracy across tournaments

456 Forecasting error was higher in the first tournament than the second tournament (see Fig.
457 4), $F(1, 889.48) = 64.59$, $P < .001$, $R^2 = .063$. We explored several possible differences between
458 the tournaments that may account for this effect. One possibility is that type of teams differed
459 between tournaments (team size, gender, number of forecasted domains, field specialization and
460 team diversity, number of PhDs on a team, prior experience with forecasting). However, the
461 difference between the tournaments remained equally pronounced when running parallel
462 analyses with team characteristics as covariates, $F(1, 847.79) = 90.45$, $P < .001$, $R^2 = .062$.

463 Another hypothesis is that forecasts for twelve months (Tournament 1) include further
464 removed data points than forecasts for more immediate six months (Tournament 2), and it is the
465 greater temporal distance between the tournament and the moment to forecast that results in
466 greater inaccuracy at Tournament 1. To test this hypothesis, we zeroed in on Tournament 1
467 inaccuracy scores for the first and the last six months, while including domain type as a control
468 dummy variable. By focusing on Tournament 1 data, we kept other characteristics such as team
469 composition as a constant. Contrary to this seemingly straightforward hypothesis, error for the
470 forecasts for the first six months was in fact significantly greater (MASE = 3.16, $SE = 0.21$,
471 95%CI [2.77, 3.60]) than for the last six months (MASE = 2.59, $SE = 0.17$, 95%CI [2.27, 2.95]),
472 $F(1, 621.41) = 29.36$, $P < .001$, $R^2 = .012$. As Supplementary Fig. 1 shows, for many domains,
473 social scientists underpredicted societal change in Tournament 1, and this difference between
474 predicted and observed values was more pronounced in the first versus last six months. This
475 suggests that for several domains social scientists anchored (19) their forecasts on the most
476 recent historical data. Figure 2 further indicates that many domains showed unusual shifts (vis-à-
477 vis prior historical data) in the first six months of the pandemic and started to return to the
478 historical baseline in the following six months. For these domains, forecasts anchored on the
479 most recent historical data were more inaccurate for the May-October 2020 forecasts compared
480 to the November 2020-April 2021 forecasts.

481 Finally, we tested whether providing teams additional six months of historical trend
482 capturing the on-set of the novel pandemic at Tournament 2 may have contributed to lower error
483 compared to Tournament 1. To this end, we compared the inaccuracy of forecasts for the six
484 months period of Nov 2020-April 2021 done in May 2020 (Tournament 1) and when provided
485 with more data in October 2020 (Tournament 2). We focused only on participants who
486 completed both tournaments to keep the number of participating teams and team characteristics
487 constant. Indeed, Tournament 1 forecasts had significantly more error (MASE $M = 2.54$, $SE =$
488 0.17 , 95%CI [2.23, 2.90]) than Tournament 2 forecasts (MASE $M = 1.99$, $SE = 0.13$, 95%CI
489 [1.74, 2.27]), $F(1, 607.79) = 31.57$, $P < .001$, $R^2 = .017$. suggesting that it was availability of new
490 (pandemic-specific) information rather than temporal distance that contributed to more accurate
491 forecasts in the second compared to the first tournament.

492 Consistency in forecasting

493 Despite variability across scientific teams, domains, and tournaments, the accuracy of
494 scientific predictions was highly systematic. Accuracy in one subset of predictions (ranking of
495 model performance across odd months) was highly correlated with accuracy in the other subset
496 (ranking of model performance across even months), first tournament: multilevel $r_{\text{across domains}}$
497 = .88, 95% CI [.85; .90], $t(357) = 34.80$, $P < .001$; domain-specific $.55 < .rs \leq .99$; second
498 tournament: multilevel $r_{\text{across domains}} = .72$, 95% CI [.67; .75], $t(544) = 23.95$, $P < .001$; domain-
499 specific $.24 < .rs \leq .96$. Further, results of a linear mixed model with MASE scores in
500 Tournament 1, domain, and their interaction predicting MASE in Tournament 2 showed that for
501 eleven out of twelve domains accuracy in Tournament 1 was associated with greater accuracy in
502 Tournament 2, $Md(\text{standardized } \beta) = .26$.

503 Moreover, the ranking of models based on performance in the initial 12-months
504 tournament corresponds to the ranking of the updated models in the follow-up 6-months
505 tournament (Fig. 4). Harder-to-predict domains in the initial tournament remained most
506 inaccurate in the second tournament. Fig. 3 shows one notable exception. Bias against African
507 Americans was easier to predict than other domains in the second tournament. Though
508 speculative, this exception appears consistent with the idea that George Floyd's death catalyzed
509 movements in racial awareness just after the first tournament (Supplementary Fig. 14 for a
510 timeline of major historical events).

511 Which strategies and team characteristics promoted accuracy?

512 Finally, we examined forecasting approaches and individual characteristics of more
513 accurate forecasters in the tournaments. In the main text, we focused on central tendencies across
514 forecasting teams, whereas in the supplementary analyses we reviewed strategies of winning
515 teams and characteristics of the top 5 performers in each domain (Supplementary Figs. 4-11). We
516 compared forecasting approaches relying on 1) no data modeling (but possible consideration of
517 theories); 2) pure data modeling (but no consideration of subject matter theories); 3) hybrid
518 approaches. Roughly half of the teams relied on data-based modeling as a basis for their
519 forecasts, whereas the other half of the teams in each tournament relied only on their intuitions or
520 theoretical considerations (Fig. 5). This pattern was similar across domains (Supplementary Fig.
521 3).

522 In both tournaments, pre-registered linear mixed model analyses with approach as a
523 factor, domain type as a control dummy variable, and MASE scores nested in forecasting teams
524 as a dependent variable revealed that forecasting approaches significantly differed in accuracy,
525 first tournament: $F(2, 149.10) = 5.47$, $P = .005$, $R^2 = .096$; second tournament: $F(2, 177.93) =$
526 5.00 , $P = .008$, $R^2 = .091$ (Fig. 5). Forecasts that considered historical data as part of the forecast
527 modeling were more accurate than models that did not, first tournament: $F(1, 56.29) = 20.38$, P
528 $< .001$, $R^2 = .096$; second tournament: $F(1, 159.11) = 8.12$, $P = .005$, $R^2 = .084$. Model
529 comparison effects were qualified by significant model type X domain interaction; first
530 tournament: $F(11, 278.67) = 4.57$, $P < .001$, $R^2 = .045$; second tournament: $F(11, 462.08) = 3.38$,
531 $P = .0002$, $R^2 = .028$. Post-hoc comparisons in Supplementary Table 4 revealed that data-
532 inclusive (data-driven and hybrid) models were significantly more accurate than data-free
533 models that did not include data in three domains (explicit and implicit racial bias against Asian-
534 Americans and implicit gender-career bias) in Tournament 1 and two domains (life satisfaction,
535 explicit gender-career bias) in Tournament 2. There were no domains where data-free models
536 were more accurate than data-inclusive models. Analyses further demonstrated that, in the first
537 tournament, data-free forecasts of social scientists were not significantly better than lay

538 estimates, $t(577) = 0.87$, $P = .385$, whereas data-inclusive models tended to perform significantly
539 better than lay estimates, $t(470) = 3.11$, $P = .006$, *Cohen's d* = 0.391.

540 To examine the incremental contribution of specific forecasting strategies and team
541 characteristics to accuracy, we pooled data from both tournaments in a linear mixed model with
542 inaccuracy (MASE) as a dependent variable. As Fig. 6 shows, we included predictors
543 representing forecasting strategies, team characteristics, domain expertise (quantified via
544 publications by team members on the topic) and forecasting expertise (quantified via prior
545 experience with forecasting tournaments). We further included domain type as a control dummy
546 variable, and nested responses in teams.

547 The full model fixed effects explained 31% of the variance in accuracy ($R^2 = .314$),
548 though much of it was accounted for by differences in accuracy between domains (non-domain
549 R^2 [partial] = .043). Consistent with prior research²², model sophistication—i.e., considering a
550 larger number of exogenous predictors, COVID-trajectory, or counterfactuals—did not
551 significantly improve accuracy (Fig. 6 and Supplementary Table 5). In fact, forecasting models
552 based on simpler procedures turned out to be significantly more accurate than complex models,
553 $B = 0.14$, $SE = 0.06$, $t(220.82) = 2.33$, $P = .021$, R^2 (partial) = .010.

554 On the one hand, experts' subjective confidence in their forecasts was not related to the
555 accuracy of their estimates. On the other, people with expertise made more accurate forecasts.
556 Teams were more accurate if they had members who published academic research on the
557 forecasted domain, $B = -0.26$, $SE = 0.09$, $t(711.64) = 3.01$, $P = .003$, R^2 (partial) = .007, and who
558 took part in prior forecasting competitions, $B = -0.35$, $SE = 0.17$, $t(56.26) = 2.02$, $P = .049$,
559 R^2 (partial) = .010 (also see Supplementary Table 5). Critically, even though some of these effects
560 were significant, only two factors – complexity of the statistical method and prior experience
561 with forecasting tournaments – showed a non-negligible partial effect size (R^2 above .009).
562 Additional testing whether inclusion of US-based scientists influenced forecasting accuracy did
563 not yield significant effects, $F(1, 106.61) < 1$, *ns*.

564 In the second tournament, we provided teams with the opportunity to compare their
565 original forecasts (Tournament 1, May 2020) to new data at a later point of time and to update
566 their predictions (22) (Tournament 2, Nov 2020). Therefore, we tested whether updating
567 improved people's predictive accuracy. Out of the initial 356 forecasts in the first tournament,
568 180 were updated in the second tournament (from 37% of teams for life satisfaction to 60% of
569 teams for implicit Asian American bias). Updated forecasts in the second tournament
570 (November) were significantly more accurate than the original forecasts in the first tournament
571 (May), $t(94.5) = 6.04$, $P < .001$, *Cohen's d* = 0.804, but so were forecasts from the 34 new teams
572 recruited in November, $t(75.9) = 6.30$, $P < .001$, *Cohen's d* = 0.816. Furthermore, updated
573 forecasts were not significantly different from forecasts provided by new teams recruited in
574 November, $t(77.8) < 0.10$, $P = .928$. This observation suggests that updating did not lead to
575 more accurate forecasts (Supplementary Table 6 reports additional analyses probing different
576 updating rationales).

577 **Discussion**

578 How accurate are social scientists' forecasts of societal change²³? Results from two
579 forecasting tournaments conducted during the first year of the COVID-19 pandemic show that
580 for most domains social scientists' predictions were no better than those from a sample of the
581 (non-specialist) general public. Further, apart from a few domains concerning racial and gender-
582 career bias, scientists' original forecasts were typically not much better than naïve statistical
583 benchmarks derived from historical averages, linear regressions, or random walks. Even when

584 confining the analysis to the top 5 forecasts by social scientists per domain, a simple linear
585 regression produced less error roughly half of the time (Supplementary Figs. 5 and 9).

586 Forecasting accuracy systematically varied across societal domains. In both tournaments,
587 positive sentiment and gender-career stereotypes were easier to forecast than other phenomena,
588 whereas negative sentiment and bias toward African Americans were most difficult to forecast.
589 Domain differences in forecasting accuracy corresponded to historical volatility in the time
590 series. Differences in the complexity of positive and negative affect are well-documented^{25,26}.
591 Moreover, racial attitudes showed more change than attitudes regarding gender during this
592 period (perhaps due to movements like Black Lives Matter).

593 Which strategies and team characteristics were associated with more effective forecasts?
594 One defining feature of more effective forecasters was that they relied on prior data rather than
595 theory alone. This observation fits with prior studies on the performance of algorithmic versus
596 intuitive human judgments²². Social scientists who relied on prior data also performed better than
597 lay crowds and were overrepresented among the winning teams (Supplementary Figs. 4 and 8).

598 Forecasting experience and subject matter expertise on a forecasted topic also
599 incrementally contributed to better performance in tournaments (R^2 [partial] = .010). This is in
600 line with some prior research on the value of subject-matter expertise for geopolitical forecasts⁶
601 and for prediction of success of behavioral science interventions²⁷. Notably, we found that
602 publication track record on a topic, rather than subjective confidence in domain expertise or
603 confidence in the forecast, contributed to greater accuracy. It is possible that subjective
604 confidence in domain expertise conflates expertise and overconfidence²⁸ (versus intellectual
605 humility). There is some evidence that overconfident forecasters are less accurate²⁹. These
606 findings, along with a lack of domain-general effect of social science expertise on performance
607 compared to the general public, invite consideration of whether what usually counts as expertise
608 in social sciences translates into greater ability to predict future real-world trends.

609 The nature of our forecasting tournaments allowed social scientists to self-select any of
610 the twelve forecasting domains, inspect three years of historical trends for each domain, and to
611 update their predictions based on feedback on their initial performance in the first tournament.
612 These features emulated typical forecasting platforms (e.g., metaculus.com). We argue that this
613 approach enhances our ability to draw externally valid and generalizable inferences from a
614 forecasting tournament. However, this approach also resulted in a complex, unbalanced design.
615 Scholars interested in isolating psychological mechanisms fostering superior forecasts may
616 benefit from a simpler design whereby all forecasting teams make forecasts for all requested
617 domains.

618 Another issue in designing forecasting tournaments involves determination of domains
619 one may want participants to forecast. In designing the present tournaments, we provided
620 participants with at least three years of monthly historical data for each forecasting domain. An
621 advantage of making the same historical data available for all forecasters is that it establishes a
622 “common task framework”^{9,16,17}, ensuring main sources of information about the forecasting
623 domains remain identical across all participants. However, this approach restricts the types of
624 social issues participants can forecast. A simpler design without inclusion of historical data
625 would have had the advantage of a greater flexibility in selecting forecasting domains.

626 Why were forecasts of societal change largely inaccurate, even though participants had
627 data-based resources and ample time to deliberate? One possibility concerns self-selection.
628 Perhaps participants in the Forecasting Collaborative were unusually bad at forecasting
629 compared to social scientists as a whole. This possibility seems unlikely. We made efforts to

630 recruit highly successful social scientists at different career stages and from different sub-
631 disciplines (see Supplementary Materials). Indeed, many of our forecasters are well-established
632 scholars. Thus, we do not expect members of the Forecasting Collaborative to be worse at
633 forecasting than other members of the social science community. Nevertheless, albeit
634 impractical, only a random sample of social scientists would have fully addressed the self-
635 selection concern.

636 Second, it is possible that social scientists were not adequately incentivized to perform
637 well in our tournaments. Though we provided reputational incentives by announcing winners and
638 ranking of participating teams, like other big-team science projects^{8,30} we did not provide
639 performance-based monetary incentives³¹, because they may not be key motivating factors for
640 intrinsically motivated social scientists³². Indeed, the drop-out rate between Tournaments 1 and
641 2 was marginal, suggesting that participating teams were motivated to continue being part of the
642 initiative. This reasoning aside, it is possible that stronger incentives for accurate forecasting
643 (whether reputation-based or monetary) could have stimulated some scientists to perform better
644 in our forecasting tournament, opening doors for future directions to address this question
645 directly.

646 Third, social scientists often deal with phenomena with small effect sizes that are
647 overestimated in the literature^{8,30,33}. Additionally, social scientists frequently study social
648 phenomena in conditions that maximize experimental control but may have little external
649 validity, and it is argued that this not only limits the generalizability of findings but in fact
650 reduces their internal validity. In the world beyond the laboratory, where more factors are at
651 play, such effects may be smaller than social scientists might think based on their lab studies, and
652 in fact, such effects may be spurious given the lack of external validity. Thus, social scientist
653 may over(and mis)estimate the impact of the effects they study in the lab on real-world
654 phenomena³⁴

655 Fourth, social scientists tend to theorize about individuals and groups and conduct
656 research at those scales. However, findings from such work may not scale up when predicting
657 phenomena on a scale of entire societies³⁹. Like other dynamical systems in economics, physics,
658 or biology, societal level processes may also be genuinely stochastic rather than deterministic. If
659 so, stochastic models will be hard to outperform.

660 Fifth, training in predictive modeling is not a requirement in many social sciences
661 programs¹⁰. Social scientists often prioritize explanations over formal predictions⁵. For
662 instance, statistical training in social sciences typically emphasizes unbiased estimation of model
663 parameters in the sample over predictive out-of-sample accuracy⁴⁰. Moreover, typical graduate
664 curricula in many areas of social science, such as social or clinical psychology, do not require
665 computational training in predictive modeling. The formal empirical study of societal change is
666 relatively uncommon in these disciplines. Most social scientists approach individual- or group-
667 level phenomena in an a-temporal fashion³⁹. Scientists may favor post-hoc explanations of
668 specific one-time events rather than the future trajectory of social phenomena. Although time is a
669 key theoretical variable for foundational theories in many subfields of social sciences, such as
670 field theory⁴¹, it has remained an elusive concept.

671 Finally, perhaps it is unreasonable to expect theories and models developed during a
672 relatively stable Post World War II period to accurately predict societal trends during a once-in-
673 a-century health crisis. Precisely for this reason, we targeted predictions in domains possessing
674 pandemic-relevant theoretical models (for instance, models about the impact of pathogen spread
675 or social isolation). In this way, we sought to provide a “stress test” of ostensibly relevant

676 theoretical models in a context (pandemic-induced crisis) where change was most likely to be
677 both meaningful and measurable. Nevertheless, the present work suggests that social scientists
678 may not be particularly accurate at forecasting societal trends in this context, though it remains
679 possible that they would perform better during more “normal” times. Considerations above
680 notwithstanding, future work should seek to address this question.

681 How can social scientists become better forecasters? Perhaps the first steps might involve
682 probing the limits of social science theories by evaluating if a given theory is suitable for making
683 societal predictions in the first place or if it is too narrow or too vague^{5,42}. Relatedly, social
684 scientists need to test their theories using representatively designed experiments. Moreover,
685 social scientists may benefit from testing whether a societal trend is deterministic and hence can
686 benefit from theory-driven components, or if it unfolds in a purely stochastic fashion. For
687 instance, one can start by decomposing a time series into the trend, autoregressive, and seasonal
688 components, examining each of them and their meaning for one’s theory and model. One can
689 further perform a unit root test to examine whether the time series is non-stationary. Training in
690 recognizing and modeling properties of time series and dynamical systems may need to become
691 more firmly integrated into graduate curricula in the field. A classic insight in the time series
692 literature is that the mean of the historical time series may be among the best multi-step ahead
693 predictor for a stationary time series⁴³. Using such insights to build predictions from the ground-
694 up can afford greater accuracy. In turn, such training can open the door to more robust models of
695 social phenomena and human behavior, with a promise of greater generalizability in the real-
696 world.

697 Given the broad societal impact of phenomena like prejudice, political polarization or
698 well-being, the ability to accurately predict trends in these variables would appear to be of
699 crucial importance for policy makers and the experts guiding them. But despite common beliefs
700 that social science experts are best equipped to accurately predict these trends compared to non-
701 experts¹, the current findings suggest that social and behavioral scientists have a lot of room for
702 growth³⁸. The good news is that forecasting skills can be improved. Consider the growing
703 accuracy in forecasting models in meteorology in the second part of the 20th century⁴⁴. Greater
704 consideration of representative experimental designs, temporal dynamics, better training in
705 forecasting methods, and more practice with formal forecasting all may improve social
706 scientists’ ability to accurately forecast societal trends going forward.

707 **Methods**

708 The study was approved by the Office of Research Ethics of the University of Waterloo
709 under protocol # 42142.

710 **Pre-registration and deviations.** Forecasts of all participating teams along with their
711 rationales were pre-registered on Open Science Framework (<https://osf.io/6wgbj/registrations>).
712 Additionally, in an a priori specific document shared with the journal in April 2020, we outlined
713 the operationalization of the key dependent variable (MASE), operationalization of covariates
714 and benchmarks (i.e., use of naive forecasting methods), along with the key analytic procedures
715 (linear mixed model and contrasts being different forecasting approaches; osf.io/7ekfm). We did
716 not pre-register the use of a Prolific sample from the general public as an additional benchmark
717 before their forecasting data was collected, though we did pre-register this benchmark in early
718 September of 2020, prior to data pre-processing or analyses. Deviating from the pre-registration,
719 to protect against inflating *p*-values, we performed a single analysis with all covariates in the
720 same model rather than performing separate analyses for each set of covariates. Further, due to
721 scale differences between domains, we chose not to feature analyses concerning absolute

722 percentage errors of each time point in the main paper (but see corresponding analyses on the
723 GitHub site of the project <https://github.com/grossmania/Forecasting-Tournament>, which
724 replicate the key effects presented in the main manuscript).

725 **Participants & recruitment.** We initially aimed for a minimum sample of 40 forecasting
726 teams in our tournament after prescreening to ensure that participants possess at minimum a
727 bachelor's degree in behavioral, social, or computer sciences. To compare groups of scientists
728 employing different forecasting strategies (e.g., data-free versus data-inclusive methods), we
729 subsequently tripled the target size of the final sample ($N = 120$), the target we accomplished by
730 the November phase of the tournament, to ensure sufficient sample for comparison of teams
731 using different strategies (see Supplementary Table 1 for demographics).

732 The Forecasting Collaborative website we used for recruitment
733 (<https://predictions.uwaterloo.ca/faq>) outlined guidelines for eligibility and approach for
734 prospective participants. We incentivized participating teams in two ways. First, prospective
735 participants had an opportunity for a co-authorship in a large-scale citizen science publication.
736 Second, we incentivized accuracy by emphasizing throughout the recruitment that we will be
737 announcing winners and will share the ranking of scientific teams in terms of performance in
738 each tournament (per domain and in total).

739 As outlined in the recruitment materials, we considered data-driven (e.g., model-based) or
740 expertise-based (e.g., general intuition, theory-based) forecasts from any field. As part of the
741 survey, participants selected which method(s) they used to generate their forecasts. Next, they
742 elaborated on how they generated their forecasts in an open-ended question. There are no
743 restrictions, though all teams were encouraged to report their education, as well as areas of
744 knowledge or expertise. Participants were recruited via large scale advertising on social media,
745 mailing lists in the behavioral and social sciences, decision sciences, and data science,
746 advertisement on academic social networks including ResearchGate, and through word of mouth.
747 To ensure broad representation across the academic spectrum of relevant disciplines, we targeted
748 groups of scientists working on computational modeling, social psychology, judgment and
749 decision-making, and data science to join the Forecasting Collaborative.

750 The Forecasting Collaborative started by the end of April 2020, during which time the U.S.
751 Institute for Health Metrics and Evaluation projected the initial peak of the COVID-19 pandemic
752 in the US. The recruitment phase continued until mid-June 2020, to ensure at least 40 teams
753 joined the initial tournament. We were able to recruit 86 teams for the initial 12-month
754 tournament ($M_{\text{age}} = 38.18$; $SD = 8.37$; 73% of forecasts made by scientists with a Doctorate
755 degree), each of which provided forecasts for at least one domain ($M = 4.17$; $SD = 3.78$). At the
756 six-month mark after 2020 US Presidential Election, we provided the initial participants with an
757 opportunity to update their forecasts (44% provided updates), while simultaneously opening the
758 tournament to new participants. This strategy allowed us to compare new forecasts against the
759 updated predictions of the original participants, resulting in 120 teams for this follow-up six-
760 month tournament ($M_{\text{age}} = 36.82$; $SD = 8.30$; 67% of forecasts made by scientists with a
761 Doctorate degree; $M_{\text{forecasted domains}} = 4.55$; $SD = 3.88$). Supplementary analyses showed that
762 updating likelihood did not significantly differ when comparing data-free and data-inclusive
763 models, $z = 0.50$, $P = .618$.

764 **Procedure.** Information for this project was available on the designated website
765 (predictions.uwaterloo.ca), which included objectives, instructions, and prior monthly data for
766 each of the 12 domains they can use for modeling. Researchers who decided to partake in the
767 tournament signed up via a Qualtrics survey, which asked them to upload their estimates for

768 forecasting domains of their choice in a pre-programmed Excel sheet that presented the historical
769 trend and automatically juxtaposed their point estimate forecasts against the historical trend on a
770 plot (see Appendix S1) and answer a set of questions about their rationale and forecasting team
771 composition. Once all data was received, de-identified responses were used to pre-register the
772 forecasted values and models on the Open Science Framework (<https://osf.io/6wgbj/>).

773 At the half-way point (i.e., at six months), participants were provided with a comparison
774 summary of their initial point estimates forecasts vs. actual data for the initial six months.
775 Subsequently, they were provided with an option to update their forecasts, provide a detailed
776 description of the updates, and answer an identical set of questions about their data model and
777 rationale for their forecasts, as well as the consideration of possible exogenous variables and
778 counterfactuals.

779 **Materials**

780 **Forecasting Domains and Data Pre-Processing.** Computational forecasting models
781 require enough prior time series data for reliable modeling. Based on prior recommendations⁴⁵,
782 in the first tournament we provided each team with 39 monthly estimates—from January 2017 to
783 March 2020—for each of the domains participating teams chose to forecast. This approach
784 enabled the teams to perform data-driven forecasting (should teams choose to do so) and to
785 establish a baseline estimate prior to the U.S. peak of the pandemic. In the second tournament,
786 conducted six months later, we provided the forecasting teams with 45 monthly timepoints—
787 from January 2017 to September 2020.

788 Because of the requirement for rich standardized data for computational approaches to
789 forecasting⁹, we limited forecasting domains to issues of broad societal significance. Our
790 domain selection was guided by the discussion of broad social consequences associated with
791 these issues at the beginning of the pandemic^{46,47}, along with general theorizing about
792 psychological and social effects of threats of infectious disease^{48,49}. Additional pragmatic
793 consideration concerning the availability of large-scale longitudinal monthly time series data for
794 a given issue. The resulting domains include affective well-being and life satisfaction, political
795 ideology and polarization, bias in explicit and implicit attitudes towards Asian Americans and
796 African Americans, as well as stereotypes regarding gender and career vs. family. To establish
797 the “common task framework”—a necessary step for the evaluation of predictions in data
798 sciences^{9,17}, we standardized methods for obtaining relevant prior data for each of these
799 domains, made the data publicly available, recruited competitor teams for a common task of
800 inferring predictions from the data, and a priori announced how the project leaders will evaluate
801 accuracy at the end of the tournament.

802 Further, each team had to 1) download and inspect the historical trends (visualized on an
803 Excel plot; example in Appendix); 2) add their forecasts in the same document, which
804 automatically visualized their forecasts against the historical trends; 3) confirm their forecasts; 4)
805 answer prompts concerning their forecasting rationale, their theoretical assumptions, models,
806 conditionals, and consideration of additional parameters in the model. This procedure ensured all
807 teams, at the minimum, considered historical trends, juxtaposed them against their forecasted
808 time series, and deliberated on their forecasting assumptions.

809 *Affective Well-being and Life Satisfaction.* We used monthly Twitter data to estimate
810 markers of affective well-being (positive and negative affect) and life satisfaction over time. We
811 rely on Twitter because no polling data for monthly well-being over the required time period
812 exists, and because prior work suggests that national estimates obtained via social media
813 language can reliably track subjective well-being⁵⁰. For each month, we used previously

814 validated predictive models of well-being, as measured by affective well-being and life
815 satisfaction scales⁵¹. Affective well-being was calculated by applying a custom lexicon⁵² to
816 message unigrams. Life satisfaction was estimated using a ridge regression model trained on
817 *latent Dirichlet allocation* topics, selected using univariate feature selection and dimensionally
818 reduced using randomized principal component analysis, to predict Cantril ladder life satisfaction
819 scores. Such twitter-based estimates closely follow nationally representative polls⁵³. We applied
820 the respective models to Twitter data from January 2017 to March 2020 to obtain estimates of
821 affective well-being and life satisfaction via language on social media.

822 *Ideological Preferences.* We approximated monthly ideological preferences via aggregated
823 weighted data from the Congressional Generic Ballot polls conducted between January 2017 and
824 March 2020 (projects.fivethirtyeight.com/congress-generic-ballot-polls), which ask
825 representative samples of Americans to indicate which party they would support in an election.
826 We weighed polls based on FiveThirtyEight pollster ratings, poll sample size, and poll
827 frequency. FiveThirtyEight pollster ratings are determined by their historical accuracy in
828 forecasting elections since 1998, participation in professional initiatives that seek to increase
829 disclosure and enforce industry best practices and inclusion of live-caller surveys to cellphones
830 and landlines. Based on this data, we then estimated monthly averages for support of Democrat
831 and Republican parties across pollsters (e.g., Marist College, NBC/Wall Street Journal, CNN,
832 YouGov/Economist).

833 *Political Polarization.* We assessed political polarization by examining differences in
834 presidential approval ratings by party identification from Gallup polls
835 (<https://news.gallup.com/poll/203198/presidential-approval-ratings-donald-trump.aspx>). We
836 obtained a difference score in % of Republican versus Democrat approval ratings and estimated
837 monthly averages for the time period of interest. The absolute value of the difference score will
838 ensure possible changes following the 2020 Presidential election will not change the direction of
839 the estimate.

840 *Explicit and Implicit Bias.* Given the natural history of the COVID-19 pandemic, we sought
841 to examine forecasted bias in attitudes towards Asian American (vs. European-Americans). To
842 further probe racial bias, we sought to examine forecasted racial bias in preferences for African
843 American (versus European-American) people. Finally, we sought to examine gender bias in
844 associations of female (vs. male) gender with family versus career. For each domain we sought
845 to obtain both estimates of explicit attitudes⁵⁴ and estimates of implicit attitudes⁵⁵. To this end,
846 we obtained data from the Project Implicit website (<http://implicit.harvard.edu>), which has
847 collected continuous data concerning explicit stereotypes and implicit associations from a
848 heterogeneous pool of volunteers (50,000 - 60,000 unique tests on each of these categories per
849 month). Further details about the website and test materials are publicly available at
850 <https://osf.io/t4bnj>. Recent work suggests that Project Implicit data can provide reliable societal
851 estimates of consequential outcomes^{56,57} and when studying cross-temporal societal shifts in
852 U.S. attitudes⁵⁸. Despite the non-representative nature of the Project Implicit data, recent
853 analyses suggest that bias scores captured by Project Implicit are highly correlated with
854 nationally representative estimates of explicit bias, $r = .75$, indicating that group aggregates of
855 the bias data from Project Implicit can reliably approximate group-level estimates⁵⁷. To further
856 correct possible non-representativeness, we applied stratified weighting to the estimates, as
857 described below.

858 Implicit attitude scores were computed using the revised scoring algorithm of the implicit
859 association test (IAT)⁵⁹. The IAT is a computerized task comparing reaction times to categorize

860 paired concepts (in this case, social groups, e.g., Asian American vs. European American) and
861 attributes (in this case, valence categories, e.g., good vs. bad). Average response latencies in
862 correct categorizations were compared across two paired blocks in which participants
863 categorized concepts and attributes with the same response keys. Faster responses in the paired
864 blocks are assumed to reflect a stronger association between those paired concepts and attributes.
865 Implicit gender-career bias was measured using the IAT with category labels of “male” and
866 “female” and attributes of “career” / “family”). In all tests, positive IAT *D* scores indicate a
867 relative preference for the typically preferred group (European-Americans) or association (men-
868 career).

869 Respondents whose scores fell outside of the conditions specified in the scoring
870 algorithm did not have a complete IAT *D* score and were therefore excluded from analyses.
871 Restricting the analyses to only complete IAT *D* scores resulted in an average retention of 92%
872 of the complete sessions across tests. The sample was further restricted to include only
873 respondents from the United States to increase shared cultural understanding of attitude
874 categories. The sample was restricted to include only respondents with complete demographic
875 information on age, gender, race/ethnicity, and political ideology.

876 For explicit attitude scores, participants provided ratings on feeling thermometers
877 towards Asian-Americans and European Americans (to assess Asian-American bias), and White
878 and Black Americans (to assess racial bias), on a 7-point scale ranging from -3 to +3. Explicit
879 gender-career bias was measured using 7-point Likert-type scales assessing the degree to which
880 an attribute was female/male, from strongly female (-3) to strongly male (+3). Two questions
881 assessed explicit stereotypes for each attribute (e.g., career with female/male, and, separately, the
882 association of family). To match the explicit bias scores with the relative nature of the IAT,
883 relative explicit stereotype scores were created by subtracting the “incongruent” association from
884 the “congruent” association (e.g., [male vs. female-career] - [male vs. female-family]). Thus, for
885 racial bias, -6 reflects a strong explicit preference for the minority over the majority (European-
886 American) group, and +6 reflects a strong explicit preference for the majority over the minority
887 (Asian American / African American) group. Similarly, for gender-career bias, counter-
888 stereotype association (e.g., male-arts/female-science), and +6 reflects a strong stereotypic
889 association (e.g., female-arts/male-science). In both cases, the midpoint of 0 represented equal
890 liking of both groups.

891 We used explicit and implicit bias data for January 2017 – March 2020 and created
892 monthly estimates for each of the explicit and implicit bias domains. Because of possible
893 selection bias among the Project Implicit participants, we adjusted population estimates by
894 weighting the monthly scores based on their representativeness of the demographic frequencies
895 in the U.S. population (age, race, gender, education; estimated biannually by the U.S. Census
896 Bureau; [https://www.census.gov/data/tables/time-series/demo/popest/2010s-national-
897 detail.html](https://www.census.gov/data/tables/time-series/demo/popest/2010s-national-detail.html)). Further, we adjusted weights based on political orientation (1 = “strongly
898 conservative;” 2 = “moderately conservative;” 3 = “slightly conservative;” 4 = “neutral;” 5 =
899 “slightly liberal;” 6 = “moderately liberal;” 7 = “strongly liberal”), using corresponding annual
900 estimates from the General Social Survey. With the weighting values for each participant, we
901 computed weighted monthly means for each attitude test. These procedures ensured that
902 weighted monthly averages approximated the demographics in the U.S. population. We cross-
903 validated this procedure by comparing weighted annual scores to nationally representative
904 estimates for feeling thermometer for African American and Asian American estimates from the
905 American National Election studies in 2017 and 2018.

906 An initial procedure was developed for computing post-stratification weights for African
907 American, Asian and gender career bias (implicit and explicit) to ensure that the sample was
908 representative of the US population at large as much as possible. Originally, we computed
909 weights for the entire year, which were then applied to each month in the year. After receiving
910 feedback from co-authors, a more optimal approach was adopted wherein weights were
911 computed on monthly, as opposed to yearly basis. This was necessary as demographic
912 characteristics varied from month to month each year. This meant that using yearly weights had
913 the potential of amplifying as opposed to reducing bias. Consequently, our new procedure
914 ensured that sample representativeness was maximized. This insight affected forecasts from
915 seven teams who provided them before the change. The teams were informed, and four teams
916 chose to provide updated estimates using newly weighted historical data.

917 For each of these domains, forecasters were provided with 39 monthly estimates in the
918 initial tournament (45 estimates in the follow-up tournament), as well as detailed explanation
919 about the origin and the calculation of respective indices. Thereby, we aim to standardize the
920 data source for the purpose of the forecasting competition⁹. See Supplementary Appendix S1 for
921 example worksheets provided to participants for submissions of their forecasts.

922 **Forecasting Justifications.** For each forecasting model submitted to the tournament,
923 participants provided detailed descriptions. They described the type of model they computed
924 (e.g., time series, game theoretic models, other algorithms), model parameters, additional
925 variables they included in their predictions (e.g., COVID-19 trajectory, presidential election
926 outcome), and underlying assumptions.

927 **Confidence in forecast.** Participants rated their confidence in their forecasted points for
928 each forecast model they submitted on a 7-point scale from 1 (not at all) to 7 (extremely).

929 **Confidence in expertise.** Participants provided ratings of their teams' expertise for a
930 particular domain by indicating extent of agreement with the statement "my team has strong
931 expertise on the research topic of [field]," on a 7-point scale from 1 (Strongly Disagree) to 7
932 (Strongly Agree).

933 **COVID-19 Conditional.** We considered the COVID-19 pandemic as a conditional of
934 interest given links between infectious disease and the target social issues selected for this
935 tournament. In Tournament 1, participants reported if they used the past or predicted trajectory of
936 the COVID-19 pandemic (as measured by number of deaths or prevalence of cases or new
937 infections) as a conditional in their model, and if so, provided their forecasted estimates for the
938 COVID-19 variable included in their model.

939 **Counterfactuals.** Counterfactuals are hypothetical alternative historic events that would
940 be thought to affect the forecast outcomes if they were to occur. Participants described the key
941 counterfactual events between December 2019 and April 2020 that they theorized would have
942 led to different forecasts (e.g., U.S.-wide implementation of social distancing practices in
943 February). Two independent coders evaluated the distinctiveness of counterfactuals (interrater κ
944 = .80). When discrepancies arose, they discussed individual cases with other members of the
945 forecasting collaborative to make the final evaluation. In primary analyses, we focus on the
946 presence of counterfactuals (yes/no).

947 **Team Expertise.** Because expertise can mean many things^{2,60}, we used a telescopic
948 approach and operationalized expertise in four ways of varying granularity. First, we examined
949 broad, domain-general expertise in social sciences by comparing social scientists' forecasts to
950 forecasts provided by the general public without the same training in social science theory and
951 methods. Second, we operationalized the prevalence of graduate training on a team as a more

952 specific marker of domain-general expertise in social sciences. To this end, we asked each
953 participating team to report how many team members have a doctorate degree in social sciences
954 and calculated the percentage of doctorates on a team. Moving to domain-specific expertise, we
955 instructed participating teams to report if any of their members had previously researched or
956 published on the topic of their forecasted variable, operationalizing domain-specific expertise
957 through this measure. Finally, moving to the most subjective level, we asked each participating
958 team to report their subjective confidence in teams' expertise in a given domain (see
959 Supplementary Information)

960 **General Public Benchmark.** In parallel to the tournament with 86 teams, on June 2-3,
961 2020, we recruited a regionally, gender- and socio-economically stratified sample of US
962 residents via the Prolific crowdworker platform (targeted $N = 1,050$ completed responses) and
963 randomly assigned them to forecast societal change for a subset of domains used in the
964 tournaments (a. wellbeing: life satisfaction, positive and negative sentiment on social media; b.
965 politics: political polarization, ideological support for Democrats and Republicans; c. Asian
966 American Bias: explicit and implicit trends; d. African American Bias: explicit and implicit
967 trends; e. Gender-career Bias: explicit and implicit trends). During recruitment, participants were
968 informed that in exchange for 3.65 GDP they have to be able to open and upload forecasts in an
969 Excel worksheet.

970 We considered responses if they provided forecasts for 12 months in at least one domain
971 and if predictions did not exceed the possible range for a given domain (e.g., polarization above
972 100%). Moreover, three coders (intercoder $\kappa = .70$ unweighted, $\kappa = .77$ weighted) reviewed all
973 submitted rationales from lay people and excluded any submissions where participants either
974 misunderstood the task or wrote bogus bot-like responses. Coder disagreements were resolved
975 via a discussion. Finally, we excluded responses if participants spent under 50s making their
976 forecasts, which included reading instructions, downloading the files, providing forecasts, and
977 re-uploading their forecasts (final $N = 802$, 1,467 forecasts; $M_{age} = 30.39$, $SD = 10.56$, 46.36%
978 female; education: 8.57% high school/GED, 28.80% some college, 62.63% college or above;
979 ethnicity: 59.52% white, 17.10% Asian American, 9.45% African American/Black, 7.43%
980 Latinx, 6.50% mixed/other; M_d annual income = \$50,000-\$75,000; residential area: 32.37%
981 urban, 57.03% suburban, 10.60% rural).

982 **Exclusions of the General Public Sample.** Supplementary Table 7 outlines exclusions by
983 category. In the initial step, we considered all submissions via the Qualtrics platform, including
984 partial submissions without any forecasting data ($N = 1,891$). Upon removing incomplete
985 responses without forecasting data, and removing duplicate submissions from the same Prolific
986 IDs, we removed 59 outliers whose data exceeded the range of possible values in a given
987 domain. Subsequently, we removed responses independent coders flagged as either
988 misunderstood ($n = 6$) or bot-like bogus responses ($n = 26$). See Supplementary Appendix S2 for
989 verbatim examples of each screening category and exact coding instructions. Finally, we
990 removed responses where participants took less than 50 seconds to provide their forecasts
991 (including reading instructions, downloading the Excel file, filling it out, re-uploading the Excel
992 worksheet, and completing additional information on their reasoning about the forecast). Finally,
993 one response was removed based on open-ended information where the participant indicated they
994 made forecasts for a different country than the US.

995 **Naïve Statistical Benchmarks.** There is evidence from data science forecasting
996 competitions that the dominant statistical benchmarks are the Theta method, ARIMA, and ETS⁷.
997 Given the socio-cultural context of our study, to avoid loss of generality we decided to employ

998 more traditional benchmarks like naïve/Random walk, historical average, as well as the basic
999 linear regression model—i.e., a method that is used more than anything else in practice and
1000 science. In short, we selected three benchmarks based on their common application in the
1001 forecasting literature (historical mean and random walk are most basic forecasting benchmarks)
1002 or the behavioral / social science literature (linear regression is the most basic statistical approach
1003 to test inferences in sciences). Furthermore, these benchmarks target distinct features of
1004 performance (historical mean speaks to the base rate sensitivity, linear regression speaks to
1005 sensitivity to the overall trend, whereas random walk captures random fluctuations and
1006 sensitivity to dependencies across consecutive time points). Each of these benchmarks may
1007 perform better in some but not in other circumstances. Consequently, to test the limits in
1008 scientists' performance, we examine if social scientists' performance is better than each of the
1009 three benchmarks. To obtain metrics of uncertainty around the naïve statistical estimates, we
1010 chose to simulate these three naïve approaches for making forecasts: 1) random resampling of
1011 historical data; 2) a naïve out-of-sample random walk based on random resampling of historical
1012 *change*; 3) extrapolation from a naïve regression based on a randomly selected interval of
1013 historical data. We describe each approach in the Supplement.

1014 **Analytic Plan**

1015 **Categorization of Forecasts.** We categorized forecasts based on modeling approaches.
1016 Two independent research associates categorize forecasts for each domain based on provided
1017 justifications: i. purely based on (a) theoretical model(s); ii. purely based on data-driven
1018 model(s); iii. a combination of theoretical and data-driven models – i.e., computational model
1019 relies on specific theoretical assumptions. See Appendix S3 for exact coding instructions and
1020 description of the classification (interrater $\kappa = .81$ unweighted, $\kappa = .90$ weighted). We further
1021 examined modelling complexity of approaches that relied on the extrapolation of time series
1022 from the data we provided (e.g., ARIMA, moving average with lags; yes/no; see Appendix S4
1023 for exact coding instructions). Disagreements between coders here (interrater $\kappa = .80$
1024 unweighted, $\kappa = .87$ weighted) and each coding task below were resolved through joint
1025 discussion with the leading author of the project.

1026 **Categorization of Additional Variables.** We tested how the presence and number of
1027 additional variables as parameters in the model impact forecasting accuracy. To this end, we
1028 ensured that additional variables are distinct from one another. Two independent coders
1029 evaluated the distinctiveness of each reported parameter (interrater $\kappa = .56$ unweighted, $\kappa = .83$
1030 weighted).

1031 **Categorization of Teams.** We next categorized teams based on compositions. First, we
1032 counted the number of team members per team. We also sorted teams based on disciplinary
1033 orientation, comparing behavioral and social scientists to teams from computer and data science.
1034 Finally, we used information teams provided concerning their objective and subjective expertise
1035 level for a given subject domain.

1036 **Forecasting Update Justifications.** Given that participants received both new data and a
1037 summary of diverse theoretical positions they can use as a basis for their updates, two
1038 independent research associates scored participants' justifications for forecasting updates on
1039 three dummy-categories: i. new six months of data we provided; ii. new theoretical insights; iii.
1040 consideration of other external events (interrater $\kappa = .63$ unweighted/weighted). See Appendix
1041 S5 for exact coding instructions.

1042 **Statistical analyses.** A priori (<https://osf.io/6wgbj/>) we specified a linear mixed model as
1043 a key analytical procedure, with MASE scores for different domains nested in participating

1044 teams as repeated measures. Prior to analyses, we inspected MASE scores to determine
1045 violations of linearity, which we corrected via log-transformation prior to performing analyses.
1046 All P values refer to two-sided t -tests. For simple effects by domain, we applied Benjamini-
1047 Hochberg false discovery rate corrections. For 95% confidence intervals by domain, we
1048 simulated a multivariate t distribution²⁰ to adjust scores for simultaneous inference of estimates
1049 for 12 domains in each tournament.

1050

1051 ***Data availability***

1052 All data used in the main text and supplementary analysis is accessible on GitHub
1053 (<https://github.com/grossmania/Forecasting-Tournament>). All prior data presented to forecasters
1054 are available on <https://predictions.uwaterloo.ca/>. Historical and ground truth markers are
1055 obtained from Project Five Thirty Eight ([https://projects.fivethirtyeight.com/polls/generic-](https://projects.fivethirtyeight.com/polls/generic-ballot)
1056 [ballot](https://projects.fivethirtyeight.com/polls/generic-ballot)), Gallup ([https://news.gallup.com/poll/203198/presidential-approval-ratings-donald-](https://news.gallup.com/poll/203198/presidential-approval-ratings-donald-trump.aspx)
1057 [trump.aspx](https://news.gallup.com/poll/203198/presidential-approval-ratings-donald-trump.aspx)), Project Implicit (see Open Science Framework website <https://osf.io/t4bnj>), and
1058 U.S. Census Bureau ([https://www.census.gov/data/tables/time-series/demo/popest/2010s-](https://www.census.gov/data/tables/time-series/demo/popest/2010s-national-detail.html)
1059 [national-detail.html](https://www.census.gov/data/tables/time-series/demo/popest/2010s-national-detail.html)).

1060 ***Code availability***

1061 Our project page at <https://github.com/grossmania/Forecasting-Tournament> displays all code
1062 from this paper. See reporting summary for R packages and their versions.

1063

1064 ***Acknowledgments***

1065 This program of research was supported by Basic Research Program at the National Research
1066 University Higher School of Economics (M. Fabrykant), John Templeton Foundation grant 62260
1067 (I.G. and P.T.), Kega 079UK-4/2021 (P.K.), National Center for Complementary & Integrative
1068 Health of the National Institutes of Health under Award Number K23AT010879 (Simon B.
1069 Goldberg), National Science Foundation RAPID Grant 2026854 (M.E.W.V.), PID2019-
1070 111512RB-I00 (M.S.), NPO Systemic Risk Institute (LX22NPO5101) (I.R.), Slovak Research
1071 and Development Agency under contract no. APVV-20-0319 (M.A.), Social Sciences and
1072 Humanities Research Council of Canada Insight Grant 435-2014-0685 (I.G.), Social Sciences and
1073 Humanities Research Council of Canada Connection Grant 611-2020-0190 (I.G.), Swiss National
1074 Science Foundation grant PP00P1_170463 (O. Strijbis). Funders played no role in the
1075 conceptualization, design, analysis, or decision to publish this research. We thank Jordan Axt for
1076 providing monthly estimates of Project Implicit data, and the members of the Forecasting
1077 Collaborative who chose to remain anonymous for their contribution to the tournaments.

1078

1079 ***Author Contributions Statement***

1080 **Conceptualization:** I.G., A.R., C.A.H., M.E.W.V., L.T., and P.E.T.

1081 **Data curation:** I.G., K.S., G.T.S., and O.J.T.

1082 **Forecasting:** S.A., M.K.D., X.E.G., M.J. Hirshberg, M.K.-Y., D.R.M., L.R., A.V., L.W., M.A.,
1083 A.A., P.A., K.B., G.B., F.B., E.B., C.B., M.B., C.K.B., D.T.B., E.M.C., R.C., B.-T.C., W.J.C.,
1084 C.W.C., L.G.C., M. Davis, M.V.D., N.A.D., J.D.D., M. Dziekan, C.T.E., S.E., M. Fabrykant, M.
1085 Firat, G.T.F., J.A.F., J.M.G., S.B.G., A.G., J.G., L.G.-V., S.D.G., S.H., A.H., M.J. Hornsey,
1086 P.D.L.H., A.I., B.J., P.K., Y.J.K., R.K., D.G.L., H.-W.L., N.M.L., V.Y.Q.L., A.W.L., A.L.L.,
1087 C.R.M., M. Maier, N.M.M., D.S.M., A.A.M., M. Misiak, K.O.R.M., J.M.N., J.N., K.N., J.O.,
1088 T.O., M.P.-C., S.P., J.P., Q.R., I.R., R.M.R., Y.R., E.R., L.S., A.S., M.S., A.T.S., O. Simonsson,

1089 M.-C.S., C.-C.T., T.T., B.A.T., D.T., D.C.K.T., J.M.T., L.U., D.V., L.V.W., H.A.V., Q.W.,
 1090 K.W., M.E.W., C.E.W., T.Y., K.Y., S.Y., V.r.A., J.R.A.-H., P.A.B., A.B., L.C., M.C., S.D.-H.,
 1091 Z.E.F., C.R.K., S.T.K., A.L.O., L.M., M.S.M., M.F.R.C.M., E.K.M., P.M., J.B.N., W.N., R.B.R.,
 1092 P.S., A.H.S., O. Strijbis, D.S., E.T., A.v.L., J.G.V., M.N.A.W., and T.W.
 1093 **Formal analysis:** I.G. and C.A.H.
 1094 **Funding acquisition:** I.G.
 1095 **Investigation:** I.G., A.R., and C.A.H.
 1096 **Methodology:** I.G., A.R., C.A.H., K.S., M.E.W.V., S.A., D.R.M., L.R., L.T., A.V., R.N.C.,
 1097 L.U., and D.V.
 1098 **Project administration:** I.G., A.R., M.E.W.V., M.K.-Y., and O.J.T.
 1099 **Resources:** I.G., A.R., J.N., and G.T.S.
 1100 **Supervision:** I.G.
 1101 **Validation:** K.S., X.E.G., and L.W.
 1102 **Visualization:** I.G. and M.K.D.
 1103 **Writing - original draft:** I.G.
 1104 **Writing - review & editing:** I.G., A.R., C.A.H., K.S., M.E.W.V., S.A., M.K.D., X.E.G., M.J.
 1105 Hirshberg, M.K.-Y., D.R.M., L.R., L.T., A.V., L.W., M.A., A.A., P.A., K.B., G.B., F.B., E.B.,
 1106 C.B., M.B., C.K.B., D.T.B., E.M.C., R.C., B.-T.C., W.J.C., R.N.C., C.W.C., L.G.C., M. Davis,
 1107 M.V.D., N.A.D., J.D.D., M. Dziekan, C.T.E., S.E., M. Fabrykant, M. Firat, G.T.F., J.A.F.,
 1108 J.M.G., S.B.G., A.G., J.G., L.G.-V., S.D.G., S.H., A.H., M.J. Hornsey, P.D.L.H., A.I., B.J., P.K.,
 1109 Y.J.K., R.K., D.G.L., H.-W.L., N.M.L., V.Y.Q.L., A.W.L., A.L.L., C.R.M., M. Maier, N.M.M.,
 1110 D.S.M., A.A.M., M. Misiak, K.O.R.M., J.M.N., K.N., J.O., T.O., M.P.-C., S.P., J.P., Q.R., I.R.,
 1111 R.M.R., Y.R., E.R., L.S., A.S., M.S., A.T.S., O. Simonsson, M.-C.S., C.-C.T., T.T., B.A.T.,
 1112 P.E.T., D.T., D.C.K.T., J.M.T., L.V.W., H.A.V., Q.W., K.W., M.E.W., C.E.W., T.Y., K.Y., and
 1113 S.Y.

1114 **Competing interests Statement**

1115 Authors declare that they have no competing interests.

1116

1117 **Tables**

1118

1119 Table 1. Contrasts of Mean-level Inaccuracy (MASE) among Lay Crowds and Social Scientists

1120

<i>Domain</i>	<i>t-ratio</i>	<i>df</i>	<i>p-value</i>	<i>Cohen's d [95% CI]</i>	<i>Bayes Factor</i>	<i>Interpretation</i>
Life Satisfaction	4.321	1725	< .001	0.93 [0.32;1.55]	22.72	Substantial ev. for difference
Explicit Gender-career Bias	3.204	1731	.006	0.90 [0.10; 1.71]	1.37	Some evidence for difference
Implicit Gender-career Bias	3.161	1747	.006	0.88 [0.09; 1.67]	2.49	Some evidence for difference
Political Polarization	2.819	1802	.015	0.71 [-0.01; 1.42]	0.77	Not enough evidence
Positive Affect	2.128	1796	.080	0.54 [-0.18; 1.26]	0.12	Substantial ev. for no difference
Exp. Asian American Bias	1.998	1789	.092	0.53 [-0.23; 1.29]	0.11	Substantial ev. for no difference
Ideology Republicans	1.650	1794	.170	0.40 [-0.29; 1.08]	0.06	Substantial ev. for no difference
Ideology Democrats	1.456	1795	.204	0.35 [-0.34; 1.04]	0.04	Substantial ev. for no difference
Imp. Asian American Bias	1.430	1802	.204	0.36 [-0.36; 1.09]	0.11	Substantial ev. for no difference
Exp. African American Bias	0.939	1747	.218	0.26 [-0.53; 1.05]	0.04	Substantial ev. for no difference
Imp. African American Bias	0.536	1780	.646	0.14 [-0.63; 0.91]	0.02	Substantial ev. for no difference
Negative Affect	-0.271	1796	.787	0.07 [-0.79; 0.65]	0.02	Substantial ev. for no difference

1121 *Note.* Scores > 1: greater accuracy of scientific forecasts. Scores < 1: greater accuracy of lay
 1122 crowds. Pairwise contrasts were obtained via *emmeans* package in R²¹, drawing on the restricted

1123 information maximum likelihood model with group (scientist or naïve crowd), domain, and their
1124 interaction as predictors of the $\log(\text{MASE})$ scores, with responses nested in participants. To
1125 avoid skew, tests are performed on log-transformed scores. Degrees of freedom were obtained
1126 via Kenward-Roger approximation. P -values are adjusted for false discovery rate. $CI =$
1127 Confidence intervals of effect size (*Cohen's d*), which are adjusted for simultaneous inference of
1128 12 domains by simulating a multivariate t distribution²⁰. For Bayesian analyses we relied on
1129 weakly informative priors for our linear mixed model (see Supplement for more detail).
1130 Interpretation of Bayes factor is in the right column. Bayes factors greater than 3 are interpreted
1131 as substantial evidence of a difference, values between 3 and 1 suggest some evidence of a
1132 difference, values between $\frac{1}{3}$ and 1 indicate that there is not enough evidence to interpret, and
1133 values $< \frac{1}{3}$ indicate substantial evidence in favor of the null hypothesis (no difference between
1134 groups).

1135 **Figure Legends/Captions**

1136
1137 Figure 1: Social scientists' average forecasting errors, compared against different benchmarks.
1138 We rank domains from least to most error in Tournament 1, assessing forecasting errors via
1139 mean absolute scaled error (MASE). Estimated means for Scientists and Naïve Crowd indicate
1140 the fixed effect coefficients of a linear mixed model with domain ($k = 12$) and group (in
1141 Tournament 1: $n_{\text{scientists}} = 86$, $n_{\text{naïve crowd}} = 802$; only scientists in Tournament 2: $n = 120$) as a
1142 predictor of forecasting error (MASE) scores nested in teams (Tournament 1 observations:
1143 $n_{\text{scientists}} = 359$, $n_{\text{naïve crowd}} = 1467$; Tournament 2 observations: $n = 546$), using restricted
1144 maximum likelihood estimation. To correct for right skew, we used log-transformed MASE
1145 scores, which are subsequently back-transformed when calculating estimated means and 95%
1146 confidence intervals. In each tournament, confidence intervals are adjusted for simultaneous
1147 inference of estimates for 12 domains in each tournament by simulating a multivariate t
1148 distribution²⁰. Benchmarks represent the naïve crowd and best performing naïve statistical
1149 benchmark (either historic mean, average random walk with an autoregressive lag of one, or
1150 linear regression). Statistical benchmarks were obtained via simulations ($k = 10,000$) with
1151 resampling (see Supplement). Scores to the left of the dotted vertical line show better
1152 performance than naïve in-sample random walk. Scores to the left of the dashed vertical line
1153 show better performance than median performance in M4 tournaments⁷.

1154 Figure 2. Forecasts and ground truth – are forecasts anchoring on the last few historical data
1155 points? Historical time series (40 months before the Tournament 1) and ground truth series (12
1156 months over the Tournament 1), along with forecasts of individual teams (light blue), lowess
1157 curve and 95% confidence interval across social scientists' forecasts (blue), and lowess curve
1158 and 95% confidence interval across naïve crowd's forecasts (salmon). For most domains,
1159 Tournament 1 forecasts of both scientists and the naïve crowd start near the last few historical
1160 data points they received prior to the tournament (January – March 2020). Note, April 2020
1161 forecast was not provided to the participants.

1162
1163 Figure 3. Ratios of Benchmarks against Scientific Forecasts. Scores >1 : greater accuracy of
1164 scientific forecasts. Scores <1 : greater accuracy of naïve benchmarks. Domains are ranked from
1165 least to most error among scientific teams in Tournament 1. Estimated means indicate the fixed
1166 effect coefficient of a linear mixed model with domain ($k = 12$) in each tournament ($n_{\text{Tournament 1}} =$
1167 86 ; $n_{\text{Tournament 2}} = 120$) as a predictor of benchmark-specific ratio scores nested in teams

1168 (observations: $n_{\text{Tournament 1}} = 359$, $n_{\text{Tournament 2}} = 546$), using restricted maximum likelihood
1169 estimation. To correct for right skew, we used square-root or log-transformed MASE scores,
1170 which were subsequently back-transformed when calculating estimated means and 95%
1171 confidence intervals. Confidence intervals are adjusted for simultaneous inference of estimates
1172 for 12 domains in each tournament by simulating a multivariate t distribution²⁰.

1173 Figure 4: Slope graph showing consistency in the ranking of domains in terms of estimated mean
1174 forecasting error across all teams in each tournament, assessed via mean absolute scaled error,
1175 from most to least inaccurate forecasts across both tournaments. Solid line = change in accuracy
1176 between tournaments is statistically significant ($P < .05$); dashed line = non-significant change.
1177 Significance is determined via pairwise comparisons of $\log(\text{MASE})$ scores for each domain,
1178 drawing on the restricted information maximum likelihood model with Tournament (first or
1179 second), domain, and their interaction as predictors of the $\log(\text{MASE})$ scores, with responses
1180 nested in scientific teams ($N_{\text{teams}} = 120$, $N_{\text{observations}} = 905$).

1181
1182 Figure 5: Forecasting errors by prediction approach. Estimated means and 95% confidence
1183 intervals are based on a restricted information maximum likelihood linear mixed effects model
1184 with model type (data-driven, hybrid or intuition/theory-based) as a fixed effects predictor of the
1185 $\log(\text{MASE})$ scores, domain as a fixed effects covariate, and responses nested in participants. We
1186 ran separate models for each tournament (first: $N_{\text{groups}} = 86$; $N_{\text{observations}} = 359$; second: $N_{\text{groups}} =$
1187 120 ; $N_{\text{observations}} = 546$). Scores below the dotted vertical line show better performance than naïve
1188 in-sample random walk. Scores below the dashed vertical line show better performance than
1189 median performance in M4 tournaments⁷.

1190
1191 Figure 6: Contribution of specific forecasting strategies (n parameters, statistical model
1192 complexity, consideration of exogenous events and counterfactuals) and team characteristics for
1193 forecasting accuracy (reversed MASE scores), ranked in terms of magnitude. Scores to the right
1194 of the dashed vertical line contribute positively to accuracy, whereas estimates to the left of the
1195 dashed vertical line contribute negatively. Analyses control for domain type. All continuous
1196 predictors are mean-centered and scaled by 2 standard deviations, to afford comparability²⁴. The
1197 reported standard errors are heteroskedasticity-robust. Thicker bands show 90% confidence
1198 interval, whereas thinner lines show at 95% confidence interval. Effects are statistically
1199 significant if the 95% confidence interval does not include zero (dashed vertical line).

1200
1201
1202

1203 **References**

- 1204 1. Hutcherson, C. *et al.* On the accuracy, media representation, and public perception of
1205 psychological scientists' judgments of societal change. *American Psychologist* (in press).
- 1206 2. Collins, H. & Evans, R. *Rethinking Expertise*. (University of Chicago Press, 2009).
- 1207 3. Fama, E. F. Efficient Capital Markets: A Review of Theory and Empirical Work. *J.*
1208 *Finance* **25**, 383 (1970).
- 1209 4. Tetlock, P. E. *Expert Political Judgement: How Good Is It?* (2005).
- 1210 5. Hofman, J. M. *et al.* Integrating explanation and prediction in computational social
1211 science. *Nature* **595**, 181–188 (2021).
- 1212 6. Mandel, D. R. & Barnes, A. Accuracy of forecasts in strategic intelligence. *Proc. Natl.*
1213 *Acad. Sci.* **111**, 10984–10989 (2014).
- 1214 7. Makridakis, S., Spiliotis, E. & Assimakopoulos, V. The M4 Competition: 100,000 time
1215 series and 61 forecasting methods. *Int. J. Forecast.* **36**, 54–74 (2020).
- 1216 8. Open Science Collaboration. Estimating the reproducibility of psychological science.
1217 *Science*. **349**, aac4716–aac4716 (2015).
- 1218 9. Hofman, J. M., Sharma, A. & Watts, D. J. Prediction and explanation in social systems.
1219 *Science*. **355**, 486–488 (2017).
- 1220 10. Yarkoni, T. & Westfall, J. Choosing prediction over explanation in psychology: Lessons
1221 from machine learning. *Perspect. Psychol. Sci.* **12**, 1100–1122 (2017).
- 1222 11. Fincher, C. L. & Thornhill, R. Parasite-stress promotes in-group assortative sociality: The
1223 cases of strong family ties and heightened religiosity. *Behav. Brain Sci.* **35**, 61–79 (2012).
- 1224 12. Varnum, M. E. W. & Grossmann, I. Pathogen prevalence is associated with cultural
1225 changes in gender equality. *Nat. Hum. Behav.* **1**, 3 (2016).
- 1226 13. Schaller, M. & Murray, D. R. Pathogens, personality, and culture: Disease prevalence
1227 predicts worldwide variability in sociosexuality, extraversion, and openness to experience.
1228 *J. Pers. Soc. Psychol.* **95**, 212–221 (2008).
- 1229 14. van Leeuwen, F., Park, J. H., Koenig, B. L. & Graham, J. Regional variation in pathogen
1230 prevalence predicts endorsement of group-focused moral concerns. *Evol. Hum. Behav.* **33**,
1231 429–437 (2012).
- 1232 15. Hawkey, L. C. & Cacioppo, J. T. Loneliness Matters: A Theoretical and Empirical
1233 Review of Consequences and Mechanisms. *Ann. Behav. Med.* **40**, 218–227 (2010).
- 1234 16. Salganik, M. J. *et al.* Measuring the predictability of life outcomes with a scientific mass
1235 collaboration. *Proc. Natl. Acad. Sci.* **117**, 8398–8403 (2020).
- 1236 17. Liberman, M. Reproducible Research and the Common Task Method. (2015).
- 1237 18. Hyndman, R. J. & Koehler, A. B. Another look at measures of forecast accuracy. *Int. J.*
1238 *Forecast.* **22**, 679–688 (2006).
- 1239 19. Eyal, P., David, R., Andrew, G., Zak, E. & Ekaterina, D. Data quality of platforms and
1240 panels for online behavioral research. *Behav. Res. Methods* 1–20 (2021).

- 1241 doi:10.3758/s13428-021-01694-3
- 1242 20. Genz, A. & Bretz, F. *Computation of Multivariate Normal and t Probabilities*. (Springer-
1243 Verlag, 2009).
- 1244 21. Lenth, R., Singmann, H., Love, J. & Maxime, H. emmeans: Estimated Marginal Means,
1245 aka Least-Squares Means. (2020).
- 1246 22. Green, K. C. & Armstrong, J. S. Simple versus complex forecasting: The evidence. *J. Bus.*
1247 *Res.* **68**, 1678–1685 (2015).
- 1248 23. Grossmann, I., Twardus, O., Varnum, M. E. W., Jayawickreme, E. & McLevey, J. Expert
1249 predictions of societal change: Insights from the world after COVID project. *American*
1250 *Psychologist* **77**, 276–290 (2022).
- 1251 24. Gelman, A. Scaling regression inputs by dividing by two standard deviations. *Stat. Med.*
1252 **27**, 2865–2873 (2008).
- 1253 25. Grossmann, I., Huynh, A. C. & Ellsworth, P. C. Emotional complexity: Clarifying
1254 definitions and cultural correlates. *J. Pers. Soc. Psychol.* **111**, 895–916 (2016).
- 1255 26. Alves, H., Koch, A. & Unkelbach, C. Why Good Is More Alike Than Bad: Processing
1256 Implications. *Trends Cogn. Sci.* **21**, 69–79 (2017).
- 1257 27. Dimant, E. *et al.* Politicizing mask-wearing: predicting the success of behavioral
1258 interventions among republicans and democrats in the U.S. *Sci. Rep.* **12**, 7575 (2022).
- 1259 28. Dunning, D., Heath, C. & Suls, J. M. Flawed Self-Assessment. *Psychol. Sci. Public*
1260 *Interes.* **5**, 69–106 (2004).
- 1261 29. Mellers, B., Tetlock, P. E. & Arkes, H. R. Forecasting tournaments, epistemic humility
1262 and attitude depolarization. *Cognition* **188**, 19–26 (2019).
- 1263 30. Klein, R. A. *et al.* Many Labs 2: Investigating Variation in Replicability Across Samples
1264 and Settings. *Adv. Methods Pract. Psychol. Sci.* **1**, 443–490 (2018).
- 1265 31. Voslinsky, A. & Azar, O. H. Incentives in experimental economics. *J. Behav. Exp. Econ.*
1266 **93**, 101706 (2021).
- 1267 32. Cerasoli, C. P., Nicklin, J. M. & Ford, M. T. Intrinsic motivation and extrinsic incentives
1268 jointly predict performance: A 40-year meta-analysis. *Psychol. Bull.* **140**, 980–1008
1269 (2014).
- 1270 33. Richard, F. D., Bond Jr., C. F. & Stokes-Zoota, J. J. One Hundred Years of Social
1271 Psychology Quantitatively Described. *Rev. Gen. Psychol.* **7**, 331–363 (2003).
- 1272 34. Cesario, J. What Can Experimental Studies of Bias Tell Us About Real-World Group
1273 Disparities? *Behav. Brain Sci.* 1–80 (2021). doi:10.1017/S0140525X21000017
- 1274 35. Henrich, J., Heine, S. J. & Norenzayan, A. The weirdest people in the world? *Behav.*
1275 *Brain Sci.* **33**, 61–83 (2010).
- 1276 36. Yarkoni, T. The generalizability crisis. *Behav. Brain Sci.* **45**, e1 (2022).
- 1277 37. IJzerman, H. *et al.* Use caution when applying behavioural science to policy. *Nat. Hum.*
1278 *Behav.* **4**, 1092–1094 (2020).

- 1279 38. Makridakis, S. & Taleb, N. Living in a world of low levels of predictability. *Int. J.*
1280 *Forecast.* **25**, 840–844 (2009).
- 1281 39. Varnum, M. E. W. & Grossmann, I. Cultural Change: The How and the Why. *Perspect.*
1282 *Psychol. Sci.* **12**, 956–972 (2017).
- 1283 40. Breiman, L. Statistical Modeling: The Two Cultures (with comments and a rejoinder by
1284 the author). *Stat. Sci.* **16**, (2001).
- 1285 41. Lewin, K. Defining the ‘field at a given time.’ *Psychol. Rev.* **50**, 292–310 (1943).
- 1286 42. Turchin, P., Currie, T. E., Turner, E. A. L. & Gavrillets, S. War, space, and the evolution
1287 of Old World complex societies. *Proc. Natl. Acad. Sci.* **110**, 16384–16389 (2013).
- 1288 43. Brockwell, P. J. & Davis, R. A. *Introduction to Time Series and Forecasting.* (Springer
1289 International Publishing, 2016). doi:10.1007/978-3-319-29854-2
- 1290 44. Hitchens, N. M., Brooks, H. E. & Kay, M. P. Objective Limits on Forecasting Skill of
1291 Rare Events. *Weather Forecast.* **28**, 525–534 (2013).
- 1292 45. Jebb, A. T., Tay, L., Wang, W. & Huang, Q. Time series analysis for psychological
1293 research: examining and forecasting change. *Front. Psychol.* **6**, (2015).
- 1294 46. Van Bavel, J. *et al.* Using social and behavioural science to support COVID-19 pandemic
1295 response. *Nat. Hum. Behav.* **4**, 460–471 (2020).
- 1296 47. Seitz, B. M. *et al.* The pandemic exposes human nature: 10 evolutionary insights. *Proc.*
1297 *Natl. Acad. Sci.* **117**, 27767–27776 (2020).
- 1298 48. Schaller, M. & Park, J. H. The behavioral immune system (and why it matters). *Curr. Dir.*
1299 *Psychol. Sci.* (2011). doi:10.1177/0963721411402596
- 1300 49. Wang, I. M., Michalak, N. M. & Ackerman, J. M. Threat of Infectious Disease. in *The*
1301 *SAGE Handbook of Personality and Individual Differences: Volume II: Origins of*
1302 *Personality and Individual Differences* (eds. Zeigler-Hill, V., Shackelford, T. K., Zeigler-
1303 Hill, V. & Shackelford, T. K.) 321–345 (2018). doi:10.4135/9781526451200.n18
- 1304 50. Luhmann, M. Using Big Data to study subjective well-being. *Curr. Opin. Behav. Sci.* **18**,
1305 28–33 (2017).
- 1306 51. Schwartz, H. A. *et al.* Predicting individual well-being through the language of social
1307 media. in *Biocomputing 2016* 516–527 (2016). doi:10.1142/9789814749411_0047
- 1308 52. Kiritchenko, S., Zhu, X. & Mohammad, S. M. Sentiment analysis of short informal texts.
1309 *J. Artif. Intell. Res.* **50**, 723–762 (2014).
- 1310 53. Witters, D. & Harter, J. *In U.S., Life Ratings Plummet to 12-Year Low.* (2020).
- 1311 54. Axt, J. R. The Best Way to Measure Explicit Racial Attitudes Is to Ask About Them. *Soc.*
1312 *Psychol. Personal. Sci.* **9**, 896–906 (2018).
- 1313 55. Nosek, B. A. *et al.* Pervasiveness and correlates of implicit attitudes and stereotypes. *Eur.*
1314 *Rev. Soc. Psychol.* **18**, 36–88 (2007).
- 1315 56. Hehman, E., Flake, J. K. & Calanchini, J. Disproportionate Use of Lethal Force in
1316 Policing Is Associated with Regional Racial Biases of Residents. *Soc. Psychol. Personal.*
1317 *Sci.* **9**, 393–401 (2018).

- 1318 57. Ofosu, E. K., Chambers, M. K., Chen, J. M. & Hehman, E. Same-sex marriage
 1319 legalization associated with reduced implicit and explicit antigay bias. *Proc. Natl. Acad.*
 1320 *Sci.* **116**, 8846–8851 (2019).
- 1321 58. Charlesworth, T. E. S. & Banaji, M. R. Patterns of Implicit and Explicit Attitudes: I.
 1322 Long-Term Change and Stability From 2007 to 2016. *Psychol. Sci.* **30**, 174–192 (2019).
- 1323 59. Greenwald, A. G., Nosek, B. A. & Banaji, M. R. Understanding and using the Implicit
 1324 Association Test: I. An improved scoring algorithm. *J. Pers. Soc. Psychol.* **85**, 197–216
 1325 (2003).
- 1326 60. Gobet, F. The Future of Expertise: The Need for a Multidisciplinary Approach. **1**, 7
 1327 (2018).

1328

1329 **The Forecasting Collaborative**

1330 Igor Grossmann^{1*}, Amanda Rotella^{2,1}, Cendri A. Hutcherson³, Konstantyn Sharpinskyi¹, Michael
 1331 E. W. Varnum⁴, Sebastian Achter⁵, Mandeep K. Dharmi⁶, Xinqi Evie Guo⁷, Mane Kara-
 1332 Yakoubian⁸, David R. Mandel^{9,10}, Louis Raes¹¹, Louis Tay¹², Aymeric Vie^{13,14}, Lisa Wagner¹⁵,
 1333 Matus Adamkovic^{16,17}, Arash Arami¹⁸, Patricia Arriaga¹⁹, Kasun Bandara²⁰, Gabriel Baník¹⁶,
 1334 František Bartoš²¹, Ernest Baskin²², Christoph Bergmeir²³, Michał Białek²⁴, Caroline K.
 1335 Børsting²⁵, Dillon T. Browne¹, Eugene M. Caruso²⁶, Rong Chen²⁷, Bin-Tzong Chie²⁸, William J.
 1336 Chopik²⁹, Robert N. Collins⁹, Chin W. Cong³⁰, Lucian G. Conway³¹, Matthew Davis³², Martin V.
 1337 Day³³, Nathan A. Dhaliwal³⁴, Justin D. Durham³⁵, Martyna Dziekan³⁶, Christian T. Elbaek²⁵,
 1338 Eric Shuman³⁷, Marharyta Fabrykant^{38,39}, Mustafa Firat⁴⁰, Geoffrey T. Fong^{1,41}, Jeremy A.
 1339 Frimer⁴², Jonathan M. Gallegos⁴³, Simon B. Goldberg⁴⁴, Anton Gollwitzer^{45,46}, Julia Goyal⁴⁷,
 1340 Lorenz Graf-Vlachy^{48,49}, Scott D. Gronlund³⁵, Sebastian Hafenbrädl⁵⁰, Andree Hartanto⁵¹,
 1341 Matthew J. Hirshberg⁵², Matthew J. Hornsey⁵³, Piers D. L. Howe⁵⁴, Anoosha Izadi⁵⁵, Bastian
 1342 Jaeger⁵⁶, Pavol Kačmár⁵⁷, Yeun Joon Kim⁵⁸, Ruslan Krenzler^{59,60}, Daniel G. Lannin⁶¹, Hung-
 1343 Wen Lin⁶², Nigel Mantou Lou^{63,64}, Verity Y. Q. Lua⁵¹, Aaron W. Lukaszewski^{65,66}, Albert L.
 1344 Ly⁶⁷, Christopher R. Madan⁶⁸, Maximilian Maier⁶⁹, Nadyanna M. Majeed⁷⁰, David S. March⁷¹,
 1345 Abigail A. Marsh⁷², Michal Misiak^{24,73}, Kristian Ove R. Myrseth⁷⁴, Jaime M. Napan⁶⁷, Jonathan
 1346 Nicholas⁷⁵, Konstantinos Nikolopoulos⁷⁶, Jiaqing O⁷⁷, Tobias Otterbring^{78,79}, Mariola Paruzel-
 1347 Czachura^{80,81}, Shiva Pauer²¹, John Protzko⁸², Quentin Raffaelli⁸³, Ivan Ropovik^{84,85}, Robert M.
 1348 Ross⁸⁶, Yefim Roth⁸⁷, Espen Røysamb⁸⁸, Landon Schnabel⁸⁹, Astrid Schütz⁹⁰, Matthias Seifert⁹¹,
 1349 A. T. Sevincer⁹², Garrick T. Sherman⁹³, Otto Simonsson^{94,95}, Ming-Chien Sung⁹⁶, Chung-Ching
 1350 Tai⁹⁶, Thomas Talhelm⁹⁷, Bethany A. Teachman⁹⁸, Philip E. Tetlock^{99,100}, Dimitrios
 1351 Thomakos¹⁰¹, Dwight C. K. Tse¹⁰², Oliver J. Twardus¹⁰³, Joshua M. Tybur⁵⁶, Lyle Ungar⁹³, Daan
 1352 Vandermeulen¹⁰⁴, Leighton Vaughan Williams¹⁰⁵, Hrag A. Vosgerichian¹⁰⁶, Qi Wang¹⁰⁷, Ke
 1353 Wang¹⁰⁸, Mark E. Whiting^{109,110}, Conny E. Wollbrant¹¹¹, Tao Yang¹¹², Kumar Yogeeswaran¹¹³,
 1354 Sangsuk Yoon¹¹⁴, Ventura r. Alves¹¹⁵, Jessica R. Andrews-Hanna^{83,116}, Paul A. Bloom⁷⁵,
 1355 Anthony Boyles¹¹⁷, Loo Charis¹¹⁸, Mingyeong Choi¹¹⁹, Sean Darling-Hammond¹²⁰, Zoe E.
 1356 Ferguson¹²¹, Cheryl R. Kaiser⁴³, Simon T. Karg¹²², Alberto López Ortega¹²³, Lori Mahoney¹²⁴,
 1357 Melvin S. Marsh¹²⁵, Marcellin F. R. C. Martinie⁵⁴, Eli K. Michaels¹²⁶, Philip Millroth¹²⁷, Jeanean
 1358 B. Naqvi¹²⁸, Weiting Ng¹²⁹, Robb B. Rutledge¹³⁰, Peter Slattery¹³¹, Adam H. Smiley⁴³, Oliver
 1359 Strijbis¹³², Daniel Sznycer¹³³, Eli Tsukayama¹³⁴, Austin van Loon¹³⁵, Jan G. Voelkel¹³⁵, Margaux
 1360 N. A. Wienk⁷⁵, Tom Wilkening¹³⁶

- 1361 ¹Department of Psychology, University of Waterloo; Waterloo, Canada.
- 1362 ²Department of Psychology, Northumbria University, UK.
- 1363 ³Department of Psychology, University of Toronto Scarborough; Toronto, Canada.
- 1364 ⁴Department of Psychology, Arizona State University, Tempe, USA.
- 1365 ⁵Institute of Management Accounting and Simulation, Hamburg University of Technology;
1366 Hamburg, Germany.
- 1367 ⁶Department of Psychology, Middlesex University London; London, UK.
- 1368 ⁷Department of Experimental Psychology, University of California San Diego; San Diego, USA.
- 1369 ⁸Department of Psychology, Toronto Metropolitan University; Toronto, Canada.
- 1370 ⁹Defence Research and Development Canada; Ottawa, Canada.
- 1371 ¹⁰Department of Psychology, York University; Toronto, Canada.
- 1372 ¹¹Department of Economics, Tilburg University, Tilburg, Netherlands.
- 1373 ¹²Department of Psychological Sciences, Purdue University; West Lafayette, USA.
- 1374 ¹³Mathematical Institute, University of Oxford; Oxford, UK.
- 1375 ¹⁴Institute of New Economic Thinking, University of Oxford, Oxford, UK.
- 1376 ¹⁵Jacobs Center for Productive Youth Development, University of Zurich; Zurich, Switzerland.
- 1377 ¹⁶Institute of Psychology, University of Prešov; Prešov, Slovakia.
- 1378 ¹⁷Institute of Social Sciences, CSPS, Slovak Academy of Sciences; Bratislava, Slovakia.
- 1379 ¹⁸Department of Mechanical and Mechatronics Engineering, University of Waterloo; Waterloo,
1380 Canada.
- 1381 ¹⁹Iscte-University Institute of Lisbon, CIS; Lisbon, Portugal.
- 1382 ²⁰Melbourne Centre for Data Science, University of Melbourne; Melbourne, Australia.
- 1383 ²¹Faculty of Social and Behavioural Sciences, University of Amsterdam; Amsterdam,
1384 Netherlands.
- 1385 ²²Department of Food Marketing, Haub School of Business, Saint Joseph's University;
1386 Philadelphia, USA.
- 1387 ²³Department of Data Science and Artificial Intelligence, Monash University; Melbourne,
1388 Australia.
- 1389 ²⁴Institute of Psychology, University of Wrocław; Wrocław, Poland.
- 1390 ²⁵Department of Management, Aarhus University; Aarhus, Denmark.
- 1391 ²⁶Anderson School of Management, UCLA; Los Angeles, USA.
- 1392 ²⁷Department of Psychology, Dominican University of California; San Rafael, USA.
- 1393 ²⁸Department of Industrial Economics, Tamkang University; New Taipei City, Taiwan.
- 1394 ²⁹Department of Psychology, Michigan State University; East Lansing, USA.

- 1395 ³⁰Department of Psychology and Counselling, Universiti Tunku Abdul Rahman; Kampar,
1396 Malaysia.
- 1397 ³¹Psychology Department, Grove City College; Grove City, USA.
- 1398 ³²Department of Economics, Siena College; Loudonville, USA.
- 1399 ³³Department of Psychology, Memorial University of Newfoundland; St. John's, Canada.
- 1400 ³⁴UBC Sauder School of Business, University of British Columbia; Vancouver, Canada.
- 1401 ³⁵Department of Psychology, University of Oklahoma; Norman, USA.
- 1402 ³⁶Faculty of Psychology and Cognitive Science, Adam Mickiewicz University; Poznań, Poland.
- 1403 ³⁷Department of Psychology, University of Groningen; Groningen, Netherlands.
- 1404 ³⁸Laboratory for Comparative Studies in Mass Consciousness, Expert Institute, HSE University;
1405 Moscow, Russia.
- 1406 ³⁹Faculty of Philosophy and Social Sciences, Belarusian State University; Minsk, Belarus.
- 1407 ⁴⁰Department of Sociology, Radboud University; Nijmegen, Netherlands.
- 1408 ⁴¹Ontario Institute for Cancer Research; Toronto, Canada.
- 1409 ⁴²Department of Psychology, University of Winnipeg; Winnipeg, Canada.
- 1410 ⁴³Department of Psychology, University of Washington; Seattle, USA.
- 1411 ⁴⁴Department of Counseling Psychology, University of Wisconsin - Madison; Madison, USA.
- 1412 ⁴⁵Department of Leadership and Organizational Behaviour, BI Norwegian Business School; Oslo,
1413 Norway.
- 1414 ⁴⁶Center for Adaptive Rationality, Max Planck Institute for Human Development; Berlin,
1415 Germany.
- 1416 ⁴⁷School of Public Health Sciences, University of Waterloo; Waterloo, Canada.
- 1417 ⁴⁸TU Dortmund University; Dortmund, Germany.
- 1418 ⁴⁹ESCP Business School; Berlin, Germany.
- 1419 ⁵⁰IESE Business School; Barcelona, Spain.
- 1420 ⁵¹School of Social Sciences, Singapore Management University; Singapore.
- 1421 ⁵²Center for Healthy Minds, University of Wisconsin-Madison; Madison, USA.
- 1422 ⁵³University of Queensland Business School; Brisbane, Australia.
- 1423 ⁵⁴Melbourne School of Psychological Sciences, University of Melbourne; Melbourne, Australia.
- 1424 ⁵⁵Department of Marketing, University of Massachusetts Dartmouth; Dartmouth, USA.
- 1425 ⁵⁶Faculty of Behavioural and Movement Sciences, Vrije Universiteit Amsterdam; Amsterdam,
1426 Netherlands.
- 1427 ⁵⁷Department of Psychology, Faculty of Arts, Pavol Jozef Šafárik University; Košice, Slovakia.
- 1428 ⁵⁸Cambridge Judge Business School, University of Cambridge; Cambridge, UK.

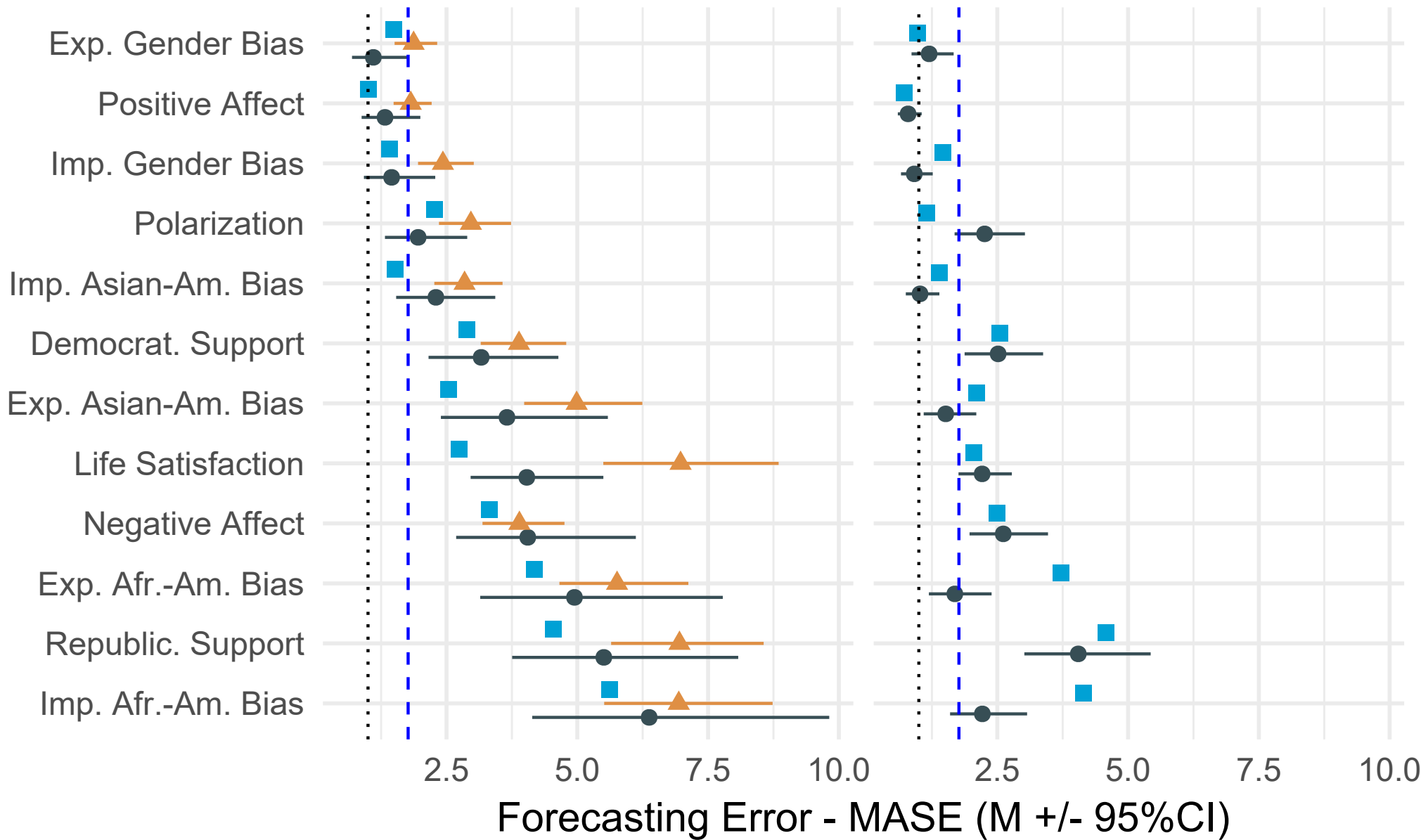
- 1429 ⁵⁹Hermes Germany GmbH; Hamburg, Germany.
- 1430 ⁶⁰University of Hamburg; Hamburg, Germany.
- 1431 ⁶¹Department of Psychology, Illinois State University; Normal, USA.
- 1432 ⁶²Department of Marketing and Logistics Management, National Penghu University of Science
1433 and Technology; Magong City, Taiwan.
- 1434 ⁶³Department of Psychology, University of Victoria; Victoria, Canada.
- 1435 ⁶⁴Centre for Youth and Society, University of Victoria; Victoria, Canada.
- 1436 ⁶⁵Department of Psychology, California State University - Fullerton; Fullerton, USA.
- 1437 ⁶⁶Center for the Study of Human Nature, California State University - Fullerton; Fullerton, USA.
- 1438 ⁶⁷Department of Psychology, Loma Linda University; Loma Linda, USA.
- 1439 ⁶⁸University of Nottingham; Nottingham, UK.
- 1440 ⁶⁹Department of Experimental Psychology, University College London; London, UK.
- 1441 ⁷⁰Singapore Management University; Singapore.
- 1442 ⁷¹Department of Psychology, Florida State University; Tallahassee, USA.
- 1443 ⁷²Department of Psychology, Georgetown University; Washington DC, USA.
- 1444 ⁷³School of Anthropology & Museum Ethnography, University of Oxford; Oxford, UK.
- 1445 ⁷⁴School for Business and Society, University of York; York, UK.
- 1446 ⁷⁵Department of Psychology, Columbia University; New York City, USA.
- 1447 ⁷⁶IHRR Forecasting Laboratory, Durham University; Durham, UK.
- 1448 ⁷⁷Department of Psychology, Aberystwyth University; Aberystwyth, UK.
- 1449 ⁷⁸School of Business and Law, Department of Management, University of Agder; Kristiansand,
1450 Norway.
- 1451 ⁷⁹Institute of Retail Economics; Stockholm, Sweden.
- 1452 ⁸⁰Institute of Psychology, University of Silesia in Katowice, Poland.
- 1453 ⁸¹ Department of Neurology, Penn Center for Neuroaesthetics, University of Pennsylvania, United
1454 States.
- 1455 ⁸²Central Connecticut State University; New Britain, USA.
- 1456 ⁸³Department of Psychology, University of Arizona; Tucson, USA.
- 1457 ⁸⁴Faculty of Education, Institute for Research and Development of Education, Charles University;
1458 Prague, Czech Republic.
- 1459 ⁸⁵Faculty of Education; University of Prešov; Prešov, Slovakia.
- 1460 ⁸⁶School of Psychology, Macquarie University; Sydney, Australia.
- 1461 ⁸⁷Department of Human Service, University of Haifa; Haifa, Israel.
- 1462 ⁸⁸Promenta Center, Department of Psychology, University of Oslo; Oslo, Norway.

1463 ⁸⁹Department of Sociology, Cornell University; Ithaca, USA.
1464 ⁹⁰Institute of Psychology, University of Bamberg; Bamberg, Germany.
1465 ⁹¹IE Business School, IE University; Madrid, Spain.
1466 ⁹²Faculty of Psychology and Human Movement Science, University of Hamburg; Hamburg,
1467 Germany.
1468 ⁹³Department of Computer and Information Science, University of Pennsylvania; Philadelphia,
1469 USA.
1470 ⁹⁴Department of Clinical Neuroscience, Karolinska Institutet; Solna, Sweden.
1471 ⁹⁵Department of Sociology, University of Oxford; Oxford, UK.
1472 ⁹⁶Department of Decision Analytics and Risk, University of Southampton; Southampton, UK.
1473 ⁹⁷University of Chicago Booth School of Business; Chicago, USA.
1474 ⁹⁸Department of Psychology, University of Virginia; Charlottesville, USA.
1475 ⁹⁹Psychology Department, University of Pennsylvania; Philadelphia, USA.
1476 ¹⁰⁰Wharton School of Business, University of Pennsylvania; Philadelphia, USA.
1477 ¹⁰¹Department of Economics, National and Kapodistrian University of Athens; Athens, Greece.
1478 ¹⁰²School of Psychological Sciences and Health, University of Strathclyde; Glasgow, UK.
1479 ¹⁰³Department of Psychology, University of Guelph; Guelph, Canada.
1480 ¹⁰⁴Psychology Department, Hebrew University of Jerusalem; Jerusalem, Israel.
1481 ¹⁰⁵Department of Economics, Nottingham Trent University; Nottingham, UK.
1482 ¹⁰⁶Department of Management and Organizations, Northwestern University; Evanston, USA.
1483 ¹⁰⁷College of Human Ecology, Cornell University, Ithaca, USA.
1484 ¹⁰⁸Harvard Kennedy School, Harvard University; Cambridge, USA.
1485 ¹⁰⁹Computer and Information Science, University of Pennsylvania; Philadelphia, USA.
1486 ¹¹⁰Operations, Information, and Decisions Department, The Wharton School, University of
1487 Pennsylvania; Philadelphia, USA.
1488 ¹¹¹School of Economics and Finance, University of St. Andrews; St. Andrews, UK.
1489 ¹¹²Department of Management, Cameron School of Business, University of North Carolina;
1490 Wilmington, USA.
1491 ¹¹³School of Psychology, Speech and Hearing, University of Canterbury; Christchurch, New
1492 Zealand, ¹¹⁴Department of Marketing, University of Dayton; Dayton, USA.
1493 ¹¹⁵ISG Universidade Lusofona; Lisbon, Portugal.
1494 ¹¹⁶Cognitive Science, University of Arizona; Tucson, USA.
1495 ¹¹⁷Ephemer AI; Atlanta, USA.
1496 ¹¹⁸Questrom School of Business, Boston University; Boston, USA.

1497 ¹¹⁹Institute of Social Science Research, Pusan National University; Busan, South
1498 Korea, ¹²⁰Fielding School of Public Health, UCLA; Los Angeles, USA.
1499 ¹²¹Psychology Department, University of Washington; Seattle, USA.
1500 ¹²²Department of Political Science, Aarhus University; Aarhus, USA.
1501 ¹²³Faculty of Social Sciences, Vrije Universiteit Amsterdam; Amsterdam, Netherlands.
1502 ¹²⁴College of Science and Mathematics, Wright State University; Fairborn, USA.
1503 ¹²⁵Department of Psychology, Georgia Southern University; Statesboro, USA.
1504 ¹²⁶Division of Epidemiology, School of Public Health, University of California, Berkeley;
1505 Berkeley, USA.
1506 ¹²⁷Department of Psychology, Uppsala University; Uppsala, Sweden.
1507 ¹²⁸Department of Psychology, Carnegie Mellon University; Pittsburgh, USA.
1508 ¹²⁹School of Humanities & Behavioral Sciences, Singapore University of Social Sciences;
1509 Singapore.
1510 ¹³⁰Department of Psychology, Yale University; New Haven, USA.
1511 ¹³¹BehaviourWorks Australia. Monash University; Melbourne, Australia.
1512 ¹³²Institute of Political Science, University of Zurich; Zurich, Switzerland.
1513 ¹³³Department of Psychology, Oklahoma State University; Stillwater, USA.
1514 ¹³⁴Department of Business Administration, University of Hawaii - West Oahu; Kapolei, USA.
1515 ¹³⁵Department of Sociology, Stanford University; Stanford, USA.
1516 ¹³⁶Department of Economics, University of Melbourne; Melbourne, Australia.

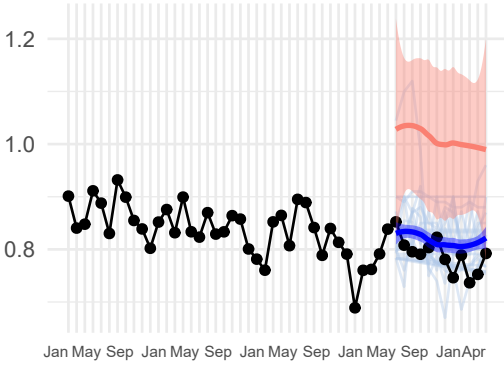
First Tournament (May 2020)

Second Tournament (Nov 2020)

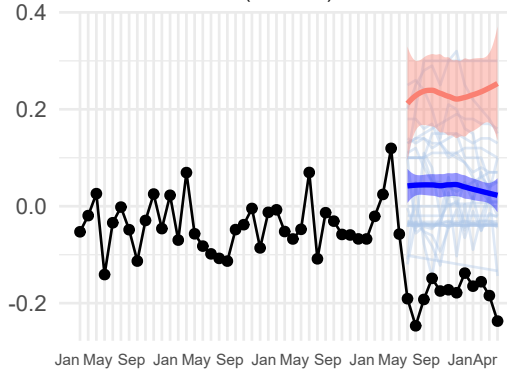


Scientists
 Naive Crowd
 Top Naive Statistic

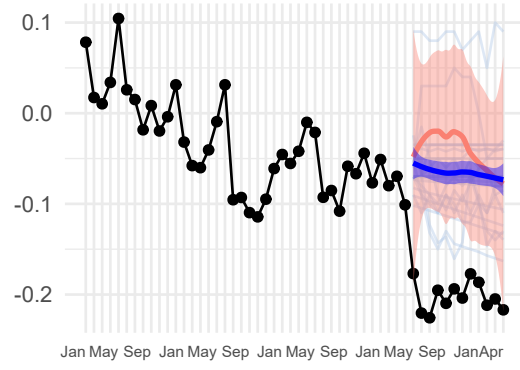
Exp. Bias Vs. Women-Career
higher=stereotype-consistent
(-3 to +3)



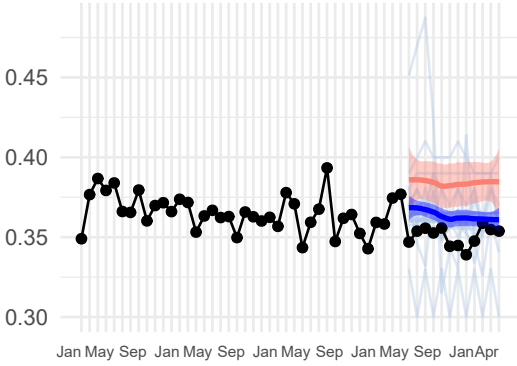
Exp. Bias Vs. Asian.-Am
higher=stereotype-consistent
(-3 to +3)



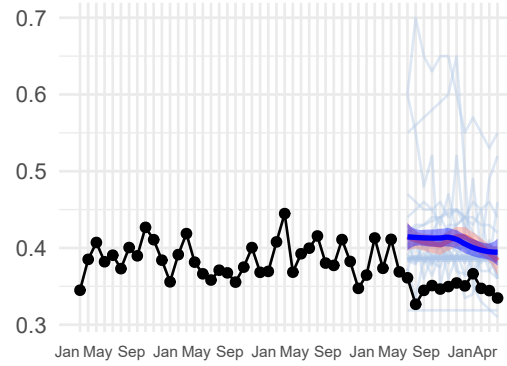
Exp. Bias Vs. Afr.-Am
higher=stereotype-consistent
(-3 to +3)



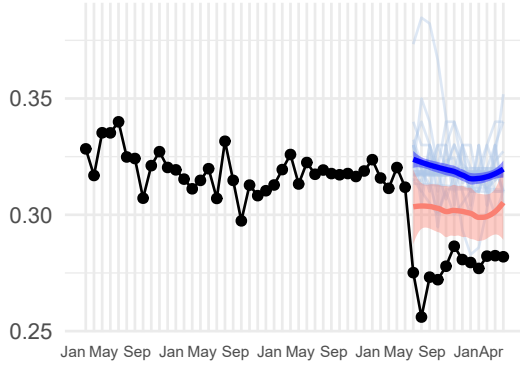
Imp. Bias Vs. Women-Career
higher=stereotype-consistent
(IAT D score)



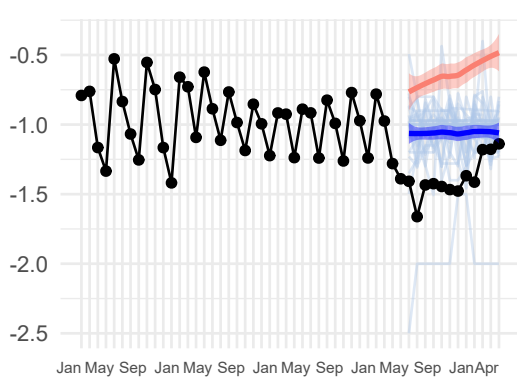
Imp. Bias Vs. Asian.-Am.
higher=stereotype-consistent
(IAT D score)



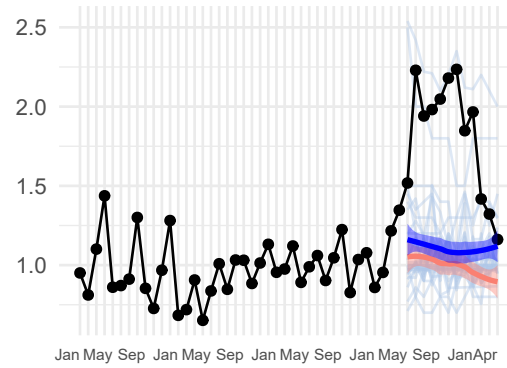
Imp. Bias Vs. Afr.-Am.
higher=stereotype-consistent
(IAT D score)



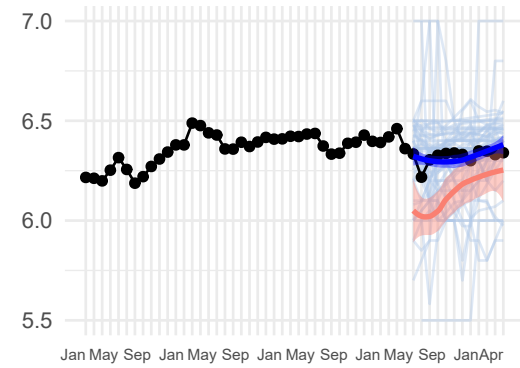
Positive Affect
(z-score)



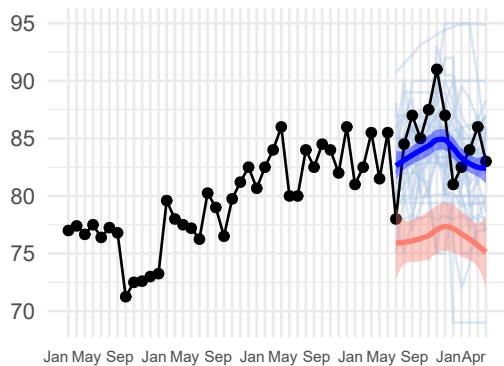
Negative Affect
standardized Vs. historical M/SD
(z-score)



Life Satisfaction
Cantril ladder
(0-10 scale)



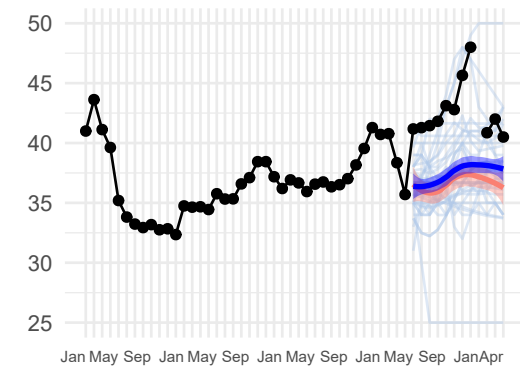
Polit. Polarization
% of Rep. Vs. Dem. approvals
(absolute difference score)



Democratic Support
(% Population)



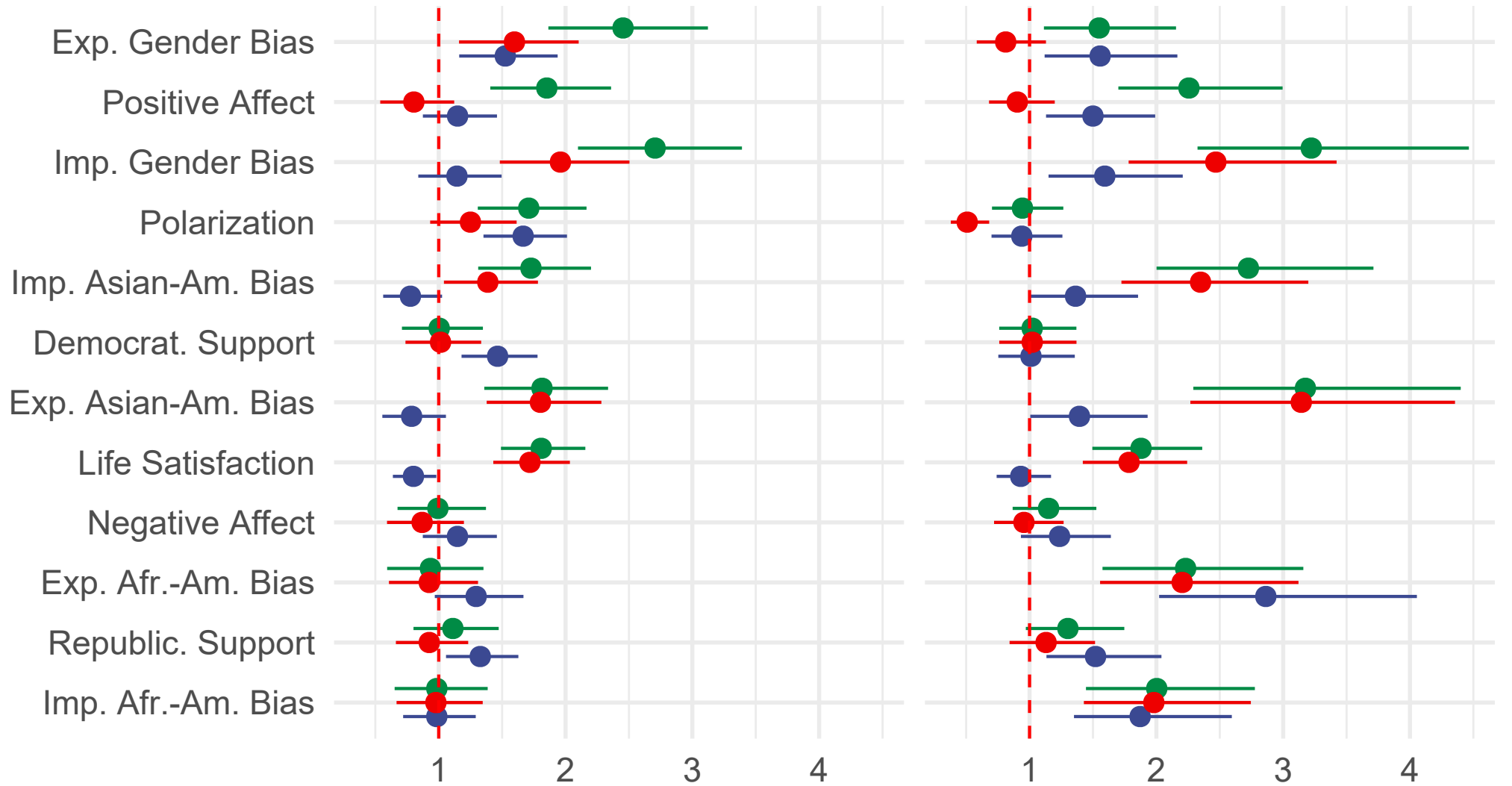
Republican Support
(% Population)



— Scientists — Naive Crowd

Tournament 1 (May 2020)

Tournament 2 (Nov 2020)



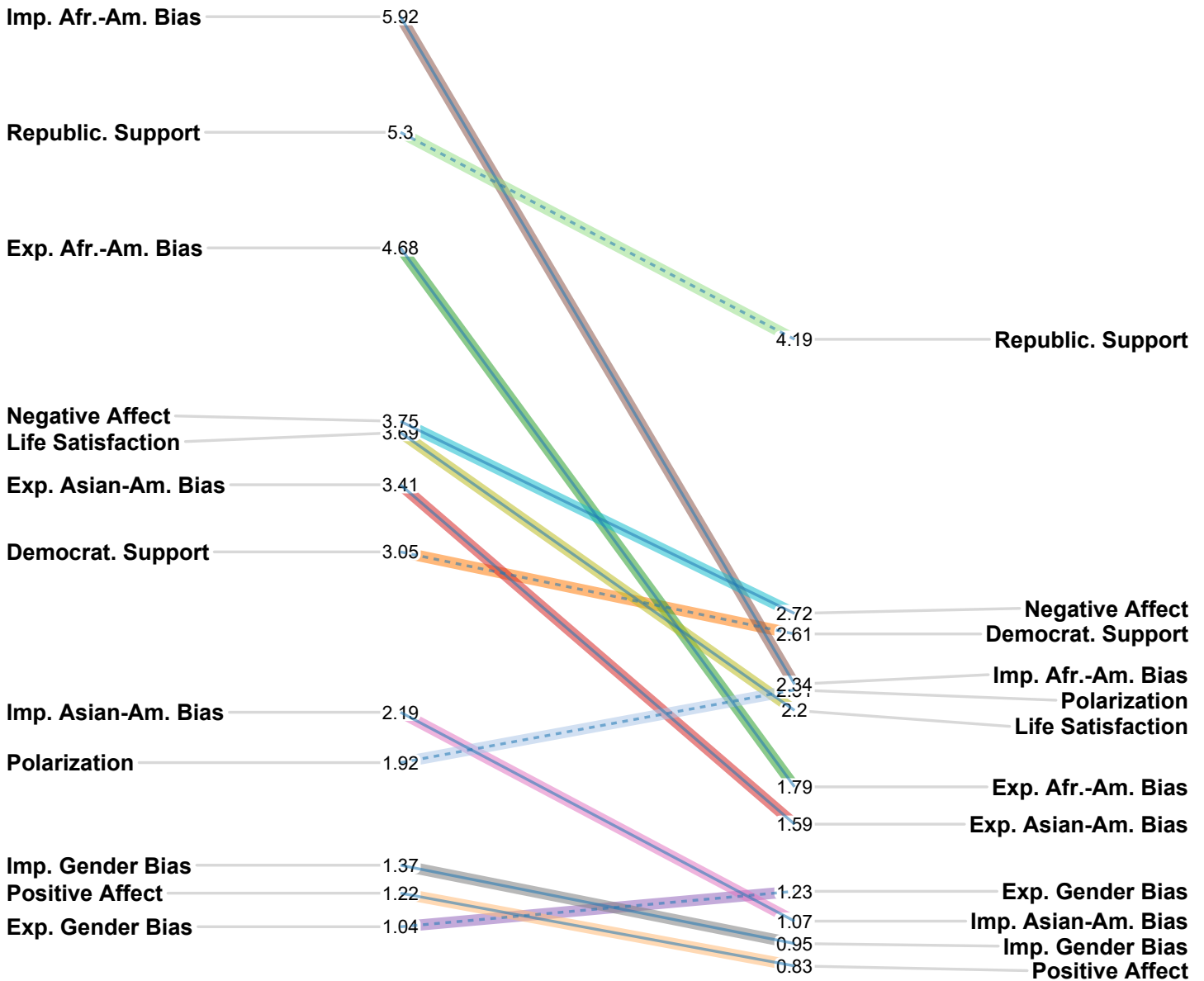
Naïve Benchmark / Scientific Forecast Error Ratio (M +/- 95%CI)

● Historical Mean
 ● Linear Regression
 ● Random Walk

Which domains are harder to predict?

First Tournament
May 2020

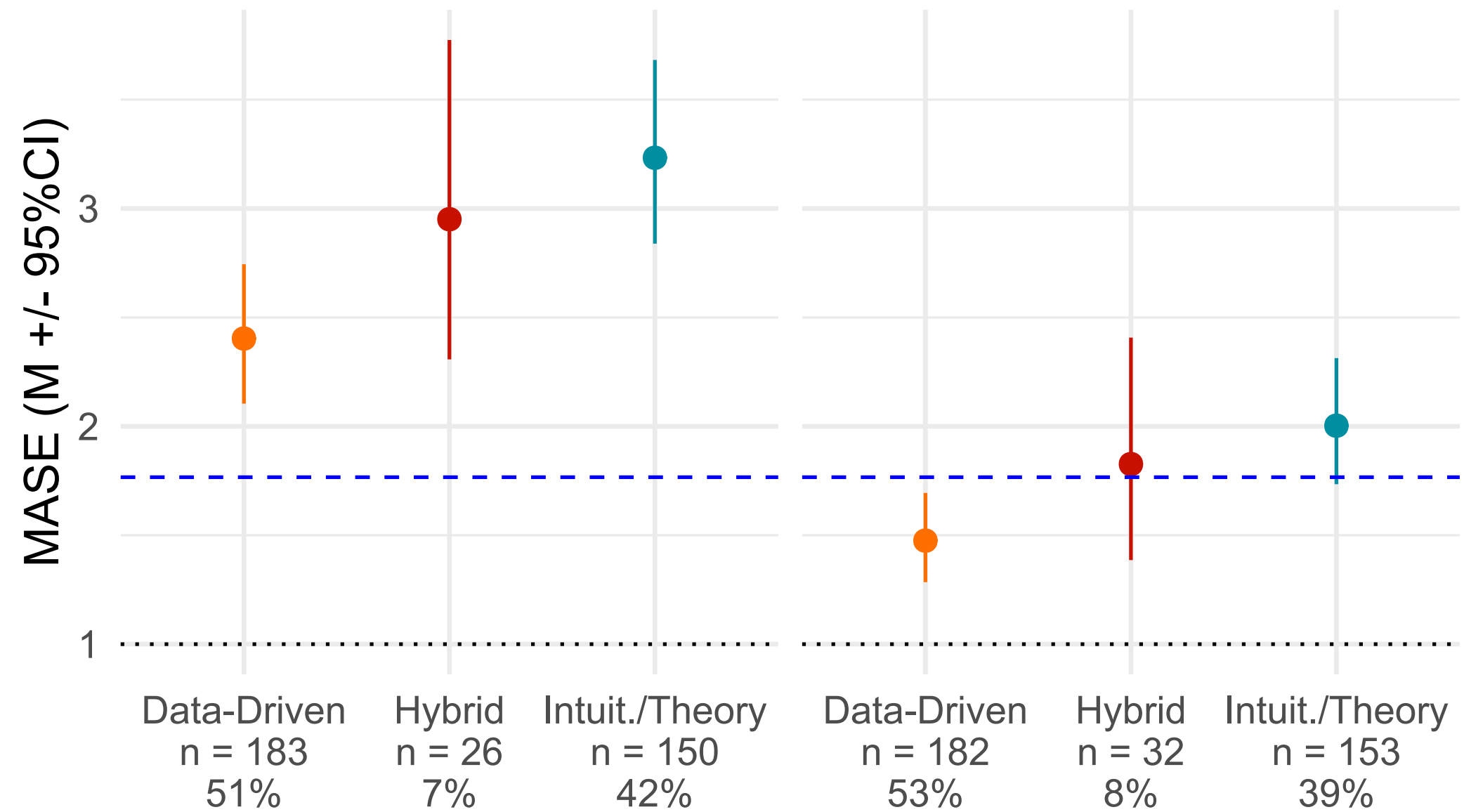
Second Tournament
Nov 2020

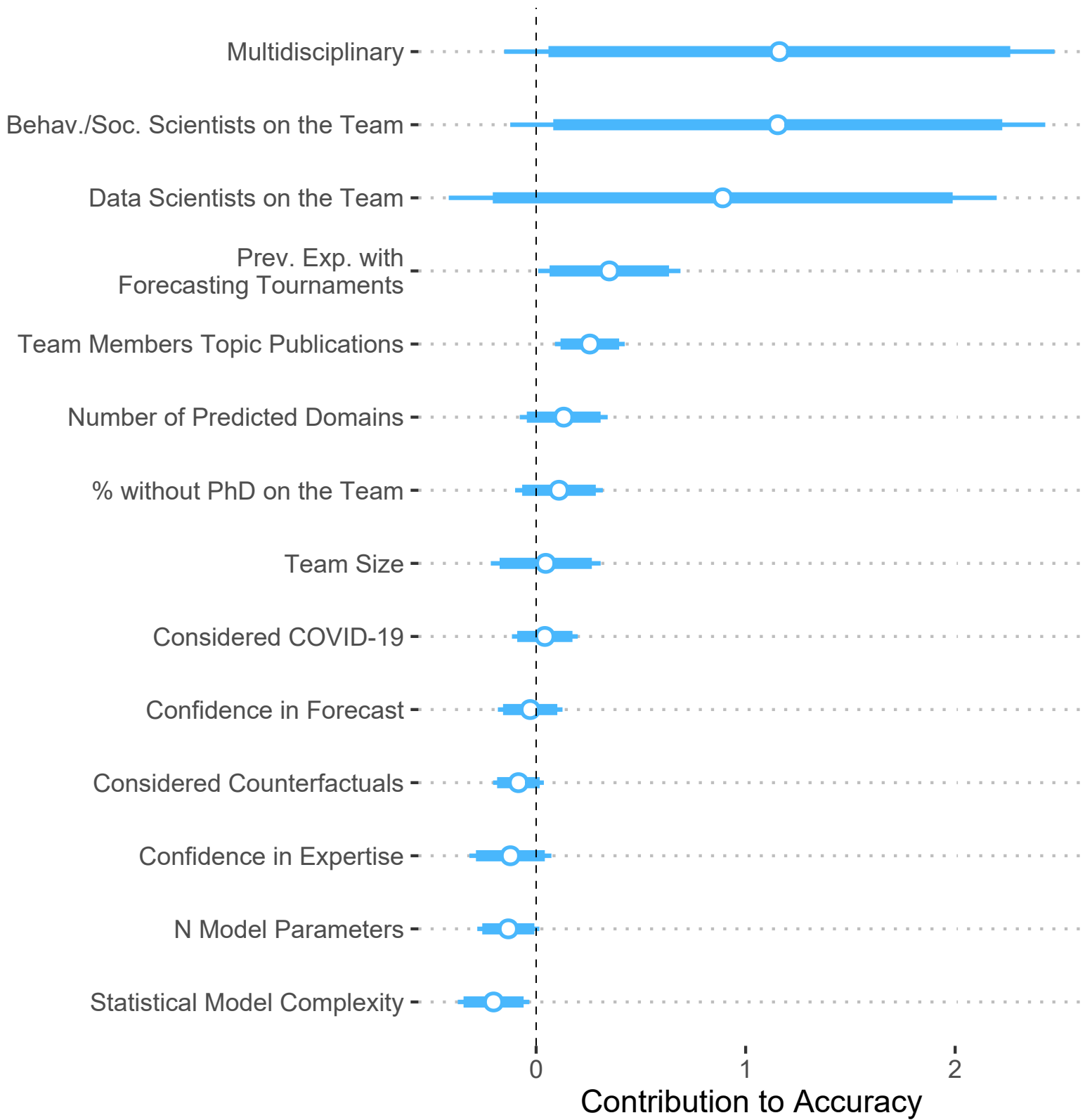


Ranking based on MASE scores per domain

First Tournament (May 2020)

Second Tournament (Nov 2020)





most negative <===== > most positive