



# City Research Online

## City St George's, University of London

**Citation:** Brodeur, A., Valenta, D., Marcoci, A., Aparicio, J. P., Mikola, D., Barbarioli, B., Alexander, R., Deer, L., Stafford, T., Vilhuber, L., et al (2026). AI-assisted teams outperform AI-led teams but not human-only teams in assessing research reproducibility in quantitative social science. *Proceedings of the National Academy of Sciences*, 123(22), e2524747123. doi: 10.1073/pnas.2524747123

This is the accepted version of the paper.

This version of the publication may differ from the final published version. To cite this item please consult the publisher's version.

**Permanent repository link:** <https://openaccess.city.ac.uk/id/eprint/37614/>

**Link to published version:** <https://doi.org/10.1073/pnas.2524747123>

**Copyright and Reuse:** Copyright and Moral Rights remain with the author(s) and/or copyright holders. Copies of full items can be used for personal research or study, educational, or not-for-profit purposes without prior permission or charge, unless otherwise indicated, provided that the authors, title and full bibliographic details are credited, a hyperlink and/or URL is given for the original metadata page and the content is not changed in any way. For full details of reuse please refer to [City Research Online policy](#).

# 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</sup>, David Valenta<sup>a,b</sup>, Alexandru Marcoci<sup>c</sup>, Juan P. Aparicio<sup>a,b</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 Motoki<sup>j,k</sup>, Mohamed Abdelhady<sup>l</sup>, Youssa Abdelmoula<sup>m</sup>, Ghina Abdul Baki<sup>n,b</sup>, Tomás Aguirre<sup>o</sup>, Sriraj Aiyer<sup>p</sup>, Shumi Akhtar<sup>q</sup>, Farida Akhtar<sup>r</sup>, Melle R. Albada<sup>s</sup>, Tobias Altman<sup>t</sup>, David Anagnostou<sup>u,v</sup>, Zahra Arjmandi Lari<sup>w</sup>, Jorge Armando De León Tejada<sup>x</sup>, David Rodríguez Arana<sup>y</sup>, Igor Asanov<sup>z</sup>, Anastasiya-Mariya Novikova<sup>aa</sup>, Rebecca Ashong<sup>ab</sup>, Tobias Auer<sup>ac</sup>, Francisco J. Bahomonde-Birk<sup>ad</sup>, Bradley J. Baker<sup>ae</sup>, Söhnke M. Bartram<sup>af,ag</sup>, Dongqi Bao<sup>ah</sup>, Lucija Batinovic<sup>ai</sup>, Tommaso Batistoni<sup>aj</sup>, Monica Beeder<sup>ak</sup>, Louis-Philippe Beland<sup>al</sup>, Carsten Gero Bienz<sup>am</sup>, Christ Billy Aryanto<sup>an</sup>, Cylcia Boilbaugh<sup>ao</sup>, Carl Bonander<sup>ap,aq</sup>, Ramiro Bravo<sup>ar</sup>, Egor Bronnikov<sup>as,at</sup>, Stephan Bruns<sup>au,av,aw</sup>, Nino Buliskeria<sup>ax</sup>, Sara Caicedo-Silva<sup>ay</sup>, Andrea Calef<sup>az</sup>, Juan Sebastian Cano Arias<sup>ba</sup>, Gustavo A. Castillo Alvarez<sup>bb</sup>, Solomon Caulker<sup>bb,bc</sup>, Simonas Cepenas<sup>bc</sup>, Arthur Chatton<sup>bd,be</sup>, Zirou Chen<sup>bf</sup>, Ngozi Chioma Ewurum<sup>bg</sup>, Anda-Bianca Ciocirlan<sup>bh</sup>, Felix J. Clouth<sup>bi</sup>, Jason Collins<sup>bj</sup>, Nikolai Cook<sup>bj</sup>, Cesar Cornejo<sup>bk,bl</sup>, João Craveiro<sup>f</sup>, Jonathan Créche<sup>bm</sup>, Jing Cui<sup>bn,bo</sup>, Niveditha Chail Vayalabron<sup>bp</sup>, Christian Czyszara<sup>bq</sup>, Carlos Daniel Bermúdez Jaramillo<sup>br</sup>, Hannes Datta<sup>bs</sup>, Lien Denoo<sup>bt</sup>, Arshia Dhalwal<sup>l</sup>, Nancy Dhameja<sup>bu</sup>, Elodie Djemai<sup>bv</sup>, Erwan Dujancourt<sup>bw,bx</sup>, Ugurcan Dündar<sup>by</sup>, Thibaut Duprey<sup>bz</sup>, Yasmine Eissa<sup>ca</sup>, Youssef El Fassi<sup>cb</sup>, Ismail El Fassi<sup>cc</sup>, Keaton Ellis<sup>cd</sup>, Ali Elminejad<sup>ax</sup>, Mahmood Elsherif<sup>ce,cf</sup>, Aysil Emirhammutoglu<sup>cg</sup>, Gülilan Etingin-Frati<sup>ch</sup>, Emeke Eze<sup>ci</sup>, Jan Fabian Dollbaum<sup>cj</sup>, Jan Feld<sup>ck</sup>, Andres Felipe Rengifo Jaramillo<sup>cl,cm</sup>, Guidon Fenig<sup>a</sup>, Victoria Fernandes<sup>bz,cn</sup>, Lenka Fiala<sup>co</sup>, Lukas Fink<sup>cp</sup>, Mojtaba Firouzjaeiangalough<sup>cq</sup>, Sara Fish<sup>ct</sup>, Jack Fitzgerald<sup>cs,ct</sup>, Rachel Forshaw<sup>cu</sup>, Alexandre Fortier-Chouinard<sup>cv</sup>, Louis Fréget<sup>tw</sup>, Joris Frese<sup>cx</sup>, Jacopo Gabani<sup>cy,cz</sup>, Sebastian Gallegos<sup>da</sup>, Max C. Gamill<sup>db</sup>, Attila Gáspár<sup>dc,dd</sup>, Romain Gauriot<sup>de</sup>, Evelina Gavrilova<sup>cg</sup>, Diogo Geraldes<sup>df,dg,dh</sup>, Giulio Giacomo Cantone<sup>di</sup>, Grant Gibson<sup>dj,dk</sup>, Dirk Goldschmidt<sup>dl</sup>, Amélie Gourdon-Kanhukamwe<sup>dm</sup>, Andrea Gregor de Varda<sup>dn</sup>, Idaliya Grigoryeva<sup>do</sup>, Alexi Gushvili<sup>dp</sup>, Aaron H.A. Fletcher<sup>dq</sup>, Florian Habermann<sup>dr,ds</sup>, Márton Hablicsek<sup>dt</sup>, Joanne Haddad<sup>du</sup>, Jonathan D. Hall<sup>dv</sup>, Olle Hammar<sup>bd,dw</sup>, Malek Hassouneh<sup>dx</sup>, Carina I. Hausladen<sup>dy</sup>, Sophie C. F. Hendrikse<sup>dz</sup>, Matthew Hepplewhite<sup>ea</sup>, Anson T. Y. Ho<sup>eb</sup>, Senan Hogan-Hennessy<sup>ec</sup>, Elliot Howley<sup>ed</sup>, Gaoyang Huang<sup>ee,ef</sup>, Hélioise Hulstaer<sup>au,eg</sup>, Zlatomira G. Ilchovska<sup>eh,ei,ed</sup>, Paola Jaimes Santamaría<sup>ej,ek</sup>, Niklas Jakobsson<sup>el</sup>, Joakim Jansson<sup>bx,em</sup>, Ewa Jarosz<sup>en</sup>, Hossein Jebeli<sup>fo</sup>, Yanchen Jiang<sup>fp</sup>, Hiba Junaid<sup>fq</sup>, Rohan Kallurya<sup>es</sup>, Sunny Karim<sup>ft</sup>, Edmund Kelly<sup>fv</sup>, Eva Kimmel<sup>fw</sup>, Soravich Kingsuwankul<sup>ew,et</sup>, Valentin Klotzbücher<sup>ex,ey,ez</sup>, Daniel Krährmer<sup>fa</sup>, Pijus Krūminas<sup>fb</sup>, Nicholas Kruus<sup>fc</sup>, Essi Kujansuu<sup>fd,fe</sup>, Christoph F. Kurzf<sup>ff</sup>, Stephan Küster<sup>fg</sup>, Blake Lee-Whiting<sup>fh</sup>, Felix Lewandowski<sup>fi</sup>, Tongzhe Li<sup>fi</sup>, Ruoxi Li<sup>fk</sup>, Dan Liu<sup>fl</sup>, Jiacheng Liu<sup>fm</sup>, Helix Lo<sup>fn</sup>, Katharina Lote<sup>fo</sup>, Felipe Macedo Dias<sup>fp</sup>, Christopher R. Madan<sup>fq</sup>, Nicolas Mäder<sup>fr</sup>, Marco Mandas<sup>fs</sup>, Cesar Mantilla<sup>ft</sup>, Jan Marcus<sup>fu</sup>, Diego Marino Fages<sup>fv</sup>, Xavier Martin<sup>fw</sup>, Ryan McWay<sup>fx</sup>, Daniel Medina-Gaspar<sup>fy</sup>, Sisi Meng<sup>gz</sup>, Lingyu Meng<sup>ga</sup>, Simon Merz<sup>gb</sup>, Alex P. Miller<sup>gc</sup>, Thibault Mirabel<sup>gd</sup>, Dibya Depta Mishra<sup>ge</sup>, Sumit Mishra<sup>gf</sup>, Belay W. Moges<sup>gg</sup>, Morteza Mohandes Mojarrad<sup>gh</sup>, Myra Mohnen<sup>hi</sup>, Louis-Philippe Morin<sup>hi</sup>, Lucija Muehlenbachs<sup>gi,gj</sup>, Gastón Mullin<sup>gk</sup>, Andrea Musulan<sup>gl,gm</sup>, Sara Muzzi<sup>gn,go</sup>, James A. C. Myers<sup>gp</sup>, Florian Neubauer<sup>gq</sup>, Tuan Nguyen<sup>au</sup>, Ali Niaz<sup>gi</sup>, Ardyn Nordstrom<sup>gr</sup>, Bartłomiej Nowak<sup>gs</sup>, Daneal O'Habb<sup>gt</sup>, Tim Ölkens<sup>gu</sup>, Justin Ong<sup>gv</sup>, Valeria Orozco Castiblanco<sup>gw,gx</sup>, Ömer Özak<sup>gy</sup>, Ali I. Ozkes<sup>gz,ha</sup>, Mikael Paaso<sup>hb</sup>, Shubham Pandey<sup>hc</sup>, Varvara Papazoglou<sup>hd</sup>, Romeo Penheiro<sup>he</sup>, Linh Pham<sup>hf</sup>, Ulrike Phier<sup>hg</sup>, Peter Pütz<sup>hh</sup>, Quan Qi<sup>hi</sup>, Jingyi Qiu<sup>hj</sup>, David A. Reinstein<sup>hk,hk</sup>, Juuso Repo<sup>hl</sup>, Nicolas Rudolf<sup>is</sup>, Shree Saha<sup>hm</sup>, Orkun Saka<sup>hn</sup>, Chiara Saponaro<sup>ho</sup>, Georg Sator<sup>hp</sup>, Martijn Schoenmakers<sup>hq</sup>, Raffaello Ser<sup>hq</sup>, Meet Shah<sup>eb</sup>, Paul Sibille<sup>hr</sup>, Christoph Siemroth<sup>hs</sup>, Vladimir Skavys<sup>ht,hu</sup>, Ben Slater<sup>hv</sup>, Wenting Song<sup>hw</sup>, Stefan Staubli<sup>si</sup>, Tobias Steindl<sup>hx</sup>, Nomwendé Steven Waonga<sup>a</sup>, Paul Stott<sup>hy,hz</sup>, Stephenson Strobel<sup>ia</sup>, Roshini Sudhakaran<sup>ib</sup>, Pu Sun<sup>ic</sup>, Scott D. Swain<sup>id</sup>, Aleksandr Talavera<sup>ie</sup>, Hanz M. Tantiangco<sup>if</sup>, George Tarasenko<sup>ig</sup>, Boyd Tarlinton<sup>ih</sup>, Mariam Tarraf<sup>ii</sup>, Ken Teoh<sup>ij</sup>, Rémi Thériault<sup>ij,ib</sup>, Bethan Thompson<sup>ik</sup>, Tonghui Tian<sup>il</sup>, Wenjie Tian<sup>ia</sup>, Manuel Tobias Rein<sup>im</sup>, Emmanuel Tolani<sup>in,io</sup>, Nicolai Borgen<sup>ip</sup>, Solveig Topstad Borgen<sup>iq</sup>, Javier Torralba<sup>bs</sup>, Carolina Velez-Ospina<sup>ir</sup>, Man Wai Mak<sup>is</sup>, Lukas Wallrich<sup>it,iu</sup>, Zeyang Wang<sup>iv</sup>, Leah Ward<sup>iw</sup>, Matthew D. Webb<sup>ix</sup>, Duncan Webb<sup>iy</sup>, Bryan S. Weber<sup>iz,ja</sup>, Christoph Weber<sup>ib</sup>, Wei-Chien Weng<sup>c</sup>, Christian Westheide<sup>d,j,e</sup>, Tom Wilkinson<sup>f</sup>, Kwong-Yu Wong<sup>is</sup>, Marcin Wroński<sup>jh</sup>, Zhuangchen Wu<sup>jj</sup>, Qixia Wu<sup>jk</sup>, Victor Y. Wu<sup>jl</sup>, Bohan Xiao<sup>o</sup>, Feihong Xu<sup>jk</sup>, Cong Xu<sup>jm</sup>, Pranav Yadav<sup>in</sup>, Yu Yang Chou<sup>dm</sup>, Luther Yap<sup>io</sup>, Myra Zubeck<sup>aj,p</sup>, Bo Yao<sup>q</sup>, Zuzanna Zagrodzka<sup>jr</sup>, Tahreen Zahra<sup>is</sup>, Mirela Zaneva<sup>jt</sup>, Xiaomeng Zhang<sup>ju</sup>, Ziwei Zhao<sup>lv,jw</sup>, Han Zhong<sup>ix</sup>, Aras Zirgulis<sup>lv</sup>, Jiacheng Zou<sup>iz</sup>, Floris Zoutman<sup>ka</sup>, Christelle Zozoungbo<sup>kb</sup>

<sup>a</sup>Department of Economics, Faculty of Social Sciences, University of Ottawa, Ottawa, Ontario K1N 6N5, Canada; <sup>b</sup>Institute for Replication, Institute for Replication, Ottawa, Ontario K1N 6N5, Canada; <sup>c</sup>Institute for Technology and Humanity, University of Cambridge, Cambridge, Cambridgeshire, CB2 1SB, UK; <sup>d</sup>Faculty of Information and Department of Statistical Sciences, University of Toronto, Toronto, Ontario M5S 3G6, Canada; <sup>e</sup>Department of Management and Marketing, University of Melbourne, Melbourne, Carlton, Victoria, Australia 3010; <sup>f</sup>School of Psychology, University of Sheffield, Sheffield, South Yorkshire, S1 4DP, UK; <sup>g</sup>Research on Research Institute, Research on Research Institute, London, WC1E 6JA, UK; <sup>h</sup>Department of Economics, Cornell University, Ithaca, NY 14853, USA; <sup>i</sup>Climate and Development Policy Division, RWI - 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 78839, USA; <sup>k</sup>Norwich Business School, Accounting and Quantitative Methods, University of East Anglia, Norwich, NR4 7TA, United Kingdom; <sup>l</sup>Department of Economics, Carleton University, Ottawa, Ontario K1S 5B6, Canada; <sup>m</sup>Statistics Canada, Statistics Canada, Ottawa, Ontario K1A 0T6, Canada; <sup>n</sup>Department of Economics, University of Ottawa, Ottawa, Ontario K1N 6N5, Canada; <sup>o</sup>Centre for the Governance of AI, Centre for the Governance of AI, London, 210 Pentonville Road, N1 9JY, UK; <sup>p</sup>Department of Experimental Psychology, University of Oxford, Oxford, Oxfordshire, OX1 3EL; <sup>q</sup>Finance Discipline, The University of Sydney Business, The University of Sydney, Sydney, NSW 2006, Australia; <sup>r</sup>Department of Actuarial Studies and Business Analytics, Macquarie Business School, Macquarie University, Sydney, NSW 2109, Australia; <sup>s</sup>Department of Socioeconomics, Vienna University of Economics and Business, Vienna, Austria, 1020; <sup>t</sup>Center for Research on Equitable and Open Scholarship, Massachusetts Institute of Technology, Cambridge, MA, 02139; <sup>u</sup>UTM School of Management, Technical University of Munich, Munich, Germany 80333; <sup>v</sup>Centre for Business Research, University of Cambridge, Cambridge, United Kingdom CB2 1QA; <sup>w</sup>Independent researcher, Independent researcher, Tehran, Iran; <sup>x</sup>Facultad de Economía, Universidad del Rosario, Bogotá, Bogotá 111711, Colombia; <sup>y</sup>School of Economics, Universidad del Rosario, Bogotá, Bogotá 111711; <sup>z</sup>International Center for Higher Education Research and Faculty of Economics, University of Kassel, Kassel, Hesse, Germany, 34125; <sup>aa</sup>INCHER, University of Kassel, Kassel, Hesse, Germany, 34125; <sup>ab</sup>Department of Economics, University of Ghana, Accra, P.O.Box LG 57, Legon Accra Ghana; <sup>ac</sup>Department of Economics, University of Basel, Basel, 4052 Basel, Switzerland; <sup>ad</sup>Tilburg School of Social and Behavioral Sciences, Tilburg University, Tilburg, 5037 AB, the Netherlands; <sup>ae</sup>Department of Sport, Tourism and Hospitality Management, Temple University, Philadelphia, PA 19132; <sup>af</sup>Warwick Business School, University of Warwick, Coventry, CV4 7AL, United Kingdom; <sup>ag</sup>Center for Economic Policy Research, CEPR, London, 2 Colindale Avenue, London EC1R 3HL, United Kingdom; <sup>ah</sup>Department of Economics, University of Zurich, Zurich, 8001 Zurich; <sup>ai</sup>Department of Behavioural Sciences and Learning, Linköping University, Linköping, Sweden 58183; <sup>aj</sup>Centre for Experimental Social Sciences, Nuffield College, University of Oxford, Oxford, OX1 1NF, United Kingdom; <sup>ak</sup>Department of Economics, University of Southampton, Southampton, SO17 1BJ United Kingdom; <sup>al</sup>Department of Economics, Carleton University, Ottawa, Ontario, K1S 5B6, Canada; <sup>am</sup>Department of Finance, Norwegian School of Economics, Bergen, Norway, 5045; <sup>an</sup>Faculty of Psychology, Atma Jaya Catholic University of Indonesia, Jakarta, Indonesia, 12500; <sup>ao</sup>Department of Education, University of York, York, United Kingdom YO10 5DD; <sup>ap</sup>Department of Economics, Maastricht Business School, Maastricht University, Maastricht, The Netherlands, 6525 XD; <sup>aq</sup>Department of Economics, Maastricht University, Maastricht, The Netherlands, 6525 XD; <sup>ar</sup>Department of Political Science, University of Chicago, Chicago, Chicago, 60637; <sup>as</sup>Department of Environmental Sciences, Hasselt University, Hasselt, Belgium, 3500; <sup>at</sup>International Center for Higher Education Research (INCHER), University of Kassel, Kassel, Hesse, Germany, 34125; <sup>au</sup>Meta-Research Innovation Center at Stanford (METRICS), Stanford University, Stanford, CA, 94305; <sup>av</sup>Department of Economics, Nazarbayev University, Astana, Kazakhstan, 010000; <sup>aw</sup>Facultad de Economía, Universidad de los Andes, Bogotá, Bogotá 111711, Colombia; <sup>ax</sup>UCL School of Management, University College London, United Kingdom, E14 5AA; <sup>ay</sup>Department of Economics, Universidad de los Andes, Bogotá, Bogotá D.C.; <sup>az</sup>United Methodist University Sierra Leone, United Methodist University Sierra Leone, Freetown, Sierra Leone; <sup>ba</sup>Department of Economics, ISM University of Management and Economics, Vilnius, Lithuania, 01103; <sup>bb</sup>Department of Medicine Social et Préventive, Université Laval, Québec, QC G1V0A6, Canada; <sup>bc</sup>Département de Médecine Sociale et Préventive, Université de Montréal, Montréal, QC H3N1J9, Canada; <sup>bd</sup>Department of Management, Marketing and Information Systems, Hong Kong Baptist University, Hong Kong, Kowloon Tong, Hong Kong, China (No postal code); <sup>be</sup>Department of Economics, College of Management Sciences (COLMANS), Michael O. Ogunniyi University of Agriculture, Umudike, Abia State, Nigeria, 440109; <sup>bf</sup>Department of Methodology and Statistics, Tilburg University, Tilburg, The Netherlands, 5037 AB; <sup>bg</sup>UTS Business School, University of Technology, Sydney, Ultimo, New South Wales, Australia, 2007; <sup>bh</sup>Department of Economics, Wilfrid Laurier University, Waterloo, Ontario, Canada, N2L 3G5; <sup>bi</sup>Department of Medical Statistics, The London School of Hygiene & Tropical Medicine, London, WC1E 7HT; <sup>bj</sup>Department of Computer Science, University College London, London, United Kingdom, WC1E 6BT; <sup>bk</sup>University of Ottawa, Department of Economics, Faculty of Social Sciences, University of Ottawa, Ottawa, Canada; <sup>bl</sup>Business School, Beijing Normal University, Beijing, Beijing, China, 100875; <sup>bm</sup>Belt and Road School, Beijing Normal University, Zhuhai, Guangdong, China, 519085; <sup>bn</sup>School of Earth and Planetary Science, National Institute of Science Education and Research, Bhubaneswar, India; <sup>bo</sup>Migration & Migrants, Netherlands Interdisciplinary Demographic Institute, The Hague, NL-2511 CV; <sup>bp</sup>Department of Economics, Universidad del Rosario, Bogotá, Colombia, 111711; <sup>bq</sup>Department of Marketing, Tilburg University, Tilburg, 5000 LE, The Netherlands; <sup>br</sup>Department of Strategy & Entrepreneurship, Tilburg University, Tilburg, 5000 LE, The Netherlands; <sup>bs</sup>Department of Economics, Binghamton University, Vestal, New York, 13902-6000; <sup>bt</sup>Economics, Université Paris-Dauphine, Paris, 75775 Paris Cedex 16, France; <sup>bu</sup>Swedish Institute for Social Research, Stockholm University, Stockholm, Sweden 103 91; <sup>bv</sup>Department of Economics and Statistics, Linnaeus University, Växjö, Sweden, 35195; <sup>bw</sup>Department of Marketing, WU Vienna, Vienna, Austria, 1020; <sup>bx</sup>Financial Stability Department, Bank of Canada, Ottawa, Ontario K1A 0G9, Canada; <sup>by</sup>Onsi Sawiris School of Business, The American University in Cairo, New Cairo, 11835, Egypt; <sup>bz</sup>Department of Accounting and Control, HEC Lausanne, Lausanne, Switzerland, 1015; <sup>ca</sup>Institute of Accounting, Control and Auditing, University of St. Gallen, St.Gallen, 9000 St.Gallen, Switzerland; <sup>cb</sup>Monash University, Monash University, Melbourne, Australia 3168; <sup>cc</sup>School of Psychology and Vision Science, University of Leicester, Leicester, LE1 7RH; <sup>cd</sup>Political Science, University of Birmingham, Birmingham, UK, B15 2TT; <sup>ce</sup>Department of Business and Management Science, NHH Norwegian School of Economics, Bergen, Norway, 5113; <sup>cf</sup>KOF Swiss Economic Institute, ETH Zurich, Zurich, Switzerland, 8092; <sup>cg</sup>College of Management Sciences, Michael Okpara University of Agriculture, Umudike, Abia State, Nigeria; <sup>ch</sup>School of Politics and International Relations, University College Dublin, Dublin, Ireland, Dublin 4; <sup>ci</sup>School of Economics and Finance, Victoria University of Wellington, Wellington, New Zealand, 6011; <sup>cj</sup>Business School, Universidad de los Andes, Bogotá, 111711, Colombia; <sup>ck</sup>Universidad de los Andes, Universidad de los Andes, Bogotá, 111711, Colombia; <sup>cl</sup>Vancouver School of Economics, University of British Columbia, Vancouver, British Columbia V6T 1Z4, Canada; <sup>cm</sup>Department of Economics, Tilburg University, Tilburg, North Brabant, 5037 AB, The Netherlands; <sup>cn</sup>Department of Economics, School of Business and Economics, Freie Universität Berlin, Berlin, Germany 14195; <sup>co</sup>Department of Business Management, Masaryk University, Brno, Czech Republic 602 00; <sup>cp</sup>School of Engineering and Applied Sciences, Harvard University, Cambridge, MA, 02139, USA; <sup>cq</sup>Department of Ethics, Governance, and Society, School of Business and Economics, Vrije Universiteit Amsterdam, Amsterdam, Noord-Holland 1081HV, Netherlands; <sup>cr</sup>Tinbergen Institute, Tinbergen Institute, Amsterdam, Noord-Holland, 1018WB, Netherlands; <sup>cs</sup>Department of Accountancy, Economics and Finance, Heriot-Watt University, Edinburgh, EH14 4AS, UK; <sup>ct</sup>Department of Political Science, Université Laval, Québec, QC G1V0A6, Canada; <sup>cu</sup>LEDA, CEPREMAP, Paris, 75014 Paris, France; <sup>cv</sup>Department of Political and Social Sciences, European University Institute, Fiesole, 50014, Italy; <sup>cw</sup>Health, Nutrition, and Population Global Practice, World Bank, Washington, DC, 20433, USA; <sup>cx</sup>Centre for Health Economics, University of York, Heslington, York, YO10 5DD, United Kingdom; <sup>cy</sup>Business School, Universidad Adolfo Ibáñez, Santiago, Chile 7910000; <sup>cz</sup>School of Chemical, Materials, and Biological Engineering, University of Sheffield, Sheffield, South Yorkshire, S10 2TN; <sup>da</sup>Institute of Economics, ELTE Centre for Economic and Regional Studies, Budapest, Budapest, Hungary, 1097; <sup>db</sup>Department of Economics, Central European University, Vienna, Austria, 1100; <sup>dc</sup>Department of Economics, Deakin University, Burwood, Victoria, 3125; <sup>dd</sup>School of Economics, University College Dublin, Dublin, D04 F8X4, Ireland; <sup>de</sup>CaBER, University of Coimbra, Coimbra, 3000-454 Coimbra, Portugal; <sup>df</sup>Behavioral Science, Geary Institute for Public Policy, Dublin, D04 P9C4, Ireland; <sup>dg</sup>Department of Law, Economics, and Society, "Magna Graecia" University of Catanzaro, Catanzaro, Italy 88100; <sup>dh</sup>Department of Economics, McMaster University, Hamilton, Ontario L8S4M4; <sup>di</sup>CRDCN, CRDCN, Hamilton, Ontario L8S4M4; <sup>dj</sup>School of Psychology, University of Sheffield, Sheffield, S10 2TN, UK; <sup>dk</sup>University of Sheffield, Sheffield, South Yorkshire, S1 4DP, UK; <sup>dl</sup>IoPPN, King's College London, London, United Kingdom SE5 9JL; <sup>dm</sup>McGovern Institute for Brain Research, Massachusetts Institute of Technology, Cambridge, Massachusetts 02139; <sup>dn</sup>Economics, UC San Diego, La Jolla, California 92092; <sup>do</sup>Department of Sociology and Human Geography, University of Oslo, Oslo, Norway, N-0851; <sup>dp</sup>School of Computer Science, University of Sheffield, Sheffield, South Yorkshire, S1 4DP, UK; <sup>dq</sup>UCD Lochnam Quinn School of Business, Economics University College Dublin, Dublin, Ireland, Dublin 4; <sup>dr</sup>Department of Accounting and Control, University of Lausanne, Lausanne, Switzerland, 1015; <sup>ds</sup>Mathematics Institute, Leiden University, Leiden, South Holland, 2333CA; <sup>dt</sup>Department of Economics and Economic History, Unit of Economic Analysis, Universidad Autónoma de Barcelona (UAB), Bellaterra, Cercanya de Valles, 08193 Bellaterra, Spain; <sup>du</sup>Department of Economics, Finance, and Legal Studies, University of Alabama, Tuscaloosa, Alabama, 35487; <sup>dv</sup>Institute for Futures Studies, Institute for Futures Studies, Stockholm, Sweden, 10131; <sup>dw</sup>Department of Economics, University of Toronto, Toronto, Ontario M5S 3G7, Canada; <sup>dx</sup>Computational Social Science, ETH Zurich, Zurich, Stampfenbachstrasse 48, 8057 Zurich; <sup>dy</sup>Department of Methodology and Statistics, Tilburg School of Social and Behavioral Sciences, Tilburg University, Tilburg, 5037 DB, North Brabant, the Netherlands; <sup>dz</sup>Department of Politics and International Relations, University of Oxford, Oxford, Oxfordshire OX1 3JU; <sup>ea</sup>Department of Real Estate Management, Ted Rogers School of Management, Toronto Metropolitan University, Toronto, Ontario, M5B 2K3, Canada; <sup>eb</sup>Department of Economics, Cornell University, Ithaca, NY 14853 USA; <sup>ec</sup>School of Psychology, University of Nottingham, Nottingham, United Kingdom, NG7 2RD; <sup>ed</sup>Department of Health Sciences and Technology, ETH Zurich, Switzerland, Zurich, 8001; <sup>ee</sup>Department of Economics, University of Zurich, Zurich, Switzerland, Zurich, 8001; <sup>ef</sup>HEC Liège, Department of Finance, University of Liège, Liège, Belgium, 4000; <sup>eg</sup>Department of Psychology, University of York, York, United Kingdom, YO10 5DD; <sup>eh</sup>School of Psychology, University of Birmingham, Birmingham, United Kingdom, B15 2SA; <sup>ei</sup>Political Science Research Institute, University of Massachusetts Amherst, Amherst, MA 01003, USA; <sup>ej</sup>Center for Economic and Policy Research, Center for Economic and Policy Research, Washington, DC 20009, USA; <sup>ek</sup>Karlstad Business School, Karlstad University, Karlstad, Sweden 651 88; <sup>el</sup>Research Institute of Industrial Economics, Research Institute of Industrial Economics, Stockholm, Sweden, 10215; <sup>em</sup>Faculty of Economic Sciences, University of Warsaw, Warsaw, 00-927 Warsaw, Poland; <sup>en</sup>Climate Analysis Team, Financial Stability Department, Bank of Canada, Ottawa, Ontario K1A 0G9, Canada; <sup>eo</sup>School of Engineering and Applied Sciences, Harvard University, Cambridge, Massachusetts, 02138, USA; <sup>ep</sup>Bart's Life Sciences, Bart's Health NHS Trust, London, United Kingdom, E14 5HU; <sup>eq</sup>Queen Mary University of London, Queen Mary University of London, London, E1 4NS, United Kingdom; <sup>er</sup>School of Electrical and Computer Engineering, Cornell University, Ithaca, NY 14853 USA; <sup>es</sup>Carleton University, Carleton University, Ottawa, ON K1S 5B6; <sup>et</sup>University of Oxford, University of Oxford, Oxford, United Kingdom OX1 2JD; <sup>eu</sup>Department of Psychology, University of York, York, United Kingdom YO10 5DD; <sup>ev</sup>Department of Management and Organization, School of Business and Economics, Vrije Universiteit Amsterdam, Amsterdam, Noord-Holland 1081 HV Amsterdam, Netherlands; <sup>ew</sup>Department of Economics, University of Freiburg, Freiburg im Breisgau, Germany, 79085; <sup>ex</sup>Department of Clinical Research, University of Basel, Basel, Switzerland, 4051; <sup>ey</sup>University Hospital Basel, University Hospital Basel, Basel, Switzerland, 4031; <sup>ez</sup>Department of Sociology, LMU Munich, Munich, Germany, 80539; <sup>fa</sup>Department of Economics, ISM University of Management and Economics, Vilnius, LT-01103 Vilnius, Lithuania; <sup>fb</sup>Department of Politics and International Relations, University of Oxford, Oxford, OX1 3JU, Oxfordshire, United Kingdom; <sup>fc</sup>Department of Economics, Faculty of Economics and Statistics, University of Innsbruck, Innsbruck, 6020 Innsbruck, Austria; <sup>fd</sup>Department of Economics, Turkey School of Economics, University of Turku, Turku, FI-20014 Turun yliopisto, Finland; <sup>fe</sup>Department of Health Economics, Ludwig-Maximilians-Universität München, Munich, DE-80539 Munich, Germany; <sup>ff</sup>School of Business and Economics, Freie Universität Berlin, Berlin, 14195 Berlin, Germany; <sup>fg</sup>Department of Political Science, University of Toronto, Toronto, Ontario M5S 3G5, Canada; <sup>fh</sup>School of Psychology, University of Nottingham, Nottingham, University Park NG7 2RD; <sup>fi</sup>Department of Food, Agricultural and Resource Economics, University of Guelph, Guelph, Ontario N1G 2W1, Canada; <sup>fj</sup>Department of Economics, Yale University, New Haven, CT 06520-8268 USA; <sup>fk</sup>Research School of Economics, Australian National University, Canberra, ACT 2600, Australia; <sup>fl</sup>School of Finance, Renmin University of China, Beijing, Beijing, China 100872; <sup>fm</sup>Department of Advanced Social and International Studies, University of Tokyo, Tokyo, 9-8-1 Komaba, Meguro City, Tokyo 153-8902, Japan; <sup>fn</sup>Department of Methodology, Tilburg School of Social and Behavioral Sciences, Tilburg University, Tilburg, 5037 DB, Netherlands; <sup>fo</sup>Dyson School of Applied Economics and Management, Cornell University, Ithaca, NY, 14850; <sup>fp</sup>Department of Psychology, University of Nottingham, Nottingham, United Kingdom NG7 2RD; <sup>fq</sup>Department of Economics, Knauass School of Business, University of San Diego, San Diego, 92110 San Diego, CA, USA; <sup>fr</sup>Department of Economics and Business Administration, University of Cagliari, Cagliari, 09124, Italy; <sup>fs</sup>Department of Economics, School of Administrative and Economic Sciences, Pontificia Universidad Javeriana, Bogotá, Colombia, 110231; <sup>ft</sup>Department of Economics, School of Business and Economics, FU Berlin, Berlin, Germany, 14624; <sup>fu</sup>Department of Economics, Durham University, Durham, United Kingdom DH1 3LB; <sup>fv</sup>School of Economics and Management, Tilburg University, Tilburg, Netherlands 5000 LE; <sup>fw</sup>Applied Economics, University of Minnesota, Saint Paul, MN, 55107; <sup>fx</sup>Department of Economics, School of Finance, Economics and Government, Universidad EAFIT, Medellín, Antioquia 05002, Colombia; <sup>fy</sup>SC Johnson School of Business, Cornell University, Ithaca, NY, 14853 USA; <sup>gz</sup>School of Economics, University of Sheffield, Sheffield, S10 2TU, United Kingdom; <sup>ga</sup>School of Economics and Business, University of Halle, Halle (Saale), Germany, 06108; <sup>gb</sup>Department of Marketing, Marshall School of Business, University of Southern California, Los Angeles, CA 90089-1424, USA; <sup>gc</sup>Equalis Capital, Equalis Capital, Paris, France, 75116; <sup>gd</sup>Department of Economics, Rice University, Houston, TX 77005, USA; <sup>ge</sup>IFMR Graduate School of Business, Krea University, Sri City, Andhra Pradesh 517646, India; <sup>gf</sup>Department of Psychology, Dilla University, Dilla, South Ethiopia, 419; <sup>gh</sup>Center PhD Students, Research Group: Econometrics, Tilburg School of Economics and Management, Tilburg

university, Tilburg, 5037 AB, Netherlands; <sup>65</sup>Department of Economics, University of Calgary, Calgary, Alberta T2N 1N4, Canada; <sup>66</sup>Resources for the Future, Resources for the Future, Washington, DC 20036; <sup>67</sup>Research Group: Econometrics, Tilburg School of Economics and Management, Tilburg University, Tilburg, 5037 AB, Netherlands; <sup>68</sup>Department of Political Science, University of Montreal, Montreal, Quebec, Canada, H3T 1J4; <sup>69</sup>IVADO, IVADO, Montreal, Quebec, Canada, H3N 1V5; <sup>70</sup>Department of Medicine and Surgery, University of Milano Bicocca, Milan, Italy, 20126; <sup>71</sup>Department of Statistics and Quantitative Methods, University of Milano Bicocca, Milan, Italy, 20126; <sup>72</sup>Department of Health Sciences, University of York, York, United Kingdom, YO10 5DD; <sup>73</sup>Climate and Development Policy Division, RWI - Leibniz Institute for Economic Research, Berlin, Berlin, Germany, 10115; <sup>74</sup>School of Public Policy and Administration, Carleton University, Ottawa, Ontario K1S 5B6, Canada; <sup>75</sup>Institute of Psychology, Cardinal Stefan Wyszyński University, Warsaw, Poland, 01-938; <sup>76</sup>Currency Department, Bank of Canada, Ottawa, Ontario K1A 0G9, Canada; <sup>77</sup>Department of Agricultural Economics and Rural Development, University of Göttingen, Göttingen, Platz der Göttinger Sieben 5, 37083 Göttingen, Niedersachsen, Germany; <sup>78</sup>School of Psychology, University of Sheffield, Sheffield, S10 2TN, United Kingdom; <sup>79</sup>IESE Business School, IESE, Barcelona, 08034 Barcelona, Spain; <sup>80</sup>Universidad de Navarra, Universidad de Navarra, Pamplona, 31009 Pamplona, Spain; <sup>81</sup>Department of Economics, Dedman College of Humanities and Sciences, Southern Methodist University - SMU, SMU, Dallas, TX, 75275-0496; <sup>82</sup>GREDEG, Université Côte d'Azur, SKEMA Business School, Lille, France, 59777; <sup>83</sup>Institute for Public Management and Governance, Vienna University of Economics and Business, Vienna, Austria, 1020; <sup>84</sup>Department of Finance, Rotterdam School of Management, Erasmus University Rotterdam, Rotterdam, 3012 CC, Zuid Holland; <sup>85</sup>Institute of Psychology, Universität Osnabrück, Osnabrück, Germany, 49078; <sup>86</sup>School of Computer Science, University of Sheffield, Sheffield, S1 4DP, United Kingdom; <sup>87</sup>Department of Psychology, University of Houston, Houston, Texas, 77204; <sup>88</sup>Department of Economics, Business and Finance (EBF), Lake Forest College, Lake Forest, IL 60048, USA; <sup>89</sup>Department of Marketing, Vienna University of Economics and Business, Vienna, 1020, Austria; <sup>90</sup>Department of Economics, Bielefeld University, Bielefeld, Germany, 33615; <sup>91</sup>Department of Economics, University of Albany, SUNY, Albany, NY 12222, USA; <sup>92</sup>School of Information, University of Michigan, Ann Arbor, MI 48109, USA; <sup>93</sup>The Unjournal, The Unjournal, Camden, Delaware 19934 USA; <sup>94</sup>INVEST Flagship Research Center, University of Turku, Turku, 20014 Turku, Finland; <sup>95</sup>Applied Economics and Management, Cornell University, Ithaca, NY 14853, USA; <sup>96</sup>Department of Economics, City St George's, University of London, London, EC1V 0HB United Kingdom; <sup>97</sup>Department of Psychology, University of Milano-Bicocca, Milan, 20126; <sup>98</sup>School of Economics, University of Nottingham, Nottingham, NG7 2RD, United Kingdom; <sup>99</sup>InSIDE Lab, DIECO, Università degli Studi dell'Insubria, Varese, Italy, 21100; <sup>100</sup>HEC Liège, Management School, University of Liège, Liège, 4000 Liège, Belgium; <sup>101</sup>Department of Economics, University of Essex, Colchester, CO4 3SQ United Kingdom; <sup>102</sup>Data and Digital Services Department, Bank of Canada, Ottawa, ON, Canada K1A 0G9; <sup>103</sup>Systems and Computer Engineering Department, Carleton University, Ottawa, ON, Canada K1S 5B6; <sup>104</sup>Leverhulme Centre for the Future of Intelligence, University of Cambridge, Cambridge, CB2 1SB United Kingdom; <sup>105</sup>Department of Economics, University of California, Davis, California, 95616; <sup>106</sup>Institute for Business Administration, University of Regensburg, Regensburg, Bavaria, 93053; <sup>107</sup>Department of Linguistics, Ghent University, Ghent, Belgium, 9000; <sup>108</sup>Department of Linguistics and English Language, University of Manchester, Manchester, Manchester, M13 9PL; <sup>109</sup>Department of Economics, McMaster University, Hamilton, Ontario L8S 4L8, Canada; <sup>110</sup>Department of Marketing, Tilburg School of Economics and Management, Tilburg University, Tilburg, 5037 AB, Netherlands; <sup>111</sup>School of Public Finance and Taxation, Dongbei University of Finance and Economics, Dalian, China, 116025; <sup>112</sup>Department of Marketing, Clemson University, Clemson, SC 29634; <sup>113</sup>Department of Economics, University of Birmingham, Birmingham, B15 2TT, UK; <sup>114</sup>School of Information, Journalism and Communication, University of Sheffield, Sheffield, S10 2AH, United Kingdom; <sup>115</sup>Department of Government, Cornell University, Ithaca, NY 14853, USA; <sup>116</sup>Agri-Science Queensland, Department of Primary Industries, Brisbane, Queensland; <sup>117</sup>Asia and Pacific Department, International Monetary Fund, Washington DC, District of Columbia, 20431; <sup>118</sup>Department of Psychology, New York University, New York, New York, 10003; <sup>119</sup>School of Natural and Social Sciences, SRUC, Edinburgh, Edinburgh, UK, EH9 3JG; <sup>120</sup>Department of Economics, Carleton University, Ottawa, Canada, K1S 5B6; <sup>121</sup>Department of Methodology, Tilburg School of Social and Behavioral Sciences, Tilburg University, Tilburg, 5037 AB, Netherlands; <sup>122</sup>Chair of Agricultural Production and Resource Economics, Technical University of Munich, Freising, Germany, 85354; <sup>123</sup>The Institute for Food and Resource Economics, University of Bonn, Bonn, Germany, 53115; <sup>124</sup>Centre for Research on Equality in Education; Department of Special Needs Education, University of Oslo, Oslo, 0318 Oslo, Norway; <sup>125</sup>Department of Sociology and Human Geography, University of Oslo, Oslo, 0317 Oslo, Norway; <sup>126</sup>World Bank, World Bank, Washington, DC 20433, USA; <sup>127</sup>Department of Economics, Carleton University, Ottawa, Canada; <sup>128</sup>Birkbeck Business School, Birkbeck, London, WC1E 7HX, United Kingdom; <sup>129</sup>University of London, University of London, London, WC1E 7HX, United Kingdom; <sup>130</sup>Department of Economics, Vanderbilt University, Nashville, TN 37235-1819, USA; <sup>131</sup>Division of Psychology & Mental Health, School of Health Sciences, University of Manchester, Manchester, M13 9PL, United Kingdom; <sup>132</sup>Department of Economics, Carleton University, Ottawa, Ontario, Canada, K1S 5B6; <sup>133</sup>Center for Health and Wellbeing, Princeton University, Princeton, NJ 08544, USA; <sup>134</sup>Department of Economics, College of Staten Island, Staten Island, NY 10314, USA; <sup>135</sup>City University of New York, CUNY, New York, NY 10017, USA; <sup>136</sup>Department Economics, Law, and Society, ESSCA School of Management, Angers, France, 49003; <sup>137</sup>Department of Economics, University of California, Davis, Davis, CA, 95616; <sup>138</sup>Stockholm Business School, Stockholm University, Stockholm, Stockholm, Sweden, 10691; <sup>139</sup>Leibniz Institute for Financial Research SAFE, Leibniz Institute for Financial Research SAFE, Frankfurt, Hesse, Germany, 60323; <sup>140</sup>School of Chemical, Materials and Biological Engineering, University of Sheffield, Sheffield, S1 3JD, United Kingdom; <sup>141</sup>Department of Economics, National University of Singapore, Singapore, 117570, Singapore; <sup>142</sup>Collegium of World Economy, SGH Warsaw School of Economics, Warsaw, Poland 02-554; <sup>143</sup>Department of Economics, University of Birmingham, Birmingham, United Kingdom B15 2TT; <sup>144</sup>Department of Political Science, Stanford University, Stanford, CA, 94305; <sup>145</sup>Department of Engineering Sciences & Applied Mathematics, Northwestern University, Evanston, IL, 60201; <sup>146</sup>Real Estate Economics, National Chengchi University, Taipei, 11605, Taiwan; <sup>147</sup>Aalto University, Aalto University, Espoo, 02150, Finland; <sup>148</sup>Department of Accounting, Tilburg School of Economics and Management, Tilburg University, Tilburg, The Netherlands, 5037 AB; <sup>149</sup>Department of Economics, Faculty of Arts and Social Sciences, National University of Singapore, Singapore, Singapore 117570; <sup>150</sup>School of Economics, Faculty of Business Economics and Law, University of Queensland, Brisbane, St Lucia 4072, Queensland, Australia; <sup>151</sup>Department of Psychology, Fyffe College, Faculty of Science and Technology, Lancaster University, Lancaster, United Kingdom LA1 4YF; <sup>152</sup>School of Biosciences, University of Sheffield, Sheffield, United Kingdom S10 2TN; <sup>153</sup>Department of Economics, Carleton University, Ottawa, ON K1S 5B6; <sup>154</sup>Christ Church College, University of Oxford, Oxford, OX1 1DP, Oxfordshire, United Kingdom; <sup>155</sup>Economics Experimental Lab, Nanjing Audit University, Nanjing, China, 210017; <sup>156</sup>HEC, Department of Finance, University of Lausanne, Lausanne, Switzerland, 1015; <sup>157</sup>Swiss Finance Institute, Swiss Finance Institute, Lausanne, Switzerland, 1015; <sup>158</sup>Department of Marketing, Rotman School of Management, University of Toronto, Toronto, ON, Canada, M5S 1A1; <sup>159</sup>Department of Economics, ISM University of Management and Economics, Vilnius, LT-01103, Lithuania; <sup>160</sup>Department of Industrial Engineering and Operations Research, Columbia University, New York, NY, 10027, USA; <sup>161</sup>Department of Business and Management Science, NHH Norwegian School of Economics, Bergen, 5045 Bergen, Norway; <sup>162</sup>Department of Economics, Penn State University, University Park, PA 16802, USA

This manuscript was compiled on April 14, 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.**

## Significance Statement

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.

To whom correspondence should be addressed. E-mail: abrodeur@uottawa.ca

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 “proto-agentic” 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 over-reliance 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 over-reliance 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).

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 non-academic organizations with a doctoral degree. 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 Tables S1 and S2. Descriptions of the coding errors identified prior to each replication game can be found in Tables S3 through 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 standard errors 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 re-used 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 (see 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 pre-determined 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 co-authorship 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 different version 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; and (8) Customizing ChatGPT; (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 (see 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 code provided by ChatGPT, though the teams could adjust file paths and their environment without the assistance of ChatGPT. During the pre-games 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, codes, 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). Table S8 shows the treatment arms are balanced across observables.

**A. 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 (completed computational reproducibility versus 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 versus 0.034). We also investigate the extent of false error detection and the share of errors not uncovered by 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 supplementary materials, 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.

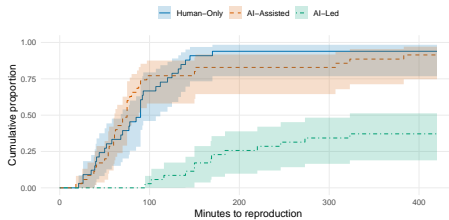
We measure differences by team type in proposing and implementing robustness tests using four measures. The first two are whether teams *proposed* one or two “good” robustness checks. The third and fourth dependent variables are whether the participants report to have *implemented* one or two of those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

## 2. Results

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

**A. 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 (see Table 1). Table 2 shows the ordinary least squares (OLS) estimates of our main regression model (see Table S10 for logit and Poisson regressions and Table S12 for coefficient estimates concerning the control variables). We find that human-only teams are about 59 percentage points more likely than 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. Figure 1 plots 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 (see Table 1). There is no statistically significant difference between human and AI-assisted teams.



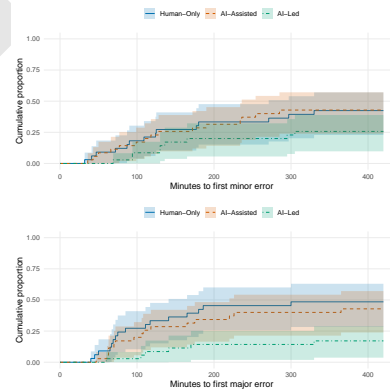
**Fig. 1.** Kaplan-Meier curves, showing the proportion of teams who computationally reproduced the paper by time  $t$  along with curve bands

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 Figure 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 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.

**B. 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 number of minor coding errors discovered 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 ( $p = 0.115$ ). The SI Appendix provide examples of errors and a discussion.

Figure 2 plots 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.



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

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 E) 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 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 over-detection.

We also provide non-causal evidence in exploratory analyses in 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, Table S15 compares human-only teams with AI-assisted teams with high versus 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 imprecisely estimated 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 (Table S13).

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

**C. 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.

**D. Additional Analyses for AI-Assisted Teams.** Table S16 presents an exploratory correlational analysis examining the relationship between AI usage (measured by total prompts) and performance in AI-assisted teams. See Figure 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, 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 over-reliance on AI support(20, 36, 37, 42–44).

**E. 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 micro tasks, for example boilerplate code and file location, while retaining conceptual control, whereas AI-led 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 the 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 suggest 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 non-trivial 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 interpreta-

tion, judgment, and more complex failures. Under this model, even partial automation can generate meaningful cost savings and efficiency gains at scale.

**A. 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.

**B. 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 over-reliance on AI tools. Finally, our analysis is based on a small, non-random 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 (see 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.

**C. 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 humans-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.

**D. 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 prompting strategies that lead to more 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 the SI Appendix.

## Materials and Methods

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), and the King's College London Research Ethics Office reviewed and approved our focus groups (MRA-24/25-48393); all participants provided informed consent. Our pre-analysis 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, pre-analysis plan, and template form available at <https://osf.io/sz2g8/> with no restrictions on sharing or re-use.

**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 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 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
p-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: Standard errors in parentheses; confidence intervals 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.

\*  $p < 0.10$ , \*\*  $p < 0.05$ , \*\*\*  $p < 0.01$

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

	(1) Pre-regression errors	(2) Regression errors	(3) Transcription/post-regression 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
p-val (AI-Assisted = AI-Led)	0.140	0.228	0.274	0.146	0.016
Obs.	103	103	103	103	103

Note: Standard errors in parentheses; confidence intervals 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.

\*  $p < 0.10$ , \*\*  $p < 0.05$ , \*\*\*  $p < 0.01$

1. Abel Brodeur, Kevin Esterling, Jörg Ankel-Peters, Natália S Bueno, Scott Desposato, Anna Dreber, Federica Genovese, Donald P Green, Matthew Hepplewhite, Fernando Hoces de la Guardia, et al. Promoting reproducibility and replicability in political science. *Research & Politics*, 11(1):20531680241233439, 2024.
2. David L Donoho, Arian Maleki, Inam Ur Rahman, Morteza Shahram, and Victoria Stodden. Reproducible research in computational harmonic analysis. *Computing in Science & Engineering*, 11(1):8–18, 2008.
3. Miloš Fišar, Ben Greiner, Christoph Huber, Elena Katok, Ali I Ozkes, and Management Science Reproducibility Collaboration. Reproducibility in management science. *Management Science*, 70(3):1343–1356, 2024.
4. Paul Gertler, Sebastian Galiani, and Mauricio Romero. How to make replication the norm. *Nature*, 554(7693):417–9, 2018.
5. Steven N Goodman, Daniele Fanelli, and John PA Ioannidis. What does research reproducibility mean? *Science Translational Medicine*, 8(341):341ps12–341ps12, 2016.
6. Marcin Milkowski, Witold M Hensel, and Mateusz Hohol. Replicability or reproducibility? on the replication crisis in computational neuroscience and sharing only relevant detail. *Journal of Computational Neuroscience*, 45(3):163–172, 2018.
7. National Academies of Sciences, Engineering, and Medicine. *Reproducibility and Replicability in Science*. National Academies Press, 2019. ISBN 978-0-309-48616-3. . URL <https://www.nap.edu/catalog/25303>.
8. Christophe Pérignon, Kamel Gadouche, Christophe Hurlin, Roxane Silberman, and Eric Debonnel. Certify reproducibility with confidential data. *Science*, 365(6449):127–128, 2019.
9. Lars Vilhuber. Report by the AEA Data Editor. *AEA Papers and Proceedings*, 112:813–23, May 2022. ISSN 2574-0768, 2574-0776. .
10. Abel Brodeur, Derek Mikola, Nikolai Cook, Thomas Bailey, 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.
11. Andrew C Chang and Phillip Li. Is economics research replicable? sixty published papers from thirteen journals say “often not”. *Critical Finance Review*, 11(1):185–206, 2022.
12. Sophia Crüwell, Deborah Apthorp, Bradley J Baker, Lincoln Colling, Malte Elson, Sandra J Geiger, Sebastian Lobentanzer, Jean Monéger, Alex Patterson, D Samuel Schwarzkopf, et al. What's in a badge? a computational reproducibility investigation of the open data badge policy in one issue of psychological science. *Psychological Science*, 34(4):512–522, 2023.
13. Open Science Collaboration. Estimating the reproducibility of psychological science. *Science*, 349(6251):aac4716, 2015.
14. Peepjn Obels, Daniel Lakens, Nicholas A Coles, Jaroslav Gottfried, and Seth A Green. Analysis of open data and computational reproducibility in registered reports in psychology. *Advances in Methods and Practices in Psychological Science*, 3(2):229–237, 2020.
15. Christophe Pérignon, Olivier Akmansoy, Christophe Hurlin, Anna Dreber, Felix Holzmeister, Jürgen Huber, Magnus Johannesson, Michael Kirchler, Albert J Menkveld, Michael Razen, et al. Computational reproducibility in finance: Evidence from 1,000 tests. *The Review of Financial Studies*, 37(11):3558–3593, 2024.
16. Victoria Stodden, Jennifer Seiler, and Zhaokun Ma. An empirical analysis of journal policy effectiveness for computational reproducibility. *Proceedings of the National Academy of Sciences*, 115(11):2584–2589, 2018.
17. Benjamin Wood, Rui Müller, and Annette Brown. Push button replication: Is open data evidence for international development verifiable?, 2018. URL <https://os.io/n7a4d/>. OSF Preprints.
18. Jean-Edouard Colliard, Christophe Hurlin, and Christophe Pérignon. The Economics of Computational Reproducibility, 2022. SSRN: <https://ssrn.com/abstract=3418896>.
19. Allison Hryciyshyn and Helen Eassom. ExplainAItions: An AI Study. Technical report, John Wiley & Sons, Hoboken, NJ, February 2025. URL <https://www.wiley.com/content/dam/wiley-com/en/pdfs/about/wiley-explanations-ai-study-february-2025vers3.pdf>.
20. Michelle Vaccaro, Abdullah Almaatouq, and Thomas Malone. When combinations of humans and ai are useful: A systematic review and meta-analysis. *Nature Human Behaviour*, pages 1–11, 2024.
21. Gagan Bansal, Tongshuang Wu, Joyce Zhou, Raymond Fok, Besmira Nushi, Ece Kamar, Marco Tulio Ribeiro, and Daniel Weld. 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*, pages 1–16, 2021.
22. Gagan Bansal, Besmira Nushi, Ece Kamar, Eric Horvitz, and Daniel S Weld. Is the most accurate ai the best teammate? optimizing ai for teamwork. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 35, pages 11405–11414, 2021.
23. Elizabeth Bondi, Raphael Koster, Hannah Sheahan, Martin Chadwick, Yoram Bachrach, Taylan Cemgil, Ulrich Paquet, and Krishnamurthy Dvijotham. Role of human-ai interaction in selective prediction. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 36, pages 5286–5294, 2022.
24. Ángel Alexander Cabrera, Adam Perer, and Jason I Hong. Improving human-ai collaboration with descriptions of ai behavior. *Proceedings of the ACM on Human-Computer Interaction*, 7(CSCW1):1–21, 2023.
25. Ethan Goh, Robert Gallo, Jason Hom, Eric Strong, Yingjie Weng, Hannah Kerman, Josephine Cool, Zahir Kanjee, Andrew S Parsons, Neera Ahuja, et al. Influence of a large language model on diagnostic reasoning: A randomized clinical vignette study. *medRxiv*, 2024.
26. Brian Koepnick, Jeff Flatten, Tamir Husain, Alex Ford, Daniel-Adriano Silva, Matthew J Bick, Aaron Bauer, Gaohua Liu, Yojiro Ishida, Alexander Boykov, et al. De novo protein design by citizen scientists. *Nature*, 570(7761):390–394, 2019.
27. Han Liu, Vivian Lai, and Chenhao Tan. Understanding the effect of out-of-distribution examples and interactive explanations on human-ai decision making. *Proceedings of the ACM on Human-Computer Interaction*, 5(CSCW2):1–45, 2021.
28. Hussein Mozannar, Jimin Lee, Dennis Wei, Prasanna Sattigeri, Subho Das, and David Sontag. Effective human-ai teams via learned natural language rules and onboarding. *Advances in Neural Information Processing Systems*, 36, 2024.
29. Shakked Noy and Whitney Zhang. Experimental evidence on the productivity effects of generative artificial intelligence. *Science*, 381(6654):187–192, 2023.
30. Carlo Reverberi, Tommaso Rigon, Aldo Solari, Cesare Hassan, Paolo Cherubini, and Andrea Cherubini. Experimental evidence of effective human-ai collaboration in medical decision-making. *Scientific Reports*, 12(1):14952, 2022.
31. Max Schemmer, Patrick Hemmer, Maximilian Nitsche, Niklas Kühn, and Michael 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*, pages 617–626, 2022.
32. Michael Henry Tessler, Michiel A Bakker, Daniel Jarrett, Hannah Sheahan, Martin J Chadwick, Raphael Koster, Georgina Evans, Lucy Campbell-Gillingham, Tatum Collins, David C Parkes, et al. Ai can help humans find common ground in democratic deliberation. *Science*, 386(6719):eadq2852, 2024.
33. Michelle Vaccaro and Jim Waldo. The effects of mixing machine learning and human judgment. *Communications of the ACM*, 62(11):104–110, 2019.
34. Bryan Wilder, Eric Horvitz, and Ece Kamar. Learning to complement humans. In *Proceedings of the 29th International Joint Conference on Artificial Intelligence*, pages 1526–1533, 2020.
35. Cherry Bekaert LLP. Report of independent auditor, 2022. ISSN 0002-8282. URL <https://pubs.aeaweb.org/doi/10.1257/aer.112.6.2083>.
36. Zana Bućina, Maja Barbara Malaya, and Krzysztof Z Gajos. To trust or to think: cognitive forcing functions can reduce overreliance on ai in ai-assisted decision-making. *Proceedings of the ACM on Human-Computer Interaction*, 5(CSCW1):1–21, 2021.
37. Linda J Skitka, Kathleen L Mosier, and Mark Burdick. Does automation bias decision-making? *International Journal of Human-Computer Studies*, 51(5):991–1006, 1999.
38. Erik Brynjolfsson, Danielle Li, and Lindsey Raymond. Generative ai at work. *Quarterly Journal of Economics*, 140:889–942, 2025.
39. Marco Del Giudice and Steven W Gangestad. A traveler's guide to the multiverse: Promises, pitfalls, and a framework for the evaluation of analytic decisions. *Advances in Methods and Practices in Psychological Science*, 4(1):2515245920954925, 2021.
40. Xun Lu and Halbert White. Robustness checks and robustness tests in applied economics. *Journal of Econometrics*, 178:194–206, 2014.
41. Michèle B Nuijten. Assessing and improving robustness of psychological research findings in four steps. In *Avoiding questionable research practices in applied psychology*, pages 379–400. Springer, 2022.
42. Vivian Lai, Han Liu, and Chenhao Tan. “why is’ chicago’deceptive?” towards building model-driven tutorials for humans. In *Proceedings of the 2020 CHI Conference on Human Factors in Computing Systems*, pages 1–13, 2020.
43. Yunfeng Zhang, Q Vera Liao, and Rachel KE 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*, pages 295–305, 2020.
44. Helena Vasconcelos, Matthew Jörke, Madeleine Grunde-McLaughlin, Tobias Gerstenberg, Michael S Bernstein, and Ranjay Krishna. Explanations can reduce overreliance on ai systems during decision-making. *Proceedings of the ACM on Human-Computer Interaction*, 7(CSCW1):1–38, 2023.
45. Yogesh K Dwivedi, Nir Kshetri, Laurie Hughes, Emma Louise Slade, Anand Jeyaraj, Arpan Kumar Kar, Abdullah M Baabdullah, Alex Koohang, Vishnupriya Raghavan, Manju Ahuja, 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. *International Journal of Information Management*, 71:102642, 2023.
46. Brady D Lund, Ting Wang, Nishith Reddy Mannuru, Bing Nie, Somipam Shimray, and Ziang Wang. Chatgpt and a new academic reality: Artificial intelligence-written research papers and the ethics of the large language models in scholarly publishing. *Journal of the Association for Information Science and Technology*, 74(5):570–581, 2023.
47. Nalin Wadhwa, Jui Pradhan, Atharv Sonwane, Surya Prakash Sahu, Nagarajan Natarajan, Aditya Kanade, Suresh Parthasarathy, and Sriram Rajamani. Core: Resolving code quality issues using llms. *Proceedings of the ACM on Software Engineering*, 1(FSE):789–811, 2024.
48. Yichi 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*, pages 683–685, 2024.
49. Daye Nam, Andrew Macvean, Vincent Helleendoorn, Bogdan Vasilescu, and Brad Myers. Using an llm to help with code understanding. In *Proceedings of the IEEE/ACM 46th International Conference on Software Engineering*, pages 1–13, 2024.

## Acknowledgments

We would like to thank Gabriel Zimmerman for research assistance.

**Funding:** This research and AI replications games were funded by Open Philanthropy 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. Shumi Akhtar acknowledges funding DP200102935 awarded by the Australian Research Council Grant.

**Author contributions:** ABr, DVa, AMa, JPap, DMi, BBa, RAl, GSa, GTa, FAk, NTBo, CCo, LFi, JFi, JFr, DRe, GGi, STo,

LPBe, MME, AMu, NMa, SMe, PSu, RSe, VSk, LYa, BWe: conception of study and revision. JPAP, ABr, DVa, GBe, CGBi, ZAr, IAs, TAu, CBo, RBr, SBr, CBo, ACh, ADh, EDu, YEi, JFi, OHa, AHo, GHu, HHu, EKe, VKI, NKr, JLi, RLi, KLo, AMANo, SMe, SMi, SMu, FNe, TNg, UPh, MRe, PSi, SSSt, BTa, MTa, OTa, DVa, CWe, VYWu, ZWu: analysis and interpretation of data. ABr, DVa, AMa, JPAP, DMi, BBA wrote the original draft, lead revision, and take responsibility for the content – while SAi, IAs, MAlt, FAK, SAK, BJBa, LBa, GBe, MBe, CGBi, CBo, YBo, RBr, SBr, ACh, JCo NCo, FCl, LDee, LDen, EDj, TDu, AEl, IElFa, GFe, JFe, LFi, LFr, JFi, JFDo, AGá, JGa, SGa, GGi, AGKa, DGo, IGr, EGZo, FHa, JHa, MHa, SHe, AHo, GHu, ZIl, JJa, EKe, EKi, SKi, NKr, EKu, SKu, BLWh, DLi, JLi, KLo, CMa, RMcW, XMa, SMe, BMo, FMo, MMo, LPMo, LMu, AMANo, FNe, AOz, OOz, SPa, UPh, PPu, QQi, NRu, OSa, DRe, JRe, MTRe, RSe, PSi, SSSt, TSt, PSu, GTa, RTh, TTi, LWa, CWe, MWe, MWr, WCWe, VYWu, BYa, LYa, MZa, HZh, XZh, AZi, ZZa, FZo contributed to editing and review through commenting on drafts of the paper. DVa: AI training. AGKa, JPAP, RAl, FBBi, RFo, ANo: focus group moderators and conception. CBAr, DAn, FCl, NCo, NDh, AFCh, GEFr, JFr, DGe, YJi, PKr, SKu, JMa, NMa, MME, CSa, OSa, WSo, TSt, GTa, AdVa, LYa: focus group participation. ABr, DVa, JPAP, DMi, RAl, LDee, TSa, LVi: Event organization. All coauthors except ABr, DVa, AMa, DMi, BBA, TSa: data acquisition through AI replication games.

**Competing interests:** The views expressed in this paper are those of the authors. No responsibility for them should be attributed to the Bank of Canada. AMa is a UKRI 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.

**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 and materials availability:** We make our (i) AI training materials and recording, (ii) data and codes, (iii) pre-analysis plan and (iv) template form available here: <https://osf.io/sz2g8/>. We declare no restrictions on sharing or re-use.

## **SI Appendix. Materials and Methods**

Figures S1 to S7

Tables S1 to S17