

## SOCIAL SCIENCE

# AI and the transformation of social science research

Careful bias management and data fidelity are key.

By Igor Grossmann<sup>1,2\*</sup>, Matthew Feinberg<sup>3</sup>, Dawn C. Parker<sup>2,4</sup>, Nicholas Christakis<sup>5</sup>, Philip E. Tetlock<sup>6</sup>, William A. Cunningham<sup>7,8</sup>

Advances in artificial intelligence (AI), particularly large language models (LLMs), are dramatically affecting social science research. These transformer-based machine-learning models pre-trained on vast amounts of text data are increasingly capable of simulating human-like responses and behaviors (1, 3), offering novel opportunities for testing theories and hypotheses about human behavior at great scale and speed. This presents urgent challenges: How can social science research practices be adapted, even reinvented, to harness the power of foundational AI? And how can this be done while ensuring transparent and replicable research?

Social sciences rely on a range of methods, including questionnaires, behavioral tests, mixed-method analyses of semi-structured responses, agent-based modeling (ABM), observational studies, and experiments. The common goal is to obtain a generalized representation of characteristics of individuals, groups, cultures, and their dynamics (3). With the advent of advanced AI systems, the landscape of data collection in social sciences may shift. LLMs take advantage of deep learning to capture complex relationships within language. Such language literacy capabilities in processing, generating, and interacting with human language in a contextually-aware and semantically-accurate fashion (1) represent a major shift from previous AI approaches, which often struggled with such nuanced aspects of language as irony, metaphor, or emotional tone. With proper conditioning (4), LLMs can more accurately simulate human behavioral responses in social science research.

LLMs might supplant human participants for data collection (Fig. 1). For example, LLMs have already demonstrated their ability to generate realistic survey responses concerning consumer behavior (3). While opinions on the

feasibility of this application vary, at a minimum, studies using simulated participants could be used to generate novel hypotheses that could then be confirmed in human populations (4, 5). The success of this approach depends on algorithmic fidelity of the trained data (4), transparency in model training, prompt engineering, and benchmark selection.

Why is this scenario plausible? Pre-trained on massive datasets, advanced AI models can represent a vast array of human experiences and perspectives, possibly giving them a higher degree of freedom to generate diverse responses than conventional human participant methods, which can help to reduce generalizability concerns in research (3). LLMs can also generate responses across a wider range of parameters than human participants due to pragmatic concerns of limited attention span, response bias, or habituation among humans, providing a less biased view of underlying latent dimensions. This makes them especially useful in high-risk projects where traditional data collection is impractical, allowing for the testing of interventions in simulated populations before real-world implementation.

LLMs could be used as surrogates in other ways (Fig. 1). They have the potential to enhance policy analysis by reproducing the views of different theoretical or ideological schools of thought. For instance, LLMs could be trained to capture nuances of complex debates, such as concerning the stability and reliability of nuclear deterrence in the face of human and technical factors (6). LLMs could be trained to capture varied perspectives, including evaluating ‘what-if’ scenarios that nearly occurred, such as the Cuban Missile Crisis in 1962, and providing assessments of how plausible these scenarios were. Once LLMs can pass the Ideological Turing test—meaning they can accurately represent opposing viewpoints in a way indistinguishable from real humans—researchers can use them to generate future scenarios. Future LLMs, appropriately trained (4), may thus outperform humans on analytic tasks such as synthesizing clashing views to generate superior forecasts and policy prescriptions.

AI could also fill the role of a “confederate” (i.e., controlled experimental partners) in social interaction research involving individuals or groups (7), potentially as components to agent-

based simulations (Fig. 1). An LLM/ABM hybrid could use LLM to derive empirically-based rules of social decision-making or behavior to simulate social interactions of individuals with specific characteristics and beliefs (5) to explore how agents with these particular characteristics influence subsequent interaction with humans, informing broader social science questions such as how misinformation spreads throughout social networks (8).

Such investigations raise questions about the limits of LLMs as human cognition and decision models. Can we “nudge” an LLM by asking it to assess the quality of a news item before sharing, replicating research with humans (8)? If so, could we use the integrated LLM/ABM model to identify interventions that would reduce the spread of misinformation through social networks? Generally, if LLM/ABMs can provide new insights on how human agents choose to share information, cooperate and compete in social dilemmas, and conform with social norms, they can provide valuable insights into both the underlying mechanisms governing human behavior and social dynamics (9) with higher fidelity than has been possible using previous human decision models.

Incorporating LLMs into ABMs introduces new challenges due to their differing operational principles. While LLMs generate and interpret language based on statistical patterns derived from vast linguistic data, traditional ABMs operate based on predefined formal rules (10) which can be generated using real-world linguistic and other qualitative data. New ABM design will be needed to take advantage of LLMs capability to simulate performance on questionnaires, behavior in ill-defined situations, or open-ended responses (3). By creating realistic initial populations for ABMs, LLMs can model subjects’ latent cognitive or affective states, surpassing traditional researchers’ capacity and opening doors for future theory generation.

LLMs’ potential future benefits include creating samples as diverse as the cultural products (3, 4) on which the models were trained, offering a more accurate portrayal of human behavior and social dynamics than conventional methods relying on typically less heterogeneous and representative convenience samples (3). Due to their population-scale

<sup>1</sup>Department of Psychology, University of Waterloo; Waterloo, Canada. <sup>2</sup>Waterloo Institute for Complexity and Innovation, University of Waterloo; Waterloo, Canada. <sup>3</sup>Rotman School of Management, University of Toronto; Toronto, Canada. <sup>4</sup>School of Planning, University of Waterloo; Waterloo, Canada. <sup>5</sup>Yale Institute for Network Science, Yale University; New Haven, USA. <sup>6</sup>Wharton School of Business, University of Pennsylvania; Philadelphia, USA. <sup>7</sup>Department of Psychology, University of Toronto; Toronto, Canada. <sup>8</sup>Schwartz Reisman Institute for Technology and Society, University of Toronto; Toronto, Canada. Email: [igorssma@uwaterloo.ca](mailto:igorssma@uwaterloo.ca)

calibration data, LLMs could help address common challenges in social science research that can lead to biased models, including generalizability and self-selection concerns (3).

### The Scientist-Humanist Dilemma

Effective AI-assisted research will depend on the AI being able to accurately mirror the perspectives of diverse demographic groups. Pre-trained models from linguistic cultural products are known to capture socio-cultural biases present in society (3, 11). When biases are recognized, a key question is their provenance: do they correctly reflect the populations, or are they artifacts of model construction (12)? Model construction bias may result from incorrect or invalid choices throughout the design and development pipeline (e.g., choosing constructs which are differentially valid across demographic groups, curating datasets that lack diversity or that encode biases of certain human annotators, selecting models that fail to capture specific patterns pertinent to minorities) or because of existing societal disparities (3).

The scientist-humanist dilemma emerges as a key issue: while scientists aim to study “pure” LLMs with embedded socio-cultural biases to simulate human behavior and trace its cultural evolution (3), ethical constraints require engineers to protect LLMs from these very biases. Already, LLM engineers have been fine-tuning pre-trained models for the world that “should be” (3) rather than the world that is, and such efforts to mitigate biases in AI training (3, 13) may thus undermine the validity of AI-assisted social science research. The proprietary “black box” nature of LLM training challenges researchers’ ability to evaluate underlying mechanisms and replicate findings. To address this, advocating for open-source LLMs, access to pre-trained but not fine-tuned models for scientific research, and transparent methodologies such as BLOOM, Cerebras-GPT, or LLaMA is essential for ensuring reliable and credible AI-driven research (3).

Overall, researchers will need to establish guidelines for the ethical use of LLMs in research, addressing concerns related to data privacy, algorithmic fairness (vs. monoculture (3)), environmental costs (3, 13) and the potential misuse of LLM-generated findings. Pragmatic concerns with data quality, fairness, and equity of access to the powerful AI systems will be substantial.

### Weighing Trade-offs and Practical Wisdom

In deciding whether to use LLMs to approximate human behavior, researchers must first validate language-mediated (latent) constructs (3). They can treat LLM-generated responses as

a “sample” of non-human participants and systematically vary prompts, akin to presenting random stimuli in traditional experiments. A crucial consideration in using LLMs for research is the tradeoff between external and internal validity. Future LLMs, trained on diverse cultural content, will offer greater external validity by simulating human-like responses and generalizing to real-world scenarios. However, their opaque nature will limit their internal validity. Conversely, lab-grown natural language processing models, built on smaller controlled datasets, will provide stronger internal validity at the expense of reduced reliability and generalizability, as the limited training data may hinder their ability to perform consistently and broadly across different contexts. Researchers should carefully choose between these approaches based on their priorities.

Researchers must also consider the context of their study. High-risk situations involving violence or situations that are plainly infeasible with large numbers of human participants may be more suitable for LLMs. For example, LLMs might be used to explore human dynamics of space travel, or create predator and victim prototypes for studies of online sexual predators, an ethically fraught realm due to potential trauma to human participants.

As AI reshapes the landscape of social science (14), researchers will diversify their expertise, embracing new roles such as model bias hunters, AI-data validators, or human-AI interactionist. In this context, maintaining conceptual clarity (3), understanding foundations of measurement (3), and adhering to ethically-grounded practical wisdom (15) for selecting an AI-assisted design that fits one’s research question will be essential. With the democratization of AI-assisted data collection, the importance of early-stage social science training and supporting quantitative methods (e.g., computation, statistics) is crucial, calling for revision of social science education programs.

Just as the prisoners in Plato’s Cave Allegory observing shadows on a wall and believing them to represent reality, LLMs rely on “shadows” of human experiences described in cultural products. These shadows offer a limited view of the true nature of the phenomena they represent, because folk psychology (3) captured in cultural products may not always reflect the mechanisms governing human behavior—a limitation essential for social scientists to acknowledge. Examining LLM’s limitations and biases also puts a mirror to common practices in many fields, be it bias in representation, sampling methods, or methodological individualism (3).

Despite these obstacles, LLMs allow social scientists to break from traditional research

methods and approach their work in innovative ways. LLM models will likely bring the downfall of online crowdsourcing platforms, the dominant source of human participant data in many social science fields, for the simple reasons of on-par performance on simple tasks, and because open-ended responses from LLM-guided bots will become indistinguishable from human participants, calling for new methods for human data verification. Social scientists must be prepared to adapt to the uncertainty (15) that comes with evolving technology while being mindful of the limitations of ongoing research practices. Only by maintaining transparency and replicability (3), can we ensure that AI-assisted social science research truly contributes to our understanding of human experience.

### REFERENCES AND NOTES

1. S. Bubeck, et al, Sparks of Artificial General Intelligence: Early experiments with GPT-4 (2023), (available at <http://arxiv.org/abs/2303.12712>).
2. J. Wei, et al, Emergent Abilities of Large Language Models (2022), (available at <http://arxiv.org/abs/2206.07682>).
3. Extended documentation of LLMs abilities, ethical challenges, and methodological concerns, along with foundational social science principles, is on Open Science Framework ([osf.io/h4e2a](https://osf.io/h4e2a)).
4. L. P. Argyle, et al, *Polit. Anal.*, 1–15 (2023).
5. J. S. Park, et al, Generative Agents: Interactive Simulacra of Human Behavior (2023), (available at <http://arxiv.org/abs/2304.03442>).
6. P. E. Tetlock, C. B. McGuire, G. Mitchell, *Annu. Rev. Psychol.* **42**, 239–276 (1991).
7. H. Shirado, N. A. Christakis, *Nature*. **545**, 370–374 (2017).
8. G. Pennycook, et al, *Psychol Sci.* **31**, 770–780 (2020).
9. M. Galesic, et al, *J. R. Soc. Interface.* **18**, 20200857 (2021).
10. P. Antosz, S. Bharwani, M. Borit, B. Edmonds, *International Journal of Social Research Methodology*. **25**, 511–515 (2022).
11. A. Abid, M. Farooqi, J. Zou, *Nat Mach Intell.* **3**, 461–463 (2021).
12. S. Fazelpour, D. Danks, Algorithmic bias: Senses, sources, solutions. *Philosophy Compass*. **16** (2021).
13. L. Weidinger, et al, “Taxonomy of Risks posed by Language Models” in *2022 ACM Conference on Fairness, Accountability, and Transparency (ACM)*, Seoul Republic of Korea, 2022; <https://dl.acm.org/doi/10.1145/3531146.3533088>, pp. 214–229.
14. J. C. Peterson, et al., *Science*. **372**, 1209–1214 (2021).
15. I. Grossmann, et al, *Psychological Inquiry*. **31**, 103–133 (2020).

**Acknowledgments:** We thank T. Charlesworth, R. Saxe, and S. Fazelpour for their feedback on earlier versions of the draft. Funding: Social Sciences and Humanities Research Council of Canada Connection grant no. 611-2020-0190 (IG); Social Sciences and Humanities Research Council of Canada Insight grant no. 435-2014-0685 (IG); John Templeton Foundation grant no. 62260 (IG and PET).

10.1126/science.adi1778



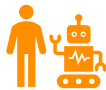
### Aid in Research

LLM helps sociologists refine surveys on social behavior norms in society.



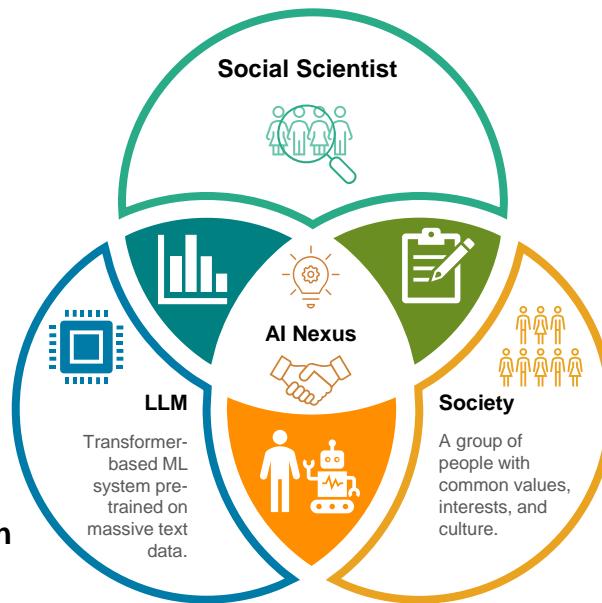
### Study of AI Surrogates

Psychologists analyze how an LLM simulates human responses in extreme isolation, mimicking the conditions of long-duration space travel.



### Society & Stand-in for Humans

In a psychological study, an LLM acts as a confederate, providing consistent responses to participants.



### AI Nexus



An interdisciplinary team of computer scientists, social scientists, and ethicists collaborate to study and address potential biases in hiring procedures due to skewed LLM outputs.

**Figure 1. Interplay of Roles Among LLMs, Social Scientists, and Society.** The diagram visualizes dynamic interactions and overlapping responsibilities among Large Language Models (LLMs), Social Scientists, and Society. Social Scientists use LLMs as research accelerators, aiding in the design and hypothesis formulation (top left). Additionally, LLMs serve as simulated AI proxies, mirroring human behavior to probe complex, elusive queries (top right) and acting as confederates in experiments to validate hypotheses (bottom left). At the heart of this relationship, the 'AI Nexus' (bottom right) illustrates a reciprocal feedback loop where each entity—LLMs, Social Scientists, and Society—mutually drives the scientific advancement in social sciences. Each box presents an example of application of LLM for transforming social science research.