



# AI-assisted teams outperform AI-led teams but not human-only teams in assessing research reproducibility in quantitative social science

Abel Brodeur<sup>a,b,1</sup> , David Valenta<sup>ab</sup> , Alexandru Marcoci<sup>c</sup> , Juan P. Aparicio<sup>ab</sup> , Derek Mikola<sup>a,b</sup> , Bruno Barbarioli<sup>a,b</sup> , Rohan Alexander<sup>d</sup> , Lachlan Deer<sup>e</sup> , Tom Stafford<sup>f,g</sup> , Lars Vilhuber<sup>h</sup> , Gunther Bensch<sup>i</sup> , Fabio Motokiki<sup>k</sup> , Mohamed Abdelhady<sup>l</sup> , Youssa Abdelmoula<sup>m</sup>, Ghina Abdul Baki<sup>a,b</sup> , Tomás Aguirre<sup>n</sup>, Sriraj Aiyer<sup>o</sup>, Shumi Akhtar<sup>p</sup> , Farida Akhtar<sup>q</sup> , Melle R. Albada<sup>r</sup> , Micah Altman<sup>s</sup>, David Angenendt<sup>t,u</sup> , Zahra Arjmandi Lari<sup>v</sup> , Jorge Armando De León Tejada<sup>w</sup>, David Rodriguez Arana<sup>x</sup> , Igor Asanov<sup>y</sup> , Anastasiya-Mariya Noha<sup>y</sup> , Rebecca Ashong<sup>z</sup> , Tobias Auer<sup>aa</sup>, Francisco J. Bahamonde-Birke<sup>bb</sup>, Bradley J. Baker<sup>cc</sup> , Söhnke M. Bartram<sup>dd,ee</sup>, Dongqi Bao<sup>ff</sup>, Lucija Batinovic<sup>gg</sup>, Tommaso Batistoni<sup>hh</sup> , Monica Beeder<sup>ii</sup> , Louis-Philippe Beland<sup>l</sup> , Carsten Gero Bienz<sup>jj</sup> , Christ Billy Aryanto<sup>kk</sup> , Cylcia Bolibaugh<sup>ll</sup> , Carl Bonander<sup>mm,nn</sup> , Ramiro Bravo<sup>oo</sup> , Egor Bronnikov<sup>pp,qq</sup>, Stephan Bruns<sup>rr,ss,tt</sup>, Nino Buliskeria<sup>uu</sup> , Sara Caicedo-Silva<sup>vv</sup>, Andrea Calef<sup>ww</sup> , Juan Sebastian Cano Arias<sup>x</sup>, Gustavo A. Castillo Alvarez<sup>xx</sup>, Solomon Caulker<sup>yy</sup>, Simonas Cepenas<sup>zz</sup> , Arthur Chatton<sup>aaa,bbb</sup> , Zirou Chen<sup>ccc</sup> , Ngozi Chioma Ewurum<sup>ddd</sup>, Anda-Bianca Ciocirlan<sup>f</sup> , Felix J. Clouth<sup>bb</sup>, Jason Collins<sup>eee</sup> , Nikolai Cook<sup>fff</sup>, Cesar Cornejo<sup>ggg,hhh</sup> , João Craveiro<sup>f</sup>, Jonathan Créchet<sup>iii</sup>, Jing Cui<sup>jjj,kkk</sup> , Niveditha Chalil Vayalabron<sup>ll</sup>, Christian Czymara<sup>mmm</sup> , Carlos Daniel Bermúdez Jaramillo<sup>nnn</sup>, Hannes Datta<sup>ooo</sup>, Lien Denoo<sup>ppp</sup> , Arshia Dhaliwal<sup>l</sup>, Nancy Dhameja<sup>qqq</sup> , Elodie Djemai<sup>rrr</sup>, Erwan Dujeancourt<sup>sss,ttt</sup> , Uğurcan Dündar<sup>uuu</sup> , Thibaut Duprey<sup>vvv</sup> , Yasmine Eissa<sup>www</sup> , Youssef El Fassi<sup>xxx</sup>, Ismail El Fassi<sup>yyy</sup> , Keaton Ellis<sup>zzz</sup>, Ali Elminejad<sup>uuu</sup> , Mahmoud Elsherif<sup>aaaa,bbbb</sup>, Aysil Emirmahmutoglu<sup>ccc</sup>, Giulian Etingin-Frati<sup>ddd</sup> , Emeke Eze<sup>eeee</sup>, Jan Fabian Dollbaum<sup>fff</sup> , Jan Feld<sup>ggg</sup> , Andres Felipe Rengifo Jaramillo<sup>hhh,iii</sup> , Guidon Fenig<sup>a</sup>, Victoria Fernandes<sup>vvv,iii</sup>, Lenka Fiala<sup>a,b,kkkk</sup> , Lukas Fink<sup>lll</sup> , Mojtaba Firouzjaeiangalouah<sup>mmmm</sup> , Sara Fish<sup>nnn</sup>, Jack Fitzgerald<sup>oooo,pppp</sup>, Rachel Forshaw<sup>qqq</sup> , Alexandre Fortier-Chouinard<sup>rrrr</sup> , Louis Fréget<sup>ssss</sup> , Joris Frese<sup>tttt</sup> , Jacopo Gabani<sup>uuuu,vvvv</sup> , Sebastian Gallegos<sup>wwww</sup> , Max C. Gamil<sup>xxxx</sup> , Attila Gáspár<sup>yyyy,zzzz</sup> , Romain Gauriot<sup>aaaaa</sup> , Evelina Gavrilova<sup>cccc</sup> , Diogo Geraldos<sup>bbbbb,cccc,dddd</sup> , Giulio Giacomo Cantone<sup>eeeee</sup> , Grant Gibson<sup>ffff,gggg</sup> , Dirk Goldschmitt<sup>f</sup> , Amélie Gourdon-Kanhukamwe<sup>hhhhh</sup>, Andrea Gregor de Varda<sup>iiii</sup>, Idaliya Grigoryeva<sup>jjjj</sup> , Alexi Gugushvili<sup>kkkkk</sup> , Aaron H. A. Fletcher<sup>llll</sup>, Florian Habermann<sup>mmmmm,nnnnn</sup> , Márton Hablicsek<sup>ooooo</sup>, Joanne Haddad<sup>ppppp</sup> , Jonathan D. Hall<sup>qqqqq</sup>, Olle Hammar<sup>tt,rrrrr</sup> , Malek Hassouneh<sup>sssss</sup>, Carina I. Hausladen<sup>ttttt</sup> , Sophie C. F. Hendrikse<sup>bb</sup>, Matthew Hepplewhite<sup>uuuuu</sup> , Anson T. Y. Ho<sup>vvvvv</sup> , Senan Hogan-Hennessy<sup>h</sup> , Elliot Howley<sup>wwwww</sup> , Gaoyang Huang<sup>xxxx,ff</sup> , Héloïse Hulstaert<sup>yyyyy</sup> , Zlatomira G. Ilchovska<sup>zzzzz,aaaaa,wwwww</sup>, Paola Jaimes Santamaria<sup>bbbbb,cccc</sup> , Niklas Jakobsson<sup>mm</sup> , Joakim Jansson<sup>ttt,dddd</sup> , Ewa Jarosz<sup>eeeee</sup>, Hossein Jebeli<sup>fffff</sup> , Yanchen Jiang<sup>nnnn</sup> , Hiba Junaid<sup>ggggg,hhhhh</sup> , Rohan Kalluraya<sup>iiiiii</sup>, Sunny Karim<sup>jjjjj</sup>, Edmund Kelly<sup>kkkkk</sup>, Eva Kimmel<sup>zzzzz</sup> , Sorraivich Kingsuwankul<sup>lllll,pppp</sup>, Valentin Klotzbücher<sup>mmmmm,nnnnn,ooooo</sup> , Daniel Krämer<sup>ppppp</sup> , Pijus Krūminas<sup>zz</sup> , Nicholas Kruus<sup>uuuuu</sup> , Essi Kujansuu<sup>qqqqq,rrrrr</sup>, Christoph F. Kurz<sup>sssss</sup>, Stephan Küster<sup>ttttt</sup> , Blake Lee-Whiting<sup>uuuuu</sup> , Felix Lewandowski<sup>wwwww</sup> , Tongzhe Li<sup>vvvvv</sup>, Ruoxi Li<sup>wwwww</sup>, Dan Liu<sup>xxxxx</sup> , Jiacheng Liu<sup>yyyyy</sup> , Helix Lo<sup>zzzzz</sup> , Katharina Loter<sup>bb</sup>, Felipe Macedo Dias<sup>aaaaa</sup>, Christopher R. Madan<sup>wwwww</sup> , Nicolas Mäder<sup>bbbbb</sup> , Marco Mandas<sup>cccc</sup> , Cesar Mantilla<sup>ddddd</sup> , Jan Marcus<sup>lll</sup> , Diego Marino Fages<sup>eeeeee</sup> , Xavier Martin<sup>fffff</sup>, Ryan McWay<sup>ggggg</sup> , Daniel Medina-Gaspar<sup>hhhhh</sup> , Sisi Meng<sup>iiiiii</sup>, Lingyu Meng<sup>jjjjj</sup> , Simon Merz<sup>kkkkk</sup> , Alex P. Miller<sup>lllll</sup> , Thibault Mirabel<sup>mmmmm</sup> , Dibya Deepta Mishra<sup>nnnnn</sup>, Sumit Mishra<sup>ooooo</sup> , Belay W. Moges<sup>ppppp</sup> , Morteza Mohandes Mojarrad<sup>qqqqq</sup>, Myra Mohnen<sup>rrrrr</sup>, Louis-Philippe Morin<sup>a</sup> , Lucija Muehlenbachs<sup>sssss,ttttt</sup>, Gastón Mullin<sup>uuuuu</sup> , Andreea Musulan<sup>wwwww,wwwww</sup> , Sara Muzzi<sup>xxxxx,yyyyy</sup>, James A. C. Myers<sup>zzzzz</sup>, Florian Neubauer<sup>aaaaa</sup> , Tuan Nguyen<sup>rr</sup> , Ali Niazi<sup>sssss</sup>, Ardyn Nordstrom<sup>bbbbb</sup>, Bartłomiej Nowak<sup>cccc</sup> , Daneal O'Habib<sup>ddddd</sup>, Tim Ölkers<sup>eeeeee</sup> , Justin Ong<sup>f</sup>, Valeria Orozco Castiblanco<sup>ffffff,ggggg</sup>, Ömer Özak<sup>hhhhh</sup> , Ali I. Ozkes<sup>iiiiii</sup> , Mikael Paaso<sup>kkkkk</sup>, Shubham Pandey<sup>lllll</sup> , Varvara Papazoglou<sup>lllll</sup> , Romeo Penheiro<sup>mmmmmm</sup> , Linh Pham<sup>nnnnn</sup>, Ulrike Phielers<sup>ooooo</sup> , Peter Pütz<sup>ppppp</sup> , Quan Qi<sup>qqqqq</sup> , Jingyi Qiu<sup>rrrrr</sup> , Manuel T. Rein<sup>bb</sup>, David A. Reinstein<sup>sssss</sup>, Juuso Repo<sup>ttttt</sup> , Nicolas Rudolf<sup>nnnnn</sup>, Shree Saha<sup>uuuuu</sup>, Orkun Saka<sup>vvvvv</sup>, Chiara Saponaro<sup>wwwww</sup> , Georg Sator<sup>xxxxx</sup> , Martijn Schoenmakers<sup>bb</sup> , Raffaello Ser<sup>yyyyy</sup> , Meet Shah<sup>vvvv</sup>, Paul Sibille<sup>yyyyy</sup>, Christoph Siemroth<sup>zzzzz</sup>, Vladimir Skavys<sup>aaaaa,bbbb</sup> , Ben Slater<sup>cccc</sup> , Wenting Song<sup>ddddd</sup>, Stefan Staubli<sup>sssss</sup>, Tobias Steindl<sup>eeeeee</sup> , Nomwendé Steven Waongo<sup>a</sup>, Paul Stott<sup>ffffff,ggggg</sup> , Stephenson Strobel<sup>ffff</sup>, Roshini Sudhaharan<sup>hhhhh</sup>, Pu Sun<sup>iiiiii</sup> , Scott D. Swain<sup>jjjjj</sup> , Oleksandr Talavera<sup>kkkkk</sup>, Hanz M. Tantiangco<sup>lllll</sup> , Georgy Tarasenko<sup>mmmm</sup> , Boyd Tarlinton<sup>nnnnn</sup> , Mariam Tarraf, Ken Teoh<sup>ooooo</sup>, Rémi Thériault<sup>ppppp</sup> , Bethan Thompson<sup>qqqqq</sup> , Tonghui Tian<sup>l</sup>, Wenjie Tian<sup>a</sup>, Emmanuel Tolani<sup>rrrrr,sssss</sup> , Nicolai Borgen<sup>ttttt</sup> , Solveig Topstad Borgen<sup>kkkk</sup> , Javier Torralba<sup>ooo</sup>, Carolina Velez-Ospina<sup>uuuuu</sup>, Man Wai Mak<sup>l</sup> , Lukas Wallrich<sup>wwwww</sup>, Zeyang Wang<sup>wwwww</sup> , Leah Ward<sup>xxxxx</sup>, Matthew D. Webb<sup>l</sup> , Duncan Webb<sup>yyyyy</sup>, Bryan S. Weber<sup>zzzzz,aaaaa</sup> , Christoph Weber<sup>bbbbb</sup> , Wei-Chien Weng<sup>cccc</sup> , Christian Westheide<sup>ddddd,eeee</sup>, Tom Wilkinson<sup>xxxx</sup> , Kwong-Yu Wong<sup>ffffff</sup> , Marcin Wroński<sup>ggggg</sup> , Zhuangchen Wu<sup>kkkkk</sup>, Qixia Wu<sup>a</sup> , Victor Y. Wu<sup>hhhhh</sup> , Bohan Xiao<sup>a</sup> , Feihong Xu<sup>iiiiii</sup> , Cong Xu<sup>jjjjj,kkkkk</sup>, Pranav Yadav<sup>lllll</sup> , Yu Yang Chou<sup>hhhh</sup> , Luther Yap<sup>mmmm</sup> , Myra Yazbeck<sup>l,nnnnn</sup> , Bo Yao<sup>ooooo</sup> , Zuzanna Zagrodzka<sup>ppppp</sup> , Tahreen Zahra<sup>l</sup> , Mirela Zaneva<sup>qqqqq</sup>, Xiaomeng Zhang<sup>rrrrr</sup>, Ziwei Zhao<sup>sssss,ttttt</sup>, Han Zhong<sup>uuuuu</sup>, Aras Zirculis<sup>zz</sup>, Jiacheng Zou<sup>wwwww</sup>, Floris Zoutman<sup>cccc</sup>, and Christelle Zozoungbo<sup>wwwww</sup>

Affiliations are included on p. 9.

Edited by Timothy D. Wilson, University of Virginia, Charlottesville, VA; received September 22, 2025; accepted March 16, 2026

Large Language Models (LLMs) such as ChatGPT are transforming how scientists conduct and validate research, offering promise as tools to improve scientific reproducibility. However, computational reproducibility and error detection remain expensive and labor-intensive. We experimentally test how collaboration between researchers and LLM assistants influences the reproduction of quantitative social science findings across different levels of AI autonomy. We randomly assigned 288 researchers to 103 teams working under three conditions: human-only, AI-assisted (using ChatGPT as a collaborative tool), or AI-led (ChatGPT operating with minimal human oversight). Teams reproduced published results from leading social science journals, detected coding errors, and proposed robustness checks. Human-only and AI-assisted teams achieved comparable reproduction rates (94% vs. 91%) and performed similarly on most outcomes, except human-only teams identified significantly more major coding errors. Both substantially outperformed AI-led teams, which achieved only a 37% reproduction rate, detected fewer errors across all categories, proposed weaker robustness checks, and required more time. This autonomous approach, however, likely represents only a lower bound of AI capabilities. Despite rapid model advances, expert human judgment currently remains indispensable for reliable empirical verification. While AI assistance did not degrade most outcomes, it provided no measurable advantages and was associated with reduced detection of major errors. However, the 37% autonomous reproduction rate indicates that AI could provide value in settings where scale or cost constraints preclude human review of papers, even though general-purpose LLMs offer no immediate advantages for human-supervised verification.

AI | reproducibility | large language models

Reproducibility is a cornerstone of robust quantitative empirical research, where complex methodologies and data handling techniques are common (1–8). Despite advancements in reproducibility protocols (9), concerns persist regarding the accuracy and reliability of published findings (10–17). Unclear reporting and methodological advances requiring expertise when evaluating quantitative studies contribute to the current reproducibility and replication crises in the behavioral and social sciences. At the same time, verifying computational reproducibility remains costly and labor-intensive (18). Even when journals require replication packages, reproducing results often involves navigating complex scripts, large datasets, and intricate empirical workflows. As empirical research becomes increasingly complex, scalable approaches to verification are needed to ensure that the reliability of published findings can be efficiently assessed.

This study investigates how artificial intelligence (AI) tools, such as Large Language Models (LLMs), could support researchers, data editors, and scientific journals in computationally reproducing research. We focus on three modes of AI and human interaction: human-only teams, human teams with AI assistance (the “AI-assisted” approach), and teams that provided only limited oversight while AI carried out reproducibility checks (the “AI-led” approach). The AI-led approach approximates a “protoagentic” system: an LLM tasked with reasoning through a reproducibility exercise with minimal human supervision. We use ChatGPT because it processes different file formats effectively for reproduction and is used most frequently by researchers (19).

This paper tests how effectively AI supports reproduction of scientific articles and works in complex cases where coding errors or methodological inconsistencies arise. We employ a randomized controlled trial design involving three treatment arms. We contribute to a large literature documenting the benefits and limitations of human–AI integration, as well as full automation (20). Evidence from human–AI decision-making suggests that performance ordering between AI alone, human alone, and human–AI teams is mixed and task-dependent, and that human–AI combinations often fail to outperform AI alone, sometimes performing worse due to miscalibrated trust and under- or overreliance on AI assistance (21–34). This is crucial for science because current methods for performing computational reproducibility and robustness checks are expensive, time consuming, and require advanced technical skills (18, 35). We also contribute to a growing body of literature documenting the potential pitfalls of integrating human and artificial intelligence, such as overreliance and expertise erosion (36, 37). This research also provides some comparative productivity measures in highly specialized intellectual tasks. This line of research mainly focuses on customer support agents and low-skill occupations, whereas we study high-skill scientific reproducibility tasks (29, 38).

Author contributions: A.B., D.V., A. Marcoci, J.P.A., D.M., B.B., R. Alexander, G.B., G.A.B., F.A., F.J.B.-B., L.-P.B., C. Cornejo, M.E., L. Fiala, J. Fitzgerald, R.F., J. Frese, G.G., A.G.-K., N.M., S. Merz, A. Musulan, A. Nordstrom, D.A.R., G.S., R. Seri, V.S., P. Sun, G.T., N.T.B., S.T.B., B.S.W., and L.Y. designed research; A.B., D.V., A. Marcoci, J.P.A., D.M., B.B., R. Alexander, L. Deer, T. Stafford, L.V., G.B., F.M., M. Abdelhady, Y.A., G.A.B., T. Aguirre, S. Ayier, S. Akhtar, F.A., M.R.A., M. Altman, D.A., Z.A.L., J.A.d.L.T., D.R.A., I.A., A.-M.N., R. Ashong, T. Auer, F.J.B.-B., B.J.B., S.M.B., D.B., L.B., T.B., M.B., L.-P.B., C.G.B., C.B.A., C. Bolibaugh, C. Bonander, R.B., E.B., S.B., N.B., S.C.-S., A. Calef, J.S.C.A., G.A.C.A., S. Caulker, S. Cepenas, A. Chatton, Z.C., N.C.E., A.-B.C., F.J.C., J. Collins, N.C., C. Cornejo, J. Craveiro, J. Créchet, J. Cui, N.C.V., C. Czymara, C.D.B.J., H.D., L. Denoo, A.D., N.D., E. Djemai, E. Dujeancourt, U.D., T.D., Y.E., Y.E.F., I.E.F., K.E., A. Elminejad, M.E., A. Emirmahmutoglu, G.E.-F., E.E., J.F.D., J. Feld, A.F.R.J., G.F., V.F., L. Fiala, L. Fink, M.F., S.F., J. Fitzgerald, R.F., A.F.-C., L. Fréget, J. Frese, J.G., S.G., M.C.G., A. Gáspár, R.G., E.G., D. Geraldes, G.G.C., G.G., D. Goldschmitt, A.G.-K., A.G.d.V., I.G., A. Gugushvili, A.H.A.F., F.H., M. Hablicsek, J.H., J.D.H., O.H., M. Hassouneh, C.I.H., S.C.F.H., M. Hepplewhite, A.T.Y.H., S.H.-H., E.H., G.H., H.H., Z.G.I., P.J.S., N.J., J.J., E.J., H. Jebeli, Y.J., H. Junaid, R.K., S. Karim, E. Kelly, E. Kimel, S. Kingsuwanukul, V.K., D.K., P.K., N.K., E. Kujansuu, C.F.K., S. Küster, B.L.-W., F.L., T.L., R.L., D.L., J.L., H.L., K.L., F.M.D., C.R.M., N.M., M. Mandas, C.M., J.M., D.M.F., X.M., R.M.W., D.M.-G., S. Meng, L. Meng, S. Merz, A.P.M., T.M., D.D.M., S. Mishra, B.W.M., M.M.M., M. Mohnen, L.-P.B., L. Muehlenbachs, G.M., A. Musulan, S. Muzzi, J.A.C.M., F.N., T.N., A. Niazi, A. Nordstrom, B.N., D.O., T.Ö., J.O., V.O.C., Ö.Ö., A.I.O., M.P., S.P., V.P., R.P., L.P., U.P., P.P., Q.Q., J.Q., D.A.R., J.R., N.R., S. Saha, O.S., C. Saponaro, G.S., M. Schoenmakers, R. Seri, M. Shah, P. Sibille, C. Siemroth, V.S., B.S., W.S., S. Staubli, T. Steindl, N.S.W., P. Stott, S. Strobel, R. Sudhaharan, P. Sun, S.D.S., O.T., H.M.T., G.T., B. Tarlinton, M.T., K.T., R.T., B. Thompson, T.T., W.T., M.T.R., E.T., N.T.B., S.T.B., J.T., C.V.-O., M.W.M., L. Wallrich, Z. Wang, L. Ward, M.D.W., D.W., B.S.W., C. Weber, W.-C.W., C. Westheide, T.W., K.-Y.W., M.W., Z. Wu, Q.W., Y.Y.W., B.X., F.X., C.X., P.Y., Y.Y.C., L.Y., M.Y., B.Y., Z. Zagrodzka, T.Z., M.Z., X.Z., Z. Zhao, H.Z., A.Z., J.Z., F.Z., and C.Z. performed research; A.B., D.V., J.P.A., G.B., Z.A.L., I.A., A.-M.N., T. Auer, C.G.B., C. Bonander, R.B., S.B., A. Chatton, A.D., E. Dujeancourt, Y.E., J. Fitzgerald, O.H., A.T.Y.H., G.H., H.H., E. Kelly, V.K., N.K., R.L., J.L., K.L., S. Merz, S. Mishra, S. Muzzi, F.N., T.N., U.P., P. Sibille, S. Staubli, O.T., B. Tarlinton, M.T., M.T.R., C. Weber, Z. Wu, and V.Y.W. analyzed data; L. Deer, T. Stafford, G.B., F.M., S. Ayier, S. Akhtar, F.A., M. Altman, I.A., A.-M.N., B.J.B., L.B., M.B., C.G.B., C. Bonander, R.B., S.B., A. Chatton, F.J.C., J. Collins, N.C., L. Denoo, E. Djemai, T.D., I.E.F., A. Elminejad, J.F.D., J. Feld, G.F., J. Fitzgerald, L. Fréget, J.G., S.G., A. Gáspár, E.G., G.G., D. Goldschmitt, A.G.-K., I.G., F.H., M. Hablicsek, J.D.H., S.C.F.H., A.T.Y.H., G.H., Z.G.I., J.J., E. Kelly, E. Kimel, S. Kingsuwanukul, N.K., E. Kujansuu, S. Küster, B.L.-W., D.L., J.L., K.L., C.M., X.M., R.M.W., S. Merz, B.W.M., M. Mohnen, L.-P.M., L. Muehlenbachs, F.N., Ö.Ö., A.I.O., S.P., U.P., P.P., Q.Q., D.A.R., J.R., N.R., O.S., R. Seri, P. Sibille, S. Staubli, P. Sun, G.T., R.T., T.T., M.T.R., L. Wallrich, M.D.W., C. Weber, W.-C.W., M.W., V.Y.W., L.Y., B.Y., Z. Zagrodzka, M.Z., X.Z., H.Z., A.Z., and F.Z. edited and reviewed the paper; A.B., D.V., J.P.A., D.M., R. Alexander, L. Deer, T. Stafford, and L.V. organized replication games; and A.B., D.V., A. Marcoci, J.P.A., D.M., B.B., G.A.B., and L. Fiala wrote the paper.

Competing interest statement: The views expressed in this paper are those of the authors. No responsibility for them should be attributed to the Bank of Canada. A. Marcoci is a UK Research and Innovation Policy Fellow seconded to the Department for Science, Innovation and Technology. The views and conclusions contained herein are those of the authors and should not be interpreted as representing the official policies, either expressed or implied, of the Department for Science, Innovation and Technology or the UK Government.

This article is a PNAS Direct Submission.

Copyright © 2026 the Author(s). Published by PNAS. This article is distributed under Creative Commons Attribution-NonCommercial-NoDerivatives License 4.0 (CC BY-NC-ND).

<sup>1</sup>To whom correspondence may be addressed. Email: abrodeur@uottawa.ca.

This article contains supporting information online at <https://www.pnas.org/lookup/suppl/doi:10.1073/pnas.2524747123/-/DCSupplemental>.

Published May 28, 2026.

We focus on three groups of outcomes across the treatment arms: 1) computational reproducibility (success rate and time required), 2) error detection capabilities, and 3) proposing and implementing quality robustness checks. Understanding the impact of the treatment on these outcomes contributes to a broader understanding of AI, and offers insights into the optimal balance of human and AI involvement in research tasks.

## 1. Procedures

The first 10 coauthors organized seven AI replication games between February and November 2024, including a pilot in February. All remaining coauthors and a few of the organizers participated in one of those games. The participating coauthors were a mix of master and PhD students, postdoctoral fellows, professors, and researchers from nonacademic organizations with a doctoral degree. *SI Appendix, Table S8* provides details on team composition. Randomization was carried out in two steps for each of the seven events. In step one, participants were randomly assigned to a team of three to evaluate the reproducibility of a quantitative social science article. The randomization in step one was conditional on the software preferences reported by participants (Stata or R) and the mode of participation (in person or virtual). In step two, each team was randomly assigned to one of three treatment arms: human-only, AI-assisted, or AI-led.

Each team was assigned a study from leading social science journals (i.e., economics, political science, or behavioral science/psychology). Each event included two studies with known coding errors (one in Stata and one in R) that had been identified by the lead authors in a prior study but were not publicly disclosed at the time of the AI replication game. Detailed information about the papers used that contained coding errors can be found in *SI Appendix, Tables S1 and S2*. Descriptions of the coding errors identified prior to each replication game can be found in *SI Appendix, Table S3* through *SI Appendix, Table S4*. Coding errors occurred when preparing data for analysis (variable definitions, incorrect merging of datasets, differing sample restrictions, not cleaning variables, missing variables) as well as when carrying out the analysis (discrepancies between code and what is written in the article). Examples of the latter include inconsistently specified SEs and control variables. Teams and local organizers had no information about the study they would be reproducing until the start of the event. Twelve studies were used in total, with a few reused for multiple events.

Relevant resources were given to the teams at 09:00 local time on the day of the event. We shared with them: the journal article and online appendix as PDFs, the original authors' replication package, and screenshots of the exhibit to reproduce from the article (*SI Appendix*). Screenshots were introduced after the pilot event to assist AI-led teams, as the AI might be better able to process tables and figures as images rather than when embedded in PDF files. Teams had seven hours to complete three tasks: i) computationally reproduce a few predetermined results, ii) detect coding errors, and iii) suggest and implement up to two robustness checks. The three tasks were independent from each other (e.g., teams did not need to fix coding errors to computationally reproduce results). However, teams were instructed to begin with reproducing the results before proceeding to specifically search for coding errors and propose robustness checks. Teams could leave before the end of the event if they believed they had completed their tasks as feasibly as possible. Upon completion, teams were asked to email the lead authors a (templated) time log documenting whether they completed computational reproducibility, with a list of all coding errors uncovered, and two robustness checks. All AI-assisted and AI-led teams used ChatGPT during the event and had to provide their AI conversation history (i.e., a transcript of all prompts and responses exchanged with ChatGPT).

Participants were offered coauthorship on this paper, independent of their team's performance or success in reproducing results. No monetary compensation or performance-based incentives were provided. While this may have led to reduced effort for some teams, it also reduced incentives for strategic behavior or protocol violations, particularly for AI-led teams who were asked not to directly examine the article, code, or data.

Access to a paid subscription of ChatGPT (powered initially by GPT-4 and subsequently by other models) was provided to all members in the AI-assisted and AI-led teams. While ChatGPT had six different versions available between February 14th 2024 (training for our pilot) and our final event on November 22nd 2024, researchers had access to the main flagship models (GPT-4, and/or GPT-4o). These models were capable of processing files, equipped with a Python environment for interpreting code and conducting data analysis, and had internet access. Additional information on the

## Significance

Verifying results of published social sciences research is essential but expensive, costing hundreds of dollars per study. With AI tools like ChatGPT becoming widespread, we tested whether they could help scientists check if research findings can be reproduced. We assigned 288 researchers to 103 teams working with no AI, with AI as an assistant, or AI leading the work with minimal human input. Human teams and AI-assisted teams performed similarly on most tasks, but humans caught more critical errors. AI working autonomously achieved a 37% reproduction rate, making it potentially useful for automated screening when human review is cost-prohibitive. These results nonetheless show that human expertise remains essential for reliable scientific validation.

different versions and use by teams at events are included in [SI Appendix](#), ChatGPT Models.

AI-assisted and AI-led teams took part in a mandatory one-hour training on the usage of ChatGPT. Participants viewed the training live or later via recording. The AI training was optional for human-only teams. The training had nine components which we outline here but describe further in [SI Appendix](#), AI Training: 1) Introduction, Overview of ChatGPT and Access; 2) Interaction with ChatGPT; 3) Sharing Chats with I4R; 4) Coding Assistance; 5) Uploading Files and Images; 6) Conducting Data Analysis Using ChatGPT; 7) ChatGPT API; 8) Customizing ChatGPT; and 9) Explanation of Differences Among ChatGPT Models. AI-led and AI-assisted teams constructed their own prompts but were given examples and best-practice guidance in the training session. Using textual analysis on all prompts, we show limited overlap in prompt wording across AI-led and AI-assisted teams ([SI Appendix](#)).

The human-only teams were not allowed to use ChatGPT or any other AI tool. The AI-assisted teams were allowed to use ChatGPT without limitation (but no other AI tool). AI-led teams had to perform the tasks only using the guidance of ChatGPT. They were not allowed to read the article or look at the data and code but could ask ChatGPT to summarize the article. They had to upload the article to ChatGPT along with an image of the table(s)/figure(s) to be reproduced, the replication code, and the data files where feasible. They were asked to first use ChatGPT's Python interpreter module to conduct the analysis. However, they were allowed to run analysis code locally (in R or Stata) when ChatGPT failed to run the analysis itself. When running code locally, the teams were not allowed to use any other code except the one provided by ChatGPT, though the teams could adjust file paths and their environment without the assistance of ChatGPT. During the pregame AI training, participants were shown examples of how to upload the article and replication files to ChatGPT and how to use the Python interpreter module. We relied on the integrity of the AI-led teams to not look at the studies, code, or files. That is, we asked them to pass everything through ChatGPT; we did not give specific guidance on how teams should operate. Teammates could work independently or jointly throughout the event.

In summary, we have 103 teams: 33 human-only teams (92 researchers), 35 AI-assisted teams (93 researchers), and 35 AI-led teams (103 researchers). [SI Appendix](#), [Table S8](#) shows the treatment arms are balanced across observables.

**1.1. Three Tasks.** Participants had three objective tasks with measurable outcomes. First, teams were asked to computationally reproduce a few selected results in the study assigned to them. The numerical results were selected by the lead author, AB, based on their relative importance to the main claims of the article. Computational reproducibility involves using the same data as the original authors and running their code. In the templated log, teams recorded the time taken to computationally reproduce the numerical result. Notably, AB, JA, and DM were able to computationally reproduce the selected results before the event, requiring only minimal adjustments (e.g., updating file paths). We have two different dependent variables for computational reproducibility: one outcome as a binary variable (completed computational reproducibility vs. did not complete), and one that is time (in minutes) from the start of the event to when teams completed a computational reproduction. A computational reproduction is defined as the successful execution

of the original authors' codes and the production of numerical results in line with those in the article.

Second, we compare how effective different team types were in finding coding errors or data irregularities. For simplicity, we refer to these as "errors." We categorize errors as major or minor based on whether they could, in theory, have an impact on the claims tested. For instance, a coding error or data irregularity that impacts the dependent or independent variables is considered a major error, as it could have an impact on the estimation results. In contrast, minor coding errors are typically easily fixed by the reproducers and do not impact the validity of the claims made by the original authors. In a set of exploratory analyses, we also categorize coding errors along three dimensions: i) whether the error occurs in preparing the data and analysis, ii) whether the error is related to the regression analysis, and iii) whether it is a transcription error (e.g., a mismatch between the coefficient reported in the article and the coefficient produced by the code, such as  $-0.034$  vs.  $0.034$ ). We also investigate the extent of false error detection and the share of errors not uncovered by the treatment arm.

Third, we asked each team to propose and perform two robustness checks. A robustness check is defined in our study as an additional statistical computation. We instructed that these robustness checks should not repeat ones already mentioned in the study or its [SI Appendix](#), that they should be feasible, and that heterogeneity analysis (e.g., comparing female and male respondents) was not considered a robustness check.

Defining what makes a robustness check "good" or "bad" is not straightforward. We define four binary criteria for evaluating the quality of robustness checks: i) clarity of purpose and execution; ii) feasibility; iii) novelty (i.e., not previously done by the original authors); and iv) relevance to the validity of the empirical strategy. Items i) through iii) are basic necessary conditions. Item iv) requires that the purpose of the robustness test is to provide evidence regarding the credibility of the empirical strategy (39–41). All four criteria must be met for a robustness check to be considered "good." Additionally, running corrected code in an attempt to correct major errors in the original paper is coded as a "good" robustness check, regardless of whether it complies with the previous criteria.

We measure differences by team type in proposing and implementing robustness tests using four measures. The first two are whether teams proposed one or two "good" robustness checks. The third and fourth dependent variables are whether the participants report to have implemented one or two of those "good" robustness checks, respectively.

## 2. Results

Our analyses were preregistered after the pilot event in Toronto. We list deviations from our preregistration in [SI Appendix](#) and note throughout whether the analysis is exploratory.

**2.1. Computational Reproducibility.** Our main finding is that computational reproducibility rates varied substantially across the groups. Most human-only (94%; 31/33) and AI-assisted (91%; 32/35) teams could computationally reproduce the results, while only 37% (13/35) of AI-led teams were able to do so ([Table 1](#)). [Table 2](#) shows the ordinary least squares (OLS) estimates of our main regression model (see [SI Appendix](#), [Table S10](#) for logit and Poisson regressions and [SI Appendix](#), [Table S12](#) for coefficient estimates concerning the control variables). We find that human-only teams are about 59 percentage points more likely than

**Table 1. Comparison of human, AI-assisted, and AI-led metrics**

Variable	Human-only	AI-assisted	AI-led	Human-only vs. AI-assisted	Human-only vs. AI-led	AI-assisted vs. AI-led
Reproduction	0.939 (0.242)	0.914 (0.284)	0.371 (0.490)	0.025 [0.697]	0.568 [<0.001]	0.543 [<0.001]
Minutes to reproduction	82.0 (39.8)	93.3 (85.4)	179.7 (68.4)	-11.3 [0.505]	-97.7 [<0.001]	-86.4 [0.002]
Number of minor errors	1.000 (1.658)	1.400 (2.488)	0.686 (1.605)	-0.400 [0.441]	0.314 [0.430]	0.714 [0.158]
Minutes to first minor error	141.6 (97.0)	139.9 (83.1)	157.6 (85.6)	1.7 [0.960]	-16.0 [0.691]	-17.7 [0.622]
Number of major errors	1.697 (2.568)	0.743 (1.120)	0.229 (0.547)	0.954 [0.049]	1.468 [0.002]	0.514 [0.017]
Minutes to first major error	110.5 (69.5)	130.3 (86.9)	152.8 (94.3)	-19.8 [0.487]	-42.3 [0.261]	-22.5 [0.606]
At least one good robustness check	1.000 (0.000)	1.000 (0.000)	0.829 (0.382)	0.000 [NA]	0.171 [0.012]	0.171 [0.010]
At least two good robustness checks	0.879 (0.331)	0.857 (0.355)	0.629 (0.490)	0.022 [0.796]	0.250 [0.017]	0.229 [0.029]
Ran at least one good robustness check	0.939 (0.242)	0.943 (0.236)	0.571 (0.502)	-0.003 [0.953]	0.368 [<0.001]	0.371 [<0.001]
Ran at least two good robustness checks	0.788 (0.415)	0.800 (0.406)	0.457 (0.505)	-0.012 [0.903]	0.331 [0.005]	0.343 [0.003]

Note: Columns 2–4 present means and SEs in parentheses for individual groups (Human-only, AI-Assisted, and AI-Led); columns 5–7 present differences in means and *P*-values in brackets for group comparisons (Human-Only vs. AI-Assisted, Human-Only vs. AI-Led, and AI-Assisted vs. AI-Led).

AI-led teams to successfully computationally reproduce the results ( $P < 0.001$ ). In contrast, there is no statistically significant difference between human-only and AI-assisted teams ( $P = 0.771$ ).

We next investigate how the distribution of time-to-computational reproduction varies across groups. Fig. 1 plots complementary Kaplan–Meier curves showing, by treatment arm, how long teams took to reproduce their paper by the end of the event. The proportion of teams that reproduce their paper does not reach 100% after seven hours in any treatment arm because all treatment arms contain some teams who could not reproduce their paper. This is especially noticeable for

AI-led teams. We find that human-only and AI-assisted teams are significantly faster than AI-led teams (Table 1). There is no statistically significant difference between human and AI-assisted teams.

In an exploratory analysis, we investigate whether AI-assisted and AI-led teams improved over time. In our setting, improvements could be due to new ChatGPT versions and increased researchers' skills over time. In *SI Appendix, Fig. S2*, we show the difference in computational reproducibility rates between the treatment groups by event. Visually, AI-led teams did not improve over time when compared to human-only teams during the first five events in 2024. We observe that the reproducibility

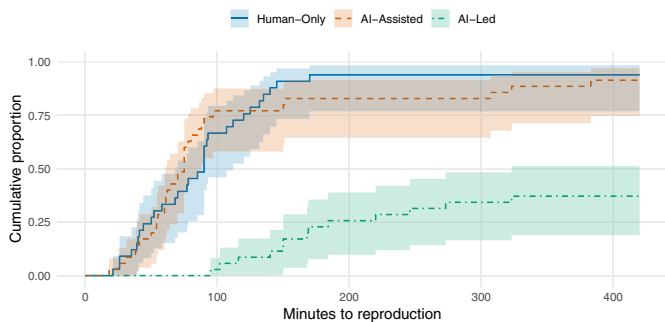
**Table 2. Causal relationship between treatment groups and reproducibility outcomes**

	(1) Reproduction	(2) Minor errors	(3) Major errors	(4) One good robustness	(5) Two good robustness	(6) Ran one robustness	(7) Ran two robustness
AI-assisted	-0.018 (0.063) [-0.144; 0.107]	0.313 (0.387) [-0.458; 1.083]	-1.022*** (0.362) [-1.743; -0.300]	-0.009 (0.027) [-0.063; 0.046]	-0.014 (0.103) [-0.220; 0.191]	-0.032 (0.061) [-0.155; 0.090]	-0.009 (0.113) [-0.233; 0.216]
AI-led	-0.593*** (0.090) [-0.773; -0.413]	-0.331 (0.350) [-1.029; 0.366]	-1.344*** (0.342) [-2.024; -0.664]	-0.167** (0.068) [-0.302; -0.031]	-0.250** (0.107) [-0.463; -0.037]	-0.323*** (0.098) [-0.518; -0.127]	-0.290** (0.126) [-0.540; -0.040]
Controls	✓	✓	✓	✓	✓	✓	✓
Mean of dep. var	0.738	1.029	0.874	0.942	0.786	0.816	0.680
<i>P</i> -val (AI-assisted = AI-led)	0.000	0.115	0.251	0.021	0.032	0.003	0.017
Obs.	103	103	103	103	103	103	103

Note: SEs in parentheses; CIs in brackets. Human-only group omitted.

Controls: number of teammates; game-by-software fixed effects; maximum and minimum position skill fixed effects; attendance fixed effects.

\* $P < 0.10$ , \*\* $P < 0.05$ , and \*\*\* $P < 0.01$ .



**Fig. 1.** Complementary Kaplan-Meier curves, showing the proportion of teams who computationally reproduced the paper by time  $t$  along with curve bands (95% CIs).

rate gap between human-only and AI-led teams was over 50 percentage points for most events in 2024. Of note, this gap had slightly narrowed by the final event of 2024.

**2.2. Coding Errors or Data Irregularities.** We have two primary dependent variables concerning coding error detection: counts of major and minor errors detected. We find that human-only teams identified on average 1.00 minor and 1.70 major errors, compared with 1.40 minor and 0.74 major errors for AI-assisted teams and 0.69 minor and 0.23 major errors for AI-led teams, respectively (Table 1). Table 2 provides OLS estimates indicating that, compared to AI-assisted and AI-led teams, human-only teams uncovered more major errors ( $P = 0.006$  and  $P < 0.001$ , respectively). The difference in the number of minor coding errors detected is, however, not significant ( $P = 0.421$  and  $P = 0.347$ , respectively). We further find that AI-assisted teams uncovered more minor errors than AI-led teams, but the estimate is not significant at any conventional level ( $P = 0.115$ ). SI Appendix provides examples of errors and a discussion.

Fig. 2 plots complementary Kaplan-Meier curves showing how long teams took to find a first minor error (Top panel) and a first major error (Bottom panel). We find that the speed at which AI-assisted teams uncover a first (minor or major) error is not statistically significantly different from that of human-only teams and that AI-led teams are statistically significantly slower than human-only teams at uncovering a first major error.

Our findings suggest that human-only teams were more effective at detecting both major and minor errors compared to AI-led teams, highlighting a challenge in AI-led teams' ability to autonomously navigate and interpret complex code and detect data irregularities.

In exploratory analyses, we explore whether AI-led and AI-assisted teams are better at uncovering different types of errors, distinguishing between coding mistakes that require substantive understanding of the paper and those that do not. Table 3 provides OLS estimates indicating that, compared to AI-led teams, human-only teams uncovered more errors that occur in preparing the data and analysis (although not significantly different,  $P = 0.165$ ), more errors related to the regression analysis ( $P = 0.007$ ) and more transcription errors ( $P = 0.034$ ). Human-only teams uncover more errors in these three categories than AI-assisted teams, but only one of the point estimates is statistically significant at the 10% level ( $P = 0.993$ ,  $P = 0.072$ , and  $P = 0.318$ ).

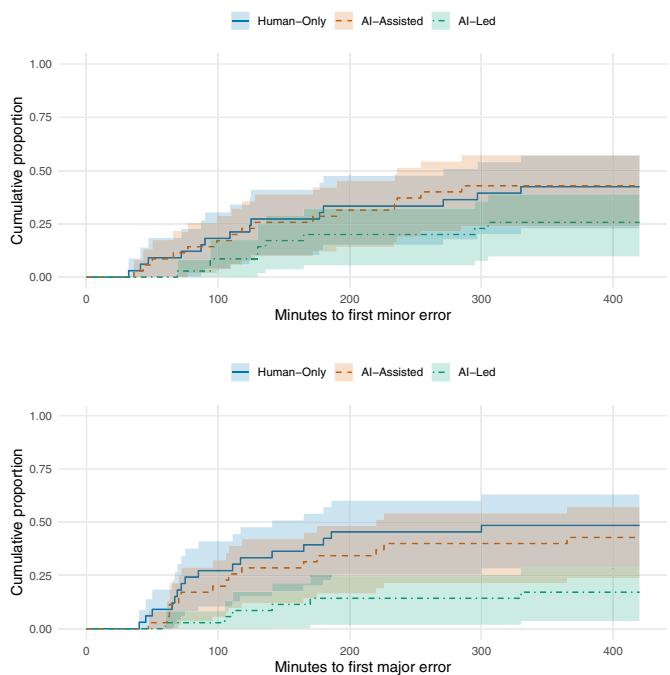
Our qualitative analysis (Section 2.5) suggests that some AI-led teams experienced prompt fatigue and hallucinated reasoning paths. This qualitative evidence motivates our exploratory analysis of whether AI-assisted and AI-led teams are more likely to

produce false error detection. We find no evidence that this is the case ( $P = 0.248$ ,  $P = 0.642$ ), suggesting that hallucinations occur at other stages of the reproduction pipeline. This previous result may mask the fact that many AI-led teams did not detect any errors. We thus investigate the proportion of errors that remain undetected by each team. We find that AI-led teams detected a significantly smaller proportion of known errors than human-only and AI-assisted teams ( $P < 0.001$ ,  $P = 0.016$ ). These results suggest that AI-led teams' primary limitation lies in error discovery rather than erroneous overdetection.

We also provide noncausal evidence in exploratory analyses in SI Appendix, Table S14 that AI-assisted teams with more AI experience uncovered coding errors faster, although these estimates are only statistically significant at the 10% level ( $P = 0.067$  and  $P = 0.070$ ). Extending this analysis, SI Appendix, Table S15 compares human-only teams with AI-assisted teams with high vs. low/medium AI experience. These comparisons should be interpreted with caution, as the number of AI-assisted teams in each subgroup is small, resulting in limited statistical power. Nonetheless, the point estimates are consistent with the hypothesis that AI experience improves the effectiveness of AI-assisted teams. AI-assisted teams with high AI experience appear to uncover coding errors faster than human-only teams and detect more minor errors on average. The magnitudes of these differences are sizable, but the estimates are imprecise and not statistically significant at conventional levels. These findings are consistent with the behavioral evidence presented in Section 3.4, which examines how AI-assisted teams used ChatGPT.

We also find that Stata teams uncovered significantly more major errors ( $P < 0.001$ ), with the human-only groups using Stata finding significantly more major errors than all other groups (SI Appendix, Table S13).

In an exploratory analysis, we investigate if the performance of AI-led teams in detecting errors improved over time. SI Appendix, Figs. S4 and S6 suggest no improvement of AI-led teams relative to human-only teams over the year 2024.



**Fig. 2.** Complementary Kaplan-Meier curves, showing the proportion of teams who found their first coding error by time  $t$  along with curve bands (95% CIs).

**Table 3. Causal relationship between treatment groups and error types**

	(1) Preregression errors	(2) Regression errors	(3) Transcription/postregression errors	(4) False error detection	(5) Share of known errors not found
AI-assisted	-0.002 (0.267) [-0.534; 0.530]	-0.604* (0.331) [-1.263; 0.055]	-0.343 (0.341) [-1.023; 0.337]	-0.359 (0.309) [-0.974; 0.256]	0.042 (0.048) [-0.053; 0.137]
AI-led	-0.407 (0.290) [-0.986; 0.171]	-0.886*** (0.321) [-1.524; -0.248]	-0.652** (0.301) [-1.251; -0.052]	0.221 (0.473) [-0.721; 1.162]	0.151*** (0.043) [0.065; 0.236]
Controls	✓	✓	✓	✓	✓
Mean of dep. var	0.786	0.660	0.583	0.680	0.846
<i>P</i> -val (AI-assisted = AI-led)	0.140	0.228	0.274	0.146	0.016
Obs.	103	103	103	103	103

Note: SEs in parentheses; CIs in brackets. Human-only group omitted.

Controls: number of teammates; game-by-software fixed effects; maximum and minimum position skill fixed effects; attendance fixed effects.

\* $P < 0.10$ , \*\* $P < 0.05$ , and \*\*\* $P < 0.01$ .

**2.3. Proposed Robustness Checks.** We find a clear, consistent performance hierarchy across both conditions: Human-only and AI-assisted teams outperform AI-led teams. We find that all human-only (33/33) and AI-assisted (35/35) teams proposed at least one good robustness check, whereas only 83% (29/35) of AI-led teams did so. Table 2 provides OLS estimates and show that the difference between AI-led groups and the other two groups is statistically significant ( $P = 0.017$  and  $P = 0.021$ , respectively). We find that 29 of 33 human-only and 30 of 35 AI-assisted teams suggested two good checks, compared with just 22 of 35 AI-led teams (Table 2,  $P = 0.022$  and  $P = 0.032$ ).

Looking at whether teams report to have implemented those checks, AI-led teams were almost 32 percentage points less likely than the other two groups to report having conducted a robustness check that was classified as “good” ( $P = 0.002$  and  $P = 0.003$ ), and six AI-led teams supplied no robustness checks evaluated as “good” at all. These six teams’ checks were judged as “bad” mostly because of a lack of clarity and duplicating analyses already run by the original authors.

Our results indicate that AI-led teams, while able to produce robustness checks with some level of quality, faced more challenges in aligning with the criteria. These difficulties may stem from omission of relevant information when describing the task to the AI or from limited ability of the AI to interpret the empirical strategy and to assess the feasibility of the checks.

**2.4. Additional Analyses for AI-Assisted Teams.** SI Appendix, Table S16 presents an exploratory correlational analysis examining the relationship between AI usage (measured by total prompts) and performance in AI-assisted teams. See SI Appendix, Fig. S7 for descriptive statistics on AI usage for AI-assisted teams. For this analysis, we divided teams into lower and higher AI-usage groups using a median split based on the total number of prompts they employed.

The findings indicate that AI-assisted teams with lower AI usage were less likely to achieve computational reproduction of the original results and uncovered less major and minor coding errors. Of note, our sample is small and none of the differences are statistically significant at the 5% level. To further explore potential mechanisms behind this heterogeneity, SI Appendix, Table S14 also reports differences in prompting behavior by AI experience. We find that AI-assisted teams with lower AI

experience tend to interact with ChatGPT more frequently, as measured by the number of prompts, but the difference is not statistically significant due to the small sample size ( $P = 0.307$ ). This pattern is consistent with the idea that less experienced teams rely more heavily on iterative prompting, which may contribute to longer task completion times. These results relate to a literature studying overreliance on AI support (20, 36, 37, 42–44).

**2.5. Focus Groups.** In additional exploratory analysis, between 18 April and 30 April 2025, we conducted six one-hour focus groups ( $n = 25$ ) involving AI-led ( $n = 8$ ), AI-assisted ( $n = 11$ ), and human-only ( $n = 6$ ) participants. The participants were aware of the headline quantitative results. While this creates risk of confirmation bias and demand characteristics, we addressed this by emphasizing process, task allocation, and failure points rather than outcomes in the discussion guide, and by treating focus group material as explanatory and triangulatory evidence rather than independent support for treatment differences. Accordingly, we use qualitative themes to illuminate mechanisms behind observed patterns, and we flag any tensions between participant claims and the experimental results. Consistent with this stance, where participant views exceeded what the quantitative results support, we report the discrepancy rather than treat it as confirmation.

Thematic analysis revealed the following patterns: AI-assisted participants reported that AI assistance sped things up, while AI-led participants reported the opposite. Human-only participants believed they were the most effective at detecting major errors, with AI-led participants trailing. AI-assisted teams strategically outsourced microtasks, for example boilerplate code and file location, while retaining conceptual control, whereas AI-led teams were required to cede entire analytic stages to ChatGPT and struggled when automation failed.

Data illuminated the practical consequences of these differences. Initial optimism about LLMs quickly gave way to prompt fatigue by participants, model’s overconfidence, and mounting frustration, especially among AI-led participants who faced hallucinated paths, truncated context windows, and prolonged debugging loops. AI-assisted teams report that human expertise remained necessary for detecting subtle errors and for arbitrating disagreements between AI output and reality. Nevertheless, when used judiciously, LLMs accelerated routine work, suggested

robustness checks, and broadened analytical ambition for less experienced coders. Therefore, our focus group findings imply that effective LLM prompting is becoming a specialized research skill and that near term gains will come from augmenting, not replacing, human judgment. Additional details on methodology and results from the focus group analysis are provided in *SI Appendix*.

### 3. Discussion

Computational reproducibility, error detection, and robustness checks are essential components of empirical research validation, but are resource-intensive tasks. Ensuring that research can be reproduced is financially demanding. Recent research suggests that, across 10 top economics journals, the average expense of reproducing a single study is about USD \$365 (18). For the American Economic Association, data-editor activities cost approximately USD \$750 per article (9). Against this backdrop, our comparative analysis of human-only, AI-assisted, and AI-led teams sheds light on how AI may be integrated into the costly reproducibility pipeline, potentially accelerating some stages of the process and reshaping how replication labor is allocated.

A key finding of our study is that AI-led teams were able to successfully computationally reproduce approximately 37% of results. This result is nontrivial and suggests that a first automated pass at computational reproducibility was already within reach for a meaningful subset of empirical work in 2024.

At the same time, our results temper expectations of immediate, widespread AI autonomy in reproducibility. While recent advances in large language models have expanded the scope for AI integration in research workflows (45, 46), AI-led and AI-assisted teams do not yet outperform human-only teams on average. Moreover, current AI deployments introduce additional costs, such as paid model subscriptions, without consistently delivering higher success rates. As a result, fully autonomous AI reproduction does not yet offer clear cost savings relative to experienced human researchers.

However, our study likely represents only a lower bound of the capabilities of a more fully developed autonomous AI replication system. In practice, an AI system could deploy more sophisticated prompting strategies, exploit parallel experimentation, and possibly be supervised by trained research assistants or undergraduate students rather than senior researchers, reducing labor costs while maintaining acceptable levels of oversight. Our findings thus imply that future iterations of AI-led reproducibility systems may achieve higher success rates without proportional increases in human effort.

This perspective reframes AI not as a replacement for human expertise, but as a tool for redistributing effort across stages of the reproducibility pipeline. AI systems may handle routine debugging, error detection, and preliminary robustness checks (47–49), while human researchers focus on interpretation, judgment, and more complex failures. Under this model, even partial automation can generate meaningful cost savings and efficiency gains at scale.

**3.1. Summary of Findings.** AI-led teams faced notable challenges compared to both AI-assisted and human-only teams. Only 37% of AI-led teams were able to successfully complete computational reproducibility, highlighting a substantial gap in the capacity of AI in 2024 to autonomously guide researchers through complex quantitative analyses. Similarly, in error detection, AI-led teams documented significantly fewer major and minor errors than either AI-assisted or human-only teams. These findings

underscore the importance of still integrating human expertise. As LLMs continue to evolve, sustained benchmarking against humans will be crucial to ensure that future AI-led efforts close and potentially surpass the existing performance gap.

**3.2. Limitations.** One limitation is our sole focus on OpenAI's ChatGPT, meaning that we cannot generalize to all current AI models. Furthermore, the limited timeframe of seven hours for study teams to complete their reproductions may not adequately reflect the conditions under which reproducibility efforts are conducted depending on the field of science. In addition, participant incentives and attribution dynamics may have encouraged some teams to minimize time or effort, potentially increasing overreliance on AI tools. Finally, our analysis is based on a small, nonrandom set of studies spanning a limited range of social science methodologies and replication difficulty levels; although we provide detailed proxies for task complexity and error types, this sample composition constrains the extent to which our findings on AI assistance generalize across papers of different difficulty and across other scientific fields (*SI Appendix, Table S5*).

We note that participant behavior may have been influenced by observation and professional identity, generating a Hawthorne-type effect. Researchers with strong coding skills or a personal stake in reproducibility may have exerted greater effort in human-only teams, while responsibility may have been partially shifted to the AI in AI-assisted or AI-led settings. While this could bias relative performance comparisons, it may also reflect real-world incentive and attribution dynamics that shape how AI tools are adopted in research practice.

**3.3. Implications for Human-AI Collaboration in Research.** Our findings support the notion that, while AI tools hold promise for aiding in reproducibility tasks, the state of technology as of late 2024 is not yet advanced enough for full autonomy in complex empirical workflows. Human expertise remains critical to navigate challenges and provide interpretative guidance for reproducibility and error detection. The AI-assisted model—where humans work alongside AI tools—did not emerge as a winner over human-only teams in overall outcomes but outperformed AI-led teams on most of our outcomes.

In scenarios where computational reproducibility, error detection, and robustness checks require in-depth understanding, domain knowledge, and flexible problem-solving, human involvement currently adds value. The ability to contextualize, interpret, and implement complex quantitative research remains a human strength, highlighting the limits of current AI in fully autonomous reproduction.

**3.4. Outlook.** Advancements in models and further optimization of AI for reproduction may soon address the limitations we reported. Future advancements in models optimized through reinforcement learning to solve reasoning problems using chain of thought could address the limitations we reported, possibly improving the model's ability to reproduce complex quantitative research through iterative, reasoning-driven processes.

Future research should consider the potential for training models specifically in social science and quantitative research contexts. Current LLMs are trained on vast datasets but may lack specificity in understanding the unique demands of empirical social science research. AI systems tailored for social science reproduction (e.g., with native support for R and Stata) could potentially improve reproducibility outcomes, reducing the

barriers AI currently faces in autonomously handling the nuances of quantitative research. Additionally, incorporating continuous feedback and learning mechanisms could allow AI-assisted and AI-led teams to improve performance over time, as AI learns from each reproduction task and adapts based on human feedback.

Future research should also focus on analyzing which prompting strategies leads to successful reproductions and which paths lead to failure, insights that could inform the development of AI systems better tailored for social science research. We make the chat transcripts publicly available and conduct an exploratory analysis of ChatGPT transcripts in *SI Appendix*.

## 4. Materials and Methods

Participants in the AI replication games experiments coauthor this study. The University of Ottawa Office of Research. Our preanalysis plan was preregistered on the Open Science Framework (OSF) on May 2nd, 2024, after our pilot event at the University of Toronto (<https://osf.io/sz2g8/>). AI-assisted and AI-led teams took part in a mandatory one-hour ChatGPT training, while the same training was optional for human-only teams; slides and recordings are available on OSF. A version-tagged copy of the code and data is permanently archived at <https://github.com/I4Replication/AI-Games>, and we make our AI training materials and recording, data and code, preanalysis plan, and template form available at <https://osf.io/sz2g8/> with no restrictions on sharing or reuse.

## 5. Research Ethics Boards

Participants in the AI replication games experiments coauthor this study. The University of Ottawa Office of Research Ethics and Integrity reviewed and approved our AI games (H-09-25-12041). The King's College London Research Ethics Office reviewed and approved our focus groups (MRA-24/25-48393). All participants provided informed consent.

**Data, Materials, and Software Availability.** We make our i) AI training materials and recording, ii) data and code, iii) preanalysis plan and iv) template form available here: <https://github.com/I4Replication/AI-Games> (50). We declare no restrictions on sharing or reuse.

**ACKNOWLEDGMENTS.** We would like to thank Gabriel Zimmerman for research assistance. This research and AI replication games were funded by Coefficient Giving project "Benchmarking LLM agents on real-world tasks: Reproducibility" and the Alfred P. Sloan Foundation Foundation grant G-2023-22326. We also benefited from funding to host games from the Universities of Toronto, Ottawa, Cornell, and Tilburg. Mahmoud Elsherif acknowledges funding from Leverhulme Early Career Research Fellowship-ECF-2022-761. S. Akhtar acknowledges funding DP200102935 awarded by the Australian Research Council Grant. R. Seri acknowledges funding from project "Dipartimento di Eccellenza 2023-2027" awarded by the Ministero dell'Università e della Ricerca.

Author affiliations: <sup>a</sup>Department of Economics, Faculty of Social Sciences, University of Ottawa, Ottawa, ON K1N 6N5, Canada; <sup>b</sup>Institute for Replication, Ottawa, ON K1N 6N5, Canada; <sup>c</sup>Institute for Technology and Humanity, University of Cambridge, Cambridge CB2 1SB, United Kingdom; <sup>d</sup>Department of Statistical Sciences, Faculty of Information, University of Toronto, Toronto, ON M5S 3G6, Canada; <sup>e</sup>Department of Management and Marketing, University of Melbourne, Melbourne, Carlton, VIC 3010, Australia; <sup>f</sup>School of Psychology, University of Sheffield, Sheffield S1 4DP, United Kingdom; <sup>g</sup>Research on Research Institute, London WC1E 6JA, United Kingdom; <sup>h</sup>Department of Economics, Cornell University, Ithaca, NY 14853; <sup>i</sup>Climate and Development Policy Division, Leibniz Institute for Economic Research - Leibniz Institute for Economic Research, Essen, NRW 45128, Germany; <sup>j</sup>Robert C. Vackar College of Business and Entrepreneurship, University of Texas Rio Grande Valley, Edinburg, TX 78539; <sup>k</sup>Norwich Business School, Accounting and Quantitative Methods, University of East Anglia, Norwich NR4 7TJ, United Kingdom; <sup>l</sup>Department of Economics, Carleton University, Ottawa, ON K1S 5B6, Canada; <sup>m</sup>Statistics Canada, Ottawa, ON K1A 0T6, Canada; <sup>n</sup>GovAI, N1 9JY, London, United Kingdom; <sup>o</sup>Department of Experimental Psychology, University of Oxford, Oxfordshire OX1 3EL, United Kingdom; <sup>p</sup>Finance Discipline, University of Sydney Business School,

The University of Sydney, Sydney, NSW 2006, Australia; <sup>q</sup>Department of Actuarial Studies and Business Analytics, Macquarie Business School, Macquarie University, Sydney, NSW 2019, Australia; <sup>r</sup>Department of Socioeconomics, Vienna University of Economics and Business, Vienna 1020, Austria; <sup>s</sup>Center for Research on Equitable and Open Scholarship, Massachusetts Institute of Technology, Cambridge, MA 02139; <sup>t</sup>Technical University of Munich School of Management, Technical University of Munich, Munich 80333, Germany; <sup>u</sup>Centre for Business Research, University of Cambridge, Cambridge CB2 1QA, United Kingdom; <sup>v</sup>Independent researcher, Shiraz, Iran; <sup>w</sup>Facultad de Economía, Universidad del Rosario, Bogotá 111711, Colombia; <sup>x</sup>School of Economics, Universidad del Rosario, Bogotá 111711, Colombia; <sup>y</sup>International Center for Higher Education Research and Faculty of Economics, University of Kassel, Kassel, Hessen 34125, Germany; <sup>z</sup>Department of Economics, University of Ghana, Legon, Accra, Ghana; <sup>aa</sup>Department of Economics, University of Basel, Basel 4052, Switzerland; <sup>ab</sup>Department of Methodology and Statistics, Tilburg School of Social and Behavioral Sciences, Tilburg University, Tilburg 5037 AB, The Netherlands; <sup>ac</sup>Department of Sport, Tourism and Hospitality Management, Temple University, Philadelphia, PA 19312; <sup>ad</sup>Warwick Business School, University of Warwick, Coventry CV4 7AL, United Kingdom; <sup>ae</sup>Center for Economic Policy Research, Coldbath Square, London EC1R 5HL, United Kingdom; <sup>af</sup>Department of Economics, University of Zurich, Zurich CH-8001, Switzerland; <sup>ag</sup>Department of Behavioural Sciences and Learning, Linköping University, Linköping 58183, Sweden; <sup>ah</sup>Centre for Experimental Social Sciences, Nuffield College, University of Oxford, Oxford OX1 1NF, United Kingdom; <sup>ai</sup>Department of Economics, University of Southampton, Southampton SO17 1BJ, United Kingdom; <sup>aj</sup>Department of Finance, Norwegian School of Economics, Bergen 5045, Norway; <sup>ak</sup>Faculty of Psychology, Atma Jaya Catholic University of Indonesia, Jakarta 12930, Indonesia; <sup>al</sup>Department of Education, University of York, York YO10 5DD, United Kingdom; <sup>am</sup>Karlstad Business School, Karlstad University, Karlstad SE-651 88, Sweden; <sup>an</sup>Center for Societal Risk Research, Karlstad University, Karlstad 65188, Sweden; <sup>ao</sup>Faculty of Biology Medicine and Health, The University of Manchester, Manchester M13 9NT, United Kingdom; <sup>ap</sup>School of Business and Economics, Maastricht University, Maastricht 6211 LM, The Netherlands; <sup>aq</sup>Department of Political Science, University of Chicago, Chicago, IL 60637; <sup>ar</sup>Centre for Environmental Sciences, Hasselt University, Hasselt 3500, Belgium; <sup>as</sup>International Center for Higher Education Research, University of Kassel, Kassel, Hessen 34125, Germany; <sup>at</sup>Meta-Research Innovation Center at Stanford, Stanford University, Stanford, CA 94305; <sup>au</sup>Department of Economics, Nazarbayev University, Astana 010000, Kazakhstan; <sup>av</sup>Facultad de Economía, Universidad de los Andes, Bogotá 111711, Colombia; <sup>aw</sup>University College London School of Management, University College London, London E14 5AA, United Kingdom; <sup>ax</sup>Department of Economics, Universidad de Los Andes, Bogotá 111711, Colombia; <sup>ay</sup>United Methodist University Sierra Leone, Freetown, Sierra Leone; <sup>az</sup>Department of Economics, International School of Management University of Management and Economics, Vilnius LT-01103, Lithuania; <sup>aaa</sup>Département de Médecine Sociale et Préventive, Université Laval, Québec, QC G1V0A6, Canada; <sup>bbb</sup>Département de Médecine Sociale et Préventive, Université de Montréal, Montréal, QC H3N1X9, Canada; <sup>ccc</sup>Department of Management, Marketing and Information Systems, Hong Kong Baptist University, Kowloon Tong, Hong Kong Special Administrative Regions of China; <sup>ddd</sup>Department of Economics, College of Management Sciences, Michael Okpara University of Agriculture, Umudike, Abia State 440109, Nigeria; <sup>eee</sup>Business School, University of Technology Sydney, Sydney, NSW 2007, Australia; <sup>fff</sup>Department of Economics, Wilfrid Laurier University, Waterloo, ON N2L 3C5, Canada; <sup>ggg</sup>Department of Medical Statistics, The London School of Hygiene & Tropical Medicine, London WC1E 7HT, United Kingdom; <sup>hhh</sup>Department of Computer Science, University College London, London WC1E 6BT, United Kingdom; <sup>iii</sup>Department of Economics, Faculty of Social Sciences, University of Ottawa, Ottawa ON, Canada; <sup>jjj</sup>Business School, Beijing Normal University, Beijing 100875, China; <sup>kkk</sup>Beijing Normal University, Belt and Road School, Zhuhai, Guangdong 519085, China; <sup>lll</sup>School of Earth and Planetary Science, National Institute of Science Education and Research, Odisha 752050, India; <sup>mmm</sup>Migration & Migrants, Netherlands Interdisciplinary Demographic Institute, The Hague NL-2511 CV, the Netherlands; <sup>nnn</sup>Department of Economics, Universidad del Rosario, Bogotá 111711, Colombia; <sup>ooo</sup>Department of Marketing, Tilburg University, Tilburg 5000 LE, The Netherlands; <sup>ppp</sup>Department of Strategy & Entrepreneurship, Tilburg University, Tilburg 5000 LE, The Netherlands; <sup>qqq</sup>Department of Economics, Binghamton University, Binghamton, NY 13902; <sup>rrr</sup>Economics, Université Paris Dauphine-Paris Sciences et Lettres, Paris CEDEX 16 75775, France; <sup>sss</sup>Swedish Institute for Social Research, Stockholm University, Stockholm 103 91, Sweden; <sup>ttt</sup>Department of Economics and Statistics, Linnaeus University, Växjö 35195, Sweden; <sup>uuu</sup>Department of Marketing, Vienna University of Economics and Business Vienna, Vienna 1020, Austria; <sup>vvv</sup>Financial Stability Department, Bank of Canada, Ottawa, ON K1A 0G9, Canada; <sup>www</sup>Onsi Sawiris School of Business, The American University in Cairo, New Cairo 11835, Egypt; <sup>xxx</sup>Department of Accounting and Control, Hautes études commerciales Lausanne, Lausanne 1015, Switzerland; <sup>yyy</sup>Institute for Accounting, Controlling and Auditing, University of St. Gallen, St. Gallen 9000, Switzerland; <sup>zzz</sup>Monash University, Melbourne, VIC 3168, Australia; <sup>aaaa</sup>School of Psychology and Vision Science, University of Leicester, Leicester LE1 7RH, United Kingdom; <sup>bbbb</sup>Psychology, University of Birmingham, Birmingham B15 2TT, United Kingdom; <sup>cccc</sup>Department of Business and Management Science, Norwegian School of Economics, Bergen 5045, Norway; <sup>dddd</sup>Konjunkturforschungsstelle Swiss Economic Institute, ETH Zurich, Zurich 8092, Switzerland; <sup>eeee</sup>Department of Economics, College of Management Sciences, Michael Okpara University of Agriculture, Umudike 440109, Abia State, Nigeria; <sup>ffff</sup>School of Politics and International Relations, University College Dublin, Dublin 4, Ireland; <sup>gggg</sup>School of Economics and Finance, Victoria University of Wellington, Wellington 6011, New Zealand; <sup>hhhh</sup>Business School, Universidad de los Andes, Bogotá 111711, Colombia; <sup>iiii</sup>Universidad de los Andes, Bogotá 111711, Colombia; <sup>jjjj</sup>Vancouver School of Economics, University of British Columbia, Vancouver, BC V6T 1Z4, Canada; <sup>kkkk</sup>Department of Economics, Tilburg University, Tilburg, North Brabant 5037 AB, The Netherlands; <sup>llll</sup>Department of Economics, School of Business and Economics, Freie Universität Berlin, Berlin 14195, Germany; <sup>mmmm</sup>Department of Business Management, Masaryk University, Brno 602 00, Czech Republic; <sup>nnnn</sup>School of Engineering and Applied Sciences, Harvard University, Cambridge, MA 02138; <sup>oooo</sup>Department of Ethics, Governance, and Society, School of Business and Economics, Vrije Universiteit Amsterdam, Noord-Holland, Amsterdam 1081HV, The Netherlands; <sup>pppp</sup>Tinbergen Institute, Amsterdam, Noord-Holland 1018WB, The Netherlands; <sup>qqqq</sup>Department of Accountancy, Economics and Finance, Heriot-Watt University, Edinburgh EH14 4AS, United Kingdom; <sup>rrrr</sup>Department of Political Science, Université Laval, Québec, QC G1V0A6, Canada; <sup>ssss</sup>laboratoire d'Économie de

Dauphine, Centre Pour la Recherche EconoMique et ses Applications, Paris 75014, France; <sup>ttttt</sup>Department of Political and Social Sciences, European University Institute, Fiesole 50014, Italy; <sup>uuuuu</sup>Health, Nutrition, and Population Global Practice, World Bank, Washington, DC 20433; <sup>vvvvv</sup>Centre for Health Economics, University of York, Heslington YO10 5DD, United Kingdom; <sup>wwwww</sup>Business School, Universidad Adolfo Ibañez, Santiago 7910000, Chile; <sup>xxxxx</sup>School of Chemical, Materials and Biological Engineering, University of Sheffield, Sheffield S1 3JD, United Kingdom; <sup>yyyyy</sup>Institute of Economics, Eötvös Loránd University Centre for Economic and Regional Studies, Budapest 1097, Hungary; <sup>zzzzz</sup>Department of Economics, Central European University, Vienna 1100, Austria; <sup>aaaaa</sup>Department of Economics, Deakin University, Burwood, VIC 3125, Australia; <sup>bbbbb</sup>School of Economics, University College Dublin, Dublin D04 F6X4, Ireland; <sup>ccccc</sup>Centre for Business and Economics Research, University of Coimbra, Coimbra 3000-145, Portugal; <sup>ddddd</sup>Behavioral Science, Gears Institute for Public Policy, Dublin D04 P9C4, Ireland; <sup>eeeee</sup>Department of Law, Economics, and Sociology, "Magna Graecia" University of Catanzaro, Catanzaro 88100, Italy; <sup>fffff</sup>Department of Economics, McMaster University, Hamilton, ON L8S 4M4, Canada; <sup>ggggg</sup>The Canadian Research Data Centre Network, Hamilton, ON L8S4M4, Canada; <sup>hhhhh</sup>Institute for Psychiatry, Psychology & Neuroscience, King's College London, London SE5 9RJ, United Kingdom; <sup>iiiiii</sup>McGovern Institute for Brain Research, Massachusetts Institute of Technology, Cambridge, MA 02139; <sup>jjjjj</sup>Economics, University of California San Diego, La Jolla, California 92092; <sup>kkkkk</sup>Department of Sociology and Human Geography, University of Oslo, Oslo N-0317, Norway; <sup>lllll</sup>School of Computer Science, University of Sheffield, Sheffield S1 4DP, United Kingdom; <sup>mmmmm</sup>University College Dublin Lochlann Quinn School of Business, University College Dublin, Dublin 4, Ireland; <sup>nnnnn</sup>Department of Accounting and Control, University of Lausanne, Lausanne 1015, Switzerland; <sup>ooooo</sup>Mathematics Institute, Leiden University, Leiden, South Holland 2333CA, The Netherlands; <sup>ppppp</sup>Department of Economics and Economic History, Unit of Economic Analysis, Universitat Autònoma de Barcelona, Bellaterra, Cerdanyola de Vallès 08193, Spain; <sup>qqqqq</sup>Department of Economics, Finance, and Legal Studies, University of Alabama, Tuscaloosa, AL 35487; <sup>rrrrr</sup>Institute for Futures Studies, Stockholm 10131, Sweden; <sup>sssss</sup>Department of Economics, University of Toronto, Toronto, ON M5S 3G7, Canada; <sup>ttttt</sup>Computational Social Science, ETH Zürich, Zurich 8057, Switzerland; <sup>uuuuu</sup>Department of Politics and International Relations, University of Oxford, Oxford OX1 3JU, United Kingdom; <sup>wwwww</sup>Department of Real Estate Management, Ted Rogers School of Management, Toronto Metropolitan University, Toronto, ON M5B 2K3, Canada; <sup>xxxxx</sup>School of Psychology, University of Nottingham, Nottingham NG7 2RD, United Kingdom; <sup>xxxxx</sup>Department of Health Sciences and Technology, ETH Zurich, Zurich 8001, Switzerland; <sup>yyyyy</sup>Management School, HEC Liège, University of Liège, Liège 4000, Belgium; <sup>zzzzz</sup>Department of Psychology, University of York, York YO10 5DD, United Kingdom; <sup>aaaaa</sup>School of Psychology, University of Birmingham, Birmingham B15 2SA, United Kingdom; <sup>bbbbb</sup>Political Economy Research Institute, University of Massachusetts Amherst, Amherst, MA 01003; <sup>ccccc</sup>Center for Economic and Policy Research, Washington, DC 20009; <sup>ddddd</sup>Research Institute of Industrial Economics, Stockholm 10215, Sweden; <sup>eeeee</sup>Faculty of Economic Sciences, University of Warsaw, Warsaw 00-927, Poland; <sup>fffff</sup>Financial Stability Department, Climate Analysis Team, Bank of Canada, Ottawa, ON K1A 0G9, Canada; <sup>ggggg</sup>Bart's Life Sciences, Bart's Health National Health Services Trust, London E14 5HJ, United Kingdom; <sup>hhhhh</sup>Queen Mary University of London, London E1 4NS, United Kingdom; <sup>iiiiii</sup>School of Electrical and Computer Engineering, Cornell University, Ithaca, NY 14853; <sup>jjjjj</sup>Carleton University, Ottawa, ON K1S 5B6; <sup>kkkkk</sup>University of Oxford, Oxford OX1 2JD, United Kingdom; <sup>lllll</sup>Department of Management and Organization, School of Business and Economics, Vrije Universiteit Amsterdam, Amsterdam, Noord-Holland 1081 HV, Netherlands; <sup>mmmmm</sup>Department of Economics, University of Freiburg, Freiburg im Breisgau 79085, Germany; <sup>nnnnn</sup>Department of Clinical Research, University of Basel, Basel 4051, Switzerland; <sup>ooooo</sup>University Hospital Basel, Basel 4031, Switzerland; <sup>ppppp</sup>Department of Sociology, Ludwig Maximilian University of Munich, Munich 80539, Germany; <sup>qqqqq</sup>Department of Economics, Faculty of Economics and Statistics, University of Innsbruck, Innsbruck 6020, Austria; <sup>rrrrr</sup>Department of Economics, Turku School of Economics, University of Turku, Turku FI-20014, Finland; <sup>sssss</sup>Department of Health Economics, Ludwig-Maximilians-Universität Munich, Munich DE-80539, Germany; <sup>ttttt</sup>School of Business and Economics, Freie Universität Berlin, Berlin 14195, Germany; <sup>uuuuu</sup>Department of Political Science, University of Toronto, Toronto, ON M5S 3G5, Canada; <sup>wwwww</sup>Department of Food, Agricultural and Resource Economics, University of Guelph, Guelph, ON N1G 2W1, Canada; <sup>xxxxx</sup>Department of Economics, Yale University, New Haven, CT 06520-8268; <sup>xxxxx</sup>Research School of Economics, Australian National University, Canberra, ACT 2600, Australia; <sup>yyyyy</sup>School of Finance, Renmin University of China, Beijing 100872, China; <sup>zzzzz</sup>Department of Advanced Social and International Studies, University of Tokyo, Meguro City, Tokyo 153-8902, Japan; <sup>aaaaa</sup>Dyson School of Applied Economics and Management, Cornell University, Ithaca, NY 14850; <sup>bbbbb</sup>Department of Economics, Knauss School of Business, University of San Diego, San Diego, CA 92110; <sup>ccccc</sup>Department of Economics and Business Administration, University of Cagliari, Cagliari 09124, Italy; <sup>ddddd</sup>Department of Economics, School of Administrative and Economic Sciences, Pontificia Universidad Javeriana, Bogotá 110231, Colombia; <sup>eeeee</sup>Department of Economics, Durham University, Durham DH1 3LB, United Kingdom; <sup>fffff</sup>School of Economics and Management, Tilburg University, Tilburg 5000 LE, The Netherlands; <sup>ggggg</sup>Applied Economics, University of Minnesota, Saint Paul, MN 55108; <sup>hhhhh</sup>Department of Economics, School of Finance, Economics and Government, Universidad Escuela Finanzas Economía y Gobierno, Antioquia, Medellín 050022, Colombia; <sup>iiiiii</sup>Samuel Curtis Johnson School of Business, Cornell University, Ithaca, NY 14853; <sup>jjjjj</sup>School of Economics, University of Sheffield, Sheffield S10 2TU, United Kingdom; <sup>kkkkk</sup>School of Economics and Business, University of Halle, Halle (Saale) 06108, Germany; <sup>lllll</sup>Department of Marketing, Marshall School of Business, University of Southern California, Los Angeles, CA 90089-1424; <sup>mmmmm</sup>Equalis Capital, Paris 75116, France; <sup>nnnnn</sup>Department of Economics, Rice University, Houston, TX 77005; <sup>ooooo</sup>Institute for Financial Management and Research Graduate School of Business, Krea University, Sri City, Andhra Pradesh 517646, India; <sup>ppppp</sup>Department of Psychology, Dilla University, Dilla, South Ethiopia 419, Ethiopia; <sup>qqqqq</sup>Center PhD Students, Research Group: Econometrics, Tilburg School of Economics and Management, Tilburg university, Tilburg 5037 AB, Netherlands; <sup>rrrrr</sup>Department of Economics, University of Ottawa, Ottawa, ON K1N 6N5, Canada; <sup>sssss</sup>Department of Economics, University of Calgary, Calgary, AB T2N 1N4, Canada; <sup>ttttt</sup>Resources for the Future, Washington, DC 20036; <sup>uuuuu</sup>Research Group: Econometrics, Tilburg School of Economics and Management, Tilburg University, Tilburg 5037 AB, The Netherlands; <sup>wwwww</sup>Department of Political Science, University of Montreal, Montreal, QC H3T 1J4, Canada; <sup>xxxxx</sup>Institut de

valorisation des données, Montreal, QC H3N 1V5, Canada; <sup>xxxxx</sup>Department of Medicine and Surgery, University of Milano-Bicocca, Milan 20126, Italy; <sup>yyyyy</sup>Department of Statistics and Quantitative Methods, University of Milano-Bicocca, Milan 20126, Italy; <sup>zzzzz</sup>Department of Health Sciences, University of York, York YO10 5DD, United Kingdom; <sup>aaaaaaa</sup>Climate and Development Policy Division, RWI - Leibniz Institute for Economic Research, Berlin 10115, Germany; <sup>bbbbb</sup>School of Public Policy and Administration, Carleton University, Ottawa, ON K1S 5B6, Canada; <sup>ccccc</sup>Institute of Psychology, Cardinal Stefan Wyszyński University, Warsaw 01-938, Poland; <sup>ddddd</sup>Currency Department, Bank of Canada, Ottawa, ON K1A 0G9, Canada; <sup>eeeee</sup>Department of Agricultural Economics and Rural Development, University of Göttingen, Göttingen, Niedersachsen 37083, Germany; <sup>fffff</sup>Instituto de Estudios Superiores de la Empresa Business School, University of Navarra, Barcelona 08034, Spain; <sup>ggggg</sup>Universidad de Navarra, Pamplona 31009, Spain; <sup>hhhhh</sup>Department of Economics, Dedman College of Humanities and Sciences, Southern Methodist University, Dallas, TX 75275-0496; <sup>iiiiii</sup>SKEMA Business School, Groupe de Recherche en Droit, Économie et Gestion, Université Côte d'Azur, Lille 59777, France; <sup>jjjjj</sup>Institute for Public Management and Governance, Vienna University of Economics and Business, Vienna 1020, Austria; <sup>kkkkkk</sup>Department of Finance, Rotterdam School of Management, Erasmus University Rotterdam, Rotterdam, Zuid Holland 3012 CC The Netherlands; <sup>llllll</sup>Institute of Psychology, Universität Osnabrück, Osnabrück 49078, Germany; <sup>mmmmmm</sup>Department of Psychology, University of Houston, Houston, TX 77204; <sup>nnnnnn</sup>Department of Economics, Business and Finance, Lake Forest College, Lake Forest, IL 60048; <sup>ooooooo</sup>Department of Marketing, Vienna University of Economics and Business, Vienna 1020, Austria; <sup>ppppppp</sup>Department of Economics, Bielefeld University, Bielefeld 33615, Germany; <sup>qqqqqqq</sup>Department of Economics, University at Albany, State University of New York, Albany, NY 12222; <sup>rrrrrrr</sup>School of Information, University of Michigan, Ann Arbor, MI 48109; <sup>sssssss</sup>The Journal, Camden, DE 19934; <sup>ttttttt</sup>Finnish Flagship Inequalities, Interventions, and New Welfare State Flagship Research Centre, University of Turku, Turku 20014, Finland; <sup>uuuuuuu</sup>Applied Economics and Management, Cornell University, Ithaca, NY 14853; <sup>wwwww</sup>Department of Economics, City St George's, University of London, London EC1V 0HB, United Kingdom; <sup>xxxxxxx</sup>Department of Psychology, University of Milano-Bicocca, Milan 20126, Italy; <sup>xxxxxxx</sup>School of Economics, University of Nottingham, Nottingham NG7 2RD, United Kingdom; <sup>yyyyyyy</sup>InsIDE Lab, Dipartimento di Economia, Università degli Studi dell'Insubria, Varese 21100, Italy; <sup>zzzzzzz</sup>Department of Economics, University of Essex, Colchester CO4 3SQ, United Kingdom; <sup>aaaaaaaaa</sup>Data and Digital Services Department, Bank of Canada, Ottawa, ON K1A 0G9, Canada; <sup>bbbbbbb</sup>Systems and Computer Engineering Department, Carleton University, Ottawa, ON K1S 5B6, Canada; <sup>ccccccc</sup>Leverhulme Centre for the Future of Intelligence, University of Cambridge, Cambridge CB2 1SB, United Kingdom; <sup>ddddddd</sup>Department of Economics, University of California, Davis, CA 95616; <sup>eeeeeee</sup>Institute for Business Administration, University of Regensburg, Regensburg, Bavaria 93053 Germany; <sup>fffffft</sup>Department of Linguistics, Ghent University, Ghent 9000, Belgium; <sup>ggggggg</sup>Department of Linguistics and English Language, University of Manchester, Manchester M13 9PL, United Kingdom; <sup>hhhhhhh</sup>Department of Marketing, Tilburg School of Economics and Management, Tilburg University, Tilburg 5037 AB, The Netherlands; <sup>iiiiiii</sup>School of Public Finance and Taxation, Dongbei University of Finance and Economics, Dalian 116025, China; <sup>jjjjjjj</sup>Department of Marketing, Clemson University, Clemson, SC 29634; <sup>kkkkkkk</sup>Department of Economics, University of Birmingham, Birmingham B15 2TT, United Kingdom; <sup>lllllll</sup>School of Information, Journalism and Communication, University of Sheffield, Sheffield S10 2AH, United Kingdom; <sup>mmmmmmm</sup>Department of Government, Cornell University, Ithaca, NY 14853; <sup>nnnnnnn</sup>Agri-Science Queensland, Department of Primary Industries, Brisbane, QLD 4102, Australia; <sup>ooooooooo</sup>Asia and Pacific Department, International Monetary Fund, Washington, DC 20431; <sup>ppppppp</sup>Department of Psychology, New York University, New York, NY 10003; <sup>qqqqqqq</sup>School of Natural and Social Sciences, Scotland's Rural College, Edinburgh EH9 3JG, United Kingdom; <sup>rrrrrrr</sup>Chair of Agricultural Production and Resource Economics, Technical University of Munich, Freising 85354, Germany; <sup>sssssss</sup>The Institute for Food and Resource Economics, University of Bonn, Bonn 53115, Germany; <sup>ttttttt</sup>Department of Special Needs Education, Centre for Research on Equality in Education, University of Oslo, Oslo 0318, Norway; <sup>uuuuuuu</sup>World Bank, Washington, DC 20433; <sup>wwwww</sup>Birkbeck Business School, Birkbeck, University of London, London WC1E 7JL, United Kingdom; <sup>xxxxxxx</sup>Department of Economics, Vanderbilt University, Nashville, TN 37235-1819; <sup>xxxxxxx</sup>Division of Psychology & Mental Health, School of Health Sciences, University of Manchester, Manchester M13 9PL, United Kingdom; <sup>yyyyyyy</sup>Center for Health and Wellbeing, Princeton University, Princeton, NJ 08544; <sup>zzzzzzz</sup>Department of Economics, College of Staten Island, Staten Island, NY 10314; <sup>aaaaaaaaa</sup>City University of New York, New York, NY 10017; <sup>bbbbbbb</sup>Department Economics, Law, and Society, Ecole supérieure des sciences commerciales School of Management, Angers 49003, France; <sup>ccccccc</sup>Department of Economics, University of California, Davis, CA 95616; <sup>ddddddd</sup>Stockholm Business School, Stockholm University, Stockholm 10691, Sweden; <sup>eeeeeee</sup>Leibniz Institute for Financial Research Sustainable Architecture for Finance in Europe, Frankfurt, Hesse 60323, Germany; <sup>fffffft</sup>Department of Economics, National University of Singapore, Singapore 117570, Singapore; <sup>ggggggg</sup>Collegium of World Economy, Szkoła Główna Handlowa Warsaw School of Economics, Warsaw 02-554, Poland; <sup>hhhhhhh</sup>Department of Political Science, Stanford University, Stanford, CA 94305; <sup>iiiiiii</sup>Department of Engineering Sciences & Applied Mathematics, Northwestern University, Evanston, IL 60201; <sup>jjjjjjj</sup>Real Estate Economics, National Chengchi University, Taipei 11605, Taiwan; <sup>kkkkkkk</sup>Aalto University, Espoo 02150, Finland; <sup>lllllll</sup>Department of Accounting, Tilburg School of Economics and Management, Tilburg University, Tilburg 5037 AB, The Netherlands; <sup>mmmmmmm</sup>Department of Economics, Faculty of Arts and Social Sciences, National University of Singapore, Singapore 117570, Singapore; <sup>nnnnnnn</sup>School of Economics, Faculty of Business Economics and Law, University of Queensland, Brisbane, St Lucia, QLD 4072, Australia; <sup>ooooooooo</sup>Department of Psychology, Faculty of Science and Technology, Fylde College, Lancaster University, Lancaster LA1 4YF, United Kingdom; <sup>ppppppp</sup>School of Biosciences, University of Sheffield, Sheffield S10 2TN, United Kingdom; <sup>qqqqqqq</sup>Christ Church College, University of Oxford, Oxfordshire OX1 1DP, United Kingdom; <sup>rrrrrrr</sup>Economics Experimental Lab, Nanjing Audit University, Nanjing 210017, China; <sup>sssssss</sup>Department of Finance, HEC, University of Lausanne, Lausanne 1015, Switzerland; <sup>ttttttt</sup>Swiss Finance Institute, Lausanne 1015, Switzerland; <sup>uuuuuuu</sup>Department of Marketing, Rotman School of Management, University of Toronto, Toronto, ON M5S 1A1, Canada; <sup>wwwww</sup>Department of Industrial Engineering and Operations Research, Columbia University, New York, NY 10027; and <sup>xxxxxxx</sup>Department of Economics, Penn State University, University Park, PA 16802

1. A. Brodeur *et al.*, Promoting reproducibility and replicability in political science. *Res. Polit.* **11**, 20531680241233439 (2024).
2. D. L. Donoho, A. Maleki, I. U. Rahman, M. Shahram, V. Stodden, Reproducible research in computational harmonic analysis. *Comput. Sci. Eng.* **11**, 8–18 (2008).
3. M. Fišar *et al.*, Reproducibility in management science. *Manag. Sci.* **70**, 1343–1356 (2024).
4. P. Gertler, S. Galiani, M. Romero, How to make replication the norm. *Nature* **554**, 417–419 (2018).
5. S. N. Goodman, D. Fanelli, J. P. Ioannidis, What does research reproducibility mean? *Sci. Transl. Med.* **8**, 341ps12–341ps12 (2016).
6. M. Milkowski, W. M. Hensel, M. Hohol, Replicability or reproducibility? On the replication crisis in computational neuroscience and sharing only relevant detail. *J. Comput. Neurosci.* **45**, 163–172 (2018).
7. National Academies of Sciences, Engineering, and Medicine, *Reproducibility and Replicability in Science* (National Academies Press, 2019).
8. C. Pérignon, K. Gadouche, C. Hurlin, R. Silberman, E. Debonnel, Certify reproducibility with confidential data. *Science* **365**, 127–128 (2019).
9. L. Villhuber, Report by the AEA data editor. *AEA Pap. Proc.* **112**, 813–823 (2022).
10. A. Brodeur *et al.*, Reproducibility and robustness of economics and political science research. *Nature* **652**, 151–156 (2026).
11. A. C. Chang, P. Li, Is economics research replicable? Sixty published papers from thirteen journals say “often not”. *Crit. Finance Rev.* **11**, 185–206 (2022).
12. S. Crüwell *et al.*, What's in a badge? A computational reproducibility investigation of the open data badge policy in one issue of psychological science. *Psychol. Sci.* **34**, 512–522 (2023).
13. Open Science Collaboration, Estimating the reproducibility of psychological science. *Science* **349**, aac4716 (2015).
14. P. Obels, D. Lakens, N. A. Coles, J. Gottfried, S. A. Green, Analysis of open data and computational reproducibility in registered reports in psychology. *Adv. Methods Pract. Psychol. Sci.* **3**, 229–237 (2020).
15. C. Pérignon *et al.*, Computational reproducibility in finance: Evidence from 1,000 tests. *Rev. Financial Stud.* **37**, 3558–3593 (2024).
16. V. Stodden, J. Seiler, Z. Ma, An empirical analysis of journal policy effectiveness for computational reproducibility. *Proc. Natl. Acad. Sci. U.S.A.* **115**, 2584–2589 (2018).
17. B. Wood, R. Müller, A. Brown, Push button replication: Is impact evaluation evidence for international development verifiable? *PLoS one* **13**, e0209416 (2018).
18. J. E. Colliard, C. Hurlin, C. Pérignon, The Economics of Computational Reproducibility (2022). <https://ssrn.com/abstract=3418896>.
19. A. Hryciyshyn, H. Eassom, “ExplanAltions: An AI study” (Tech. Rep. 1, John Wiley & Sons, Hoboken, NJ, 2025).
20. M. Vaccaro, A. Almaatouq, T. Malone, When combinations of humans and ai are useful: A systematic review and meta-analysis. *Nat. Hum. Behav.* **8**, 1–11 (2024).
21. G. Bansal *et al.*, “Does the whole exceed its parts? The effect of AI explanations on complementary team performance” in *Proceedings of the 2021 CHI conference on human factors in computing systems* (2021), pp. 1–16.
22. G. Bansal, B. Nushi, E. Kamar, E. Horvitz, D. S. Weld, “Is the most accurate AI the best teammate? Optimizing AI for teamwork” in *Proceedings of the AAAI Conference on Artificial Intelligence* (PKP Publishing Services Network, 2021), vol. 35, pp. 11405–11414.
23. E. Bondi *et al.*, “Role of human-AI interaction in selective prediction” in *Proceedings of the AAAI Conference on Artificial Intelligence* (PKP Publishing Services Network, 2022), vol. 36, pp. 5286–5294.
24. Á. A. Cabrera, A. Perer, J. I. Hong, Improving human-AI collaboration with descriptions of AI behavior. *Proc. ACM Hum. Comput. Interact.* **7**, 1–21 (2023).
25. E. Goh *et al.*, Influence of a large language model on diagnostic reasoning: A randomized clinical vignette study. medRxiv [Preprint] (2024). <https://doi.org/10.1101/2024.03.12.24303785> (Accessed 20 February 2026).
26. B. Koepnick *et al.*, De novo protein design by citizen scientists. *Nature* **570**, 390–394 (2019).
27. H. Liu, V. Lai, C. Tan, Understanding the effect of out-of-distribution examples and interactive explanations on human-AI decision making. *Proc. ACM Hum. Comput. Interact.* **5**, 1–45 (2021).
28. H. Mozannar *et al.*, Effective human-AI teams via learned natural language rules and onboarding. *Adv. Neural Inf. Process. Syst.* **36**, e01007 (2024).
29. S. Noy, W. Zhang, Experimental evidence on the productivity effects of generative artificial intelligence. *Science* **381**, 187–192 (2023).
30. C. Reverberi *et al.*, Experimental evidence of effective human-AI collaboration in medical decision-making. *Sci. Rep.* **12**, 14952 (2022).
31. M. Schemmer, P. Hemmer, M. Nitsche, N. Kühl, M. Vössing, “A meta-analysis of the utility of explainable artificial intelligence in human-AI decision-making” in *Proceedings of the 2022 AAAI/ACM Conference on AI, Ethics, and Society* (2022), pp. 617–626.
32. M. H. Tessler *et al.*, AI can help humans find common ground in democratic deliberation. *Science* **386**, eadq2852 (2024).
33. M. Vaccaro, J. Waldo, The effects of mixing machine learning and human judgment. *Commun. ACM* **62**, 104–110 (2019).
34. B. Wilder, E. Horvitz, E. Kamar, “Learning to complement humans” in *Proceedings of the 29th International Joint Conference on Artificial Intelligence* (ACM Digital Library, 2020), pp. 1526–1533.
35. Cherry Bekaert LLP, Report of independent auditor. *Am. Econ. Rev.* **112**, 2083–2098 (2022).
36. Z. Bućinca, M. B. Malaya, K. Z. Gajos, To trust or to think: Cognitive forcing functions can reduce overreliance on AI in AI-assisted decision-making. *Proc. ACM Hum. Comput. Interact.* **5**, 1–21 (2021).
37. L. J. Skitka, K. L. Mosier, M. Burdick, Does automation bias decision-making? *Int. J. Hum. Comput. Stud.* **51**, 991–1006 (1999).
38. E. Brynjolfsson, D. Li, L. Raymond, Generative AI at work. *Q. J. Econ.* **140**, 889–942 (2025).
39. M. Del Giudice, S. W. Gangestad, A traveler's guide to the multiverse: Promises, pitfalls, and a framework for the evaluation of analytic decisions. *Adv. Methods Pract. Psychol. Sci.* **4**, 2515245920954925 (2021).
40. X. Lu, H. White, Robustness checks and robustness tests in applied economics. *J. Econom.* **178**, 194–206 (2014).
41. M. B. Nuijten, “Assessing and improving robustness of psychological research findings in four steps” in *Avoiding Questionable Research Practices in Applied Psychology* (Springer, 2022), pp. 379–400.
42. V. Lai, H. Liu, C. Tan, “why is' chicago' deceptive? Towards building model-driven tutorials for human” in *Proceedings of the 2020 CHI Conference on Human Factors in Computing Systems* (ACM Digital Library, 2020), pp. 1–13.
43. Y. Zhang, Q. V. Liao, R. K. Bellamy, “Effect of confidence and explanation on accuracy and trust calibration in AI-assisted decision making” in *Proceedings of the 2020 Conference on Fairness, Accountability, and Transparency* (ACM Digital Library, 2020), pp. 295–305.
44. H. Vasconcelos *et al.*, Explanations can reduce overreliance on AI systems during decision-making. *Proc. ACM Hum. Comput. Interact.* **7**, 1–38 (2023).
45. Y. K. Dwivedi *et al.*, Opinion paper: “So what if chatgpt wrote it?” Multidisciplinary perspectives on opportunities, challenges and implications of generative conversational ai for research, practice and policy. *Int. J. Inf. Manag.* **71**, 102642 (2023).
46. B. D. Lund *et al.*, Chatgpt and a new academic reality: Artificial intelligence-written research papers and the ethics of the large language models in scholarly publishing. *J. Assoc. Inf. Sci. Technol.* **74**, 570–581 (2023).
47. N. Wadhwa *et al.*, Core: Resolving code quality issues using LLMs. *Proc. ACM Softw. Eng.* **1**, 789–811 (2024).
48. Y. Zhang, “Detecting code comment inconsistencies using LLM and program analysis” in *Companion Proceedings of the 32nd ACM International Conference on the Foundations of Software Engineering* (ACM Digital Library, 2024), pp. 683–685.
49. D. Nam, A. Macvean, V. Hellendoorn, B. Vasilescu, B. Myers, “Using an LLM to help with code understanding” in *Proceedings of the IEEE/ACM 46th International Conference on Software Engineering* (ACM Digital Library, 2024), pp. 1–13.
50. A. Brodeur *et al.*, AI Replication Games. <https://github.com/I4Replication/AI-Games>. GitHub. Accessed 2 May 2026.

**AI-Assisted Teams Outperform AI-Led Teams but Not Human-Only Teams in Assessing Research Reproducibility in Quantitative Social Science**

Abel Brodeur (Corresponding author), David Valenta, Alexandru Marcoci, Juan P. Aparicio, Derek Mikola, Bruno Barbarioli, Rohan Alexander, Lachlan Deer, Tom Stafford, Lars Vilhuber, Gunther Bensch, Fabio Motoki, Mohamed Abdelhady, Yousra Abdelmoula, Ghina Abdul Baki, Tomás Aguirre, Sriraj Aiyer, Shumi Akhtar, Farida Akhtar, Melle R. Albada, Micah Altman, David Angenendt, Zahra Arjmandi Lari, Jorge Armando De León Tejada, David Rodriguez Arana, Igor Asanov, Anastasiya-Mariya Noha, Rebecca Ashong, Tobias Auer, Francisco J. Bahamonde-Birke, Bradley J. Baker, Söhnke M. Bartram, Dongqi Bao, Lucija Batinovic, Tommaso Batistoni, Monica Beeder, Louis-Philippe Beland, Carsten Gero Bienz, Christ Billy Aryanto, Cylcia Bolibaugh, Carl Bonander, Ramiro Bravo, Egor Bronnikov, Stephan Bruns, Nino Buliskeria, Sara Caicedo-Silva, Andrea Calef, Juan Sebastian Cano Arias, Gustavo A. Castillo Alvarez, Solomon Caulker, Simonas Cepenas, Arthur Chatton, Zirou Chen, Ngozi Chioma Ewurum, Anda-Bianca Ciocîrlan, Felix J. Clouth, Jason Collins, Nikolai Cook, Cesar Cornejo, João Craveiro, Jonathan Créchet, Jing Cui, Niveditha Chalil Vayalabron, Christian Czymara, Carlos Daniel Bermúdez Jaramillo, Hannes Datta, Lien Denoo, Arshia Dhaliwal, Nancy Dhameja, Elodie Djemai, Erwan Dujeancourt, Uğurcan Dündar, Thibaut Duprey, Yasmine Eissa, Youssef El Fassi, Ismail El Fassi, Keaton Ellis, Ali Elminejad, Mahmoud Elsherif, Aysil Emirmahmutoglu, Giulian Etingin-Frati, Emeka Eze, Jan Fabian Dollbaum, Jan Feld, Andres Felipe Rengifo Jaramillo, Guidon Fenig, Victoria Fernandes, Lenka Fiala, Lukas Fink, Mojtaba Firouzjaeiangalougah, Sara Fish, Jack Fitzgerald, Rachel Forshaw, Alexandre Fortier-Chouinard, Louis Fréget, Joris Frese, Jacopo Gabani, Sebastian Gallegos, Max C. Gamill, Attila Gáspár, Romain Gauriot, Evelina Gavrilova, Diogo Geraldos, Giulio Giacomo Cantone, Grant Gibson, Dirk Goldschmitt, Amélie Gourdon-Kanhukamwe, Andrea Gregor de Varda, Idaliya Grigoryeva, Alexi Gugushvili, Aaron H.A. Fletcher, Florian Habermann, Márton Hablicsek, Joanne Haddad, Jonathan D. Hall, Olle Hammar, Malek Hassouneh, Carina I. Hausladen, Sophie C. F. Hendrikse, Matthew Hepplewhite, Anson T. Y. Ho, Senan Hogan-Hennessy, Elliot Howley, Gaoyang Huang, Héloïse Hulstaert, Zlatomira G. Ilchovska, Paola Jaimes Santamaria, Niklas Jakobsson, Joakim Jansson, Ewa Jarosz, Hossein Jebeli, Yanchen Jiang, Hiba Junaid, Rohan Kalluraya, Sunny Karim, Edmund Kelly, Eva Kimel, Sorravich Kingsuwankul, Valentin Klotzbücher, Daniel Krähmer, Pijus Krūminas, Nicholas Kruus, Essi Kujansuu, Christoph F. Kurz, Stephan Küster, Blake Lee-Whiting, Felix Lewandowski, Tongzhe Li, Ruoxi Li, Dan Liu, Jiacheng Liu, Helix Lo, Katharina Loter, Felipe Macedo Dias, Christopher R. Madan, Nicolas Mäder, Marco Mandas, Cesar Mantilla, Jan Marcus, Diego Marino Fages, Xavier Martin, Ryan McWay, Daniel Medina-Gaspar, Sisi Meng, Lingyu Meng, Simon Merz, Alex P. Miller, Thibault Mirabel, Dibya Deepta Mishra, Sumit Mishra, Belay W. Moges, Morteza Mohandes Mojarrad, Myra Mohnen, Louis-Philippe Morin, Lucija Muehlenbachs, Gastón Mullin, Andreea Musulan, Sara Muzzi, James A. C. Myers, Florian Neubauer, Tuan Nguyen, Ali Niazi, Ardyn Nordstrom, Bartłomiej Nowak, Daneal O’Habib, Tim Ölkens, Justin Ong, Valeria Orozco Castiblanco, Ömer Özak, Ali I. Ozkes, Mikael Paaso, Shubham Pandey, Varvara Papazoglou, Romeo Penheiro, Linh Pham, Ulrike Phieler, Peter Pütz, Quan Qi, Jingyi Qiu, David A. Reinstein, Juuso Repo, Nicolas Rudolf, Shree Saha, Orkun Saka, Chiara Saponaro, Georg Sator, Martijn Schoenmakers, Raffaello Seri, Meet Shah, Paul Sibille, Christoph Siemroth, Vladimir Skavysh, Ben Slater, Wenting Song, Stefan Staubli, Tobias Steindl, Nomwendé Steven Waongo, Paul Stott, Stephenson Strobel, Roshini Sudhaharan, Pu Sun, Scott D. Swain, Oleksandr Talavera, Hanz M. Tantiangco, Georgy Tarasenko, Boyd Tarlinton, Mariam Tarraf, Ken Teoh, Rémi Thériault, Bethan Thompson, Tonghui Tian, Wenjie Tian, Manuel Tobias Rein, Emmanuel Tolani, Nicolai Borgen, Solveig Topstad Borgen, Javier Torralba, Carolina Velez-Ospina, Man Wai Mak, Lukas Wallrich, Zeyang Wang, Leah Ward, Matthew D. Webb, Duncan Webb, Bryan S. Weber, Christoph Weber, Wei-Chien Weng, Christian Westheide, Tom Wilkinson, Kwong-Yu Wong, Marcin Wroński, Zhuangchen Wu, Qixia Wu, Victor Y. Wu, Bohan Xiao, Feihong Xu, Cong Xu, Pranav Yadav, Yu Yang Chou, Luther Yap, Myra Yazbeck, Bo Yao, Zuzanna Zagrodzka, Tahreen Zahra, Mirela Zaneva, Xiaomeng Zhang, Ziwei Zhao, Han Zhong, Aras Zirculis, Jiacheng Zou, Floris Zoutman, Christelle Zozoungbo.

**This PDF file includes:** Materials and Methods

Figures S1 to S7

Tables S1 to S17

**Materials and Methods.** A version-tagged copy of the code and data is permanently archived at <https://github.com/I4Replication/Al-Games>. All our materials are available here: <https://osf.io/sz2g8/>.

**Research Ethics Boards.** Participants in the AI replication games experiments coauthor this study. The University of Ottawa Office of Research Ethics and Integrity reviewed and approved our AI games (H-09-25-12041). The King's College London Research Ethics Office reviewed and approved our focus groups (MRA-24/25-48393). All participants provided informed consent.

**Pre-Registration.** Our pre-analysis plan was preregistered on the Open Science Framework (OSF) on May 2nd, 2024: <https://osf.io/sz2g8/>. The preregistration was done after our pilot event at the University of Toronto.

Of note, the pre-analysis plan refers to AI-assisted teams as cyborg teams and AI-led teams as machine teams. AI-led teams are also referred to as 'machines with limited human assistance.'

We note the following deviations from the pre-analysis plan:

- The pre-analysis plan mentions the Ottawa, Sheffield, Cornell, Cambridge, and Tilburg games. We ended up not organizing games at the University of Cambridge and replaced those games with the Bogota games. We also mentioned in the pre-analysis plan that we were hoping to have at least one more event in 2024/2025. We added a 2024 fully virtual game session on November 22nd, 2024.
- The pre-analysis plan mentions three dependent variables for the robustness checks. We added a fourth dependent variable: the replicators managed to implement their two robustness checks. The other three dependent variables are all preregistered.
- We did not conduct exploratory analysis of improvement over time in proposing/implementing robustness checks among AI-led and AI-assisted teams.
- We revised our classification of coding errors. This coding was done by DVa, JPAP, DMi, and LFi, and disagreements were resolved through discussion rather than voting.
- In our revised classification of coding errors, missing packages/paths or versioning issues are not considered minor errors; they are not considered coding errors at all.
- In our revised classification of coding errors, we distinguished between minor and major errors as pre-specified, but we also included a category for false-positives (i.e., issues raised by teams that are not actual errors), and we coded at which stage of data processing each error occurred.
- We included additional secondary analyses in the manuscript: (i) OLS regressions testing differences for R and Stata teams, (ii) heterogeneity analysis by number of prompts, (iii) heterogeneity analysis by experience.
- We also included paper metrics to compare the papers we selected for the AI games to other papers published in the same journals in the same years.
- Qualitative analyses were not preregistered.
- We rely on Kaplan-Meier curves.

**Research Questions.** Here are the primary research questions that were preregistered:

1. Do AI-led teams computationally reproduce more results than AI-assisted and human-only teams?
2. Are AI-led teams faster to computationally reproduce results than AI-assisted and human-only teams?
3. Do AI-led teams detect more major and minor coding errors or data irregularities than AI-assisted and human-only teams?
4. Are AI-led teams faster at detecting major and minor coding errors or data irregularities than AI-assisted and human-only teams?
5. Do AI-led teams propose better robustness checks than AI-assisted and human-only teams?
6. Are AI-led teams more capable of implementing robustness checks than AI-assisted and human-only teams?
7. Do AI-assisted teams computationally reproduce more results than human-only teams?
8. Are AI-assisted teams faster to computationally reproduce results than human-only teams?
9. Do AI-assisted teams detect more major and minor coding errors or data irregularities than human-only teams?
10. Are AI-assisted teams faster at detecting major and minor coding errors or data irregularities than human-only teams?
11. Do AI-assisted teams propose better robustness checks than human-only teams?
12. Are AI-assisted teams more capable of implementing robustness checks than human-only teams?

We also explored the following exploratory (preregistered) research questions:

13. Are AI-led teams improving their performance over time at computationally reproducing results, detecting coding errors or data irregularities, and providing good robustness checks?"
14. Are AI-assisted teams improving their performance over time at computationally reproducing results, detecting coding

errors or data irregularities, and providing good robustness checks?

We also tackle an exploratory research question that was not preregistered in the article:

15. Do AI-assisted teams over-rely or under-rely on AI?

**AI Replication Games Advertisement.** The Institute for Replication advertised the AI replication games through social media (Bluesky and X) and emails. Events were also promoted on the Institute’s webpage (<https://i4replication.org/games.html>).

The typical social media posts included the following information:

“This is a one-day event that brings researchers together to collaborate on reproducing quantitative results published in high-ranking social science journals. You will have the opportunity to network with fellow researchers and develop your coding and AI skills.

Open to all researchers: faculty, post-docs, and graduate students. Knowledge of Python or R or Stata is essential. Participants will be randomly assigned to one of three teams: Machine with restricted human assistance, Cyborg or Human.

All participants will get coauthorship on a meta-research journal paper which combines the work of all teams.

Register here: (Link to Registration Form).”

**Participation, Incentives and Participants Exclusion.** Our AI Replication Games were open to graduate students, postdoctoral fellows, professors and researchers from non-academic organizations with a PhD. All participants were offered coauthorship on this paper conditional on participation, independent of their team’s performance or success in reproducing results. No monetary compensation or performance-based rewards were provided, which may have influenced how participants allocated effort across teams and conditions. Furthermore, participants did not have access to preliminary results at the time of participation.

We did not accept registration from participants with no knowledge of Stata nor R. We also excluded from participating a very small number of researchers with no knowledge of R and who did not have a Stata license.

As noted in the main text, a few organizers participated in one of the games. They did not know about the papers to be reproduced at their respective event.

**Randomization.** As mentioned in the main text, randomization was carried out in two steps for each of the seven events. In step one, coauthors were randomly assigned to a team of three, conditional on the software preferences reported by participants (Stata or R) and the mode of participation (in person or virtual). In step two, each team was randomly assigned to one of three treatment arms.

Each team was assigned a study from leading social science journals (i.e., economics, political science, or behavioral science/psychology). All team members within a team were not necessarily from the same field. Studies were assigned based on knowledge of R and Stata. In practice, most psychologists are more comfortable with R, so they were assigned a behavioral or political science study, whereas economists are more comfortable with Stata and were thus much more likely to be assigned an economics study.

**Documents to be Filled During the Games.** During the event, each team filled out an Excel sheet documenting their outcomes. See the Excel document here: <https://osf.io/sz2g8/>. The document “Template Time Stamp” includes 3 sheets to be filled by each team. The first sheet is for computational reproducibility. Teams need to fill out the time that they have computationally reproduced the exhibit. The second sheet documents coding errors detected. Teams need to add a row for each coding error and data irregularity and enter the time they have detected them. In the last sheet, teams need to provide a description of their two robustness checks and provide estimates if they could implement those. Researchers in the AI-assisted and AI-led groups are also asked to share their prompts/conversations at the end of the event.

**Information Provided.** The materials provided to researchers participating in our events are available here: <https://osf.io/sz2g8/overview>. All researchers took part in a one-hour-long training. AI-led and AI-assisted teams had an additional presentation on general ChatGPT use. This was also offered to No-AI teams but was voluntary. Our slides and recordings of the training sessions are available on OSF.

Teams were made aware of the article that they needed to reproduce at the beginning of the event. A note was added in the PDF of the article clarifying which exhibits needed to be reproduced: “First, computationally reproduce Table X, columns Y and Z. Second, detect coding errors/data irregularities. Last, propose and try to implement two robustness check.”

Our pre-games slides provided no rigorous definition of a coding error, a discrepancy or a data irregularity. The slides can be found here: <https://osf.io/sz2g8/files/7aguk>. That said, our pre-analysis plan provided a brief definition of Coding Errors or Data Irregularities (<https://osf.io/sz2g8/files/b42ue>) and we provided to participants Our Template Time Stamp (<https://osf.io/sz2g8/files/p7hd9>), which provide a very simple example: “we uncovered a coding error in the do-file entitled ‘XYZ.do’, line 45. The error involves forgetting a control variable or mispecifying something or miscoding a variable, etc.”

**Selection of Papers Used During Events.** For each event, two studies published in leading social science journals are selected by AB. Table S1 lists the 12 different articles given to teams for the AI Games. The studies are published in a journal with a data and code availability policy. One study is coded in Stata; the other is coded in R. The studies have all been reproduced by the Institute for Replication before the AI replication games. The Institute for Replication runs about two “regular” replication games, in contrast to these AI events, each month. At every such event, teams of researchers try to reproduce results from peer-reviewed publications. They then prepare reports of their findings, which are subsequently shared with the original authors

and made public on average six months following an event. Importantly, this means the Institute for Replication had over 20 published studies with known reproduction results *that have not yet been made publicly available* to choose from at any point in 2024. We could not take studies with publicly known coding issues since ChatGPT may be able to “know” coding errors or data discrepancies without “finding” them. This list of papers with known reproduction results that had not yet been made public is the corpus we sampled from for each of the AI replication games. Tables S3 and S4 then provide a summary of the coding errors known to the Institute for Replication.

The sampling cannot be random or blind for a few reasons. First, the variation in reproduction packages (sometimes called *replication packages* in the social sciences) is too large. In both scenarios, when folders reproduce studies perfectly or when folders cannot be deciphered at all, would yield no variation in at least one of our outcomes. (No coding errors exist in the former; we cannot evaluate the correctness of the code in the latter.) Second, studies need to rely on publicly available data and codes, or the exercise is futile. Third, we need to match the software abilities of participants to each study. Within this corpus, we selected studies known to have coding errors or data irregularities. All teams were told that they needed to uncover coding errors or data irregularities.

Of note, the Institute for Replication uncovers coding errors for about 25% of studies, with some studies containing multiple errors (1).

Table S5 reports summary statistics on the 12 articles and their associated replication packages, comparing AI Games papers to a comparison group. The comparison group is composed of 120 articles with a replication package randomly chosen from the same journals and publication years as the 12 AI Games articles.

On average, articles in our sample contain about 25 pages, with AI Games articles being slightly longer (25.1 pages) than comparison articles (23.7 pages). The difference is not statistically significant at conventional levels ( $p = 0.718$ ). Replication packages for AI Games papers are substantially larger, averaging 994 MB, compared to 631 MB for comparison papers. The difference is not statistically significant ( $p = 0.449$ ), possibly due to the small sample size. Documentation and structure are more prevalent among AI Games packages: 58% include a README file compared to 45% in the comparison group ( $p = 0.407$ ).

Overall, the table highlights that AI Games replication packages are larger, more complex, and more software-intensive than those of the comparison articles (although the differences are mostly not statistically significant), while article length remains broadly similar across groups.

DRAFT

**Table S1. Papers Assigned to Teams for AI Games**

Article Title	Year Published	Journal	Programming Language	Computationally Reproducible?	Any Errors in Codes?	Games Used
<b>Economics Papers</b>						
Major Reforms in Electricity Pricing: Evidence from a Quasi-Experiment	2022	The Economic Journal	Stata	Yes	Yes	Toronto, Ottawa
Sorting or Steering: The Effects of Housing Discrimination on Neighborhood Choice	2022	Journal of Political Economy	R, Stata	Yes	Yes	Toronto
Taste-Based Gender Favouritism In High-Stake Decisions: Evidence from the Price is Right	2023	The Economic Journal	R	Yes	Yes	Sheffield
The Heterogeneous Tax Pass-Through Under Different Vertical Relationships	2022	The Economic Journal	Stata	Yes	Yes	Sheffield, Cornell
The Interplay Among Savings Accounts and Network-Based Financial Arrangements: Evidence from a Field Experiment	2023	The Economic Journal	Stata	Yes	Yes	Bogota
How the Other Half Died: Immigration and Mortality in U.S. Cities	2024	The Review of Economic Studies	R	Yes	Yes	Virtual (Europe)
Gambling, Saving, and Lumpy Liquidity Needs	2021	American Economic Journal: Applied Economics	Stata	Yes	Yes	Virtual (Europe), Virtual (North America)
<b>Political Science Papers</b>						
Acquiescence Bias Inflates Estimates of Conspiratorial Beliefs and Political Misperceptions	2023	Political Analysis	R	Yes	Yes	Cornell, Bogota
Do Policy Makers Listen to Experts? Evidence from a National Survey of Local and State Policy Makers	2022	American Political Science Review	R	Yes	Yes	Tilburg
<b>Papers Published in Behavioural Science Journals</b>						
Arrests and Convictions but Not Sentence Length Deter Terrorism in 28 European Union Member States	2023	Nature Human Behaviour	R	Yes	Yes	Ottawa
Examining Inequality in the Time Cost of Waiting	2023	Nature Human Behaviour	Stata	Yes	Yes	Tilburg
Mindful-Gratitude Practice Reduces Prejudice at High Levels of Collective Narcissism	2024	Psychological Science	R	Yes	Yes	Virtual (North America)

A broad description of the 12 different articles given to teams for the AI Games. *Article Title* corresponds the paper provided to different teams across our events. The column for *Computationally Reproducible?* means the replication packages had previously been run and could reproduce the results. *Any Errors in Codes?* indicates whether there were errors in the article/replication package. Commas separating the entries in the *Games Used* column mean the papers were used for two different games.

**Table S2. Information Regarding the Paper's Used in AI Games Events**

Article Title	DOI for Paper	Replication Folder
<b>Economics Papers</b>		
Major Reforms in Electricity Pricing: Evidence from a Quasi-Experiment	10.1093/ej/ueab076	<a href="https://zenodo.org/records/5423782">https://zenodo.org/records/5423782</a>
Sorting or Steering: The Effects of Housing Discrimination on Neighborhood Choice	10.1086/720140	<a href="https://github.com/peterchristensen/Sorting-or-Steering">https://github.com/peterchristensen/Sorting-or-Steering</a>
Taste-Based Gender Favouritism In High-Stake Decisions: Evidence from the Price is Right	10.1093/ej/uead087	<a href="https://doi.org/10.5281/zenodo.8372384">https://doi.org/10.5281/zenodo.8372384</a>
The Heterogeneous Tax Pass-Through Under Different Vertical Relationships	10.1093/ej/ueac007	<a href="https://doi.org/10.5281/zenodo.5824590">https://doi.org/10.5281/zenodo.5824590</a>
The Interplay Among Savings Accounts and Network-Based Financial Arrangements: Evidence from a Field Experiment	10.1093/ej/ueac053	<a href="https://doi.org/10.5281/zenodo.6985683">https://doi.org/10.5281/zenodo.6985683</a>
How the Other Half Died: Immigration and Mortality in U.S. Cities	10.1093/restud/rdad035	<a href="https://dx.doi.org/10.5281/zenodo.7506459">https://dx.doi.org/10.5281/zenodo.7506459</a>
Gambling, Saving, and Lumpy Liquidity Needs	10.1257/app.20180177	<a href="https://www.openicpsr.org/openicpsr/project/115162/version/V1/view">https://www.openicpsr.org/openicpsr/project/115162/version/V1/view</a>
<b>Political Science Papers</b>		
Acquiescence Bias Inflates Estimates of Conspiratorial Beliefs and Political Misperceptions	10.1017/pan.2022.2	<a href="https://doi.org/10.7910/DVN/TVJCTX">https://doi.org/10.7910/DVN/TVJCTX</a>
Do Policy Makers Listen to Experts? Evidence from a National Survey of Local and State Policy Makers	10.1017/S0003055421000800	<a href="https://doi.org/10.7910/DVN/S2SNOT">https://doi.org/10.7910/DVN/S2SNOT</a>
<b>Papers Published in Behavioural Science Journals</b>		
Arrests and Convictions but Not Sentence Length Deter Terrorism in 28 European Union Member States	10.1038/s41562-023-01695-6	<a href="https://doi.org/10.5281/zenodo.8196717">https://doi.org/10.5281/zenodo.8196717</a>
Examining Inequality in the Time Cost of Waiting	10.1038/s41562-023-01524-w	<a href="https://github.com/stevebholt/waiting-time">https://github.com/stevebholt/waiting-time</a>
Mindful-Gratitude Practice Reduces Prejudice at High Levels of Collective Narcissism	10.1177/09567976231220902	<a href="https://osf.io/t7kxa/overview">https://osf.io/t7kxa/overview</a>

*Article Title* corresponds to the paper provided to different teams across our events. *DOI* (column 2) and *Replication Folder* reference what the Institute for Replication accessed prior to each of the AI Games.

**Table S3. Summary of Found Coding Errors in Published Economics Papers**

Article Title	Description of the Found Coding Errors
<b>Economics Papers</b>	
Major Reforms in Electricity Pricing: Evidence from a Quasi-Experiment	There are two discrepancies across four estimated models. First, the authors did not cluster their standard errors in Table two, columns four and eight of their regressions (where the manuscript suggest that they intended to). Second, the manuscript suggests that quarter dummies should have been included in columns two and six of Table 2 but were not included.
Sorting or Steering: The Effects of Housing Discrimination on Neighborhood Choice	Two coding errors: the first where city names are not correctly cleaned leading to too many fixed effects in their main specification. Second, their discrimination variable is incorrectly coded.
Taste-Based Gender Favouritism In High-Stake Decisions: Evidence from the Price is Right	There were three concrete discrepancies noted in this paper. First, their main variable of interest in Table 3 is defined in absolute distance as opposed to the relative measure as stated in the table notes. Second, the original authors use all ratings (four) to build their main variable as opposed to the last two ratings as claimed in the manuscript. Finally, there was an issue with merging, the replicators who originally worked on this replication folder noted approximately 7% different samples in rating scores when correcting this error.
The Heterogeneous Tax Pass-Through Under Different Vertical Relationships	Small inconsistencies in Figure 1 where the number of observations differ from what is stated in the published manuscript.
The Interplay Among Savings Accounts and Network-Based Financial Arrangements: Evidence from a Field Experiment	There are two stated discrepancies by the team of replicators. First, there are discrepancies in some reported p-values (TC of Table 1, CTd X Vi of Table 3, etc.). Second, the model specified in the manuscript includes time-fixed effects while the code and data that implements that model does not include time-fixed effects.
How the Other Half Died: Immigration and Mortality in U.S. Cities	There is a major issue with the merging of the database before analysis. Instead of merging one-to-one, there is a many-to-many merge, which duplicates observations.
Gambling, Saving, and Lumpy Liquidity Needs	There are two inconsistencies between the code provided in the replication folder and the manuscript. First, table 6 is supposed to drop those who already had the money necessary to make a purchase (manuscript) but does not appear to do so in the code. Second, code for Table 7 drops those who already had the money necessary to make a purchase but the manuscript does not suggest they should be excluded. Finally, the main writing in the manuscript does not suggest survey round fixed effects are included in the model for Table 5 but both the table notes and the code include survey round fixed effects.

*Article Title* corresponds to the paper provided to different teams across our events. We then provide a *Description of the Coding Errors* which had been identified by prior to the experiment.

**Table S4. Summary of Found Coding Errors in Published Political Science and Psychology Articles**

Article Title	Description of the Found Coding Errors
<b>Political Science Papers</b>	
Acquiescence Bias Inflates Estimates of Conspiratorial Beliefs and Political Misperceptions	The main variable (educational outcome) which is supposed to be an 8-category variable, contains outliers (values of -3105) for 14 respondents.
Do Policy Makers Listen to Experts? Evidence from a National Survey of Local and State Policy Makers	There are two concerns with this paper. First, the scripts in the replication folder repeatedly try to use a variable ( <i>bias</i> ) which does not exist in the data. To get this to work, one must create a variable for <i>bias</i> to try and reproduce Figure 4. Second, attempting to reproduce Table A22 contains counts and proportions which did not match the manuscript.
<b>Papers Published in Behavioural Science Journals</b>	
Arrests and Convictions but Not Sentence Length De-ter Terrorism in 28 European Union Member States	There are two concerns with this paper. First, there are 292 different values of a main dependent variable (attack rates) when using an inverse hyperbolic sine transformation when the raw value is zero. Second, variables contain impossible values or undisclosed imputations.
Examining Inequality in the Time Cost of Waiting	There are two concerns with this paper: first, Table 1, Panel (A), Column (6), incorrectly condition their sample to have positive waiting time which is inconsistent with what is written in the paper. Second, survey weights are not included (but were stated to have been used) in Table 1, Column 6.
Mindful-Gratitude Practice Reduces Prejudice at High Levels of Collective Narcissism	There were multiple coding errors in this study. The sample size on page 140 was slightly different than what was observed in the replication folder when keeping the sample stated in the text. There were claims of bootstrapping in the text which were not apparent in the code. Small deviations of confidence intervals and p-values in Table 1, with some statistics not reported at all. Table 2 could not be reproduced. There were general coding errors for Study 2.

*Article Title* corresponds to the paper provided to different teams across our events. We then provide a *Description of the Coding Errors* which had been identified by prior to the experiment.

**Table S5. Balance of paper metrics (AI games vs comparison)**

Metric	AI Games	Comparison	Difference
Size (MB)	994.179 (430.428)	631.185 (182.024)	362.994 [0.449]
Files in replication package	647.583 (484.275)	86.492 (46.105)	561.092 [0.273]
Script files	15.750 (7.164)	5.417 (0.964)	10.333 [0.180]
Data files	136.167 (97.244)	63.217 (44.115)	72.950 [0.504]
Script-to-data ratio	1.042 (0.459)	0.868 (0.278)	0.174 [0.749]
README present	0.583 (0.149)	0.450 (0.046)	0.133 [0.407]
Main file present	0.167 (0.112)	0.242 (0.039)	-0.075 [0.539]
Declared package count	9.417 (2.254)	3.292 (0.573)	6.125 [0.021]
Implied package count	11.750 (2.125)	2.975 (0.352)	8.775 [0.002]
Total pages	25.083 (3.498)	23.717 (1.249)	1.367 [0.718]
Main pages	22.583 (3.827)	23.717 (1.249)	-1.133 [0.783]

*Note:* Columns 2–3 report means and standard errors in parentheses for AI games (N=12) and comparison papers (N=120). The difference column shows AI minus comparison means with Welch *t*-test *p*-values in brackets. Binary variables are reported as proportions. Size is reported in MB.

**AI Training.** Researchers took part in a one-hour-long training on the usage of ChatGPT. This training was mandatory for the researchers in AI-assisted and AI-led groups.

Recordings and materials are publicly available here: <https://osf.io/sz2g8/>. The training included the following topics:

### 1. Introduction, Overview of ChatGPT, and Access

- Introduction to the capabilities of ChatGPT and its applications in reproducing scientific studies, coding, and data analysis.
- Instructions on accessing ChatGPT, creating an account, and accessing the Institute for Replication workspace/team subscription.
- Explanation of subscription tiers, model capabilities, limitations on message usage, and privacy settings.

### 2. Interaction with ChatGPT

- Techniques for optimizing prompts and ChatGPT's responses, such as providing contextual information.
- Strategies to manage randomness in outputs or when the model gets “stuck,” such as opening new chats and regenerating answers.

### 3. Sharing Chats with I4R

- Information on how to generate shareable links to sessions and manage privacy, including restrictions on who can access shared chats.
- Explanation on how to save chats as a webpage when the chat cannot be shared as a link (e.g., when the chat includes images).

### 4. Coding Assistance

- Explanation of how ChatGPT can assist with coding, including practical examples such as writing code for converting data formats (e.g., R's .rds to Stata's .dta) and debugging code.

## 5. Uploading Files and Images

- Introduction to ChatGPT's ability to process uploaded files.
- Overview of supported file types (e.g., PDFs, Word documents, CSVs, Excel files) and limitations regarding file size.
- Example of uploading an academic article to inquire about research questions, identification strategy, and robustness checks.
- Explanation of the potential benefits of uploading an image of a results table/figure instead of only the PDF file.
- Example of uploading an image of a results table from a study and inquiring about it.
- The image upload was not mentioned or demonstrated during training for the Toronto event.

## 6. Conducting Data Analysis Using ChatGPT

- Introduction to using ChatGPT's Data Analysis Module for executing Python code and performing data analysis.
- Example of uploading a replication package of an article and replicating regression analyses using the Python module.
- AI-led teams were instructed to first attempt to run the authors' codes/scripts using the data analysis capabilities of ChatGPT. If this analysis failed, teams were instructed to run the code in their local environment by following instructions provided by ChatGPT, as introduced in the Coding Assistance example.

## 7. ChatGPT API

- Explanation of the ChatGPT API for automating repetitive tasks and integrating AI capabilities into code.
- Example code shown for connecting to the ChatGPT API in R.

## 8. Customizing ChatGPT

- Information on setting up personalized models with custom instructions for specific needs.
- Mention of ChatGPT's memory feature that retains information across sessions, and how information that should not be retained can be deleted. The ChatGPT memory feature was not mentioned during the Toronto training session.

## 9. Explanation of Differences Among ChatGPT Models

- Differences between ChatGPT 4 and 4o were first discussed during the Sheffield event's training.
- Introduction of GPT-o1-preview and GPT-o1-mini models was first provided during the Bogota event's training.
- Capabilities of ChatGPT 4o with canvas were introduced during the last event's training.

**ChatGPT Models.** Researchers in the AI-assisted and AI-led groups were provided with access to ChatGPT Team. Table S6 presents an overview of the ChatGPT models available to these researchers during each event. Table S7 provides details about the capabilities of these models. Throughout all events, researchers had access to the main flagship model, GPT-4, and/or GPT-4o. These models were capable of processing files, equipped with a Python environment for interpreting code and conducting data analysis, and had internet access. OpenAI's web user interface and the desktop app do not allow for the execution of any code in a program language other than Python. That is still the case as of December 2025.

The file upload was limited to a maximum of 512MB per file, and further limited to two million tokens for text files, approximately 50MB for CSV files and spreadsheets, and 20MB per image for images. A user file size is capped at 10GB and organization at 100GB.<sup>(2)</sup> However, the practical limitations based on the Python environment's capabilities were likely lower.

Only researchers in the Bogota, Tilburg, and virtual-only events had access to the GPT-o1-preview and GPT-o1-mini models. These models were trained using reinforcement learning to perform complex reasoning and, unlike the 4/4o models, can produce an internal chain of thought before responding to users.<sup>(3)</sup> These models were not capable of processing files at the time of the events and thus were of very limited use to the AI-led teams.

Usage limits for certain models were applied by OpenAI. During the Toronto and Ottawa events, these limits were explicitly stated, with the Team subscription limit set at 100 messages per three hours per user. Researchers were instructed to collaborate with their teammates if the limit was reached or use the unlimited GPT-3.5 model. For the remaining events, usage limits for the GPT-4/4o models were no longer explicitly mentioned by OpenAI but were likely higher. The GPT-o1-preview model was limited to 50 queries per week, while GPT-o1-mini was limited to 50 queries per day.

Games	Date	Training Date	Image <sup>*</sup>	ChatGPT versions available					
				3.5	4	4o	4o-mini	o1-preview and o1-mini	4o with canvas
Toronto	Feb 20	Feb 14	No	Yes	Yes				
Ottawa	May 3	Apr 26	Yes	Yes	Yes				
Sheffield	Jun 17	Jun 12	Yes	Yes	Yes	Yes			
Cornell	Aug 12	Jul 31	Yes		Yes	Yes	Yes		
Bogota	Oct 4	Sep 23 <sup>†</sup>	Yes		Yes	Yes	Yes	Yes	Yes <sup>‡</sup>
Tilburg	Oct 18	Sep 30	Yes		Yes	Yes	Yes	Yes	Yes <sup>‡</sup>
Virtual	Nov 22	Nov 8	Yes		Yes	Yes	Yes	Yes	Yes

\* Image upload trained as part of the pre-games training and screenshots of relevant results from the studies provided to researchers.

<sup>†</sup> Training using recording of the Cornell training + o1-preview model slide added to presentation

<sup>‡</sup> While GPT-4o with canvas was available for the Bogota and Tilburg events, it was not mentioned during the training.

**Table S6. ChatGPT models available by event**

Model	Date Introduced	File Upload	Python Code Interpreter	Web Browsing	Reasoning
GPT-3.5	Before 1st event	No	No	No	No
GPT-4	Before 1st event	Yes	Yes	Yes	No
GPT-4o	May 13, 2024 <sup>1</sup>	Yes	Yes	Yes	No
GPT-4o-mini	July 18, 2024 <sup>2</sup>	Yes <sup>*</sup>	Yes <sup>*</sup>	Yes <sup>*</sup>	No
o1-preview	September 12, 2024 <sup>3</sup>	No	No	No	Yes
o1-mini	September 12, 2024 <sup>3</sup>	No	No	No	Yes
GPT-4o with canvas	October 3, 2024 <sup>4</sup>	Yes	Yes	No	No

\* While 4o-mini supported these functions at the time of the last training it did not necessarily at the time of introduction.

[1] <https://openai.com/index/hello-gpt-4o/>

[2] <https://openai.com/index/gpt-4o-mini-advancing-cost-efficient-intelligence/>

[3] <https://openai.com/index/introducing-openai-o1-preview/>

[4] <https://openai.com/index/introducing-canvas/>

**Table S7. ChatGPT capabilities by model**

**ChatGPT Prompts.** We conducted an exploratory analysis of ChatGPT transcripts. This revealed meaningful variation in prompting styles: in iteration and delegation. These patterns could form a basis for a future typology of prompting behaviors or failure modes. Preliminary review of AI-led transcript logs suggests that common failure modes included: (i) reliance on default package behavior without validation, (ii) misunderstanding file structure or model specifications, (iii) hallucination of non-existent variables or steps, (iv) limited ability to debug code execution errors, and (v) lack of prompt iteration or self-critique.

As a case study, we investigate the prompts from two teams of three researchers.

In the first team, each user brought value, but in distinct ways. Teammate 1 excelled in questioning the correctness of the figure generation process and raised crucial flags about the interpretability of marginal effects – though they sometimes required nudging to resolve issues independently. Teammate 2, while technically competent, engaged less rigorously with the underlying model logic and accepted ChatGPT’s answers too readily, potentially missing subtle but important discrepancies. Teammate 3 was the most thorough and experimental, implementing multiple robustness checks (cross-validation, outlier drops, subgroup splits) and producing interpretable, well-supported diagnostic output. Teammate 3’s iterative prompting was both strategic and technically sound, positioning ChatGPT more as a collaborator than a calculator.

Takeaway: To maximize ChatGPT’s potential for replication, teams benefit most from strategic prompting (Teammate 3), skeptical review of outputs (Teammate 1), and a clear sense of why each analysis step is done. Future participants might learn from Teammate 3’s modular prompting, Teammate 1’s critical eye, and Teammate 2’s clarity of implementation—ideally combining all three.

The second team showed strong consistency in core tasks: all three members successfully reproduced Figure 4 and applied the core logistic model. The first and second teammates shared overlapping concerns around code reliability, variable scaling, and indexing accuracy. However, only the third teammate conducted a wide array of robustness tests (cross-validation, subgroup splits, clustering, outlier removal). The second teammate completed reproduction but did not explore model sensitivity or robustness.

The biggest missed opportunity was a lack of follow-up on proposed robustness ideas (standardization, TF-IDF alternatives, placebo tests), despite their clear feasibility. A more integrated approach would have helped cross-pollinate promising ideas across team members.

DRAFT

**Coding Errors and Data Irregularities.** All studies to be reproduced had known coding errors that had been identified as part of a prior study but were not publicly disclosed at the time of the AI replication game. Of note, participants during the AI replication games identified additional errors.

Participants were not informed that the assigned papers necessarily contained at least one coding error or data irregularity, nor were they told how many such issues might exist or of what type. Participants were instructed to attempt a computational reproduction of the assigned study under conditions intended to mirror standard reproducibility exercises, in which researchers typically do not know ex ante whether errors are present. This design choice was made to avoid priming participants to search specifically for mistakes and to better approximate real-world reproduction settings.

Prior to the games, participants were not provided with a formal or exhaustive definition of what constituted a “coding error,” “discrepancy,” or “data irregularity.” The pre-games slides (<https://osf.io/sz2g8/files/7aguk>) focused on task logistics rather than conceptual definitions. However, our pre-analysis plan did include an operational definition of coding errors and data irregularities (<https://osf.io/sz2g8/files/b42ue>), which guided our ex post evaluation. In addition, participants were given a standardized Template Time Stamp (<https://osf.io/sz2g8/files/p7hd9>), which provide a very simple example: Example: we uncovered a coding error in the do-file entitled “XYZ.do”, line 45. The error involves forgetting a control variable or misspecifying something or miscoding a variable, etc.

We define coding errors as minor or major depending on whether the coding error could, in theory, have an impact on the claims tested. Additionally, we define a false positive to be an issue reported by a team that is not an actual error. DVA, JJPAp, DMi, and LFi discussed all errors uncovered during the AI games and classified coding errors as major, minor, or false positives. Minor coding errors uncovered were for example reporting the wrong p-value without impacting significance (e.g., 0.007 instead of 0.006). Major coding errors included mis-coding of the dependent variable or main independent variable or conducting a many-to-many merge instead of a many-to-one merge.

We did not keep track of the exact number of times that the four of us discussed hard cases (i.e., whether a reported error was not an error, a minor error or a major error), but it was rare. We always agreed on whether a reported error was truly an error (e.g., we always agreed that missing packages/paths or versioning issues were not errors), and all agreed on whether errors uncovered by teams were major or minor.

In what follows, we provide concrete examples of major coding errors and data irregularities. In the article entitled “Arrests and Convictions but Not Sentence Length Deter Terrorism in 28 European Union Member States,” one of the major coding errors is in the coding of the dependent variable. The authors state in the article that the terrorism rate used as their dependent variable is the inverse hyperbolic sine (IHS) of the per capita rate of terrorist attacks. But the code reveals that the dependent variable takes impossible values and is thus not the IHS of the per capita rate of terrorist attacks. For instance, countries with zero terror attacks are assigned strictly positive values, which is impossible. Another major coding error is that some European countries were imputed as experiencing zero terror attacks when the number of terror attacks was missing in the raw data. This article was retracted on January 8th, 2026. The retraction occurred as a result of a Matters Arising submission by one of our reproducers (4).

In the article “Sorting or Steering: The Effects of Housing Discrimination on Neighborhood Choice,” one of the major coding errors involved assigning a value of zero for the variable ‘of color’ to both individuals identified as ‘white’ and as ‘other’ in the raw data. A major data irregularity is the inclusion of fixed effects for the string variable ‘city’. The raw variable is case sensitive and has many spelling mistakes.

In the main article, we document that AI-led groups identified significantly fewer coding errors and data irregularities. AI-led groups likely uncovered fewer major coding errors due to the nuanced and contextual nature of these errors. It is plausible that AI-led groups struggled to identify technically correct but conceptually flawed code errors. These errors, such as a many-to-many merge instead of a many-to-one merge, produce duplicate entries without causing a runtime error. While the code executes without issues, the underlying conceptual mistake leads to incorrect data handling. This type of error is particularly challenging for AI to detect, as it requires an understanding of the conceptual intent behind the code rather than just its syntactic correctness.

More generally, many coding mistakes involve subtle misapplications of statistical transformations, such as assigning incorrect values or mishandling missing data, which often require domain expertise and a deep understanding of the data’s structure. AI tools, while efficient at automating tasks, may struggle with interpreting complex logical relationships, ambiguous data definitions, or recognizing implausible outcomes without explicit programming. In contrast, human-led groups are better equipped to identify errors that hinge on contextual reasoning, such as the incorrect coding of dependent variables or mis-assignments due to case-sensitive inconsistencies in datasets.

Finally, we also categorized coding errors along three dimensions: (i) whether the error occurs in preparing the data and analysis, (ii) whether the error is related to the regression analysis and (iii) whether it is a transcription error. We also created two additional variables: (i) the extent of false error detection and (ii) the share of errors not uncovered.

**Robustness Checks.** We propose four different binary measures which we believe qualify a good robustness check: (i) clarity (not vague) regarding purpose and execution; (ii) feasible, (iii) not previously done by the original author(s); and (iv) focuses on the validity of the empirical strategy.

In addition, we classify any corrections to major errors and rerunning the script as a “good” robustness check, although not complying with all the previous criteria. In the event that at least one of the above categories was hard to classify, we discussed and classified together (ABr, JJPAp and DMi). The classification was done by JJPAp across all events to avoid measurement error, and hard cases were discussed with ABr and DMi.

**Clarity (not vague) regarding purpose and execution:** It is possible that teams will not adequately describe their robustness check. This could be due to ChatGPT not sufficiently describing what they are doing, or from their own explanation. An example of a vague robustness check would be “adding control variables.” In contrast, a clear robustness check would be to precisely document which variable should be added as a control.

**Feasible:** For teams that are able to implement the robustness check, we categorize the robustness check as feasible. For teams that do not implement a robustness check, the question we ask is whether or not, with more time but the same resources, it could be implemented.

**Not previously done by the original author(s):** All teams of reproducers— independent of which type of team they are— have access to the original study, online appendix, and the replication packages. All teams must verify that their recommended robustness check was not previously done. AB, JA, and DM verified with each study whether the proposed robustness checks were included in the article or the appendix.

**Validity:** While robustness checks can serve multiple purposes, we view them as alternative specifications that test the main conclusion(s) of a study. A valid robustness check tests the reliability and stability of the results. Examples of invalid robustness checks include: using bad controls, misspecified models (bad instrument), etc.

**Textual and Sentiment Analysis.** We perform a textual analysis on the full set of prompts (4,111 prompts in total) for 34 AI-led (3,101 prompts) and 32 AI-assisted (1,010 prompts) teams. We do not have prompts available from one AI-led and three AI-assisted teams. Of note, we do not have prompts available from one AI-led and three AI-assisted teams. These prompts exclude scripts from programming languages. We find that AI-led teams interacted more with ChatGPT than AI-assisted teams both in the number of conversations and the number of prompts. Panel A of Table S17 shows that AI-led teams have on average 3.136 ( $p = 0.006$ ) conversations more than AI-assisted teams. A similar pattern emerges for the total number of prompts, where AI-led teams use on average 59.643 more prompts than AI-assisted teams ( $p < 0.001$ ). Together, these results indicate that AI-led teams rely more heavily on repeated and sustained interactions with the AI model.

We then conduct two types of sentiment analyses: dictionary-based and machine-based (5, 6). The main difference between these approaches lies in how sentiment is identified and quantified. Dictionary-based methods rely on a pre-defined lexicon (i.e., dictionary) that assigns sentiment or emotional categories to individual words, allowing sentiment to be measured through transparent and interpretable word-level matches. In contrast, machine-based methods use trained models to infer sentiment from the broader linguistic context, capturing complex patterns such as negation, word order, and semantic nuance that are not explicitly encoded in lexicons. While the dictionary-based analyses provide clear interpretability and facilitate direct comparisons across texts, machine-based analyses prioritize contextual understanding and predictive accuracy (7).

**Dictionary-based Sentiment Analysis:** We conducted this analysis based on the National Research Council (NRC) dictionary (8).

**Data and Sample Construction:** The analytical sample is restricted to prompts that did not contain embedded code, programming outputs, or direct quotes from the assigned papers. This restriction ensures that sentiment measures reflect natural language content rather than programming or paper syntax. Each remaining prompt was assigned a sequential numeric identifier to serve as the primary unit of analysis.

**Text Pre-processing:** Prompt text is normalized prior to sentiment analysis to reduce superficial variation. Specifically, all text was converted to lowercase, non-alphabetic characters and punctuation were removed, and excess whitespace was trimmed. The cleaned text is then tokenized into words, enabling direct matching with lexicon-based sentiment dictionaries.

**NRC Emotion Lexicon Analysis:** Emotional content is measured using the NRC Emotion Lexicon. This lexicon associates individual words with two general sentiment categories (positive and negative) and eight discrete emotions: anger, sadness, joy, fear, trust, disgust, anticipation, and surprise. Words may be associated with multiple emotion categories simultaneously.

Tokenized words are matched to the NRC Emotion lexicon, and the number of words associated with each emotion and sentiment category is counted for each prompt. These counts are aggregated at the prompt level, resulting in one row per prompt, with separate columns for each emotion and sentiment category. Prompts with no matching words in a given category are assigned a value of zero.

**Machine-Based Sentiment Analysis:** In addition to dictionary-based approaches, we conduct machine-based sentiment analyses using transformer-based language models(9, 10). These models are pre-trained on large corpora (i.e., texts) to infer sentiment and emotional content from the contextual meaning of text (rather than relying on predefined word lists like dictionary-based sentiment analysis).

We implement two distinct machine-based sentiment analysis methods: emotional classification and multiclass sentiment classification. Although both approaches rely on contextual language models, they differ fundamentally in their analytical objective, output structure, and interpretive focus. The former method identifies discrete emotional states, while the latter classifies overall sentiment polarity. Together, these approaches provide a context-sensitive measure of emotional tone that mirrors the lexicon-based NRC emotion classification

**Emotion Classification Using a Transformer Model:** Emotional content is analyzed using a pre-trained DistilRoBERTa (Distilled Robustly Optimized BERT Approach) transformer model fine-tuned for emotion classification. BERT (Bidirectional Encoder Representations from Transformers) is a transformer-based language model that learns contextual word representations by jointly considering both left and right contexts. The model assigns one of several discrete emotion labels (joy, anger, sadness, fear, disgust, surprise, and neutral) to each text, along with a confidence score reflecting the model’s certainty.

Each prompt is passed individually to the emotion classification pipeline. The model outputs a single predicted emotion label and an associated probability score for each prompt. The predicted emotion represents the emotion with the highest

posterior probability given the input text.

**Multi-Class Sentiment Classification Using a Transformer Model:** Overall sentiment polarity is also assessed using a pretrained RoBERTa (Robustly Optimized BERT Approach)-based transformer model designed for three-class sentiment classification. This model assigns one of three sentiment categories (negative, neutral, or positive) to each prompt and returns a confidence score indicating prediction certainty.

Each prompt was analyzed independently. The model evaluates the full textual context, allowing it to account for linguistic features such as negation, intensifiers, and compositional meaning that are not captured by dictionary-based methods.

**Results from Sentiment Analyses:** AI-led teams expressed higher levels of trust and positive affect, alongside lower levels of fear when compared to AI-assisted teams (Panels B and C, Table S17). AI-led teams also exhibited higher levels of anger in user-generated prompts. In contrast, AI-assisted teams interacted with ChatGPT more selectively and exhibited greater fear in their prompts. Together, these patterns suggest more cautious engagement and sustained human oversight from AI-assisted teams while AI-led teams may have been more frustrated with perceived inaccuracies of model outputs.

Using the NRC Emotion Lexicon (Table S17, Panel B), we find that AI-assisted teams express significantly more anger and fear, and significantly less trust, positive sentiment, and joy relative to AI-led teams. Prompts generated by AI-assisted teams contain higher levels of anger (mean difference = 0.031,  $p < 0.001$ ) and fear (mean difference = 0.024,  $p = 0.056$ ). In contrast, expressions of trust, positive sentiment, and joy are all significantly lower for AI-assisted teams, as shown in Table A2.

To address limitations of dictionary-based methods, we complement the NRC analysis with a BERT-based sentiment classifier (Table S17, Panel C) that incorporates contextual meaning and assigns each prompt to a single dominant sentiment category. Under this approach, AI-assisted teams exhibit less anger, more fear, and a higher prevalence of neutral sentiment relative to AI-led teams, as revealed in Table A3.

Given that the BERT classifier accounts for semantic context and enforces mutually exclusive sentiment assignments, we interpret these results as providing a more reliable characterization of emotional expression in prompts. Notably, under this approach, anger appears more dominant among AI-led teams rather than AI-assisted teams.

**Focus Groups: Corpus and Analytic Approach.** We applied reflexive thematic analysis in the sense of Braun and Clarke (2006)(11). Before coding began, we built a shared codebook seeded with theory-driven categories linked to our research questions (e.g., “task delegation,” “error detection”) and left space for inductive labels to emerge. Five analysts then read each transcript in full and independently tagged text segments with the evolving codebook. The coders iteratively reconciled discrepancies, refined code definitions, and updated the codebook; all individual codes remained visible in a common spreadsheet so that divergent interpretations stayed traceable. Finally, we overlaid the coding layers, examined convergent and divergent patterns across treatment arms, and iteratively collapsed related codes into higher-order themes. This reflexive yet transparent workflow preserved the nuance of individual readings while yielding a coherent thematic structure—crucial for triangulating the qualitative insights with the quantitative performance metrics reported below.

**Focus Groups: Additional Secondary Results.** Teams began with high hopes for AI: “we were both quite optimistic about ChatGPT and what it can do” (FG3-A-P3, 29 Apr 2025). Hopes faded quickly: “I dampened down that enthusiasm because I saw that ... the work was getting done better without it” (FG4-A-P2, 30 Apr 2025) and “something that should be push-button goes from a task that should take minutes into a task that takes hours” (FG1-L-P3, 18 Apr 2025).

Many AI-assisted teams treated ChatGPT as a tool of last resort. One recalled “we started with the assumption of going through the code by ourselves” (FG3-A-P4), another “only copy-pasted the error message ... if we were desperate” (FG2-A-P4), and a third “just ask quite a lot of questions honestly ... only when we found an error” (FG2-A-P1). In contrast, AI-led teams were tasked with fully automated runs, hoping to “put as little input as possible from myself” (FG3-L-P3). When archives overflowed the context window they switched to piecemeal prompting. Thus, speed hinged on where teams drew the *delegation threshold*: fast AI-assisted teams outsourced micro-tasks yet kept strategic control; AI-led teams ceded entire analytic stages and struggled to reclaim them.

Human expertise remained essential. One AI-led participant reflected, “Let’s approach the analysis as researchers ... build it up from the bottom up” (FG3-L-P3). Yet they still felt hopeless at error-hunting: “We could have done this for ten more hours—we would not find the error. No way.” (FG1-L-P3). Others blamed AI’s blind spots: “[ChatGPT] is very convenient for summarising stuff but then there’s this five percent it drops which is crucial” (FG1-A-P5) and “some issues are so particular—say institutional details relevant for an identification strategy” (FG3-H-P2).

Across AI arms, a pragmatic repertoire emerged: specify software version, feed minimal code snippets, and request pseudocode before executable commands. Still, prompting was effortful: “Eventually I just ran out of ideas on how else to prompt it” (FG6-L-P3). Over-confidence compounded matters; once ChatGPT raised no alarms, a teammate admitted “our second trial by ourselves was quite light ... we didn’t try very hard, to be honest” (FG3-A-P5). Some even deferred responsibility: “If the replication is wrong ... it’s GPT-4’s problem” (FG5-A-P1). The resulting pipelines were fast but shallow—mirroring AI-assisted teams’ lower major-error counts.

AI-led transcripts brim with exasperation: “ChatGPT was like an undergrad who hadn’t used econometrics ... unapologetically arrogant when it’s wrong” (FG6-L-P3). Another recalled “hours of nothing ... talking to ChatGPT, trying to find an error and not getting anywhere” (FG1-L-P4). Repetitive copy-paste loops and hallucinated file paths eroded trust: “I gave it the zipped folder and it said it couldn’t open it—though it had done so in the previous chat” (FG6-L-P1). AI-assisted frustration was milder—“that is wrong ... and then I have to ask it again ... it’s just *slow*” (FG4-A-P2)—because they

reverted to manual coding once cost–benefit turned negative: “The focus shifted to manual coding because the first attempt with ChatGPT failed, so I thought: let’s not waste too much time there” (FG4-A-P2).

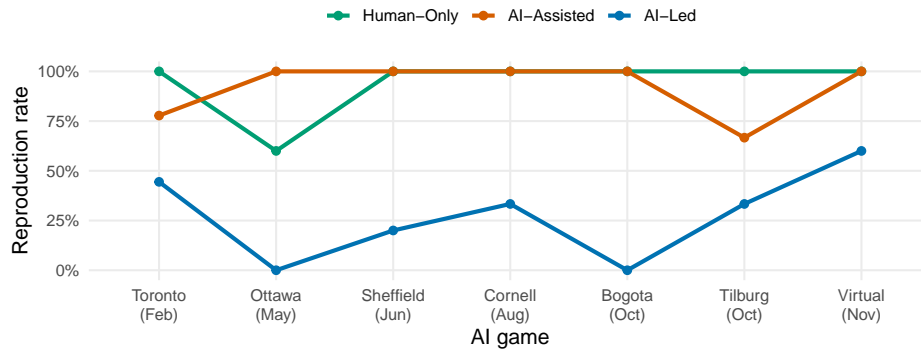
AI reshaped collaboration. AI-assisted members often kept separate ChatGPT tabs: “both . . . separately working one-on-one with ChatGPT, which was a little bit sad” (FG3-A-P6). Coordination costs mounted: “we were circling back discussing “Is this a fair prompt?”” (FG3-L-P3). Duplicated effort surfaced: “by the time ChatGPT produced the table, the team had already figured it out” (FG4-A-P2). Disagreement between AI and humans triggered arbitration: “ChatGPT flagged something; I tried to replicate it and couldn’t” (FG4-L-P1). AI-led teams felt “the human-human connection was lower than in a standard team” (FG4-L-P3). Human-only groups relied on dense pair-programming; reading papers together was “something only humans can understand” (FG2-H-P2), a practice that surfaced subtle errors invisible to models.

AI-assisted teams accelerated routine steps yet—owing to prompt fatigue and overconfidence—missed errors that diligent human-only teams caught. AI-led teams, mired in technical bottlenecks and diminished peer oversight, achieved neither speed nor depth. As one AI-led participant put it: “We could reproduce that number super-fast . . . and then it was *hours of nothing*” (FG6-L-P4). Conversely, an AI-assisted analyst conceded: “ChatGPT just sped things up . . . we were much faster. Yeah, I don’t think we were better, we were just faster” (FG4-A-P3).

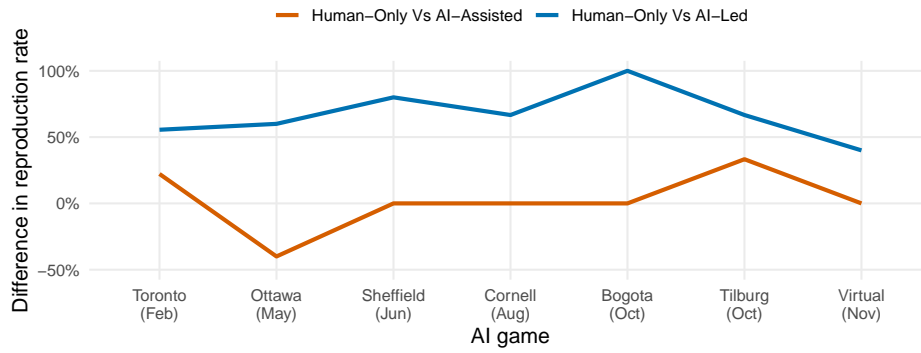
Despite frustrations, participants acknowledged AI’s benefits when used as a complement. AI suggested robustness checks: “some really interesting suggestions we didn’t think of ourselves” (FG3-L-P3). It helped locate code snippets: “it very quickly directed us to the relevant files” (FG3-L-P3). Less-experienced coders valued boilerplate help: “GPT-4 can definitely code better than me” (FG5-A-P1). Human-arm participants admitted the absence of AI nudged them toward “low-hanging fruit,” implying AI can broaden analytical ambition.

Participants also praised rapid LLM progress but warned that fully automated replication remains distant: “Right now it just seems too big of a task” (FG4-L-P1).

**Focus Groups: Implications for future human–AI workflows.** Effective LLM use is now a specialized skill on par with fluency in *Stata* or *R*. Tool designers should prioritize larger context windows, transparent code provenance, and archive-level ingestion to cut the friction that crippled AI-led teams. Pedagogically, prompt engineering and AI-assisted debugging must join the *reproduction toolkit*. Most importantly, our findings caution against premature automation: April-2025 LLMs excel as accelerators for routine sub-tasks but falter as autonomous reproducers. Collective, deliberative human judgment currently remains indispensable for detecting the subtle conceptual and coding errors that determine whether published science withstands scrutiny.



**Fig. S1.** Computational reproducibility rates across events and treatment groups



**Fig. S2.** Difference in average computational reproducibility rate by groups across AI replication game

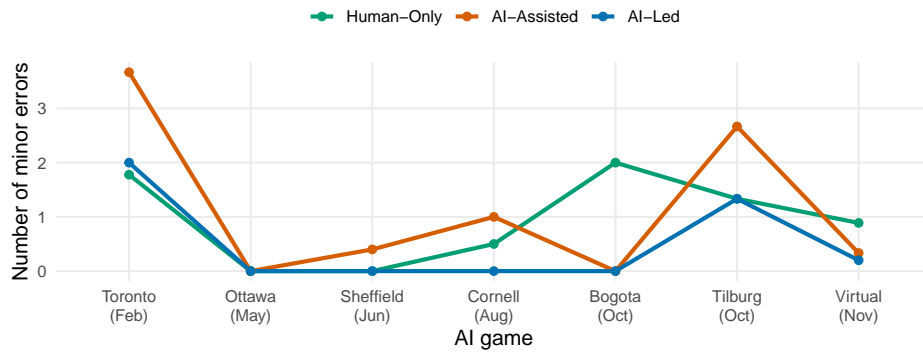


Fig. S3. Number of minor errors detected across events and treatment groups

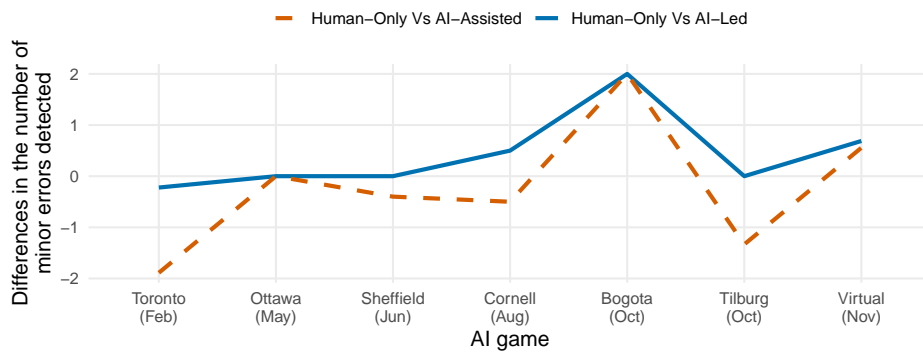


Fig. S4. Difference in average number of minor errors discovered by groups across AI replication games

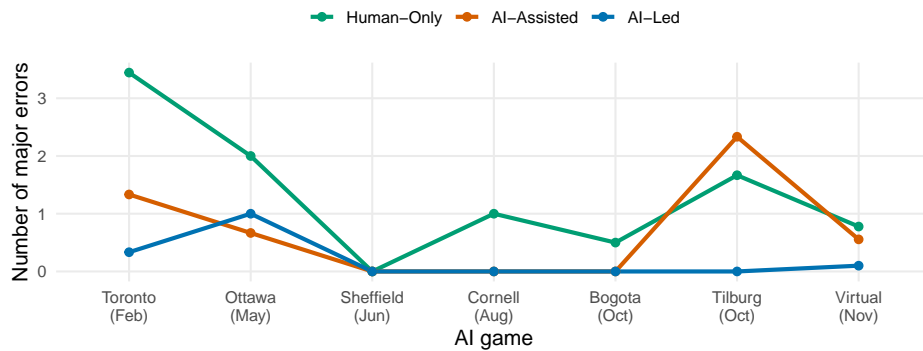


Fig. S5. Number of major errors detected across events and treatment groups

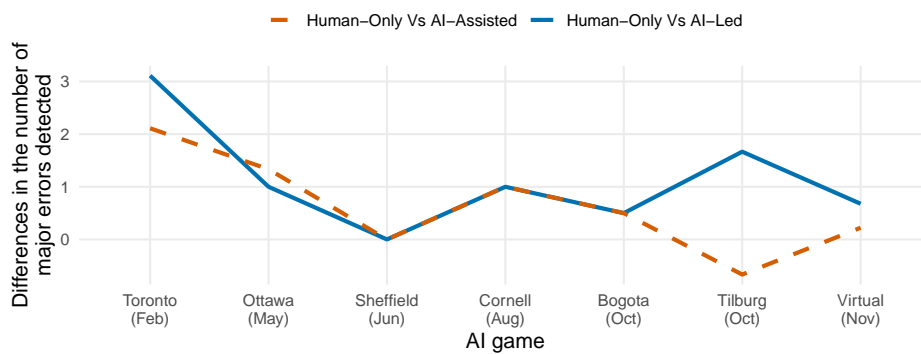


Fig. S6. Difference in average number of major errors discovered by groups across AI replication games

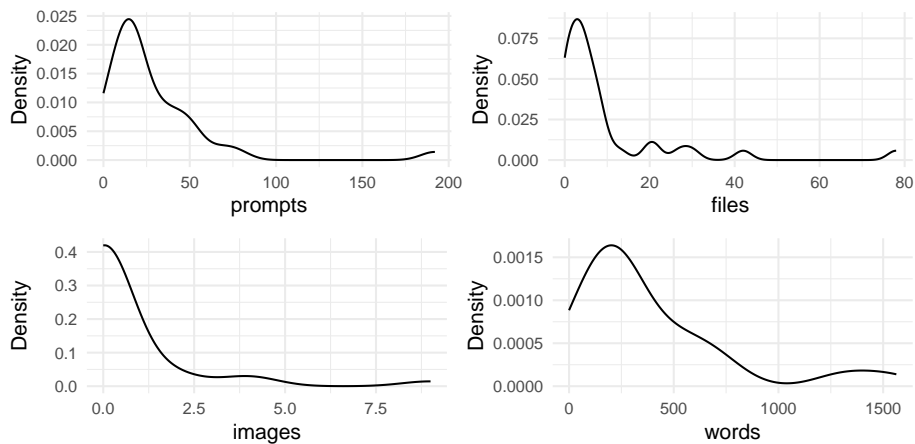


Fig. S7. Distribution of number of prompts, files, images, and words used between AI-Assisted teams and ChatGPT

**Table S8. Balance of Team-Level Characteristics by Group**

Variable	Human-Only	AI-Assisted	AI-Led	Human-Only vs AI-Assisted	Human-Only vs AI-Led	AI-Assisted vs AI-Led
Number of teammates	2.606 (0.496)	2.429 (0.655)	2.829 (0.568)	0.177 [0.214]	-0.223 [0.091]	-0.400 [0.008]
Minimum academic level: Professor	0.091 (0.292)	0.086 (0.284)	0.086 (0.284)	0.005 [0.941]	0.005 [0.941]	0.000 [1.000]
Minimum academic level: Postdoc	0.030 (0.174)	0.114 (0.323)	0.057 (0.236)	-0.084 [0.190]	-0.027 [0.597]	0.057 [0.400]
Minimum academic level: Researcher	0.152 (0.364)	0.171 (0.382)	0.029 (0.169)	-0.020 [0.827]	0.123 [0.076]	0.143 [0.047]
Minimum academic level: Student	0.727 (0.452)	0.629 (0.490)	0.829 (0.382)	0.099 [0.392]	-0.101 [0.321]	-0.200 [0.061]
Maximum academic level: Professor	0.576 (0.502)	0.514 (0.507)	0.686 (0.471)	0.061 [0.617]	-0.110 [0.355]	-0.171 [0.147]
Maximum academic level: Postdoc	0.152 (0.364)	0.257 (0.443)	0.143 (0.355)	-0.106 [0.289]	0.009 [0.921]	0.114 [0.238]
Maximum academic level: Researcher	0.091 (0.292)	0.057 (0.236)	0.000 (0.000)	0.034 [0.600]	0.091 [0.070]	0.057 [0.156]
Maximum academic level: Student	0.182 (0.392)	0.171 (0.382)	0.171 (0.382)	0.010 [0.912]	0.010 [0.912]	-0.000 [1.000]
Average years of coding experience	9.000 (4.484)	8.267 (3.060)	9.740 (3.365)	0.733 [0.431]	-0.740 [0.442]	-1.474 [0.059]
Min ChatGPT level: Never	0.303 (0.467)	0.143 (0.355)	0.257 (0.443)	0.160 [0.115]	0.046 [0.679]	-0.114 [0.238]
Min ChatGPT level: Beginner	0.485 (0.508)	0.629 (0.490)	0.571 (0.502)	-0.144 [0.239]	-0.087 [0.482]	0.057 [0.632]
Min ChatGPT level: Intermediate	0.152 (0.364)	0.200 (0.406)	0.143 (0.355)	-0.048 [0.607]	0.009 [0.921]	0.057 [0.533]
Min ChatGPT level: Advanced	0.061 (0.242)	0.000 (0.000)	0.029 (0.169)	0.061 [0.144]	0.032 [0.527]	-0.029 [0.321]
Max ChatGPT level: Never	0.000 (0.000)	0.029 (0.169)	0.029 (0.169)	-0.029 [0.335]	-0.029 [0.335]	-0.000 [1.000]
Max ChatGPT level: Beginner	0.152 (0.364)	0.143 (0.355)	0.086 (0.284)	0.009 [0.921]	0.066 [0.408]	0.057 [0.460]
Max ChatGPT level: Intermediate	0.515 (0.508)	0.514 (0.507)	0.629 (0.490)	0.001 [0.994]	-0.113 [0.352]	-0.114 [0.341]
Max ChatGPT level: Advanced	0.333 (0.479)	0.286 (0.458)	0.257 (0.443)	0.048 [0.677]	0.076 [0.498]	0.029 [0.792]

*Note:* Columns 2–4 present means and standard errors in parentheses for individual groups (Human-only, AI-Assisted, and AI-Led); the difference columns show mean differences and *p*-values in brackets for the indicated group comparisons.

**Table S9. Comparison of Error Type Metrics by Group**

Variable	Human-Only	AI-Assisted	AI-Led	Human-Only vs	Human-Only vs	AI-Assisted vs
				AI-Assisted	AI-Led	AI-Led
Pre-regression errors	0.909 (1.284)	0.971 (1.424)	0.486 (1.040)	-0.062 [0.851]	0.423 [0.139]	0.486 [0.108]
Regression errors	1.182 (2.068)	0.657 (1.136)	0.171 (0.618)	0.525 [0.196]	1.010 [0.007]	0.486 [0.030]
Transcription/post-regression errors	0.909 (1.508)	0.571 (1.267)	0.286 (0.572)	0.338 [0.320]	0.623 [0.026]	0.286 [0.228]
False error detection	0.727 (1.701)	0.257 (0.505)	1.057 (1.830)	0.470 [0.122]	-0.330 [0.445]	-0.800 [0.015]
Share of known errors not found	0.784 (0.231)	0.816 (0.215)	0.936 (0.101)	-0.032 [0.560]	-0.152 [<0.001]	-0.120 [0.004]

Note: Columns 2–4 present means and standard errors in parentheses for individual groups (Human-only, AI-Assisted, and AI-Led); columns 5–7 present differences in means and p-values in brackets for group comparisons (Human-Only vs AI-Assisted, Human-Only vs AI-Led, and AI-Assisted vs AI-Led).

**Table S10. Causal relationship between treatment groups and reproducibility outcomes using Logit and Poisson regressions**

	(1)	(2)	(3)	(4)	(5)	(6)
	Reproduction	Minor errors	Major errors	Two good robustness	Ran one robustness	Ran two robustness
AI-Assisted	-0.724 ( 1.124) [ -2.928; 1.479]	0.267 ( 0.389) [ -0.495; 1.030]	-0.572** ( 0.273) [ -1.106; -0.038]	-0.388 ( 1.202) [ -2.743; 1.968]	-0.959 ( 1.758) [ -4.405; 2.487]	-0.025 ( 0.949) [ -1.886; 1.835]
AI-Led	-6.478*** ( 1.648) [ -9.709; -3.247]	-0.449 ( 0.420) [ -1.272; 0.374]	-1.543*** ( 0.451) [ -2.428; -0.659]	-2.480** ( 1.227) [ -4.885; -0.075]	-3.701*** ( 1.280) [ -6.209; -1.192]	-1.777* ( 0.926) [ -3.591; 0.037]
Model	Logit	Poisson	Poisson	Logit	Logit	Logit
Controls	✓	✓	✓	✓	✓	✓
Mean of dep. var	0.738	0.963	0.881	0.786	0.816	0.680
p-val (AI-Assisted vs. AI-Led)	0.002	0.112	0.030	0.049	0.129	0.024
Obs.	103	82	84	103	103	103

Note: Standard errors in parentheses, confidence intervals in brackets; human-only group omitted. Logit coefficients reported.

Controls include number of teammates; game-by-software fixed effects; maximum and minimum position skill fixed effects; attendance fixed effects.

\*  $p < 0.1$ , \*\*  $p < 0.05$ , \*\*\*  $p < 0.01$

**Table S11. Causal relationship between treatment groups and error types using Poisson regressions**

	(1)	(2)	(3)	(4)
	Pre-regression errors	Regression errors	Transcription/post-regression errors	False error detection
AI-Assisted	0.004 ( 0.274) [ -0.534; 0.542]	-0.392 ( 0.333) [ -1.045; 0.262]	-0.660 ( 0.472) [ -1.584; 0.264]	-0.822 ( 0.583) [ -1.966; 0.321]
AI-Led	-0.587 ( 0.367) [ -1.306; 0.132]	-1.487** ( 0.689) [ -2.837; -0.136]	-1.270*** ( 0.454) [ -2.160; -0.379]	0.399 ( 0.563) [ -0.705; 1.503]
Model	Poisson	Poisson	Poisson	Poisson
Controls	✓	✓	✓	✓
Mean of dep. var	0.773	0.517	0.529	0.750
p-val (AI-Assisted vs. AI-Led)	0.116	0.119	0.231	0.031
Obs.	88	60	70	92

Note: Standard errors in parentheses, confidence intervals in brackets; human-only group omitted. Poisson for count outcomes.

Controls include number of teammates; game-by-software fixed effects; maximum and minimum position skill fixed effects; attendance fixed effects.

\*  $p < 0.1$ , \*\*  $p < 0.05$ , \*\*\*  $p < 0.01$

Table S12. Estimates for the control variables

	(1)	(2)	(3)	(4)	(5)	(6)	(7)
	Reproduction	Minor errors	Major errors	One good robustness	Two good robustness	Ran one robustness	Ran two robustness
Branch: AI-Assisted	-0.018 (0.063)	-0.455 (0.263)	-0.646** (0.254)	-0.009 (0.027)	-0.014 (0.103)	-0.032 (0.061)	-0.009 (0.113)
Branch: AI-Led	[-0.142; 0.105]	[-0.970; 0.061]	[-1.145; -0.147]	[-0.062; 0.045]	[-0.217; 0.188]	[-0.153; 0.088]	[-0.229; 0.212]
Number of teammates	-0.593*** (0.090)	-1.029*** (0.255)	-1.136** (0.235)	-0.167*** (0.068)	-0.250** (0.107)	-0.323*** (0.098)	-0.290** (0.126)
Game: Ottawa	[-0.770; -0.416]	[-1.528; -0.529]	[-1.597; -0.674]	[-0.300; -0.033]	[-0.460; -0.040]	[-0.515; -0.130]	[-0.536; -0.044]
Game: Sheffield	0.052 (0.069)	0.368* (0.202)	0.224 (0.196)	0.050 (0.050)	0.091 (0.083)	-0.038 (0.077)	0.059 (0.094)
Game: Cornell	[-0.083; 0.187]	[-0.028; 0.765]	[-0.161; 0.609]	[-0.124; 0.071]	[-0.071; 0.253]	[-0.190; 0.114]	[-0.125; 0.242]
Game: Bogota	-0.086 (0.155)	-1.410** (0.692)	-0.243 (0.460)	-0.026 (0.172)	0.098 (0.217)	-0.351* (0.177)	0.069 (0.205)
Game: Tilburg	[-0.389; 0.217]	[-2.767; -0.053]	[-1.144; 0.659]	[-0.386; 0.288]	[-0.326; 0.523]	[-0.698; -0.003]	[-0.471; 0.334]
Game: Virtual Europe	-0.180 (0.259)	-1.610** (0.648)	-0.653 (0.464)	0.159 (0.136)	-0.006 (0.400)	-0.336 (0.310)	-0.010 (0.360)
Game: Virtual North America	[-0.688; 0.328]	[-2.879; -0.341]	[-1.563; 0.257]	[-0.106; 0.425]	[-0.791; 0.779]	[-0.944; 0.271]	[-0.715; 0.696]
Game: Bogota	0.276 (0.190)	-1.539** (0.540)	-1.075 (0.499)	0.061 (0.120)	0.317 (0.201)	0.055 (0.180)	0.285 (0.233)
Game: Tilburg	[-0.097; 0.649]	[-2.597; -0.482]	[-2.053; -0.097]	[-0.176; 0.296]	[-0.077; 0.711]	[-0.299; 0.409]	[-0.176; 0.746]
Game: Virtual Europe	0.014 (0.175)	-1.101** (0.998)	-1.465** (1.021)	0.016 (0.136)	0.071 (0.350)	-0.189 (0.178)	0.170 (0.354)
Game: Virtual North America	[-0.328; 0.357]	[-3.058; 0.855]	[-3.467; 0.536]	[-0.250; 0.283]	[-0.615; 0.757]	[-0.537; 0.160]	[-0.864; 0.524]
Software: R	0.231 (0.173)	-2.190*** (0.701)	0.532 (0.720)	0.067 (0.112)	0.398 (0.168)	-0.076 (0.175)	0.230 (0.184)
Maximum academic level: Researcher	[-0.109; 0.570]	[-3.563; -0.816]	[-0.878; 1.943]	[-0.151; 0.286]	[-0.069; 0.728]	[-0.418; 0.266]	[-0.130; 0.591]
Maximum academic level: Postdoc	0.004 (0.161)	-1.980*** (0.620)	0.818 (0.587)	0.098 (0.110)	0.351* (0.174)	0.092 (0.128)	0.245 (0.250)
Maximum academic level: Professor	[-0.311; 0.319]	[-3.196; -0.763]	[-1.970; 0.333]	[-0.117; 0.313]	[-0.010; 0.691]	[-0.158; 0.343]	[-0.246; 0.735]
Minimum academic level: Researcher	0.123 (0.180)	-1.544*** (0.468)	0.583 (0.458)	0.088 (0.112)	0.327 (0.163)	0.130 (0.179)	0.092 (0.205)
Minimum academic level: Postdoc	[-0.229; 0.476]	[-2.460; -0.627]	[-1.481; 0.315]	[-0.131; 0.308]	[-0.009; 0.646]	[-0.481; 0.222]	[-0.309; 0.494]
Minimum academic level: Professor	-0.169 (0.154)	0.926** (0.619)	0.249 (0.513)	0.006 (0.121)	0.118 (0.183)	-0.014 (0.123)	0.096 (0.184)
Attendance: In-Person	[-0.471; 0.133]	[-0.286; 2.138]	[-0.756; 1.255]	[-0.232; 0.244]	[-0.240; 0.476]	[-0.256; 0.228]	[-0.265; 0.457]
Game: Ottawa × Software: R	0.148 (0.180)	-1.535** (0.626)	0.336 (1.559)	-0.015 (0.091)	-0.213 (0.195)	0.105 (0.159)	-0.053 (0.219)
Game: Sheffield × Software: R	[-0.205; 0.501]	[-2.761; -0.308]	[-2.720; 3.391]	[-0.194; 0.165]	[-0.595; 0.168]	[-0.207; 0.416]	[-0.483; 0.377]
Game: Cornell × Software: R	0.110 (0.177)	0.042 (0.332)	0.275 (0.318)	0.032 (0.082)	-0.090 (0.145)	0.314** (0.146)	0.196 (0.172)
Game: Bogota × Software: R	[-0.237; 0.457]	[-0.609; 0.693]	[-0.348; 0.898]	[-0.129; 0.193]	[-0.375; 0.194]	[-0.027; 0.601]	[-0.141; 0.533]
Game: Tilburg × Software: R	0.030 (0.140)	0.027 (0.243)	0.340 (0.262)	-0.043 (0.090)	-0.165 (0.128)	0.107 (0.145)	-0.008 (0.147)
Game: Virtual Europe × Software: R	[-0.244; 0.305]	[-0.449; 0.503]	[-0.174; 0.854]	[-0.220; 0.134]	[-0.415; 0.085]	[-0.177; 0.391]	[-0.296; 0.279]
Game: Virtual North America × Software: R	-0.140 (0.091)	-0.112 (0.348)	0.285 (0.519)	-0.012 (0.066)	0.184 (0.116)	0.078 (0.108)	0.269 (0.137)
Game: Ottawa × Software: R	[-0.319; 0.039]	[-0.793; 0.570]	[-0.732; 1.302]	[-0.140; 0.117]	[-0.044; 0.411]	[-0.134; 0.289]	[-0.000; 0.537]
Game: Sheffield × Software: R	-0.080 (0.214)	0.121 (0.838)	-0.150 (0.435)	-0.127 (0.134)	-0.082 (0.183)	0.033 (0.123)	0.005 (0.195)
Game: Cornell × Software: R	[-0.500; 0.339]	[-1.522; 1.763]	[-1.001; 0.702]	[-0.390; 0.136]	[-0.441; 0.278]	[-0.209; 0.275]	[-0.377; 0.388]
Game: Bogota × Software: R	-0.094 (0.150)	0.698* (0.430)	0.535 (0.414)	0.001 (0.051)	-0.081 (0.187)	0.006 (0.141)	-0.076 (0.223)
Game: Tilburg × Software: R	[-0.388; 0.199]	[-0.145; 1.541]	[-0.277; 1.347]	[-0.099; 0.101]	[-0.446; 0.285]	[-0.270; 0.282]	[-0.513; 0.362]
Game: Virtual Europe × Software: R	-0.124 (0.120)	0.365 (0.362)	0.205 (0.245)	-0.012 (0.085)	-0.072 (0.121)	0.227** (0.120)	0.196 (0.127)
Game: Virtual North America × Software: R	[-0.358; 0.111]	[-0.346; 1.075]	[-0.274; 0.685]	[-0.178; 0.153]	[-0.309; 0.165]	[-0.009; 0.463]	[-0.053; 0.446]
Game: Ottawa × Software: R	-0.281 (0.285)	-1.463** (0.678)	-0.306 (0.647)	-0.094 (0.250)	-0.066 (0.307)	-0.216 (0.324)	-0.352 (0.353)
Game: Sheffield × Software: R	[-0.840; 0.278]	[-2.791; -0.134]	[-1.573; 0.962]	[-0.585; 0.397]	[-0.668; 0.535]	[-0.851; 0.419]	[-1.044; 0.341]
Game: Cornell × Software: R	0.322 (0.316)	-1.369 (0.639)	-0.751 (0.621)	-0.276 (0.200)	-0.260 (0.458)	0.178 (0.341)	-0.321 (0.430)
Game: Bogota × Software: R	[-0.297; 0.942]	[-2.621; -0.116]	[-1.968; 0.465]	[-0.668; 0.116]	[-1.157; 0.637]	[-0.490; 0.846]	[-1.164; 0.522]
Game: Tilburg × Software: R	-0.265 (0.244)	-1.320 (0.614)	0.424 (0.742)	-0.012 (0.151)	-0.115 (0.239)	-0.449 (0.227)	-0.557 (0.303)
Game: Virtual Europe × Software: R	[-0.743; 0.212]	[-2.523; -0.117]	[-1.031; 1.879]	[-0.308; 0.284]	[-0.584; 0.354]	[-0.894; -0.005]	[-1.150; 0.037]
Game: Virtual North America × Software: R	0.112 (0.319)	-1.838* (1.103)	0.727 (1.032)	0.010 (0.165)	0.103 (0.384)	0.054 (0.259)	0.044 (0.420)
Game: Ottawa × Software: R	[-0.514; 0.737]	[-3.999; 0.323]	[-1.296; 2.751]	[-0.314; 0.334]	[-0.649; 0.855]	[-0.562; 0.454]	[-0.778; 0.867]
Game: Sheffield × Software: R	-0.364 (0.252)	-0.229 (0.805)	-1.553* (0.898)	-0.026 (0.134)	-0.697** (0.294)	-0.250 (0.310)	-0.915** (0.270)
Game: Cornell × Software: R	[-0.858; 0.129]	[-1.807; 1.349]	[-3.312; 0.207]	[-0.290; 0.237]	[-1.274; -0.119]	[-0.857; 0.357]	[-1.445; -0.386]
Game: Bogota × Software: R	0.234 (0.222)	-0.836 (0.773)	-0.597 (0.745)	-0.024 (0.140)	-0.424 (0.287)	-0.111 (0.189)	-0.287 (0.341)
Game: Tilburg × Software: R	[-0.201; 0.669]	[-2.350; 0.679]	[-2.057; 0.864]	[-0.298; 0.250]	[-0.986; 0.138]	[-0.480; 0.259]	[-0.955; 0.381]
Game: Virtual Europe × Software: R	0.076 (0.285)	-1.148 (0.722)	-0.216 (0.836)	0.002 (0.138)	-0.312 (0.316)	0.110 (0.275)	-0.210 (0.377)
Game: Virtual North America × Software: R	[-0.482; 0.634]	[-2.563; 0.268]	[-1.855; 1.423]	[-0.268; 0.272]	[-0.931; 0.307]	[-0.429; 0.649]	[-0.949; 0.528]
Controls	✓	✓	✓	✓	✓	✓	✓
Mean of dep. var	0.738	0.942	0.728	0.942	0.786	0.816	0.680
Obs.	103	103	103	103	103	103	103

Note: Standard errors in parentheses, confidence intervals in brackets; human-only group omitted.

Controls include number of teammates; game-by-software fixed effects; maximum and minimum position skill fixed effects; attendance fixed effects.

\*  $p < 0.1$ , \*\*  $p < 0.05$ , \*\*\*  $p < 0.01$

**Table S13. Causal relationship between treatment groups and reproducibility outcomes for different softwares**

	(1)	(2)	(3)	(4)	(5)	(6)	(7)
	Reproduction	Minor errors	Major errors	One good robustness	Two good robustness	Ran one robustness	Ran two robustness
AI-Assisted	-0.019 ( 0.062) [-0.143; 0.105]	-0.667 ( 0.645) [-1.950; 0.617]	-2.175*** ( 0.703) [-3.573; -0.776]	0.000 ( 0.035) [-0.070; 0.071]	0.148 ( 0.107) [-0.066; 0.361]	0.024 ( 0.059) [-0.092; 0.141]	0.245* ( 0.130) [-0.013; 0.504]
AI-Led	-0.554*** ( 0.146) [-0.844; -0.265]	-0.820 ( 0.724) [-2.260; 0.619]	-2.676*** ( 0.637) [-3.945; -1.408]	-0.141 ( 0.099) [-0.339; 0.056]	-0.192 ( 0.160) [-0.511; 0.126]	-0.365*** ( 0.136) [-0.636; -0.094]	-0.185 ( 0.192) [-0.567; 0.197]
R	-0.146* ( 0.087) [-0.320; 0.027]	0.278 ( 0.694) [-1.103; 1.659]	-2.522*** ( 0.640) [-3.796; -1.248]	-0.012 ( 0.039) [-0.090; 0.067]	0.051 ( 0.131) [-0.210; 0.312]	-0.074 ( 0.086) [-0.246; 0.097]	0.062 ( 0.168) [-0.271; 0.396]
AI-Assisted × R	-0.022 ( 0.138) [-0.296; 0.252]	2.134** ( 1.011) [ 0.123; 4.145]	2.164*** ( 0.781) [ 0.610; 3.717]	-0.011 ( 0.051) [-0.113; 0.092]	-0.322 ( 0.204) [-0.728; 0.084]	-0.128 ( 0.125) [-0.376; 0.121]	-0.511** ( 0.226) [-0.961; -0.060]
AI-Led × R	-0.072 ( 0.194) [-0.458; 0.313]	0.810 ( 0.840) [-0.861; 2.482]	2.310*** ( 0.695) [ 0.927; 3.693]	-0.037 ( 0.126) [-0.288; 0.213]	-0.106 ( 0.210) [-0.525; 0.312]	0.052 ( 0.184) [-0.313; 0.418]	-0.192 ( 0.247) [-0.683; 0.300]
Controls	✓	✓	✓	✓	✓	✓	✓
Mean of dep. var	0.738	1.029	0.874	0.942	0.786	0.816	0.680
Obs.	103	103	103	103	103	103	103

Note: Standard errors in parentheses, confidence intervals in brackets; human-only group omitted; Stata papers omitted. Controls include number of teammates; game and software fixed effects; maximum and minimum position skill fixed effects; attendance fixed effects. \*  $p < 0.1$ , \*\*  $p < 0.05$ , \*\*\*  $p < 0.01$

**Table S14. AI-Assisted and AI-Led Metrics by Experience Level**

Variable	AI-Assisted high experience (n=10)	AI-Assisted low/medium experience (n=24)	AI-Led high experience (n=9)	AI-Led low/medium experience (n=26)	AI-Assisted High vs Low	AI-Led High vs Low
Reproduction	1.000 (0.000)	0.875 (0.338)	0.444 (0.527)	0.346 (0.485)	0.125 [0.083]	0.098 [0.631]
Minutes to reproduction	89.2 (85.5)	81.4 (60.6)	178.2 (69.2)	180.3 (72.3)	7.8 [0.800]	-2.1 [0.962]
Number of minor errors	1.500 (3.064)	1.417 (2.320)	0.889 (2.028)	0.615 (1.472)	0.083 [0.940]	0.274 [0.717]
Minutes to first minor error	71.3 (49.6)	157.0 (82.0)	94.0 (0.0)	175.7 (89.6)	-85.7 [0.067]	-81.7 [0.052]
Number of major errors	1.100 (1.595)	0.625 (0.875)	0.444 (0.726)	0.154 (0.464)	0.475 [0.393]	0.291 [0.287]
Minutes to first major error	85.8 (25.4)	152.6 (99.0)	102.0 (40.6)	203.7 (113.3)	-66.8 [0.070]	-101.7 [0.257]
At least one appropriate robustness check	1.000 (0.000)	1.000 (0.000)	0.778 (0.441)	0.846 (0.368)	0.000 [-]	-0.068 [0.684]
At least two appropriate robustness checks	0.900 (0.316)	0.833 (0.381)	0.667 (0.500)	0.615 (0.496)	0.067 [0.604]	0.051 [0.794]
Ran at least one appropriate robustness check	1.000 (0.000)	0.958 (0.204)	0.556 (0.527)	0.577 (0.504)	0.042 [0.328]	-0.021 [0.917]
Ran at least two appropriate robustness check	0.900 (0.316)	0.792 (0.415)	0.444 (0.527)	0.462 (0.508)	0.108 [0.417]	-0.017 [0.934]
Number of prompts	21.200 (17.139)	31.250 (39.344)	86.222 (69.009)	89.423 (63.338)	-10.050 [0.307]	-3.201 [0.904]

Note: This table reports exploratory analyses examining heterogeneity among AI-assisted and AI-led teams by level of AI experience. AI-assisted teams are divided into higher and lower/medium AI-experience groups based on self-reported familiarity with AI tools prior to the event.

**Table S15. Human-Only and AI-Assisted Metrics by Experience Level**

Variable	Human-Only (n=33)	AI-Assisted high experience (n=10)	AI-Assisted low/medium experience (n=24)	Human-Only vs AI-Assisted High	Human-Only vs AI-Assisted Low/Medium
Reproduction	0.939 (0.242)	1.000 (0.000)	0.875 (0.338)	-0.061 [0.160]	0.064 [0.430]
Minutes to reproduction	82.0 (39.8)	89.2 (85.5)	81.4 (60.6)	-7.2 [0.801]	0.5 [0.972]
Number of minor errors	1.000 (1.658)	1.500 (3.064)	1.417 (2.320)	-0.500 [0.631]	-0.417 [0.457]
Minutes to first minor error	141.6 (97.0)	71.3 (49.6)	157.0 (82.0)	70.2 [0.119]	-15.4 [0.664]
Number of major errors	1.697 (2.568)	1.100 (1.595)	0.625 (0.875)	0.597 [0.384]	1.072 [0.031]
Minutes to first major error	110.5 (69.5)	85.8 (25.4)	152.6 (99.0)	24.7 [0.249]	-42.1 [0.259]
At least one appropriate robustness check	1.000 (0.000)	1.000 (0.000)	1.000 (0.000)	0.000 [-]	0.000 [-]
At least two appropriate robustness checks	0.879 (0.331)	0.900 (0.316)	0.833 (0.381)	-0.021 [0.857]	0.045 [0.641]
Ran at least one appropriate robustness check	0.939 (0.242)	1.000 (0.000)	0.958 (0.204)	-0.061 [0.160]	-0.019 [0.751]
Ran at least two appropriate robustness check	0.788 (0.415)	0.900 (0.316)	0.792 (0.415)	-0.112 [0.375]	-0.004 [0.973]

*Note:* This table reports exploratory analyses examining heterogeneity among AI-assisted teams by level of AI experience in comparison to human-only teams. AI-assisted teams are divided into higher and lower/medium AI-experience groups based on self-reported familiarity with AI tools prior to the event.

**Table S16. Comparison of Key Metrics by Prompt Levels within AI-Assisted Group**

Variable	Above median (n=16)	Below/equal to median (n=19)	Difference
Reproduction	1.000 (0.000)	0.842 (0.375)	0.158 [0.083]
Minutes to reproduction	115.7 (115.1)	70.9 (27.7)	44.8 [0.149]
Number of minor errors	1.188 (2.562)	1.579 (2.479)	-0.391 [0.651]
Minutes to first minor error	154.2 (81.0)	130.3 (87.8)	23.8 [0.600]
Number of major errors	0.938 (1.436)	0.579 (0.769)	0.359 [0.380]
Minutes to first major error	136.9 (67.5)	124.6 (105.4)	12.2 [0.791]
At least one appropriate robustness check	1.000 (0.000)	1.000 (0.000)	0.000 [-]
At least two appropriate robustness checks	0.750 (0.447)	0.947 (0.229)	-0.197 [0.125]
Ran at least one appropriate robustness check	0.938 (0.250)	0.947 (0.229)	-0.010 [0.905]
Ran at least two appropriate robustness check	0.688 (0.479)	0.895 (0.315)	-0.207 [0.151]

*Note:* Columns 2–3 present means and standard errors in parentheses for individual groups (Human-only, AI-Assisted, and AI-Led); column 4 shows mean differences and *p*-values in brackets for the indicated group comparison. Groups are defined by the median number of prompts ( 19 ) in the AI-Assisted sample.

**Table S17. Sentiment in Prompts across AI-Led and AI-Assisted Teams**

Variable	AI-Assisted Mean	AI-Led Mean	Difference	Std. Err.	95% CI [LB; UB]
<i>Panel A: Conversation and Prompt Summary Statistics</i>					
Number of Conversations	3.188	6.324	3.136***	(1.091)	[0.956; 5.316]
Number of Prompts	31.563	91.206	59.643***	(12.581)	[34.511; 84.776]
<i>Panel B: Dictionary-Based Sentiment (NRC Lexicon)</i>					
Positive	0.048	0.609	0.127**	(0.059)	[0.011; 0.243]
Negative	0.259	0.240	-0.019	(0.029)	[-0.076; 0.037]
Anger	0.055	0.024	-0.031***	(0.007)	[-0.045; 0.017]
Fear	0.093	0.069	-0.024*	(0.024)	[-0.049; 0.001]
Disgust	0.022	0.019	-0.003	(0.005)	[-0.013; 0.007]
Joy	0.081	0.107	0.026*	(0.015)	[-0.004; 0.055]
Sadness	0.129	0.113	-0.016	(0.016)	[-0.049; 0.017]
Surprise	0.085	0.076	-0.008	(0.011)	[-0.031; 0.014]
Trust	0.245	0.331	0.086**	(0.042)	[0.005; 0.168]
Anticipation	0.115	0.142	0.027	(0.020)	[-0.012; 0.067]
<i>Panel C: Machine-Based Sentiment (Emotion)</i>					
Anger	0.001	0.008	0.007**	(0.003)	[0.002; 0.013]
Fear	0.006	0.002	-0.004**	(0.002)	[-0.008; 0.000]
Disgust	0.020	0.003	0.001	(0.002)	[-0.003; 0.005]
Joy	0.007	0.007	-0.000	(0.003)	[-0.006; 0.006]
Sadness	0.024	0.021	-0.003	(0.005)	[-0.013; 0.008]
Surprise	0.036	0.039	0.003	(0.007)	[-0.010; 0.017]
Neutral	0.925	0.920	-0.005	(0.010)	[-0.024; 0.014]
<i>Panel D: Machine-Based Sentiment Model (Multi-Class Sentiment)</i>					
Positive	0.023	0.037	0.009	(0.007)	[-0.004; 0.022]
Negative	0.091	0.107	0.017	(0.011)	[-0.005; 0.038]
Neutral	0.882	0.856	-0.026**	(0.012)	[-0.050; 0.001]

*Notes:* Prompts were unavailable for one AI-led and three AI-assisted teams. There were 32 different teams who were AI-assisted compared to 34 different teams which were AI-Led. Difference = Mean(AI-Led) – Mean(AI-Assisted). Two-sample t-tests with equal variances. \*\*\*  $p < 0.01$ , \*\*  $p < 0.05$ , \*  $p < 0.10$ .

1. Abel Brodeur, Derek Mikola, Nikolai Cook, Thomas Brailey, Ryan Briggs, Alexandra de Gendre, Yannick Dupraz, Lenka Fiala, Jacopo Gabani, Romain Gauriot, et al. Mass reproducibility and replicability: A new hope, 2024. Institute for Replication Discussion Paper 107.
2. Open AI. File uploads faq, 2024. <https://help.openai.com/en/articles/8555545-file-uploads-faq/> [Accessed: November 28, 2024].
3. Open AI. Learning to reason with llms, 2024. <https://openai.com/index/learning-to-reason-with-llms/> [Accessed: November 18, 2024].
4. Jack Fitzgerald. Imputations, inverse hyperbolic sines and impossible values. *Nature Human Behaviour*, pages 1–4, 2026.
5. Tim Loughran and Bill McDonald. Textual analysis in accounting and finance: A survey. *Journal of Accounting Research*, 54(4):1187–1230, 2016.
6. Bo Pang, Lillian Lee, et al. Opinion mining and sentiment analysis. *Foundations and Trends in Information Retrieval*, 2(1–2):1–135, 2008.
7. Matthew Gentzkow, Bryan Kelly, and Matt Taddy. Text as data. *Journal of Economic Literature*, 57(3):535–574, 2019.
8. Julia Silge, David Robinson, and David Robinson. *Text mining with R: A tidy approach*. O'reilly Boston (MA), 2017.
9. Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. *Advances in Neural Information Processing Systems*, 30, 2017.
10. Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. Bert: Pre-training of deep bidirectional transformers for language understanding. In *Proceedings of the 2019 conference of the North American chapter of the association for computational linguistics: human language technologies, volume 1 (long and short papers)*, pages 4171–4186, 2019.
11. Virginia Braun and Victoria Clarke. Using thematic analysis in psychology. *Qualitative research in psychology*, 3(2):77–101, 2006.

DRAFT