

CAMBRIDGE WORKING PAPERS IN ECONOMICS

JANEWAY INSTITUTE WORKING PAPERS

Non-Standard Errors

Albert J. Menkveld Vrije Universiteit Amsterdam	Felix Holzmeister University of Innsbruck	Magnus Johannesson Stockholm School of Economics	Sebastian Neusüss Aalto University	Utz Weitzel Vrije Universiteit Amsterdam
Anna Dreber Stockholm School of Economics	Juergen Huber University of Innsbruck	Michael Kirchler University of Innsbruck	Michael Razen University of Innsbruck	Oliver Linton University of Cambridge

Abstract

In statistics, samples are drawn from a population in a data generating process (DGP). Standard errors measure the uncertainty in sample estimates of population parameters. In science, evidence is generated to test hypotheses in an evidence generating process (EGP). We claim that EGP variation across researchers adds uncertainty: non-standard errors. To study them, we let 164 teams test six hypotheses on the same sample. We find that non-standard errors are sizeable, on par with standard errors. Their size (i) co-varies only weakly with team merits, reproducibility, or peer rating, (ii) declines significantly after peer-feedback, and (iii) is underestimated by participants.

Online appendix available at <https://bit.ly/3DIQKrB>

Please note a full list of authors is available in the working paper

Reference Details

2182 2021/12	Cambridge Working Papers in Economics Janeway Institute Working Paper Series
Published	25 November 2021
Key Words JEL Codes	Market Efficiency, P-hacking, Publication bias A14, C10, C12, C59, C90, G14, G40
Websites	www.econ.cam.ac.uk/cwpe www.janeway.econ.cam.ac.uk/working-papers

Non-Standard Errors*

Albert J. Menkveld^{200,121}, *Anna Dreber*¹¹³, *Felix Holzmeister*¹⁴⁹,
*Juergen Huber*¹⁴⁹, *Magnus Johannesson*¹¹³, *Michael Kirchler*¹⁴⁹,
*Sebastian Neusüss*¹, *Michael Razen*¹⁴⁹, *Utz Weitzel*^{200,100}, David
Abad-Díaz¹³³, Menachem Abudy¹², Tobias Adrian⁵⁸, Yacine
Ait-Sahalia⁹⁵, Olivier Akmansoy^{17,21}, Jamie T. Alcock¹⁶⁹, Vitali
Alexeev¹⁷⁹, Arash Aloosh⁸¹, Livia Amato¹³⁹, Diego Amaya²⁰⁴, James
J. Angel⁴⁴, Alejandro T. Avetikian⁹³, Amadeus Bach¹⁵⁴, Edwin
Baidoo¹¹⁸, Gaetan Bakalli⁶, Li Bao¹²², Andrea Barbon¹⁷⁴, Oksana
Bashchenko¹⁰⁴, Parampreet C. Bindra¹⁴⁹, Geir H. Bjønnes⁸, Jeffrey
R. Black¹⁵⁸, Bernard S. Black⁸⁵, Dimitar Bogoev³², Santiago
Bohorquez Correa¹²⁹, Oleg Bondarenko¹⁴⁷, Charles S. Bos²⁰⁰, Ciril
Bosch-Rosa¹¹⁶, Elie Bouri⁶⁹, Christian Brownlees¹³⁰, Anna
Calamia¹¹⁵, Viet Nga Cao⁷⁸, Gunther Capelle-Blancard¹³¹, Laura M.
Capera Romero²⁰⁰, Massimiliano Caporin¹⁷⁰, Allen Carrion¹⁵⁸, Tolga
Caskurlu¹³⁴, Bidisha Chakrabarty¹⁰⁷, Jian Chen⁹⁷, Mikhail
Chernov¹²⁴, William Cheung²⁰², Ludwig B. Chincarini¹⁷², Tarun
Chordia³⁷, Sheung Chi Chow⁷, Benjamin Clapham⁴⁶, Jean-Edouard
Colliard⁴⁹, Carole Comerton-Forde¹⁵⁷, Edward Curran⁷⁴, Thong
Dao⁸⁷, Wale Dare⁴⁷, Ryan J. Davies⁹, Riccardo De Blasis⁶⁷, Gianluca
F. De Nard¹⁹³, Fany Declerck¹²², Oleg Deev⁷⁵, Hans Degryse⁶¹,
Solomon Y. Deku⁸⁷, Christophe Desagre¹²⁵, Mathijs A. van Dijk³⁸,
Chukwuma Dim⁴³, Thomas Dimpfl¹⁴⁶, Yun Jiang Dong⁹⁷, Philip A.
Drummond⁷⁸, Tom Dudda¹¹⁷, Teodor Duevski⁴⁹, Ariadna
Dumitrescu³⁶, Teodor Dyakov³³, Anne Haubo Dyhrberg¹⁷⁸, Michał
Dzieliński¹¹⁴, Asli Eksi¹⁰⁸, Izidin El Kalak²⁰, Saskia ter Ellen²²,
Nicolas Eugster¹⁷¹, Martin D. D. Evans⁴⁴, Michael Farrell¹⁸⁹, Ester
Felez-Vinas¹⁷⁹, Gerardo Ferrara¹¹, El Mehdi Ferrouhi⁵⁶, Andrea
Flori⁹², Jonathan T. Fluharty²⁰³, Sean D. V. Foley⁷⁴, Kingsley Y. L.
Fong¹⁶⁵, Thierry Foucault⁴⁹, Tatiana Franus¹⁴, Francesco
Franzoni¹⁹⁵, Bart Frijns⁸⁹, Michael Frömmel⁴⁵, Servanna M. Fu¹⁴¹,
Sascha C. Füllbrunn¹⁰⁰, Baoqing Gan¹⁷⁹, Ge Gao¹³⁵, Thomas P.
Gehrig¹⁸⁸, Roland Gemayel⁶⁵, Dirk Gerritsen¹⁹⁷, Javier
Gil-Bazo^{130,13}, Dudley Gilder²⁰, Lawrence R. Glosten²⁵, Thomas

Gomez¹⁹⁷, Arseny Gorbenko⁷⁸, Joachim Grammig¹⁸⁴, Vincent Grégoire⁴⁸, Ufuk Güçbilmez¹⁴³, Björn Hagströmer¹¹⁴, Julien Hambuckers⁴⁷, Erik Hapnes¹, Jeffrey H. Harris³, Lawrence Harris¹⁷³, Simon Hartmann¹⁹⁹, Jean-Baptiste Hasse², Nikolaus Hautsch¹⁸⁸, Xue-Zhong (Tony) He²⁰⁵, Davidson Heath¹⁸⁵, Simon Hediger¹⁹³, Terrence Hendershott¹²³, Ann Marie Hibbert²⁰³, Erik Hjalmarsson¹⁴⁴, Seth Hoelscher⁷⁷, Peter Hoffmann³⁹, Craig W. Holden⁵⁷, Alex R. Horenstein¹⁵⁹, Wenqian Huang¹⁰, Da Huang¹⁸⁵, Christophe Hurlin^{168,21}, Konrad Ilczuk¹, Alexey Ivashchenko²⁰⁰, Subramanian R. Iyer¹⁶⁴, Hossein Jahanshahloo²⁰, Naji P. Jalkh¹⁰⁶, Charles M. Jones²⁴, Simon Jurkatis¹¹, Petri Jylhä¹, Andreas T. Kaeck¹⁷⁷, Gabriel Kaiser¹⁵², Arzé Karam³⁰, Egle Karmaziene²⁰⁰, Bernhard Kassner¹⁶², Markku Kaustia¹, Ekaterina Kazak¹⁵³, Fearghal Kearney⁹⁸, Vincent van Kervel⁹⁴, Saad A. Khan⁴⁸, Marta K. Khomyn¹⁷⁹, Tony Klein⁹⁸, Olga Klein¹⁹⁰, Alexander Klos⁶³, Michael Koetter⁵⁰, Aleksey Kolokolov¹⁵³, Robert A. Korajczyk⁸⁵, Roman Kozhan¹⁹⁰, Jan P. Krahnen⁴⁶, Paul Kuhle¹²⁷, Amy Kwan¹⁷⁸, Quentin Lajaunie⁹¹, F. Y. Eric C. Lam⁵¹, Marie Lambert⁴⁷, Hugues Langlois⁴⁹, Jens Lausen⁴⁶, Tobias Lauter⁷¹, Markus Leippold¹⁹³, Vladimir Levin¹⁵², Yijie Li¹⁰⁹, Hui Li⁶⁸, Chee Yoong Liew¹²⁶, Thomas Lindner²⁰¹, Oliver Linton¹³⁸, Jiacheng Liu⁹⁶, Anqi Liu¹⁷⁸, Guillermo Llorente¹²⁷, Matthijs Lof¹, Ariel Lohr⁴, Francis Longstaff¹²⁴, Alejandro Lopez-Lira¹⁴², Shawn Mankad²⁷, Nicola Mano¹⁰², Alexis Marchal³⁵, Charles Martineau¹⁸², Francesco Mazzola³⁸, Debrah Meloso¹¹⁵, Michael G. Mi¹⁷⁸, Roxana Mihet¹⁰¹, Vijay Mohan⁹⁹, Sophie Moinas¹²², David Moore⁷², Liangyi Mu¹²⁰, Dmitriy Muravyev⁷⁶, Dermot Murphy¹⁴⁷, Gabor Neszveda⁵⁹, Christian Neumeier⁶⁰, Ulf Nielsson²⁶, Mahendrarajah Nimalendran¹⁴², Sven Nolte¹⁰⁰, Lars L. Norden¹¹⁴, Peter W. O'Neill⁴², Khaled Obaid¹⁸, Bernt A. Ødegaard¹⁷⁵, Per Östberg¹⁹³, Emiliano Pagnotta¹¹⁰, Marcus Painter¹⁰⁷, Stefan Palan¹⁴⁵, Imon J. Palit⁹⁹, Andreas Park¹⁸³, Roberto Pascual¹⁹⁴, Paolo Pasquariello¹⁶⁰, Lubos Pastor¹³⁹, Vinay Patel¹⁷⁹, Andrew J. Patton²⁹, Neil D. Pearson^{148,19}, Lorian Pelizzon⁴⁶, Michele Pelli¹⁰⁵, Matthias Pelster⁹⁰, Christophe Pérignon^{49,21}, Cameron Pffiffer¹⁶⁷, Richard Philip¹⁷⁸, Tomáš Plíhal⁷⁵, Puneet Prakash⁷⁷, Oliver-Alexander Press²⁶, Tina Prodromou¹⁹², Marcel Prokopczuk⁷¹, Talis Putnins¹⁷⁹, Ya Qian¹, Gaurav Raizada⁵³, David Rakowski¹⁸⁰, Angelo Ranaldo¹⁷⁴, Luca Regis¹⁸¹, Stefan Reitz⁶⁴, Thomas Renault¹⁹⁶, Rex W. Renjie²⁰⁰, Roberto Reno¹⁸⁶, Steven J. Riddiough¹⁸², Kalle

Rinne¹⁵², Paul J. Rintamäki¹, Ryan Riordan⁹⁷, Thomas Rittmannsberger¹⁴⁹, Iñaki Rodríguez Longarela¹¹⁴, Dominik Roesch¹¹², Lavinia Rognone¹⁵³, Brian Roseman⁸⁸, Ioanid Rosu⁴⁹, Saurabh Roy¹⁵⁶, Nicolas Rudolf¹⁵¹, Stephen R. Rush¹⁵, Khaladdin Rzayev^{140,66}, Aleksandra A. Rzeźnik²⁰⁶, Anthony Sanford¹⁵⁵, Harikumar Sankaran⁸², Asani Sarkar⁴¹, Lucio Sarno¹³⁸, Olivier Scaillet¹⁰³, Stefan Scharnowski¹⁵⁴, Klaus R. Schenk-Hoppé¹⁵³, Andrea Schertler¹⁴⁵, Michael Schneider^{28,70}, Florian Schroeder⁷⁴, Norman Schürhoff¹⁰⁴, Philipp Schuster¹⁷⁶, Marco A. Schwarz^{31,16}, Mark S. Seasholes⁴, Norman J. Seeger²⁰⁰, Or Shachar⁴¹, Andriy Shkilko²⁰⁴, Jessica Shui⁴⁰, Mario Sikic¹⁹³, Giorgia Simion²⁰¹, Lee A. Smales¹⁹¹, Paul Söderlind¹⁷⁴, Elvira Sojli¹⁶⁵, Konstantin Sokolov¹⁵⁸, Jantje Sönksen¹⁸⁴, Laima Spokeviciute²⁰, Denitsa Stefanova¹⁵², Marti G. Subrahmanyam^{80,79}, Barnabas Szaszi³⁴, Oleksandr Talavera¹³⁵, Yuehua Tang¹⁴², Nick Taylor¹³⁷, Wing Wah Tham¹⁶⁵, Erik Theissen¹⁵⁴, Julian Thimme⁶², Ian Tonks¹³⁷, Hai Tran⁷², Luca Trapin¹³⁶, Anders B. Trolle²⁶, M. Andreea Vaduva¹²⁸, Giorgio Valente⁵², Robert A. Van Ness¹⁶¹, Aurelio Vasquez⁵⁵, Thanos Verousis¹⁴¹, Patrick Verwijmeren³⁸, Anders Vilhelmsson⁷³, Grigory Vilkov⁴³, Vladimir Vladimirov¹³⁴, Sebastian Vogel³⁸, Stefan Voigt¹⁵⁰, Wolf Wagner³⁸, Thomas Walther¹⁹⁷, Patrick Weiss¹⁹⁸, Michel van der Wel³⁸, Ingrid M. Werner¹¹⁹, Joakim Westerholm¹⁷⁸, Christian Westheide¹⁸⁸, Hans C. Wika⁸⁴, Evert Wipplinger²⁰⁰, Michael Wolf¹⁹³, Christian C. P. Wolff¹⁵², Leonard Wolk²⁰⁰, Wing-Keung Wong⁵, Jan Wrampelmeyer²⁰⁰, Zhen-Xing Wu¹, Shuo Xia⁵⁰, Dacheng Xiu¹³⁹, Ke Xu¹⁸⁷, Caihong Xu¹¹⁴, Pradeep K. Yadav¹⁶⁶, José Yagüe¹⁶³, Cheng Yan¹⁴¹, Antti Yang³⁸, Woongsun Yoo²³, Wenjia Yu¹, Yihe Yu¹³², Shihao Yu²⁰⁰, Bart Z. Yueshen⁵⁴, Darya Yuferova⁸⁶, Marcin Zamojski¹⁴⁴, Abalfazl Zareei¹¹⁴, Stefan M. Zeisberger¹⁰⁰, Lu Zhang¹⁵², S. Sarah Zhang¹⁵³, Xiaoyu Zhang²⁰⁰, Lu Zhao¹¹¹, Zhuo Zhong¹⁵⁷, Zeyang (Ivy) Zhou¹⁹², Chen Zhou³⁸, Xingyu S. Zhu¹¹³, Marius Zoican¹⁸³, and Remco Zwinkels²⁰⁰

¹Aalto University, ²Aix-Marseille University, ³American University, ⁴Arizona State University, ⁵Asia University, ⁶Auburn University, ⁷Australian National University, ⁸BI Norwegian Business School, ⁹Babson College, ¹⁰Bank for International Settlements, ¹¹Bank of England, ¹²Bar-Ilan University, ¹³Barcelona School of Economics, ¹⁴Bayes Business School, ¹⁵Bowling Green State University, ¹⁶CESifo, ¹⁷CNRS, ¹⁸California State University - East Bay,

¹⁹Canadian Derivatives Institute, ²⁰Cardiff University, ²¹Cascad, ²²Central Bank of Norway, ²³Central Michigan University, ²⁴Columbia Business School, ²⁵Columbia University, ²⁶Copenhagen Business School, ²⁷Cornell University, ²⁸Deutsche Bundesbank, ²⁹Duke University, ³⁰Durham University, ³¹Düsseldorf Institute for Competition Economics, ³²EDF Energy London, ³³EDHEC Business School, ³⁴ELTE, Eotvos Lorand University, ³⁵EPFL, ³⁶ESADE Business School, Univ. Ramon Llull, ³⁷Emory University, ³⁸Erasmus University Rotterdam, ³⁹European Central Bank, ⁴⁰Federal Housing Finance Agency, ⁴¹Federal Reserve Bank of New York, ⁴²Financial Conduct Authority, ⁴³Frankfurt School of Finance and Management, ⁴⁴Georgetown University, ⁴⁵Ghent University, ⁴⁶Goethe University Frankfurt, ⁴⁷HEC Liège - University of Liège, ⁴⁸HEC Montréal, ⁴⁹HEC Paris, ⁵⁰Halle Institute for Economic Research, ⁵¹Hong Kong Institute for Monetary and Financial Research, ⁵²Hong Kong Monetary Authority, ⁵³IIM Ahmedabad, ⁵⁴INSEAD, ⁵⁵ITAM, ⁵⁶Ibn Tofail University, ⁵⁷Indiana University, ⁵⁸International Monetary Fund, ⁵⁹John von Neumann University, ⁶⁰Justus-Liebig University, ⁶¹KU Leuven, ⁶²Karlsruhe Institute of Technology, ⁶³Kiel University, ⁶⁴Kiel university, ⁶⁵King's College London, ⁶⁶Koç University, ⁶⁷LUM University, ⁶⁸La Trobe University, ⁶⁹Lebanese American University, ⁷⁰Leibniz Institute for Financial Research SAFE, ⁷¹Leibniz University Hannover, ⁷²Loyola Marymount University, ⁷³Lund University, ⁷⁴Macquarie University, ⁷⁵Masaryk University, ⁷⁶Michigan State University, ⁷⁷Missouri State University, ⁷⁸Monash University, ⁷⁹NYU Shanghai, ⁸⁰NYU Stern, ⁸¹Neoma Business School, ⁸²New Mexico state University, ⁸³None, ⁸⁴Norges Bank, ⁸⁵Northwestern University, ⁸⁶Norwegian School of Economics (NHH), ⁸⁷Nottingham Trent University, ⁸⁸Oklahoma State University, ⁸⁹Open Universiteit, ⁹⁰Paderborn University, ⁹¹Paris Dauphine University, ⁹²Politecnico di Milano, ⁹³Pontificia Universidad Católica de Chile, ⁹⁴Pontifical University of Chile, ⁹⁵Princeton University, ⁹⁶Purdue University, ⁹⁷Queen's University, ⁹⁸Queen's University Belfast, ⁹⁹RMIT University, ¹⁰⁰Radboud University, ¹⁰¹SFI at HEC Lausanne, ¹⁰²SFI at USI Lugano, ¹⁰³SFI at University of Geneva, ¹⁰⁴SFI at University of Lausanne, ¹⁰⁵SFI at University of Zurich, ¹⁰⁶Saint Joseph University, ¹⁰⁷Saint Louis University, ¹⁰⁸Salisbury University, ¹⁰⁹SandP Global Ratings, ¹¹⁰Singapore Management University, ¹¹¹Southwestern University of Finance and Economics, ¹¹²State University of New York at Buffalo, ¹¹³Stockholm School of Economics, ¹¹⁴Stockholm University, ¹¹⁵TBS Education, ¹¹⁶Technische Universität Berlin, ¹¹⁷Technische Universität Dresden, ¹¹⁸Tennessee Technological University, ¹¹⁹The Ohio State University, ¹²⁰The University of Manchester, ¹²¹Tinbergen Institute, ¹²²Toulouse 1 Capitole University, ¹²³UC Berkeley, ¹²⁴UCLA, ¹²⁵UCLouvain, ¹²⁶UCSI University, ¹²⁷Universidad Autónoma de Madrid, ¹²⁸Universidad Carlos III de Madrid, ¹²⁹Universidad EAFIT, ¹³⁰Universitat Pompeu Fabra, ¹³¹University Paris 1 Pantheon-Sorbonne, ¹³²University at Buffalo, ¹³³University of Alicante, ¹³⁴University of Amsterdam, ¹³⁵University of Birmingham, ¹³⁶University of Bologna, ¹³⁷University of Bristol, ¹³⁸University of Cambridge, ¹³⁹University of Chicago Booth School of Business, ¹⁴⁰University of Edinburgh, ¹⁴¹University of Essex, ¹⁴²University of Florida, ¹⁴³University of Glasgow, ¹⁴⁴University of Gothenburg, ¹⁴⁵University of Graz, ¹⁴⁶University of Hohenheim, ¹⁴⁷University of

Illinois at Chicago, ¹⁴⁸University of Illinois at Urbana-Champaign, ¹⁴⁹University of Innsbruck, ¹⁵⁰University of København, ¹⁵¹University of Lausanne, ¹⁵²University of Luxembourg, ¹⁵³University of Manchester, ¹⁵⁴University of Mannheim, ¹⁵⁵University of Maryland, ¹⁵⁶University of Massachusetts, Amherst, ¹⁵⁷University of Melbourne, ¹⁵⁸University of Memphis, ¹⁵⁹University of Miami, ¹⁶⁰University of Michigan, ¹⁶¹University of Mississippi, ¹⁶²University of Munich (LMU), ¹⁶³University of Murcia, ¹⁶⁴University of New Mexico, ¹⁶⁵University of New South Wales, ¹⁶⁶University of Oklahoma, ¹⁶⁷University of Oregon, ¹⁶⁸University of Orléans, ¹⁶⁹University of Oxford, ¹⁷⁰University of Padova, ¹⁷¹University of Queensland, ¹⁷²University of San Francisco, ¹⁷³University of Southern California, ¹⁷⁴University of St. Gallen, ¹⁷⁵University of Stavanger, ¹⁷⁶University of Stuttgart, ¹⁷⁷University of Sussex, ¹⁷⁸University of Sydney, ¹⁷⁹University of Technology Sydney, ¹⁸⁰University of Texas at Arlington, ¹⁸¹University of Torino, ¹⁸²University of Toronto, ¹⁸³University of Toronto Mississauga, ¹⁸⁴University of Tübingen, ¹⁸⁵University of Utah, ¹⁸⁶University of Verona, ¹⁸⁷University of Victoria, ¹⁸⁸University of Vienna, ¹⁸⁹University of Virginia, ¹⁹⁰University of Warwick, ¹⁹¹University of Western Australia, ¹⁹²University of Wollongong, ¹⁹³University of Zurich, ¹⁹⁴University of the Balearic Islands, ¹⁹⁵Università della Svizzera italiana, ¹⁹⁶Université Paris 1 Panthéon-Sorbonne, ¹⁹⁷Utrecht University, ¹⁹⁸Vienna Graduate School of Finance, ¹⁹⁹Vienna University of Economics and Business, ²⁰⁰Vrije Universiteit Amsterdam, ²⁰¹WU Vienna University of Economics and Business, ²⁰²Waseda University, ²⁰³West Virginia University, ²⁰⁴Wilfrid Laurier University, ²⁰⁵Xi'an Jiaotong-Liverpool University, ²⁰⁶York University, ²⁰⁷Zhongnan University of Economics and Law

November 23, 2021

*The first nine authors in italics are the project coordinators. They designed the project, managed it, and wrote the manuscript. Any errors are therefore their sole responsibility. The other authors all significantly contributed to the project by participating either as a member of a research team, or as a peer evaluator. This manuscript represents the views of the authors and does not reflect the views of any of the institutions authors are affiliated with or receive financing from. The coordinators thank Lucas Saru for valuable comments. They further thank Adam Gill, Eugénie de Jong, and Elmar Nijkamp for research assistance. The coordinators are grateful for financial support from (Dreber) the Knut and Alice Wallenberg Foundation, the Marianne, Marcus Wallenberg Foundation, the Jan Wallander, Tom Hedelius Foundation, (Huber) an FWF grant P29362, (Huber and Kirchler) FWF SFB F63, and (Menkveld) NWO-Vici.

Non-Standard Errors

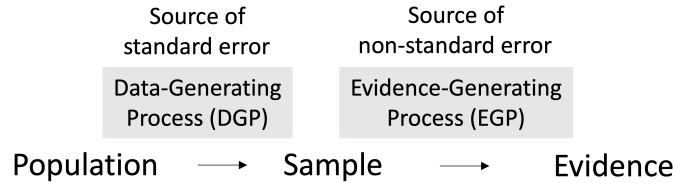
Abstract

In statistics, samples are drawn from a population in a data-generating process (DGP). Standard errors measure the uncertainty in sample estimates of population parameters. In science, evidence is generated to test hypotheses in an evidence-generating process (EGP). We claim that EGP variation across researchers adds uncertainty: *non-standard errors*. To study them, we let 164 teams test six hypotheses on the same sample. We find that non-standard errors are sizeable, on par with standard errors. Their size (i) co-varies only weakly with team merits, reproducibility, or peer rating, (ii) declines significantly after peer-feedback, and (iii) is underestimated by participants.

Online appendix available at <https://bit.ly/3DIQKrB>.

1 Introduction

Academic research recognizes randomness in data samples by computing standard errors (SEs) for parameter estimates. It, however, does *not* recognize the randomness that is in the research process itself. We believe that such randomness is the cause of, what we will call, non-standard errors (NSEs).



The above schema depicts the overarching idea of non-standard errors. Statisticians use the term data-generating process (DGP) to convey the idea that samples are random draws from a population. These samples are used to estimate population parameters, with error: standard error. Using the same language, one could say that scientists collectively engage in an evidence-generating process (EGP). This EGP exposes variation across researchers, adding further error, which we refer to as *non-standard error*. Can this additional uncertainty be safely ignored, or is it sizeable? Does it decline with peer feedback? This has been our core motivation when designing the *Finance Crowd Analysis Project* (#fincap).¹

Our working definition of non-standard errors is the standard deviation across researchers for the results they report when independently testing the same hypotheses on the same sample. For example, in #fincap we provide researchers with identical data and then ask them to independently propose a measure for market efficiency, estimate its average per-year change, and report it. The NSE for this reported change is simply the standard deviation across researchers of their reported changes.

Non-standard errors exist for a variety of reasons. They might, of course, be due to glitches in computer code or due to machine

¹To increase the credibility of the findings and address concerns of specification searching, we wrote a pre-analysis plan (PAP) and filed it with the Open Science Foundation at <https://osf.io/h82aj/>. This was done *before* distributing the sample with instructions to the #fincap participants. We follow the PAP throughout unless stated otherwise.

precision. A more intriguing source of NSEs is the different routes that researchers travel through the “garden of forking paths.” [Gelman and Loken \(2014\)](#) use this metaphor to describe the non-trivial set of choices that researchers have to make when generating evidence. For example, they have to specify an appropriate econometric model, they have to pre-process the sample to ready it for estimation (e.g., purge outliers), they have to pick a programming language, et cetera. Error is therefore to be understood in the sense of erratic rather than erroneous.

Our objective is to measure and explain the size of non-standard errors. Are they large? And, how large are they compared to standard errors?² The four questions that we focus on are the following:

1. How large are non-standard errors for research in finance?
2. Can they be “explained” in the cross-section of researchers?
Are they smaller
 - (a) for papers by higher quality teams?
 - (b) for papers with more easily reproducible results?
 - (c) for papers that score higher in peer evaluations?
3. Does peer feedback reduce these non-standard errors?
4. Are researchers accurately aware of the size of non-standard errors?

Answering these questions is extremely costly in terms of human resources. The core structure of an ideal experiment involves two sizeable sets of representative researchers. A first set of researchers independently tests the same set of hypotheses on the same sample and writes a short paper presenting the results. A second, non-overlapping set of researchers obtains these papers, evaluates them, and provides feedback.

We believe that #fincap is close to this “ideal experiment” for three main reasons. First, Deutsche Börse supported the project by offering exclusive access to 720 million trade records spanning 17

²Note that the relative size of NSEs in total uncertainty is likely to grow in view of the trend of ever larger datasets. NSEs are invariant to sample size whereas SEs decline with sample size.

years of trading in Europe’s most actively traded instrument: the EuroStoxx 50 index futures. It enables researchers to test important hypotheses on how the market changes when migrating to super-human speeds. Second, #fincap being the first crowd-sourced empirical paper in finance might have pushed the hesitant few over the line (in addition to us arguing in the invitation that this is an opportunity to forego future regret).³ 164 research teams and 34 peer evaluators participated in #fincap. A back-of-the-envelope calculation shows that this effort alone cost 27 full-time equivalent person years ($164 \times 2 \text{ months} + 34 \times 2 \text{ days} \approx 27 \text{ years}$). Third, we believe that #fincap participants did, in fact, exert serious effort for a variety of reasons, such as, strong on-time delivery statistics and a high average rating of papers (more in Section 3.3).

Summary of our findings. First off, we show that the group of #fincap participants is representative of the academic community that studies empirical finance/liquidity. The hypotheses to be tested on the Deutsche Börse sample are tailored to this group. To distinguish them from the hypotheses we test in our (meta) study, we henceforth refer to them as RT-hypotheses. About a third of the 164 research teams (RTs) have at least one member with publications in the top-three finance or the top-five economics journals.⁴ For the group of peer evaluators (PEs), this share is 85%. 52% of RTs consist of at least one associate or full professor. For PEs, this is 88%. On a scale from 1 (low) to 10, the average self-ranked score on experience with empirical-finance research is 8.1 for RTs and 8.4 for PEs. For experience with market-liquidity research, it is 6.9 for RTs and 7.8 for PEs.

The dispersion in results across research teams is sizeable. All six hypotheses had to be tested by proposing a measure and computing the average per-year percentage change. The first RT-hypothesis, for example, was “Market efficiency has not changed over time.” The across-RT dispersion in estimates is enormous, but is mostly due to

³The #fincap was presented to all involved by means of a dedicated website (<https://fincap.academy>) and a short video (<https://youtu.be/HPtnus0Yu-o>).

⁴Economics: *American Economic Review*, *Econometrica*, *Journal of Political Economy*, *Quarterly Journal of Economics*, and *Review of Economic Studies*. Finance: *Journal of Finance*, *Journal of Financial Economics*, and *Review of Financial Studies*.

the presence of extreme values. One RT reported a 74,491% increase, which explains why the simple across-RT mean and standard deviation (SD) are 446.3% and 5,817.5%, respectively. If, however, following a standard practice in finance,⁵ outliers are treated by winsorizing the sample using 2.5 and 97.5 percentiles, the mean and SD become -7.5% and 20.6%, respectively.

An average 7.5% per-year decline in efficiency from 2002 to 2018 for a leading exchange, is remarkable. The really remarkable result given *our* purposes, however, is the across-RT dispersion. An SD of 20.6% for a *winsorized* sample is striking. But, how does this non-standard error compare to the standard error? The across-RT mean SE for the winsorized sample is 13.2%. The NSE-SE ratio therefore is 1.6. The NSE, therefore, is non-trivial in a relative sense as well. The NSE-SE ratio ranges between 0.6 and 2.1 across RT-hypotheses, including some that were deliberately non-fuzzy (e.g., one on the share of client volume in total volume). Not surprisingly, the non-fuzzy ones tend to show lower NSE-SE ratios.

We further find that the size of non-standard errors is hard to explain in the cross-section of researchers. We formally test three hypotheses that relate team quality, work-flow quality, and paper quality to dispersion, or more specifically, the size of “errors.” These errors are defined as the difference between a particular RT result, and the across-RT mean result. A standard approach to model heteroskedasticity of errors regresses log squared error on explanatory factors (Harvey, 1976). Although almost all coefficients are negative consistent with a high-quality small-error association, none of them are statistically significant.⁶ The only exceptions are team quality and work-flow quality in the 2.5%-97.5% winsorized sample. The latter is proxied by verifying how easy it is to reproduce the results of a particular RT with the code that the RT provided. For work-flow quality, we find a statistically significant negative sign. For team quality we find

⁵Adams et al. (2019), for example, review how outliers are dealt with for all ordinary least-squares (OLS) regressions in articles that appeared in premier finance journals in the period from 2008 through 2017: *Journal of Finance*, *Journal of Financial Economics*, *Review of Financial Studies*, and *Journal of Financial and Quantitative Analysis*. We like to note that outlier treatment was not pre-registered.

⁶We use the conservative significance levels advocated by Benjamin et al. (2018): 0.5% for significance and 5% for weak significance. The latter is referred to as “suggestive evidence.”

only a suggestive (weakly significant) negative sign. A one standard-deviation increase in work-flow quality reduces non-standard error by 12%. For team quality, a one SD increase reduces it by 8%.

Peer feedback significantly reduces non-standard errors. The peer-feedback process involves multiple stages. By and large, we find that each stage reduces NSEs. The overall reduction across all four stages is 8.5% for the raw sample, and 53.5% for the 2.5%-97.5% winsorized sample. This stark difference is driven by extreme-result teams largely staying put in the feedback process.

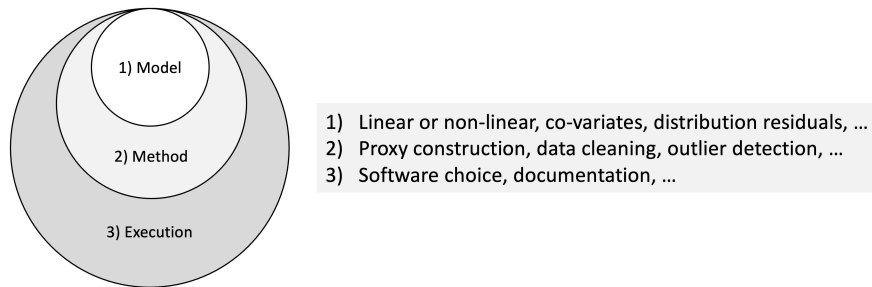
RTs mostly underestimate non-standard errors as tested in an incentivized belief survey. The extent to which the average team underestimates the realized across-RT standard deviation ranges between 9.0% and 99.5% across the six hypotheses. It is not merely due to the outliers that they did not foresee as the vast majority also underestimated the NSE for the 2.5%-97.5% winsorized sample. Such underestimation might well be the reason why non-standard errors never attracted much attention, until recently.

These findings on dispersion in percentage-change estimates carry over to the associated t -values. In a sense, t -values represent the statistical strength of a positive, or a negative percentage-change finding. The non-standard error of t -values therefore measures the extent of agreement across RTs on statistical strength. Our finding of non-trivial NSEs also for t -values therefore implies large dispersion in findings of statistical strength. Most telling, perhaps, is that for the first RT-hypothesis on efficiency, 23.8% of RTs find a significant decline, 8.5% find a significant increase, and 67.7% find no significance if their t -values are evaluated at a conventional significance level (in finance) of 5%. We find a similar pattern for all other RT-hypotheses.

Finally, we like to point out that the sizeable non-standard errors are *not* driven by the presence of poor-quality results. This could already be concluded from the very weak traction obtained from quality variables in the error-size regressions (discussed above). To further corroborate this point, we reduce the sample to research teams who score high on *all* quality variables. Doing so yields a sample of only nine out of 164 teams. Non-standard errors for this subsample remain large, including implied t -values that for some RTs are above 1.96 and for others are below -1.96 for the *same* RT-hypothesis.

Contribution to the literature. The issue of variability in the research process is not new.⁷ The extant literature is best discussed by adding some structure. The following chart depicts the various ways in which the existing literature touches on non-standard errors.

Non-standard error sourced from all stages of EGP



1. *Model.* One source of variability across researchers is model specification. Is the empirical model linear or non-linear? Which co-variates should be included? Et cetera. Economists are cognizant of specification errors and they typically do *robustness analysis* to show that their main results do not critically depend on specification choices. They have even endogenized the awareness specification error in their theoretical models. Agents in these models are assumed to make decisions in full awareness of potential model-specification errors (e.g., [Hansen and Sargent, 2001](#), for economics or [Garlappi, Uppal, and Wang, 2006](#), for finance).
2. *Method.* An additional source of variability is in the empirical method. What are reasonable empirical proxies for the parameters of interest? What filters should be applied to the sample (e.g., how to deal with outliers)? What are appropriate test statistics? It has long been known that not only is this a source of variability, it can actually produce misleading inference. In finance, [Lo and MacKinlay \(1990\)](#) were the first to warn of “data snooping” in which properties of the data are used to construct test statistics. More recently, [Mitton \(2021\)](#) provides di-

⁷[Leamer \(1983\)](#), for example, was troubled by the “fumes which leak from our computing centers.” He called for studying “fragility in a much more systematic way.”

rect evidence on substantial methodological variation by analyzing 604 empirical corporate-finance studies published in the top-three finance journals (mentioned in footnote 4).

3. *Execution.* Variability remains, *even* if one fixes model and method. This is the domain of computational reproducibility studies that go far back in economics and finance. In economics, the *American Economic Review* (AER) was the first to introduce a *Data Availability Policy* in 2005 after the AER had published two studies that illustrated how hard it was to reproduce empirical studies (Dewald, Thursby, and Anderson, 1986; McCullough and Vinod, 2003). Reproducibility issues persisted as, *at most*, half of the empirical studies published in the top economics journals could be computationally reproduced (Glandon, 2011; Chang and Li, 2017; Gertler, Galiani, and Romero, 2018). Reproducibility has become an issue for finance journals as well as exemplified by the *Journal of Finance* retracting a published article for the first time in its history (Rampini, Viswanathan, and Vuillemeys, 2021).

The presence of sizeable non-standard errors potentially corrupts the research process through a phenomenon known as “*p*-hacking.” The *p*-value of a statistical test refers to the probability (*p*) that the reported effect is solely due to chance (under the null hypothesis of there being no effect). Large non-standard errors create opportunities for researchers to pick a model, a method, and an execution that produces a low *p*-value and write a report accordingly (i.e., they are hacking the *p*-value).

Replication studies typically find much weaker effects and less statistical strength,⁸ suggestive of *p*-hacking. We caution that these results are not necessarily due to ill intent on the side of researchers, because they might be demand- rather than supply-driven. That is, journals might selectively publish papers with low *p*-values. Munafò

⁸See, for example, Open Science Collaboration (2015) and Camerer et al. (2016, 2018). Replication is an issue in finance as well (Harvey, 2017) and replication studies started to appear recently. Hou, Xue, and Zhang (2018), for example, study 452 asset-pricing anomalies, including “the bulk of published anomalies in finance and accounting” (p. 2024). They conclude that most anomalies fail to replicate. Black et al. (2021) try to replicate the results of four top finance/accounting/economics publications. They were not able to reproduce any of them.

et al. (2017) survey the various threats to credible empirical science and propose several fixes.

Our study on the magnitude of non-standard errors, therefore, not only measures the additional uncertainty in empirical results, it also reveals the scope for p -hacking. #fincap was carefully designed to minimize any p -hacking dynamics so as to obtain unbiased estimates of NSEs. For example, researchers signed contracts that clarified that their results would remain anonymous in all public communication about #fincap.

We are the first to study how large these errors are in finance, but not the first in science. Silberzahn et al. (2018) pioneered the multi-analyst study by letting multiple research teams test whether soccer referees are more likely to draw red cards for players with darker skin color. Other examples are Botvinik-Nezer et al. (2020) for neuroscience, Huntington-Klein et al. (2021) for economics, and Breznau et al. (2021) and Schweinsberg et al. (2021) for sociology. Our contribution relative to these studies is that we not only document the size of NSEs,⁹ we also try to explain them in the cross-section, and to study convergence in a peer-feedback process. Furthermore, we ascertain whether research teams themselves are accurately aware of the size of NSEs. A further strength of our study is the large sample size: $N=164$. It is at least twice and up to 20 times larger than other multi-analyst samples.

The remainder of the paper is organized as follows. Section 2 provides an in-depth discussion of the project design.¹⁰ It also presents the hypotheses that we will test and how we test them. Section 3 presents our results and Section 4 concludes.

2 Project design and hypotheses

This section starts with presenting the details of #fincap. With this in mind, we then translate our overall objectives into testable hypotheses.

⁹We like to note that documenting the size of NSEs has been done before, but, to the best of our knowledge, only in Huntington-Klein et al. (2021). Their sample size is seven RTs, ours is 164 RTs.

¹⁰The design of #fincap follows the guidelines for multi-analyst studies proposed by Aczel et al. (2021).

2.1 Project design

The core of #fincap is letting multiple research teams (RTs) independently test the same six hypotheses on the same Deutsche Börse sample. We refer to these hypotheses as RT-hypotheses and to this sample as RT-sample. This is to avoid potential confusion with the hypotheses that *we* will test based on sample generated by RTs and peer evaluators (PEs).¹¹

The RT-sample is a plain-vanilla trade dataset for the EuroStoxx 50 index futures with, added to it, an agency/principal flag.¹² For each buy and sell, we therefore know whether the exchange member traded for its own account or for a client. The sample runs from 2002 through 2018 and contains 720 million trade records. These index futures are among the world’s most actively traded index derivatives. They give investors exposure to Europe, or, more precisely, to a basket of euro-area blue-chip equities. With the exception of over-the-counter activity, all trading is done through an electronic limit-order book (see, e.g., [Parlour and Seppi, 2008](#)).

The RT-hypotheses are all statements about trends in the following market characteristics (with the null being no change):

RT-H1 market efficiency,

RT-H2 the realized bid-ask spread,

RT-H3 the share of client volume in total volume,

RT-H4 the realized spread on client orders,

RT-H5 the share of market orders in all client orders, and

RT-H6 the gross trading revenue of clients.

Appendix [A](#) discusses these RT-hypotheses in detail. For the purpose of our analysis, it suffices to know that these market

¹¹RTs and PEs have been recruited mostly by alerting appropriate candidates through suitable channels (e.g., the <https://microstructure.exchange/>). To inform them about #fincap, we created an online repository: <https://fincap.academy>. The repository remains largely unaltered (except for, e.g., adding FAQs).

¹²Trade records contain the following fields: Datetime, expiration, buy-sell indicator, size, price, aggressor flag, principal-agent flag, and a full- or partial-execution flag. More details on the sample are in Figure [OA.1](#) of the Online Appendix.

characteristics should not surprise a researcher familiar with empirical-finance/liquidity. There is, purposefully, considerable variation across RT-hypotheses in the level of abstraction. Or, in other words, the garden of forking paths is much larger for some than for others. RT-H1, for example, is on the relatively abstract notion of market efficiency. RT-H3, on the other hand, is on the market share of client volume. Such share should be a relatively straightforward to calculate because, in the RT-sample, each buy and sell trade is flagged client or proprietary.

RTs are asked to test these RT-hypotheses by estimating an average yearly change for a self-proposed measure.¹³ In the interest of readability, we refer to this estimate as the “effect size.” They are further asked to report standard errors for these estimates. We compute the ratio of the two, which we refer to as the implied *t*-value, or *t*-value for short. Collectively, we refer to these results as RT results.

RTs write a short academic paper in which they present and discuss their findings. These papers are evaluated by PEs who were recruited outside the set of researchers who registered as RTs. RT papers were randomly and evenly assigned to PEs in such a way that each paper was evaluated twice, and each PE evaluated nine or ten papers. PEs scored the papers in single-blinded process: PEs saw the names of RTs but not vice versa. This was made clear to all ex-ante.¹⁴ PEs scored the papers both at the RT-hypothesis level, and at the aggregate, paper level. They motivated their scores in a feedback form and were encouraged to add constructive feedback. RTs received this feedback unabridged and were allowed to update their results based on it. Importantly, the design of #fincap was common knowledge to all as it was communicated ex-ante via a dedicated website (see footnote 3). More specifically, #fincap consisted of the following four stages:

Stage 1 (Jan 11 - Mar 23, 2021.) RTs receive the detailed instructions along with access to the RT-sample. They conduct their

¹³RTs are asked to express their results in annualized terms. To some it was not clear. We therefore notified everyone of the following clarification that we added to the FAQ section on <https://fincap.academy>: “Research teams are asked to report annualized effect size estimates (and the corresponding standard errors); research teams are not required, however, to consider only annualized data.”

¹⁴We picked a single-blinded process instead of double-blinded process in order to incentivize RTs to exercise maximum effort.

analysis and hand in their results (short paper plus code). We emphasized in our emails and on the project website that RTs should work in *absolute secrecy* so as to ensure independence across RTs.

Stage 2 (May 10 - May 28, 2021.) RTs receive feedback from two anonymous PEs and are allowed to update their analysis based on it. They are asked to report their findings in the same way they did in Stage 1.

Stage 3 (May 31 - June 18, 2021.) RTs receive the five best papers based on the average raw PE score. The names of the authors of these five papers were removed before distributing the papers.¹⁵ Similar to Stage 2, all RTs are allowed to update their analysis and resubmit their results.

Stage 4 (June 20 - June 28, 2021.) RTs report their final results, this time not constrained by delivering code that produces them. In other words, RTs are allowed to Bayesian update their results (i.e., effect sizes and standard errors) taking in all the information that has become available to them, in particular the five best papers. They could, for example, echo the results of one of these papers, simply because of an econometric approach that they believe is superior but that is beyond their capacity to code. This stage was added to remove all constraints and see how far the RT community can get in terms of reaching consensus.

The stages subsequent to the first one mimic the feedback researchers get from various interactions with peer researchers in the research process *before* a first journal submission (e.g., feedback from colleagues over lunch or at the water cooler, during seminars, or in the coffee breaks after conferences, but also hearing about and seeing similar papers endorsed by others). The dynamics at play in a refereeing process at a scientific journal are out of scope.¹⁶

¹⁵If two papers were tied in terms of their average score, then, following the pre-analysis plan, we picked the one that had highest reproducibility score provided by the Cascad. For more information on Cascad, see the statement of H2 in Section 2.2.

¹⁶Testing the dynamics in a refereeing process requires a different experiment that involves “publishing” papers, *including* the names of the authors. Note that

2.2 Hypotheses

The overall objectives outlined in the introduction were translated into a set of pre-registered hypotheses (see footnote 1 for pre-registration details). The hypotheses all center on across-RT dispersion. We use variance as our dispersion measure because it is additive and the *de facto* standard dispersion measure in statistics. Let us therefore define dispersion at stage t for RT-hypothesis j as:

$$\text{var}(y_{jt}) = \frac{1}{n-1} \sum (y_{ijt} - \bar{y}_j)^2, \quad (1)$$

where

- $i \in \{1, \dots, n\}$ indexes RTs,
- $j \in \{1, \dots, 6\}$ indexes RT-hypotheses,
- $t \in \{1, \dots, 4\}$ indexes group stages, and

where \bar{y}_j is the overall average for RT-hypothesis j . To study the extent of dispersion requires modeling variance, which is known in econometrics as modeling heteroskedasticity. The idea is that the size of “errors” are related to co-variates. In our application these errors are defined as:¹⁷

$$\hat{u}_{ijt} = y_{ijt} - \bar{y}_j. \quad (2)$$

With all the groundwork done, we can now translate our overall objectives into three sets of hypotheses. These hypotheses are tested for the dispersion in two types of RT-results:

we do reveal the best five papers (according to PEs) to all RTs in Stage 4, but the authors of these papers remain hidden. Our focus is narrowly on the pure findings and beliefs of the RTs, avoiding any possible corruption by “the publication game.” In other words, a refereeing system has two components. The first is anonymous feedback from peers/experts and the second is an incentive to revise according to this feedback in order for one’s paper to become published. We capture the first component by asking PEs to rate the paper, but not the second one.

¹⁷The term error is used in a statistical sense of how observed data differ from a population mean. It should *not* be confused with errors necessarily being mistakes. Jumping ahead a little, we will model heteroskedasticity following [Harvey \(1976\)](#) (Section 2.3). \hat{u}_{ijt} in (2) denotes the OLS residual that results from estimating [Harvey \(1976, Equation \(2\)\)](#) with only intercepts pertaining to the six dummies for the RT-hypotheses.

- The effect size (i.e., the point estimate).
- The t -value (i.e., the point estimate divided by the standard error).

We believe that both are of interest. The effect size is of intrinsic interest as it is about the outcome in the variable of interest. The t -value captures the statistical strength of the reported effect in the sense larger magnitudes are less likely to occur under the null of there being no effect. All hypotheses on the dispersion in effect sizes and t -values are stated as null hypotheses. The hypotheses tests will be two-sided.

The first set of three hypotheses focuses on whether error variance relates to various quality measures:

H1 Team quality does not explain the size of errors in Stage 1 results. Team quality is proxied by the largest common factor in various candidate proxies for team quality. We prefer an appropriately weighted average over simply adding all proxies to maximize statistical power in the regressions. More specifically, we define team quality as the first principal component of standardized series that include:¹⁸

- Top publications*: The team has at least one top-five publication in economics or top-three publication in finance (0/1) (see footnote 4).
- Expertise in the field*: Average of self-assessed experience in market liquidity and empirical finance (scale from 0 to 10).
- Experience with big data*: The team has worked with datasets that are at least as big as the Deutsche Börse dataset analyzed in #fincap (0/1).
- Academic seniority*: At least one team member holds an associate or a full professorship (0/1).
- Team size*: The team size attains its maximum of two members (0/1).

¹⁸An important advantage of a principal-component analysis (PCA) is that the weighting is data-driven, thus avoiding subjective weights. Note that even the five proxies that enter were picked ex-ante in the pre-analysis plan filed at OSF. The PCA results will be discussed in Section 3.2.1.

H2 Work flow quality does not co-vary with the size of errors in Stage 1 results. We proxy for work-flow quality with a reproducibility score per $RT \times RT$ -hypothesis. The score measures to what extent an RT result is reproducible from the RT code. The scoring was done by the Certification Agency for Scientific Code and Data ([Cascad](#)). Cascad is a non-profit certification agency created by academics with the support of the French National Science Foundation (CNRS) and a consortium of French research institutions. The goal of this agency is to provide researchers with an innovative tool allowing them to signal the reproducibility of their research (used by, for example, the *American Economic Review*).¹⁹

H3 Paper quality as judged by the average score of PEs (hypothesis level) does not explain the size of errors in the first submission (scale from 0 to 10). To remove a possible PE fixed effect we use demeaned PE scores in all of our analysis.

The second set of hypotheses are about convergence of RT results across the four stages, and about convergence from the first to the final stage to have a test with maximum power.²⁰

H4 The error variance does not change from first to second stage.

H5 The error variance does not change from second to third stage.

H6 The error variance does not change from third to fourth stage.

H7 The error variance does not change from first to fourth stage.

The final hypothesis focuses on RT *beliefs* about the dispersion in results across RTs.

H8 The average belief of RTs on the across-RT dispersion in results, is correct. The dispersion beliefs are solicited in terms of the standard deviation measure.

¹⁹Cascad rates reproducibility on a five-category scale: RRR (perfectly reproducible), RR (practically perfect), R (minor discrepancies), D (potentially serious discrepancies), and DD (serious discrepancies). For #fincap, Cascad converted their standard categorical rating to an equal-distance numeric one: RRR, RR, R, D, and DD become 100, 75, 50, 25, 0, respectively.

²⁰The decline across consecutive stages might be small, but after summing them the total decline might be sizeable. It is in this sense that the first-to-final stage is expected to be statistically most powerful.

2.3 Statistical approach

Heterogeneity in error variance is known as heteroskedasticity in econometrics. A standard textbook in econometrics, [Greene \(2007, Ch. 9.7\)](#), refers to [Harvey \(1976\)](#) as a common approach to modeling heteroskedasticity. We adapt this approach to our setting and follow Harvey's notation for ease of comparison. Error variance is modeled as:

$$\sigma_{ijt}^2 = \text{var}(u_{ijt}) = \exp(z'_{ijt}\alpha), \quad (3)$$

where u_{ijt} is the unobserved disturbance term, and z_{ijt} contains the set co-variables which, instead of an intercept, includes six dummies corresponding to the six RT-hypotheses to account for possible heteroskedasticity across hypotheses. The parameter vector α measures the marginal *relative* change in variance when the associated co-variate is increased by one unit. For example, if the coefficient for co-variate X is 0.1, then increasing X by one unit raises variance by ten percent. It raises standard deviation by approximately half that amount: five percent.²¹

Estimation. An attractive feature of the Harvey model is that it can be estimated with ordinary least squares (OLS). More specifically, the OLS model is:²²

$$\log \hat{u}_{ijt}^2 = z'_{ijt}\alpha + w_{ijt}. \quad (4)$$

Let $\tilde{\alpha}$ denote the resulting parameter estimate. Harvey shows that $\tilde{\alpha}$ estimates α consistently, except for the intercept coefficients. If,

²¹This follows directly from a first-order Taylor approximation of $f(x) = \sqrt{x}$ around μ : $f(x) \approx \sqrt{\mu} + \frac{1}{2\sqrt{\mu}}(x - \mu)$. If the predicted change in variance is ten percent, then it follows that $f(x) \approx \sqrt{\mu} + \frac{1}{2\sqrt{\mu}}0.10\mu = 1.05\sqrt{\mu}$.

²²It is tempting to drop the natural logarithm on the dependent variable and simply use squared errors as the dependent variable in an OLS regression. [Harvey \(1976, p.6\)](#) mentions three reasons for why his model is more attractive than this alternative one: "Firstly, the likelihood function is bounded and no problems arise due to estimated variances being negative or zero. Secondly, the error terms in the two-step equation [...] are (asymptotically) homoscedastic and so the estimated covariance matrix of the two-step estimator, $\tilde{\alpha}$, is consistent. Finally, the likelihood ratio test has a much simpler form in the multiplicative model."

however, w_{ijt} is assumed to be Gaussian, then a consistent estimator for the full α vector becomes (where the subscript k refers the k th element of the vector):

$$\hat{\alpha}_k = \begin{cases} \tilde{\alpha}_k & \text{for } k > 6, \\ \tilde{\alpha}_k - \psi\left(\frac{1}{2}\right) + \log\left(\frac{1}{2}\right) & \text{for } k \leq 6, \end{cases} \quad (5)$$

where $\psi(\nu)$ is the psi (digamma) function defined as $d \log \Gamma(\nu) / d\nu$ where $\Gamma(\nu)$ is the Gamma function with ν degrees of freedom. This follows from [Harvey \(1976, Equation \(6\)\)](#). The predicted variance for RT i 's error for RT-hypothesis j in stage t becomes:

$$\hat{\sigma}_{ijt}^2 = \exp\left(z'_{ijt}\hat{\alpha}\right), \quad (6)$$

Importantly, the non-intercept elements of the parameter vector α are estimated consistently, even if w_{ijt} is non-Gaussian. The tests of our hypotheses all pertain to non-intercept elements, and therefore remain valid in the absence of normality.

Statistical inference. All hypotheses are tested using standard errors that cluster residuals w_{ijt} by RTs.²³ The R-squared of the heteroskedasticity regression in (4) is a useful measure of how much of the total dispersion can be explained. More precisely, the R-squared captures how much of the variance of log squared errors is explained by the set of explanatory variables.

Balanced sample. We will implement our tests on a balanced sample. We therefore include only observations for RTs that participated in *all* of the four stages. This makes the results in all four stages comparable as they pertain to the same set of RTs. The balanced sample, however, is not that different from the unbalanced one as only four of the 168 RTs did not complete all rounds. The balanced sample therefore consists of 164 RTs.

²³The clustering is needed to account for possible non-zero correlation in residuals w_{ijt} per RT, across RT-hypotheses and stages. For example, if the error is large for a particular RT on a particular RT-hypothesis, then it is likely to be large for the other RT-hypotheses as well. We do not cluster on RT-hypotheses because the model in (4) includes fixed effects for RT-hypotheses, which, we believe, removes most of the commonality. We do not cluster on PEs as a possible PE fixed effect was removed by demeaning PE scores in all of our analysis.

Significance levels. The hypothesis tests are two-sided and tested at significance levels of 0.005 and 0.05. Following the recommendations of [Benjamin et al. \(2018\)](#), we refer to results with a p -value smaller than 0.005 as statistically significant. Results with a p -value smaller than 0.05 are referred to as suggestive (weakly significant) evidence.

3 Results

This section summarizes all our findings. The first subsection presents various summary statistics to familiarize with the sample. The second subsection presents and discusses the test results for our hypotheses.

3.1 Summary statistics

(Insert Table 1 about here.)

Table 1 summarizes our sample by means of three sets of statistics, organized in three panels. Panel (a) summarizes the qualities of the #fincap community. It consists of 164 research teams (RTs) and 34 peer evaluators (PEs). The maximum RT size is two members, which is the size of 79% of RTs.

The statistics testify to the high quality of the #fincap community. 31% of RTs have at least one top publication in finance or economics (see footnote 4 for the list of journals). For PEs this is 85%. Those who provide feedback are therefore better published, which probably mirrors reality as the feedback flow is most likely to come from more senior or successful scholars. This finding is echoed in the percentage of RTs that have at least one member who is tenured at the associate or full professor level. This fraction is 52% for RTs and 88% for PEs.

(Insert Figure 1 about here.)

Figure 1 illustrates how RT members and PEs cover the global academic-finance community reasonably well. RT members reside in 34 countries with most (51 out of 293) residing in the United States. PEs reside in 13 countries with, again, most (13 out of 34) residing in the United States. The stronger skew towards the United States (US) is not surprising given that the more senior, well-published finance scholars are predominantly affiliated with US universities.

Most RTs and PEs seem to have the appropriate background for testing the RT-hypotheses on the RT-sample. Their average self-reported scores on having experience in the field of empirical finance is 8.1 for RTs and 8.4 for PEs on a scale from 0 (low) to 10. For experience in market-liquidity research these average scores are 6.9 for RTs and 7.8 for PEs. There is considerable variation around these averages as the across-RT standard deviation for these scores range from 1.7 to 2.4. When it comes to working with the large RT-sample (720 million trade records), again, most RTs and PEs seem up to it. 52% of RTs have worked with samples of similar size or larger. For PEs, this percentage is 88%.

Panel (b) of Table 1 shows that the average quality of the RT analysis is solid, and the dispersion is large. The average Cascad reproducibility score is 64.5 on a scale from 0 (low) to 100 (see footnote 19). This is high when benchmarked against other studies on reproducibility surveyed in [Colliard, Hurlin, and Pérignon \(2021\)](#). The across-RT standard deviation is 43.7, which implies extreme variation across RTs (most code either reproduces perfectly or not at all). The paper-quality score of PEs show a similar pattern, albeit with considerably lower dispersion. The average score across RTs is 6.2 on a scale from 0 (low) to 10. The across-RT standard deviation is 2.0.

Panel (c) provides descriptive statistics on non-standard errors: the dispersion in results across RTs. It does so, by hypothesis, and by type of result: Estimate of the effect size, standard error, and t -value. Our focus is on dispersion across RTs, which is why we relegate a discussion of RT means to Appendix A. More specifically, this appendix discusses the RT-hypotheses in-depth and summarizes what RTs, as a group, seem to find with a focus on the across-RT mean instead of the across-RT standard deviation.

(Insert Figure 2 about here.)

Perhaps the most salient feature of the extensive Panel (c) is that there is substantial across-RT variation in all hypotheses and for all results. For RT-H1 on efficiency, for example, the across-RT mean annual change is 446.3% with an across-RT standard deviation (NSE) of 5,817.5%. These extraordinary numbers are intimately linked to an extremely large value of 74,491.1% for a particular RT. If, as is common in finance ([Adams et al., 2019](#)), the RT-sample is winsorized

at a 2.5%-97.5% level (which is the default level for the remainder of this text), then the mean becomes -7.4% with standard deviation of 20.6%. Note that such NSE of 20.6% is similar in magnitude to the mean reported SE of 13.2%. The NSE-SE ratio therefore is 1.6. The dispersion is not much lower for RT-H3, where the corresponding numbers are a mean change is -2.6% with an NSE of 1.4% and a mean SE of 1.3%, leading to an NSE-SE ratio of 1.3. This pattern emerges for all RT-hypotheses with NSE-SE ratios ranging from 0.6 to 2.1. Figure 2 illustrates the substantial across-RT dispersion with boxplots.

The panel further shows that t -values also exhibit sizeable dispersion across RTs. For RT-H1, for example, the mean and standard deviation for t -value across RTs is -3.6 and 28.4, respectively. For the winsorized sample, these values are -1.4 and 5.2, respectively. At conventional threshold levels for 95% significance, -1.96 and 1.96, almost a third of RTs find statistically significant results. Interestingly, there is no agreement in this group of RTs as the 32.3% decomposes into 23.8% who find a significant decline and 8.5% who find a significant increase. This pattern is common across RT-hypotheses. Even for the relatively straightforward RT-H3, 49.4% report an (implied) t -value that is larger than 1.96 in absolute value. There is more agreement on the sign of the significance, as 45.7% finds significant decline, whereas 3.7% finds a significant increase. The 3.7% is not a single RT, it corresponds to *six* RTs who find a significant increase whereas 75 RTs (45.7%) find a significant decline.

Overall, the summary statistics show that there is substantial dispersion across RTs at various levels. There is substantial dispersion in team quality (panel (a)), the quality of their analysis (panel (b)), and in the results they report (panel (c)). These findings are promising for the hypotheses tests that we will turn to next. The sizeable dispersion in results implies that NSEs are large and therefore worthy of study. The strong variation in both team quality and analysis quality creates the statistical power needed to test whether they are significant covariates for the heteroskedasticity in results (H1-3). Are errors larger in magnitude for lower quality teams or for teams with lower quality analysis? This is what will be tested in the next subsection.

3.2 Hypotheses tests

The results on the three sets of hypotheses are discussed in the next three subsections.

3.2.1 Co-variates for stage-1 dispersion (H1-3)

(Insert Table 2 about here.)

The first set of hypotheses aims to measure whether various quality variables significantly co-vary with the size of errors. The first hypothesis centers on RT quality, which we measure by picking the first principal component of five standardized proxies. Table 2 summarizes the principal component analysis (PCA). Panel (a) presents the simple correlation matrix that enters the PCA. Note that this matrix is the covariance matrix since the five proxies are standardized. It is reassuring that the lion share of these correlations are positive. The only exception is that big-data experience is negatively related to experience. Since both are positively correlated with the rest, this finding suggests that the sample could be cut in two subgroups: those who are relatively more experienced in the field and those who are relatively more experienced in big-data analysis.

Panel (b) and (c) document that the largest principal component explains 38.3% of all variance and loads positively on all variables. It loads strongest on publications and weakest on big-data but, importantly, it loads positively on all of them. This is the principal component that serves as the team-quality measure in hypothesis test. Noteworthy is that the second principal component picks up $23.6/(100-38.3)=38.2\%$ of the remaining variance. The strongest loadings are 0.79 on big-data and -0.55 on experience. This component neatly picks up the experience/big-data distinction in the group of RTs that we commented on earlier.

(Insert Table 3 about here.)

Table 3 summarizes our findings when regressing log-squared errors on the various quality co-variates. These regressions were done for the raw sample as pre-registered in the PAP. We, however, had not foreseen the types of extreme values that some RTs reported (as evidenced, for example, by one RT reporting an average per-year increase in efficiency of 74,491.1%). Following standard practice in

finance, we treat these “outliers” by winsorizing and trimming the sample with cut-off levels at the 1% and 99% quantile and at the 2.5% and 97.5% quantile (Adams et al., 2019). Winsorizing replaces values beyond these levels by the quantile level, whereas trimming simply removes them.

For effect-size estimates, the regression results show that log-squared errors seem unrelated to the various quality co-variates, for the most part. Most coefficients are negative, and thus consistent with higher quality being associated with smaller errors. Loosely speaking, higher quality is in the center of the distribution of RT results and lower quality in the periphery. There is, however, no statistical significance with two exceptions for the 2.5%-97.5% winsorized sample.

First, we find suggestive evidence that higher team quality coincides with smaller errors. A one standard-deviation (SD) increase in team quality coincides with a 16% drop in error-variance, which, by approximation,²⁴ implies an SD drop of $0.5 \times 16 = 8\%$.

Second, we find significant evidence that better reproducibility is associated with smaller errors. A one-SD increase in reproducibility coincides with an approximate 12% drop in SD. Phrased differently, it coincides with a 12% drop in NSE. The findings for dispersion in t -values are similar with, again, statistical significance only for reproducibility, with higher reproducibility being associated with smaller errors.

In summary, the evidence is such that the first three hypotheses cannot be firmly rejected for both effect-size estimates and t -values. In other words, team quality (H1), reproducibility (H2), and paper quality (H3) seem unrelated to the size of errors. The only exception is that for the 2.5%-97.5% winsorized sample we do reject the null of no effect for H1 and H2.

(Insert Figure 3 about here.)

The almost absent statistical significance is unlikely to be due to low statistical power. The sample contains almost one thousand observations: 164 RTs times six hypotheses yields 984 observations. The result suggests that NSEs are sizeable, *also* for high-quality teams and/or high-quality results. The large sample allows us to make this

²⁴This approximation is derived in footnote 21.

point in a more straightforward way: Pull out a sub-sample of RTs that score high on all quality measures. We note that this additional analysis was not pre-registered.

Figure 3 illustrates the result. Nine RTs score highest on all measures: all RT-quality proxies, reproducibility, and the average peer-evaluator rating of their paper.²⁵ The figure shows that NSEs remain large. For the effect-size estimates depicted in the top plot, the NSE for RT-H1 is 10.1% (reported on top of the plot) and for RT-H3 it is 0.4%. For the winsorized full sample, these values were 20.6% and 1.8%, respectively (see Table 1). The bottom plot shows that also for t -values the across-RT dispersion is large. It is so large, in fact, that for RT-H4, two RTs report values larger than 1.96 in magnitude, but with *opposite* signs. If interpreted conventionally, the one RT would conclude that the realized spread on client order increased significantly in the course of the sample, whereas the other RT would conclude it significantly decreased.

3.2.2 Convergence across stages? (H4-7)

(Insert Table 4 about here.)

Whereas initial dispersion is large and seemingly only weakly related to quality, the dispersion does decline with peer feedback. Table 4 shows that error variance declines significantly when comparing the first to the last stage.

(Insert Figure 4 about here.)

Panel (a) presents the results for the effect-size estimates. In the raw sample, error variance declines by 17% and the standard deviation, therefore, by approximately half that amount: 8.5%. Figure 4 illustrates this finding by showing box plots for all four stages. One salient pattern in the figure is that extreme values seem to stay put or move only marginally. Most of the convergence seems to be in the boxes that depict the interquartile range. This observation explains why for

²⁵For the binary variables, the conditioning is trivial. For non-binary variables, we picked 7.5 for the scales from zero to 10, and 75 for the scales from 0 to 100.

the winsorized and trimmed samples the convergence across stages is much stronger.²⁶

For the winsorized sample, the standard deviation across RT estimates decline by approximately 53.5%. The decline is distributed rather evenly across the three stages. After receiving written feedback on their paper by two PEs in the second stage, it declines by 14.5%. It declines by another 20% in the third stage after RTs see the best five papers. There is a final 19% reduction in the final stage when RTs hand in their final estimates, this time not constrained by having to hand in the code to support it. All these changes are statistically significant, except for first change where the significance is weak and the result is therefore only suggestive.

The pattern for t -values is somewhat different for this winsorized sample. The SD also declines significantly across all stages. The decline, however, is smaller in magnitude: 20.5% instead of 53.5%. The most salient difference is that the SD *increases* significantly by 14% from Stage 1 to Stage 2. The subsequent decline of 4% is insignificant. The real convergence is in the final stage where it drops by 30.5%. This maybe due to the nature of this final stage where RTs might simply be Bayesian updating with less than full weight on their own paper. This necessarily leads to convergence across RTs.

This evidence makes us reject most of the null hypotheses on convergence. H7 is most firmly rejected. For both effect-size estimates and t -values do we find rejection in all stages in favor of the alternative that error variance declines. The other three hypotheses that pertain to changes across the three consecutive stages, the evidence is weaker but largely rejecting the null of no changes.

3.2.3 Are RT-beliefs about dispersion accurate? (H8)

The eighth and final hypothesis is on whether RTs are accurately aware of the size of non-standard errors: the standard deviation of results across RTs. Beliefs were solicited in an incentivized way. If their beliefs turn out to be within 50% of the realized standard devi-

²⁶The winsorization and trimming was done per hypothesis, per stage. To retain a balanced panel, the trimming procedure removes an RT for all stages in case it was removed in a single stage. The balanced-panel condition is particularly important when studying convergence across stages, because one wants to compare a fixed set of RTs across stages.

ation, they can earn a monetary reward of \$300. The details of the reward scheme are in the instruction sheet they obtained prior to reporting their beliefs (see Figure OA.10 in the Online Appendix). The hypothesis pertains to Stage 1 results because they were solicited only for this stage.

As H8 is stated in terms of the average belief being correct, testing it requires a test on the equality of means: the mean belief about standard deviation in results across RTs and the standard deviation of these results in the population. Let us define the test statistic D that measures relative distance between beliefs and realizations:

$$D = \frac{1}{6n} \sum_{i,j} \left(\frac{BeliefOnStDev_{ij} - RealizationOfStDev_j}{RealizationOfStDev_j} \right), \quad (7)$$

where $BeliefOnStDev_{ij}$ is the belief of team i on the standard deviation across RTs for hypothesis j and $RealizationOfStDev_j$ is the realized standard deviation for this hypothesis in the raw sample.²⁷ The distribution of D under the null of equal means is obtained by bootstrapping. For details on the bootstrap procedure we refer to Appendix C.

(Insert Table 5 about here.)

Table 5 presents the test results which show that non-standard errors are severely underestimated. Effect-size estimates are significantly underestimated by 71.7% and t -values by a significant 70.6%. Redoing the test by RT-hypothesis shows that there is strong heterogeneity but, where significant, there is underestimation. For effect-size estimates, it seems that for the RT-hypotheses that are relatively straightforward to test, the null of accurate beliefs cannot be rejected. An example is RT-H3, which is on the market share of client volume. The average RT-belief for this RT-hypothesis is not significantly different from the realization. For RT-hypotheses that require more creative effort to test, NSEs are significantly underestimated. For RT-H1 on market efficiency, for example, the underestimation is 99.5%. In

²⁷The benefit of a relative measure as opposed to an absolute measure is that (i) it is easy to interpret as it allows for statements of RTs over- or underestimating by some percentage and (ii) it accounts for level differences across hypotheses (e.g., under the null accurate beliefs, a uniform distribution of beliefs on the support 0.09 to 0.11 will exhibit the same dispersion as a uniform distribution of beliefs on 900 to 1100).

summary, the vast majority of tests show significant underestimation and we therefore firmly reject the null of an accurate average belief.

(Insert Figure 5 about here.)

Figure 5 plots the entire distribution of reported beliefs along with the realized values that are depicted by red dots. It illustrates that the vast majority of RTs underestimated dispersion. The interquartile range denoted by the boxes is consistently below the red dot. One might think that RTs simply overlooked the extreme values that prop up SDs. This, however, does not seem to be true since even if one trims the RT-results (which is stronger than winsorizing in this case), the boxes stay below the realized value which, for the trimmed sample, are depicted by orange dots. The only exception is RT-H3 and effect size, for which the orange dot is just within the top of the box. 75% of the RTs therefore underestimate the NSEs for 11 out of the 12 cases.

3.3 Alternative explanations

After having presented all our results, it is useful to discuss alternative explanations. Might the sizeable non-standard errors be due to the presence of inexperienced researchers testing unsuitable hypotheses with little effort? We believe this is unlikely to be the case for the following reasons.

Experience. Aware of this potential pitfall, we selectively approached research teams and peer evaluators who we knew were sufficiently experienced in the field. When signing up, they ticked a box that they understood that participating in #fincap requires research expertise and experience in empirical finance and the analysis of large datasets. Ticking the box further meant that they acknowledge that one of the team members held a PhD in finance or economics. After ticking the box, researchers had to motivate in open text box why they believe they meet these requirements. We parsed the content of this box to make sure that the team qualifies before accepting them into #fincap (see Figure OA.2 in the Online Appendix for the sign-up sheet).

Hypotheses. We proceeded with care when designing RT-hypotheses. Early versions were shared with senior scholars, and their feedback helped us fine-tune the RT-hypotheses. We therefore feel comfortable that the RT-hypotheses are suitable and well motivated hypotheses to test with the RT-sample (see Figure OA.6 in the Online Appendix for the RT instruction sheet which shows how RT-hypotheses were presented to RTs).

Related to the suitability question, one might wonder whether vagueness of an RT-hypothesis might be a viable alternative explanation for sizeable NSEs. To address this concern, we included a very precise RT-hypothesis: RT-H3 on client volume share. The results for RT-H3 show that NSEs are sizeable for relatively precise hypotheses as well. It is true, however, that NSEs tend to be lower for the more precise RT-hypotheses.

Effort. We incentivized research teams to exert effort by informing them about the following before signing up: the deadlines of the various stages so that they could plan for it; their *non-anonymized* paper would be evaluated by senior peer reviewers; the top-five (anonymized) papers would be announced to all others;²⁸ and, only those who completed all stages would become co-authors. In addition to these incentives, we believe that most scientists are propelled by an intrinsic motivation to do good research.

Looking back, we have various reasons to believe that researchers did indeed exert serious effort. First, only four out of 168 research teams did not complete all stages. 123 out of 168 teams (73.2%) handed in their Stage 1 report at least a day early, and none of the teams seriously breached any deadline. The average reproducibility score was 64.5 on a scale from 0 (low) to 100, which is high in comparison to what has been reported in other reproducibility studies (Colliard, Hurlin, and Pérignon, 2021). Finally, the average paper quality was 6.2 on a scale from zero (low) to 10.

²⁸Individuals obtain “ego utility” from positive views about their ability to do well and they exert more effort (or take more risks) when they are informed about their rank in non-incentivized competitions (Köszegi, 2006; Tran and Zeckhauser, 2012; Kirchler, Lindner, and Weitzel, 2018).

4 Conclusion

Testing an hypothesis on a data sample involves many decisions on the side of the researcher. He needs to find an appropriate econometric model, clean the data, choose suitable software for estimation, et cetera. These choices inherently generate some dispersion in results across researchers. There is no right or wrong here, there simply are many roads that lead to Rome. We propose the across-researcher standard deviation in results as a measure of such dispersion and refer to it as the non-standard error.

Studying NSEs involves letting a representative group of researchers *independently* test the same set of hypotheses on the same sample. #fincap did exactly this. 164 research teams tested relatively standard types of hypotheses on novel and proprietary trade data.

We compute NSEs and show that they are sizeable, similar in magnitude to SEs. They, for the most part, do not co-vary with various quality variables: research-team quality, work-flow quality, and paper quality. The last one is assessed by 34 external peer evaluators. PEs also provided written feedback, so that each research team received anonymized feedback from two PEs. After receiving this feedback, teams were allowed to update their analysis and post new results. NSEs declined weakly significantly in this stage. They declined significantly further in two subsequent stages, in which teams received new information. The total decline across all stages is substantial, up to 50%. Finally, we find that teams grossly underestimate NSEs, and more so for more abstract hypotheses.

In sum, we believe that our study shows that NSEs add non-negligible uncertainty to the outcome of hypotheses tests in finance. We document that NSE uncertainty is similar in magnitude as SE uncertainty. This is particularly worrisome as it creates substantial space for *p*-hacking (which had been scoped out here). An encouraging result, however, is that peer feedback reduces NSEs. We therefore believe that the profession should further encourage researchers to make data available, to post a pre-registration plan, to consider a multi-analyst approach, and to stimulate interaction among researchers. This, hopefully, does not leave us lost in the larger metropolitan area, but leads us right to the Forum Romanum.

Appendices

A RT-sample, RT-hypotheses, and results

This appendix presents the RT-hypotheses in detail and the test results of #fincap RTs as a group. We start by providing a motivating context.

A.1 Context

Electronic order matching systems (automated exchanges) and electronic order generation systems (algorithms) have changed financial markets over time. Investors used to trade through broker-dealers by paying the dealers' quoted ask prices when buying, and accepting their bid prices when selling. The wedge between dealer bid and ask prices, the bid-ask spread, was a useful measure of trading cost, and often still is.

Now, investors more commonly trade in electronic limit-order markets (as is the case for EuroStoxx 50 futures). They still trade at bid and ask prices. They do so by submitting so-called market orders and marketable limit orders. However, investors now also can quote bid and ask prices themselves by submitting (non-marketable) standing limit orders. Increasingly, investors now also use agency algorithms to automate their trades. Concurrently, exchanges have been continuously upgrading their systems to better serve their clients. Has market quality improved, in particular when taking the viewpoint of non-exchange members: (end-user) clients?

A.2 RT-hypotheses and test results

The RT-hypotheses and results are discussed based on RT-results in the final stage of the project (Table OA.3 in the Online Appendix). We therefore base our discussion on the results that RTs settled on after receiving feedback. What do RTs find after having shown some convergence across the stages? And, consistent with the main text, the presence of extreme values makes us prefer to analyze the 2.5%-97.5% winsorized sample.

(The first two hypotheses focus on all trades.)

RT-H1. Assuming that informationally-efficient prices follow a random walk, did market efficiency change over time?

Null hypothesis: Market efficiency has not changed over time.

Findings. RTs predominantly reject the hypothesis in favor of a decline of market efficiency. In the final stage, 61.6% of RTs report an implied t -value that is larger than 1.96 in absolute value, and therefore significant at a conventional 5% significance level. Of these, $51.8/61.6 = 84.1\%$ report a significant decline. The decline

seems modest as the across-RT mean²⁹ is -1.7% per year. The small changes add up, though, to a total decline in the 2002-2018 sample of $(0.983^{17} - 1) = 25.3\%$. This might reflect a trend of declining depth in the market, possibly due to new regulation in the aftermath of the global financial crisis of 2007-2008. The regulation constrained the supply of liquidity by sell-side banks (e.g. [Bao, O'Hara, and Zhou, 2018](#); [Jovanovic and Menkveld, 2021](#)). If these banks incur higher inventory costs as a result, then, in equilibrium, one observes larger transitory price pressures thus reducing market efficiency (e.g., [Pastor and Stambaugh, 2003](#); [Hendershott and Menkveld, 2014](#)). In the interest of brevity, we discuss all remaining hypotheses in the same way.

RT-H2. Did the (realized) bid-ask spread paid on market orders change over time? The realized spread could be thought of as the gross-profit component of the spread as earned by the limit-order submitter.

Null hypothesis: The realized spread on market orders has not changed over time.

Findings. The majority of RTs (53.7%) find statistical significance, mostly (90.9%) in the direction of a decline in the realized spread. The decline is 2.1% per year, which implies a 30.3% decline over the full sample. This trend might be due to the arrival of high-frequency market makers who operate at low costs. They do not have the deep pockets that sell-side banks have, but they will offer liquidity for regular small trades by posting near the inside of the market. Their arrival is typically associated with a tighter bid-ask spread, but not necessarily with better liquidity supply for large orders (e.g., [Jones, 2013](#); [Angel, Harris, and Spatt, 2015](#); [Menkveld, 2016](#)).

(The remaining hypotheses focus on agency trades only.)

RT-H3. Did the share of client volume in total volume change over time?

Null hypothesis: Client share volume as a fraction of total volume has not changed over time.

Findings. Almost all RTs (86.6%) find statistical significance, almost exclusively (97.9%) pointing towards a decline in the share of client volume. The average decline is 2.7% per year, which implies a total decline of 37.2% for the full sample. Intermediation, therefore, seems to have increased which should surprise those who believe that the arrival of agency algorithms enables investors to execute optimally themselves, thus reducing the need for intermediation.³⁰

²⁹The across-RT mean includes all RTs, thus also those who report insignificant results.

³⁰We verified with Deutsche Börse that this change is not purely mechanical in the sense that, in the sample period, many institutions became an exchange member and, with it, their volume changes status from agency to principal.

RT-H4. On their market orders and marketable limit orders, did the realized bid-ask spread that clients paid, change over time?

Null hypothesis: Client realized spreads have not changed over time.

Findings. 32.3% of RTs would reject the null hypothesis, with the majority (71.8%) in favor of a decline in realized spread. The average decline is 0.7% per year and 11.3% for the full sample. The decline in client realized spread is therefore only about a third of the total realized spread decline, which suggests that market orders of intermediaries benefited most from a general realized-spread decline.

RT-H5. Realized spread is a standard cost measure for market orders, but to what extent do investors continue to use market and marketable limit orders (as opposed to non-marketable limit orders)?

Null hypothesis: The fraction of client trades executed via market orders and marketable limit orders has not changed over time.

Findings. 31.1% of RTs would reject the null, but this time, interestingly, about half find a significant decline, whereas the other half finds a significant increase. The average per-year change is -0.2% which adds up to -3.3% for the full sample. The results seem to suggest that clients neither increased their share of market orders, nor did they decrease it. One might have expected the latter because an increased use of agency algorithms should allow them to execute more through limit orders³¹ as opposed to market orders. The benefit of execution via limit order is that one earns half the bid-ask spread rather than pays for it.

RT-H6. A measure that does not rely on the classic limit- or market-order distinction is gross trading revenue (GTR). Investor GTR for a particular trading day can be computed by assuming a zero position at the start of the day and evaluating an end-of-day position at an appropriate reference price. Relative investor GTR can then be defined as this GTR divided by the investor's total (euro) volume for that trading day. This relative GTR is, in a sense, a realized spread. It reveals what various groups of market participants pay in aggregate for (or earn on) their trading. It transcends market structure as it can be meaningfully computed for any type of trading in any type of market (be it trading through limit-orders only, through market-orders only, through a mix of both, or in a completely different market structure).

Null hypothesis: Relative gross trading revenue (GTR) for clients has not changed over time.

Findings. Only 12.8% of RTs find significance with about half in favor of GTR decline and the other half in favor of an increase. The average change is a 1.5% increase, which for the full sample implies a total change of 28.8%. Since the

³¹By executing via a limit order we mean, submitting a limit order that cannot be executed immediately and, therefore, enters the book and eventually gets matched with an incoming market order.

average client GTR is most likely negative, this result implies that it became more negative thus suggesting worse overall execution quality. We caution that this result is very weak in terms of significance and in terms of agreement across RTs on the sign. If, however, overall execution quality did indeed decline, then we would come full circle with the RT-H1 finding of a market efficiency decline, potentially due to a decline in market depth.

B Explanatory variables for error variance

B.1 Team quality

The quality measures for research teams are based on the survey that participants filled out upon registration (see Figure OA.2 in the Online Appendix). To keep the regression model both concise and meaningful, we reduce the ordinal variable “current position” and the logarithmic interval-based variable “size of largest dataset worked with” to binary variables. The academic position variable is one if a researcher is either associate or full professor. The dataset variable is one if the researcher has worked with datasets that are contained at least 100 million observations, because the #fincap sample contains 720 million observations. We aggregate these binary variables to research team level by taking the maximum across the team members.

As for self-assessed experience, we asked for both empirical finance and market liquidity, which we deem equally relevant for testing the RT-hypotheses. Thus, and because of the anticipated high correlation, we use the average of these two measures to obtain the individual score. And, in the interest of consistency, we again aggregate to the team level by taking the maximum across the team members.

B.2 Workflow quality

We proxy for workflow quality with an objectively obtained score of code quality provided by Cascad (see footnote 19). The scale ranges from 0 (serious discrepancies) to 100 (perfect reproducibility).

B.3 Paper quality

Papers are rated by an external group of peer evaluators (PEs). They rate the analyses associated with each RT-hypothesis individually, but also the paper in its entirety (see Figure OA.11 in the Online Appendix). The ratings range from 0 (very weak) to 10 (excellent). Each paper is rated by two PEs and the paper rating is the average of the two (after removing a PE fixed effect as discussed in Section 2.1).

C Bootstrap procedure for belief statistic D

The distribution of D under the null of equal means is obtained by bootstrapping as follows. For each RT-hypothesis, we subtract the difference between the average belief on standard deviation and the observed standard deviation, from the beliefs:

$$AdjBeliefOnStDev_{ij} = BeliefOnStDev_{ij} - \left[\left(\frac{1}{n} \sum_i BeliefOnStDev_{ij} \right) - RealizationOfStDev_j \right] \quad (8)$$

In this new sample with adjusted beliefs, the average belief about dispersion equals the observed dispersion, by construction. This sample is input to the bootstrapping procedure which iterates through the following steps 10,000 times:

1. As we have n RTs, in each iteration we draw n times from the new sample, with replacement. Each draw picks a particular RT and stores its beliefs and its results for all of the six RT-hypotheses. The result of these n draws therefore is a simulated sample that has the same size as the original sample.
2. The simulated sample is used to compute the test statistic D in (7). This statistic for iteration k , a scalar, is stored as D_k .

The bootstrap procedure yields 10,000 observations of the test statistic under the null. For a significance level of 0.005, the statistic observed in the #fincap sample is statistically significant if it lands below the 25th lowest simulated statistic or above the 25th highest simulated statistic. Its p -value is:³²

$$2 \min(EmpiricalQuantileFincapStatistic, 1 - EmpiricalQuantileFincapStatistic). \quad (9)$$

References

Aczel, Balazs, Barnabas Szaszi, Gustav Nilsson, Olmo R. van den Akker, Casper J. Albers, Marcel A.L.M. van Assen, Jojanneke A. Bastiaansen, Dan Benjamin, Udo Boehm, Rotem Botvinik-Nezer, Laura F. Bringmann, Niko A. Busch, Emmanuel Caruyer, Andrea M. Cataldo, Nelson Cowan, Andrew Delios, Noah N.N. van Dongen, Chris Donkin, Johnny B. van Doorn, Anna Dreber, Gilles Dutilh, Gary F. Egan, Morton Ann Gernsbacher, Rink Hoekstra, Sabine Hoffmann, Felix Holzmeister, Juergen Huber, Magnus Johannesson, Kai J. Jonas, Alexander T. Kindel, Michael Kirchler, Yoram K. Kunkels,

³²Note that the procedure accounts for within-RT correlations (i.e., including possible non-zero correlations among a particular RT's results and the beliefs that it reports). The reason the procedure accounts for these correlations is that the bootstrap uses block-sampling where, when an RT is drawn, all of its beliefs and all of its estimates are drawn. One therefore only assumes independence across RTs which holds by construction given the design of #fincap.

- D. Stephen Lindsay, Jean-Francois Mangin, Dora Matzke, Marcus R Munafò, Ben R. Newell, Brian A. Nosek, Russell A Poldrack, Don van Ravenzwaaij, Jörg Rieskamp, Matthew J. Salganik, Alexandra Sarafoglou, Tom Schonberg, Martin Schweinsberg, David Shanks, Raphael Silberzahn, Daniel J. Simons, Barbara A. Spellman, Samuel St-Jean, Jeffrey J. Starns, Eric L. Uhlmann, Jelte Wicherts, and Eric-Jan Wagenmakers. 2021. “Consensus-Based Guidance for Conducting and Reporting Multi-Analyst Studies.” *eLife* (forthcoming) .
- Adams, John, Darren Hayunga, Sattar Mansi, David Reeb, and Vincenzo Verardi. 2019. “Identifying and Treating Outliers in Finance.” *Financial Management* 48:345–384.
- Angel, James J., Lawrence E. Harris, and Chester S. Spatt. 2015. “Equity Trading in the 21st Century: An Update.” *Quarterly Journal of Finance* 5:1–39.
- Bao, Jack, Maureen O’Hara, and Xing (Alex) Zhou. 2018. “The Volcker Rule and Corporate Bond Market Making in Times of Stress.” *Journal of Financial Economics* 130:95–113.
- Benjamin, Daniel J., James O. Berger, Magnus Johannesson, Brian A. Nosek, E.-J. Wagenmakers, Richard Berk, Kenneth A. Bollen, Björn Brembs, Lawrence Brown, Colin Camerer, David Cesarini, Christopher D. Chambers, Merlise Clyde, Thomas D. Cook, Paul De Boeck, Zoltan Dienes, Anna Dreber, Kenny Easwaran, Charles Efferson, Ernst Fehr, Fiona Fidler, Andy P. Field, Malcolm Forster, Edward I. George, Richard Gonzalez, Steven Goodman, Edwin Green, Donald P. Green, Anthony G. Greenwald, Jarrod D. Hadfield, Larry V. Hedges, Leonhard Held, Teck Hua Ho, Herbert Hoijtink, Daniel J. Hruschka, Kosuke Imai, Guido Imbens, John P.A. Ioannidis, Minjeong Jeon, James Holland Jones, Michael Kirchler, David Laibson, John List, Roderick Little, Arthur Lupia, Edouard Machery, Scott E. Maxwell, Michael McCarthy, Don A. Moore, Stephen L. Morgan, Marcus Munafó, Shinichi Nakagawa, Brendan Nyhan, Timothy H. Parker, Luis Pericchi, Marco Perugini, Jeff Rouder, Judith Rousseau, Victoria Savalei, Felix D. Schönbrodt, Thomas Sellke, Betsy Sinclair, Dustin Tingley, Trisha Van Zandt, Simine Vazire, Duncan J. Watts, Christopher Winship, Robert L. Wolpert, Yu Xie, Cristobal Young, Jonathan Zinman, and Valen E. Johnson. 2018. “Redefine Statistical Significance.” *Nature Human Behavior* 2:6–10.
- Black, Bernard S., Hemang Desai, Kate Litvak, Woongsun Yoo, and Jeff Jiewei Yu. 2021. “Specification Choice in Randomized and Natural Experiments: Lessons from the Regulation SHO Experiment.” Manuscript, Northwestern University.
- Botvinik-Nezer, Rotem, Felix Holzmeister, Colin F. Camerer, Anna Dreber, Juer-gen Huber, Magnus Johannesson, Michael Kirchler, Roni Iwanir, Jeanette A. Mumford, R. Alison Adcock, Paolo Avesani, Blazej M. Baczowski, Aahana Bajracharya, Leah Bakst, Sheryl Ball, Marco Barilari, Nadège Bault, Derek Beaton, Julia Beitner, Roland G. Benoit, Ruud M.W.J. Berkers, Jamil P. Bhanji,

Bharat B. Biswal, Sebastian Bobadilla-Suarez, Tiago Bortolin, Katherine L. Bottenhorn, Alexander Bowring, Senne Braem, Hayley R. Brooks, Emily G. Brudner, Cristian B. Calderon, Julia A. Camilleri, Jaime J. Castellon, Luca Cecchetti, Edna C. Cieslik, Zachary J. Cole, Olivier Collignon, Robert W. Cox, William A. Cunningham, Stefan Czoschke, Kamalaker Dadi, Charles P. Davis, Alberto De Lucas, Mauricio R. Delgado, Lysia Demetriou, Jeffrey B. Dennison, Xin Di, Erin W. Dickie, Ekaterina Dobryakova, Claire L. Donnat, Juer-gen Dukart, Niall W. Duncan, Joke Durnez, Amr Eed, Simon B. Eickhoff, Andrew Erhart, Laura Fontanesi, G. Matthew Fricke, Shiguang Fu, Adriana Galván, Remi Gau, Sarah Genon, Tristan Glatard, Enrico Glerean, Jelle J. Goeman, Sergej A.E. Golowin, Carlos González-García, Krzysztof J. Gorgolewski, Cheryl L. Grady, Mikella A. Green, Joao F. Guassi Moreira, Olivia Guest, Shabnam Hakimi, J. Paul Hamilton, Roeland Hancock, Giacomo Handjaras, Bronson B. Harry, Colin Hawco, Peer Herholz, Gabrielle Herman, Stephan Heunis, Felix Hoffstaedter, Jeremy Hogeveen, Susan Holmes, Chuan-Peng Hu, Scott A. Huettel, Matthew E. Hughes, Vittorio Iacovella, Alexandru D. Iordan, Peder M. Isager, Ayse I. Isik, Andrew Jahn, Matthew R. Johnson, Tom Johnstone, Michael J.E. Joseph, Anthony C. Juliano, Joseph W. Kable, Michalis Kassinos, Cemal Koba, Xiang-Zhen Kong, Timothy R. Koscik, Nuri Erkut Kucukboyaci, Brice A. Kuhl, Sebastian Kupek, Angela R. Laird, Claus Lamm, Robert Langner, Nina Lauharatanahirun, Hongmi Lee, Sangil Lee, Alexander Leemans, Andrea Leo, Elise Lesage, Flora Li, Monica Y.C. Li, Phui Cheng Lim, Evan N. Lintz, Schuyler W. Liphardt, Annabel B. Losecaat Vermeer, Bradley C. Love, Michael L. Mack, Norberto Malpica, Theo Marins, Camille Maumet, Kelsey McDonald, Joseph T. McGuire, Helena Melero, Adriana S. Méndez Leal, Benjamin Meyer, Kristin N. Meyer, Glad Mihai, Georgios D. Mitsis, Jorge Moll, Dylan M. Nielson, Gustav Nilsson, Michael P. Notter, Emanuele Olivetti, Adrian I. Onicas, Paolo Papale, Kaustubh R. Patil, Jonathan E. Peelle, Alexandre Pérez, Doris Pischke, Jean-Baptiste Poline, Yanina Prystauka, Shruti Ray, Patricia A. Reuter-Lorenz, Richard C. Reynolds, Emiliano Ricciardi, Jenny R. Rieck, Anais M. Rodriguez-Thompson, Anthony Romyn, Taylor Salo, Gregory R. Samanez-Larkin, Emilio Sanz-Morales, Margaret L. Schlichting, Douglas H. Schultz, Qiang Shen, Margaret A. Sheridan, Jennifer A. Silvers, Kenny Skagerlund, Alec Smith, David V. Smith, Peter Sokol-Hessner, Simon R. Steinkamp, Sarah M. Tashjian, Bertrand Thirion, John N. Thorp, Gustav Tinghög, Loreen Tisdall, Steven H. Thompson, Claudio Toro-Serey, Juan Jesus Torre Tresols, Leonardo Tozzi, Vuong Truong, Luca Turella, Anna van der Veer, Tom Verguts, Jean M. Vettel, Sagana Vijayarajah, Khoi Vo, Matthew B. Wall, Wouter D. Weeda, Susanne Weis, David J. White, David Wisniewski, Alba Xifra-Porxas, Emily A. Yearling, Sangsuk Yoon, Rui Yuan, Kenneth S. L. Yuen, Lei Zhang, Xu Zhang, Joshua E. Zosky, Thomas E. Nichols, Russell A. Poldrack, and Tom Schonberg. 2020. "Variability in the Analysis of a Single Neuroimaging Dataset by Many Teams." *Nature* 582:84–88.

Breznau, Nate, Eike Mark Rinke, Alexander Wuttke, Muna Adem, Jule Adriaans, Amalia Alvarez-Benjumea, Henrik K. Andersen, Daniel Auer, Flavio Azevedo, Oke Bahnsen, Dave Balzer, Gerrit Bauer, Paul C. Bauer, Markus Baumann,

Sharon Baute, Verena Benoit, Julian Bernauer, Carl Berning, Anna Berthold, Felix S. Bethke, Thomas Biegert, Katharina Blinzler, Johannes N. Blumenberg, Licia Bobzien, Andrea Bohman, Thijs Bol, Amie Bostic, Zuzanna Brzozowska, Katharina Burgdorf, Kaspar Burger, Kathrin Busch, Juan Carlos-Castillo, Nathan Chan, Pablo Christmann, Roxanne Connelly, Christian S. Czymara, Elena Damian, Alejandro Ecker, Achim Edelmann, Maureen A. Eger, Simon Ellerbrock, Anna Forke, Andrea Forster, Chris Gaasendam, Konstantin Gavras, Vernon Gayle, Theresa Gessler, Timo Gnambs, Amélie Godefroidt, Max Grömping, Martin Groß, Stefan Gruber, Tobias Gummer, Andreas Hajar, Jan Paul Heisig, Sebastian Hellmeier, Stefanie Heyne, Magdalena Hirsch, Mikael Hjerm, Oshrat Hochman, Andreas Hövermann, Sophia Hunger, Christian Hunkler, Nora Huth, Zsófia S. Ignácz, Laura Jacobs, Jannes Jacobsen, Bastian Jaeger, Sebastian Jungkunz, Nils Jungmann, Mathias Kauff, Manuel Kleinert, Julia Klinger, Jan-Philipp Kolb, Marta Kołczyńska, John Kuk, Katharina Kunißen, Dafina Kurti Sinatra, Alexander Langenkamp, Philipp M. Lersch, Lea-Maria Löbel, Philipp Lutscher, Matthias Mader, Joan E. Madia, Natalia Malancu, Luis Maldonado, Helge-Johannes Marahrens, Nicole Martin, Paul Martinez, Jochen Mayerl, Oscar J. Mayorga, Patricia McManus, Kyle McWagner, Cecil Meeusen, Daniel Meierrieks, Jonathan Mellon, Friedolin Merhout, Samuel Merk, Daniel Meyer, Leticia Micheli, Jonathan Mijs, Cristóbal Moya, Marcel Neunhoffer, Daniel Nüst, Olav Nygård, Fabian Ochsenfeld, Gunnar Otte, Anna Pechenkina, Christopher Prosser, Louis Raes, Kevin Ralston, Miguel Ramos, Arne Roets, Jonathan Rogers, Guido Ropers, Robin Samuel, Gregor Sand, Ariela Schachter, Merlin Schaeffer, David Schieferdecker, Elmar Schlueter, Regine Schmidt, Katja M. Schmidt, Alexander Schmidt-Catran, Claudia Schmiedeberg, Jürgen Schneider, Martijn Schoonvelde, Julia Schulte-Cloos, Sandy Schumann, Reinhard Schunck, Jürgen Schupp, Julian Seuring, Henning Silber, Willem Slegers, Nico Sonntag, Alexander Staudt, Nadia Steiber, Nils Steiner, Sebastian Sternberg, Dieter Stiers, Dragana Stojmenovska, Nora Storz, Erich Striessnig, Anne-Kathrin Stroppe, Janna Teltemann, Andrey Tibajev, Brian Tung, Giacomo Vagni, Jasper Van Assche, Meta van der Linden, Jolanda van der Noll, Arno Van Hootegeem, Stefan Vogtenhuber, Bogdan Voicu, Fieke Wagemans, Nadja Wehl, Hannah Werner, Brenton M. Wiernik, Fabian Winter, Christof Wolf, Yuki Yamada, Nan Zhang, Conrad Ziller, Stefan Zins, Tomasz Żółtak, and Hung H.V. Nguyen. 2021. “Observing Many Researchers Using the Same Data and Hypothesis Reveals a Hidden Universe of Uncertainty.” Manuscript, University of Bremen.

Camerer, Colin F., Anna Dreber, Eskil Forsell, Teck-Hua Ho, Jürgen Huber, Magnus Johannesson, Michael Kirchler, Johan Almenberg, Adam Altmeld, Taizan Chan, Emma Heikensten, Felix Holzmeister, Taisuke Imai, Siri Isaksson, Gideon Nave, Thomas Pfeiffer, Michael Razen, and Hang Wu. 2016. “Evaluating Replicability of Laboratory Experiments in Economics.” *Science* 351:1433–1436.

Camerer, Colin F., Anna Dreber, Felix Holzmeister, Teck-Hua Ho, Jürgen Huber, Magnus Johannesson, Michael Kirchler, Gideon Nave, Brian A. Nosek, Thomas Pfeiffer, Adam Altmeld, Nick Buttrick, Taizan Chan, Yiling Chen, Eskil Forsell,

- Anup Gampa, Emma Heikensten, Lily Hummer, Taisuke Imai, Siri Isaksson, Dylan Manfredi, Julia Rose, Eric-Jan Wagenmakers, and Hang Wu. 2018. "Evaluating the Replicability of Social Science Experiments in Nature and Science." *Nature Human Behaviour* 2:637–644.
- Chang, Andrew C. and Phillip Li. 2017. "A Preanalysis Plan to Replicate Sixty Economics Research Papers That Worked Half of the Time." *American Economic Review: Papers & Proceedings* 107:60–64.
- Colliard, Jean-Edouard, Christophe Hurlin, and Christophe Pérignon. 2021. "The Economics of Research Reproducibility." Manuscript, HEC Paris.
- Dewald, William G., Jerry G. Thursby, and Richard G. Anderson. 1986. "Replication in Empirical Economics: The Journal of Money, Credit and Banking Project." *American Economic Review* 76:587–603.
- Garlappi, Lorenzo, Raman Uppal, and Tan Wang. 2006. "Portfolio Selection with Parameter and Model Uncertainty: A Multi-Prior Approach." *Review of Financial Studies* 20:41–81.
- Gelman, Andrew and Eric Loken. 2014. "The Statistical Crisis in Science." *American Scientist* 102:460–465.
- Gertler, Paul, Sebastian Galiani, and Mauricio Romero. 2018. "How to Make Replication the Norm." *Nature* 554:417–419.
- Glandon, Philip J. 2011. "Appendix to the Report of the Editor: Report on the American Economic Review Data Availability Compliance Project." *American Economic Review: Papers & Proceedings* 101:695–699.
- Greene, W.H. 2007. *Econometric Analysis*. London: Prentice Hall.
- Hansen, Lars Peter and Thomas J. Sargent. 2001. "Robust Control and Model Uncertainty." *AEA Papers and Proceedings* 91:60–66.
- Harvey, Andrew C. 1976. "Estimating Regression Models with Multiplicative Heteroscedasticity." *Econometrica* 44:461–465.
- Harvey, Campbell R. 2017. "Presidential Address: The Scientific Outlook in Financial Economics." *Journal of Finance* 72:1399–1440.
- Hendershott, Terrence and Albert J. Menkveld. 2014. "Price Pressures." *Journal of Financial Economics* 114:405–423.
- Hou, Kewei, Chen Xue, and Lu Zhang. 2018. "Replicating Anomalies." *Review of Financial Studies* 33:2019–2133.

- Huntington-Klein, Nick, Andreu Arenas, Emily Beam, Marco Bertoni, Jeffrey R. Bloem, Pralhad Burli, , Naibin Chen, Paul Grieco, Godwin Ekpe, Todd Pugatch, Martin Saavedra, and Yani Stopnitzky. 2021. "The Influence of Hidden Researcher Decisions in Applied Microeconomics." *Economic Inquiry* 59:944–960.
- Jones, Charles M. 2013. "What Do We Know About High-Frequency Trading?" Manuscript, Columbia University.
- Jovanovic, Boyan and Albert J. Menkveld. 2021. "Equilibrium Bid-Price Dispersion." *Journal of Political Economy* (forthcoming) .
- Kirchler, Michael, Florian Lindner, and Utz Weitzel. 2018. "Rankings and Risk-Taking in the Finance Industry." *Journal of Finance* 73:2271–2302.
- Köszegi, Botond. 2006. "Ego Utility, Overconfidence, and Task Choice." *Journal of the European Economic Association* 4:673–707.
- Leamer, Edward E. 1983. "Let's Take the Con Out of Econometrics." *American Economic Review* 73:31–43.
- Lo, Andrew W. and A. Craig MacKinlay. 1990. "Data-Snooping Biases in Tests of Financial Asset Pricing Models." *Review of Financial Studies* 3:431–467.
- Open Science Collaboration. 2015. "Estimating the Reproducibility of Psychological Science." *Science* 349 (6251).
- McCullough, Bruce D. and Hrishikesh D. Vinod. 2003. "Verifying the Solution from a Nonlinear Solver: A Case Study." *American Economic Review* 93:873–892.
- Menkveld, Albert J. 2016. "The Economics of High-Frequency Trading: Taking Stock." *Annual Review of Financial Economics* 8:1–24.
- Mitton, Todd. 2021. "Methodological Variation in Empirical Corporate Finance." *Review of Financial Studies* (forthcoming) .
- Munafò, Marcus R., Brian A. Nosek, Dorothy V.M. Bishop, Katherine S. Button, Christopher D. Chambers, Nathalie Percie du Sert, Uri Simonsohn, Eric-Jan Wagenmakers, Jennifer J. Ware, and John P.A. Ioannidis. 2017. "A Manifesto for Reproducible Science." *Nature Human Behaviour* 1:1–9.
- Parlour, Christine A. and Duane J. Seppi. 2008. "Limit Order Markets: A Survey." In *Handbook of Financial Intermediation and Banking*, edited by Arnoud W.A. Boot and Anjan V. Thakor. Amsterdam, Netherlands: Elsevier Publishing.
- Pastor, L. and R.F. Stambaugh. 2003. "Liquidity Risk and Expected Returns." *Journal of Political Economy* 111:642–685.
- Rampini, Adriano A., S. Viswanathan, and Guillaume Vuillemy. 2021. "Retracted: Risk Management in Financial Institutions." *Journal of Finance* 76:2709–2709.

Schweinsberg, Martin, Michael Feldman, Nicola Staub, Olmo R. van den Akker, Robbie C.M. van Aert, Marcel A.L.M. van Assen, Yang Liu, Tim Althoff, Jeffrey Heer, Alex Kale, Zainab Mohamed, Hashem Amireh, Vaishali Venkatesh Prasad, Abraham Bernstein, Emily Robinson, Kaisa Snellman, S. Amy Sommer, Sarah M.G. Otner, David Robinson, Nikhil Madan, Raphael Silberzahn, Pavel Goldstein, Warren Tierney, Toshio Murase, Benjamin Mandl, Domenico Viganola, Carolin Strobl, Catherine B.C. Schaumans, Stijn Kelchtermans, Chan Naseeb, S. Mason Garrison, Tal Yarkoni, C.S. Richard Chan, Prestone Adie, Paulius Alaburda, Casper Albers, Sara Alspaugh, Jeff Alstott, Andrew A. Nelson, Eduardo Ariño de la Rubia, Adbi Arzi, Štěpán Bahník, Jason Baik, Laura Winther Balling, Sachin Banker, David AA Baranger, Dale J. Barr, Brenda Barros-Rivera, Matt Bauer, Enuh Blaise, Lisa Boelen, Katerina Bohle Carbonell, Robert A. Briers, Oliver Burkhard, Miguel-Angel Canela, Laura Castillo, Timothy Catlett, Olivia Chen, Michael Clark, Brent Cohn, Alex Coppock, Natàlia Cugueró-Escofet, Paul G. Curran, Wilson Cyrus-Lai, David Dai, Giulio Valentino Dalla Riva, Henrik Danielsson, Rosaria de F.S.M. Russo, Niko de Silva, Curdin Derungs, Frank Dondelinger, Carolina Duarte de Souza, B. Tyson Dube, Marina Dubova, Ben Mark Dunn, Peter Adriaan Edelsbrunner, Sara Finley, Nick Fox, Timo Gnambs, Yuanyuan Gong, Erin Grand, Brandon Greenawalt, Dan Han, Paul H.P. Hanel, Antony B. Hong, David Hood, Justin Hsueh, Lilian Huang, Kent N. Hui, Keith A. Hultman, Azka Javaid, Lily Ji Jiang, Jonathan Jong, Jash Kamdar, David Kane, Gregor Kappler, Erikson Kaszubowski, Christopher M. Kavanagh, Madian Khabisa, Bennett Kleinberg, Jens Kouros, Heather Krause, Angelos-Miltiadis Kryptos, Dejan Lavbič, Rui Ling Lee, Timothy Leffel, Wei Yang Lim, Silvia Liverani, Bianca Loh, Dorte Lønsmann, Jia Wei Low, Alton Lu, Kyle MacDonald, Christopher R. Madan, Lasse Hjorth Madsen, Christina Maimone, Alexandra Mangold, Adrienne Marshall, Helena Ester Matskewich, Kimia Mavon, Katherine L. McLain, Amelia A. McNamara, Mhairi McNeill, Ulf Mertens, David Miller, Ben Moore, Andrew Moore, Eric Nantz, Ziauddin Nasrullah, Valentina Nejkovic, Colleen S. Nell, Andrew Arthur Nelson, Gustav Nilsson, Rory Nolan, Christopher E. O'Brien, Patrick O'Neill, Kieran O'Shea, Toto Olita, Jahna Otterbacher, Diana Palsea, Bianca Pereira, Ivan Pozdniakov, John Protzko, Jean-Nicolas Rey, Travis Riddle, Amal (Akmal) Ridhwan Omar Ali, Ivan Ropovik, Joshua M. Rosenberg, Stephane Rothen, Michael Schulte-Mecklenbeck, Nirek Sharma, Gordon Shotwell, Martin Skarzynski, William Stedden, Victoria Stodden, Martin A. Stoffel, Scott Stoltzman, Subashini Subbaiah, Rachael Tatman, Paul H. Thibodeau, Sabina Tomkins, Ana Valdivia, Gerrieke B. Druiff-van de Woestijne, Laura Viana, Florence Villesèche, W. Duncan Wadsworth, Florian Wanders, Krista Watts, Jason D Wells, Christopher E. Whelpley, Andy Won, Lawrence Wu, Arthur Yip, Casey Youngflesh, Ju-Chi Yu, Arash Zandian, Leilei Zhang, Chava Zibman, and Eric Luis Uhlmann. 2021. "Same Data, Different Conclusions: Radical Dispersion in Empirical Results when Independent Analysts Operationalize and Test the Same Hypothesis." *Organizational Behavior and Human Decision Processes* 165:228–249.

Silberzahn, Raphael, Eric L. Uhlmann, Dan P. Martin, Pasquale Anselmi, Frederik

Aust, Eli Awtrey, Štěpán Bahník, Feng Bai, Colin Bannard, Evelina Bonnier, Rickard Carlsson, Felix Cheung, Garret Christensen, Russ Clay, Maureen A. Craig, Anna Dalla Rosa, Lammertjan Dam, Mathew H. Evans, Ismael Flores Cervantes, Nathan Fong, Monica Gamez-Djokic, Andreas Glenz, Shauna Gordon-McKeon, Tim J. Heaton, Karin Hederos, Moritz Heene, Alicia J. Hofelich Mohr, Fabia Högden, Kent Hui, Magnus Johannesson, Jonathan Kalodimos, Erikson Kaszubowski, Deanna M. Kennedy, Ryan Lei, Thomas A. Lindsay, Silvia Liverani, Christopher R. Madan, Daniel Molden, Eric Molleman, Richard D. Morey, Laetitia B. Mulder, Bernard R. Nijstad, Nolan G. Pope, Bryson Pope, Jason M. Prenoveau, Floor Rink, Egidio Robusto, Hadiya Roderique, Anna Sandberg, Elmar Schlüter, Felix D. Schönbrodt, Martin F. Sherman, S. Amy Sommer, Kristin Sotak, Seth Spain, Christoph Spörlein, Tom Stafford, Luca Stefanutti, Susanne Tauber, Johannes Ullrich, Michelangelo Vianello, Eric-Jan Wagenmakers, Maciej Witkowiak, Sangsuk Yoon, and Brian A. Nosek. 2018. “Many Analysts, One Data Set: Making Transparent How Variations in Analytic Choices Affect Results.” *Advances in Methods and Practices in Psychological Science* 1:337–356.

Tran, Anh and Richard Zeckhauser. 2012. “Rank as an Inherent Incentive: Evidence from a Field Experiment.” *Journal of Public Economics* 96:645–650.

Table 1: Summary statistics

This table presents summary statistics. Standard deviations are in parentheses.

Panel (a): Quality of the #fincap community

	Research teams	Peer evaluators
Fraction with top econ/finance publications (see footnote 4)	0.31	0.85
Fraction including at least associate/full professor	0.52	0.88
Experience empirical-finance research (low-high, 1-10)	8.1 (1.7)	8.4 (1.8)
Experience market-liquidity research (low-high, 1-10)	6.9 (2.4)	7.8 (2.3)
Relevant experience (average of the above two items)	7.5 (1.3)	8.1 (1.7)
Fraction with "big data" experience (>#fincap sample)	0.65	0.88
Fraction teams consisting of two members (maximum team size)	0.79	
Number of observations	164	34

Panel (b): Quality of the analysis of research teams

	Research teams
Reproducibility score according to Cascad (low-high, 0-100)	64.5 (43.7)
Paper quality as judged by peer evaluators (low-high, 0-10)	6.2 (2.0)

(continued on next page)

(continued from previous page)

Panel (c): Dispersion across teams of estimates, SEs, and *t*-values

	RT-H1 Efficiency	RT-H2 RSpread	RT-H3 Client Volume	RT-H4 Client RSpread	RT-H5 Client MOrders	RT-H6 Client GTR
<i>Estimate effect size</i>						
Mean	446.3	-1,093.4	-3.5	-38,276.1	-3.5	-87.1
Mean (wins.) ^a	-7.4	9.0	-2.6	-2.1	-0.3	-27.6
Mean (trim) ^b	-6.2	5.4	-2.7	0.6	-0.2	-22.1
SD	5,817.5	14,537.2	9.4	490,024.2	37.6	728.5
SD (wins.) ^a	20.6	49.5	1.8	46.8	1.1	187.8
SD (trim) ^b	16.6	29.0	1.5	24.9	0.9	128.7
Min	-171.1	-186,074.5	-117.5	-6,275,383.0	-452.9	-8,254.5
Q(0.10)	-23.7	-6.9	-3.8	-6.7	-1.6	-192.1
Q(0.25)	-6.2	-3.6	-3.5	-2.1	-0.6	-18.2
Median	-1.1	-0.0	-3.3	0.1	-0.0	0.0
Q(0.75)	0.5	3.9	-2.4	3.8	0.2	3.2
Q(0.90)	3.7	21.5	-0.1	20.4	1.0	56.5
Max	74,491.1	4,124.0	8.7	870.2	69.5	1,119.0
<i>Standard error</i>						
Mean	468.7	1,195.3	3.7	38,302.0	6.2	148.2
Mean (wins.) ^a	13.2	23.5	1.4	26.9	1.7	104.8
Mean (trim) ^b	10.8	16.4	1.3	18.8	1.4	80.7
SD	5,810.6	14,711.9	29.5	489,929.5	40.1	526.0
SD (wins.) ^a	27.0	60.8	1.6	75.3	2.6	237.5
SD (trim) ^b	21.2	38.8	1.3	53.1	1.7	174.1
Min	0.0	0.0	0.0	0.0	0.0	0.0
Q(0.10)	0.1	0.2	0.1	0.2	0.1	0.0
Q(0.25)	0.5	1.1	0.3	1.2	0.2	0.7
Median	2.5	5.0	1.4	4.4	1.0	9.7
Q(0.75)	9.3	13.9	2.0	14.3	2.4	77.1
Q(0.90)	44.7	39.6	2.2	31.2	3.1	235.4
Max	74,425.5	188,404.1	378.8	6,274,203.0	463.7	4,836.2
<i>t-value</i>						
Mean	-3.6	35.3	-47.1	24.3	-5.7	-2.0
Mean (wins.) ^a	-1.4	-1.3	-13.3	-0.4	0.2	-0.4
Mean (trim) ^b	-1.1	-1.0	-11.6	-0.3	0.0	-0.2
SD	28.4	541.2	269.9	406.0	60.1	21.2
SD (wins.) ^a	5.2	4.1	25.4	3.0	4.2	2.0
SD (trim) ^b	3.4	3.3	22.0	2.4	2.6	1.4
Min	-322.3	-764.6	-2,770.6	-852.6	-631.6	-191.7
Q(0.10)	-4.7	-5.7	-37.4	-3.5	-2.3	-1.7
Q(0.25)	-1.9	-1.5	-11.5	-1.0	-0.6	-1.0
Median	-0.7	-0.1	-1.8	0.1	0.0	0.0
Q(0.75)	0.3	0.8	-1.6	1.0	0.8	0.7
Q(0.90)	1.7	1.5	-0.3	1.6	1.7	1.2
Max	51.6	6,880.5	29.5	5,119.5	89.6	100.6
<i>More t-value statistics</i>						
<i>t</i> < -1.96	23.8%	22.6%	45.7%	17.7%	11.0%	9.8%
<i>t</i> > 1.96	8.5%	6.1%	3.7%	9.1%	9.8%	3.0%
<i>t</i> > 1.96	32.3%	28.7%	49.4%	26.8%	20.7%	12.8%
<i>Relative size NSE (wins.)^c</i>						
NSE/SE ratio	1.6	2.1	1.3	1.7	0.6	1.8

^a: Winsorized at 2.5%-97.5%. ^b: Trimmed at 2.5%-97.5%. ^c: The non-standard error of effect size is compared to the mean standard error of effect size for the winsorized sample.

Table 2: Principal component analysis team quality

This table presents the results of a principal component analysis of the standardized team quality variables.

Panel (a): Correlation team quality measures

	Publications	Experience	Big Data	Position	#Members
Publications		0.34	0.10	0.54	0.30
Experience			-0.18	0.25	0.12
Big Data				0.14	0.14
Position					0.16

Panel (b): Fraction of variance explained

	PC1	PC2	PC3	PC4	PC5
Variance explained	38.3%	23.6%	17.1%	12.4%	8.6%

Panel (c): Loading of principal components on variables

	Publications	Experience	Big Data	Position	#Members
PC1	0.61	0.40	0.13	0.55	0.37
PC2	-0.01	-0.55	0.79	0.05	0.26
PC3	-0.10	0.06	-0.21	-0.46	0.86
PC4	-0.20	0.71	0.56	-0.35	-0.12
PC5	-0.76	0.14	-0.02	0.60	0.22

Table 3: Stage-1 error-variance regressions

This table presents the results of stage-1 error-variance regressions. The standard errors in parentheses are based on clustering on research teams. The incremental R^2 measure how much more is explained beyond a model with simple dummies for RT-hypotheses. */** correspond to significance at the 5/0.5% level, respectively.

Panel (a): Estimates

	Raw	Winsorized		Trimmed	
		1% to 99%	2.5% to 97.5%	1% to 99%	2.5% to 97.5%
Team quality (standardized)	0.01 (0.06)	-0.13 (0.07)	-0.16* (0.08)	-0.12 (0.06)	-0.12 (0.08)
Reproducibility score (standardized)	-0.13 (0.07)	-0.15 (0.09)	-0.24** (0.08)	-0.13 (0.09)	-0.08 (0.08)
Average rating (standardized)	-0.12 (0.07)	-0.14 (0.10)	-0.18 (0.09)	-0.17 (0.11)	-0.17 (0.09)
RT-hypotheses dummies	Yes	Yes	Yes	Yes	Yes
R^2	0.95	0.73	0.65	0.72	0.63
#Observations	984	984	984	972	936

Panel (b): t -values

	Raw	Winsorized		Trimmed	
		1% to 99%	2.5% to 97.5%	1% to 99%	2.5% to 97.5%
Team quality (standardized)	-0.00 (0.08)	-0.01 (0.09)	0.00 (0.11)	-0.05 (0.09)	0.05 (0.10)
Reproducibility score (standardized)	-0.12 (0.09)	-0.18 (0.09)	-0.23* (0.11)	-0.21* (0.08)	-0.09 (0.10)
Average rating (standardized)	-0.15 (0.12)	-0.21 (0.13)	-0.09 (0.12)	-0.01 (0.10)	0.02 (0.11)
RT-hypotheses dummies	Yes	Yes	Yes	Yes	Yes
R^2	0.71	0.49	0.37	0.49	0.32
#Observations	984	984	984	972	936

Table 4: All-stages error-variance regressions

This table presents the results of error-variance regressions including all stages. The standard errors in parentheses are based on clustering on research teams. */** correspond to significance at the 5/0.5% level, respectively.

Panel (a): Estimates

	Raw	Winsorized		Trimmed	
		1% to 99%	2.5% to 97.5%	1% to 99%	2.5% to 97.5%
Dummy Stage 2 - Dummy Stage 1	-0.11* (0.05)	-0.24* (0.09)	-0.29* (0.11)	-0.30** (0.10)	-0.21* (0.11)
Dummy Stage 3 - Dummy Stage 2	0.04 (0.03)	-0.26** (0.07)	-0.40** (0.09)	-0.21* (0.09)	-0.40** (0.09)
Dummy Stage 4 - Dummy Stage 3	-0.10** (0.03)	-0.26** (0.05)	-0.38** (0.06)	-0.25** (0.05)	-0.39** (0.06)
Dummy Stage 4 - Dummy Stage 1	-0.17* (0.06)	-0.75** (0.10)	-1.07** (0.12)	-0.76** (0.10)	-1.00** (0.12)
RT-hypotheses dummies	Yes	Yes	Yes	Yes	Yes
R ²	0.96	0.58	0.51	0.54	0.50
#Observations	3,936	3,936	3,936	3,888	3,744

Panel (b): *t*-values

	Raw	Winsorized		Trimmed	
		1% to 99%	2.5% to 97.5%	1% to 99%	2.5% to 97.5%
Dummy Stage 2 - Dummy Stage 1	-0.01 (0.07)	0.04 (0.09)	0.28** (0.09)	0.10 (0.08)	0.24* (0.08)
Dummy Stage 3 - Dummy Stage 2	-0.20** (0.06)	-0.23* (0.09)	-0.08 (0.08)	-0.14 (0.09)	0.14 (0.07)
Dummy Stage 4 - Dummy Stage 3	-0.06 (0.03)	-0.32** (0.05)	-0.61** (0.08)	-0.43** (0.07)	-0.73** (0.09)
Dummy Stage 4 - Dummy Stage 1	-0.26** (0.09)	-0.51** (0.11)	-0.41** (0.12)	-0.46** (0.11)	-0.36** (0.12)
RT-hypotheses dummies	Yes	Yes	Yes	Yes	Yes
R ²	0.76	0.57	0.42	0.47	0.39
#Observations	3,936	3,936	3,936	3,888	3,744

Table 5: Dispersion on research team beliefs

This presents test statistics on whether the beliefs of research teams about dispersion across research teams matches realized dispersion. More precisely, the test statistic is defined as the difference between the average belief on the standard deviation across teams minus the realized standard deviation, divided by the latter. The p-values in parentheses are obtained through bootstrapping. This analysis is the only one that uses the unwinsorized sample as beliefs were solicited for the raw data. * / ** correspond to significance at the 5/0.5% level, respectively.

	RT-H1 Efficiency	RT-H2 RSpread	RT-H3 Client Volume	RT-H4 Client RSpread	RT-H5 Client MOrders	RT-H6 Client GTR	All
Estimate	-99.5%** (0.00)	-95.4%** (0.00)	-9.0% (0.64)	-97.5%** (0.00)	-45.3% (0.50)	-83.3%** (0.00)	-71.7%** (0.00)
t-value	15.9% (0.19)	-96.0%** (0.00)	-92.4%** (0.00)	-97.1%** (0.00)	-86.0%** (0.00)	-68.2%** (0.00)	-70.6%** (0.00)

Figure 1: Countries of #fincap community

This plot illustrates how the #fincap community is dispersed around the globe. The top plot depicts how the members of the research teams are dispersed across countries. The bottom plot does the same for the peer evaluators.

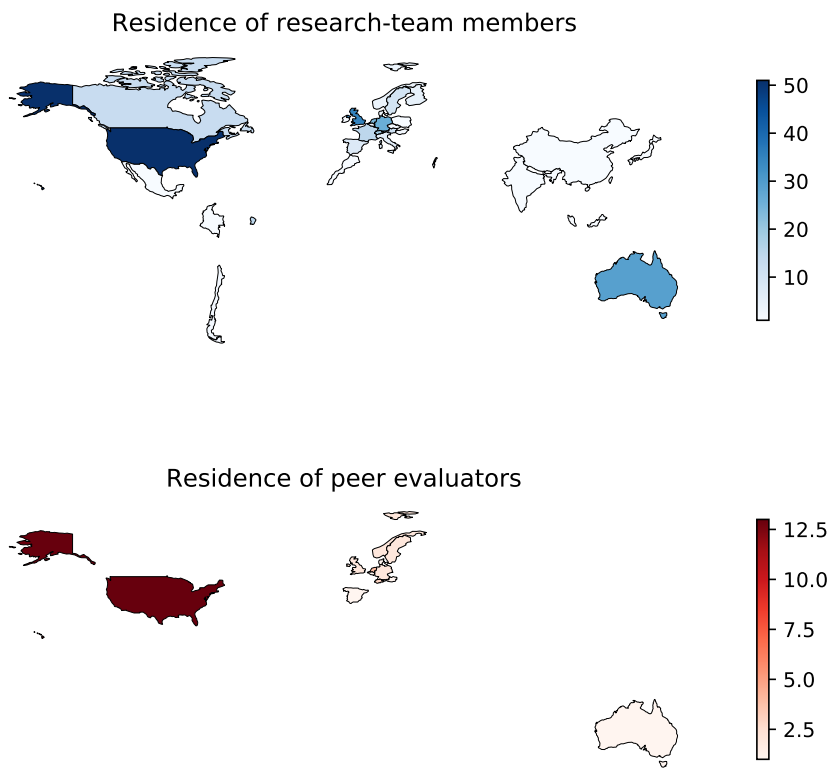


Figure 2: Dispersion in research-team estimates

This plot illustrates the dispersion in the stage-1 estimates reported by the research teams for all six hypotheses. It is based on the raw sample. The boxes span the first to the third quartile with the median as the interior horizontal line. The whiskers span 95% of the observations, starting from the 2.5% quantile to the 97.5% quantile.

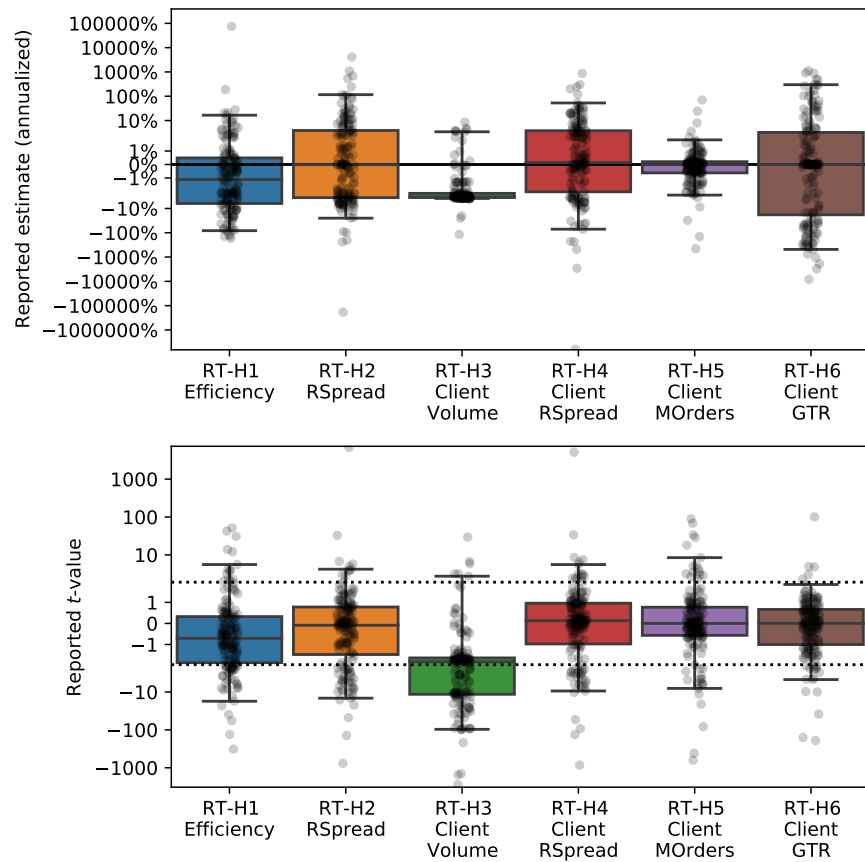


Figure 3: Dispersion research-team estimates: A high-quality subset

This plot mirrors Figure 2 but reports dispersion only for the subset of the highest-quality research teams with highest-quality results (N=9). The teams score one on *all* of the binary team-quality measures used in Table 3 and at least 7.5 in the experience field (on a 1-10 scale). The reproducibility score of their analysis is at least 75 (out of 100) and their average peer evaluator rating is at least 7.5 (out of 10). The boxes span the the first to the third quartile with the median as the interior horizontal line. The whiskers span 95% of the the observations, starting from the 2.5% quantile to the 97.5% quantile.

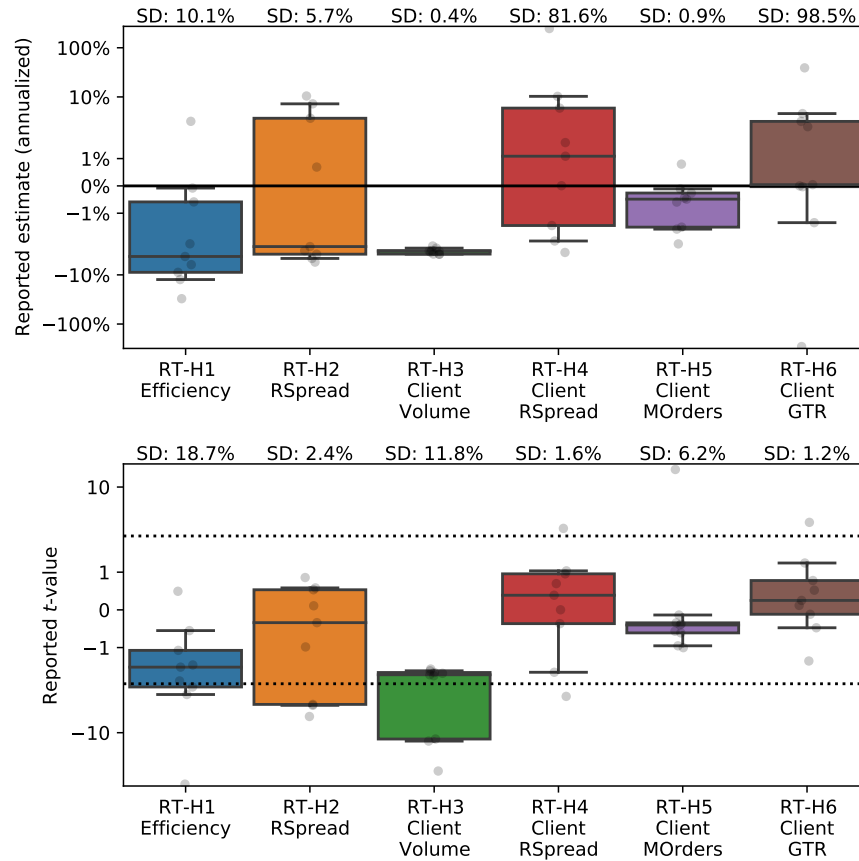


Figure 4: Dispersion research team estimates across all stages

This plot illustrates the dispersion in the estimates reported by the research teams for all six hypotheses. The boxes span the first to the third quartile with the median as the interior horizontal line. The whiskers span 95% of the observations, starting from the 2.5% quantile to the 97.5% quantile.

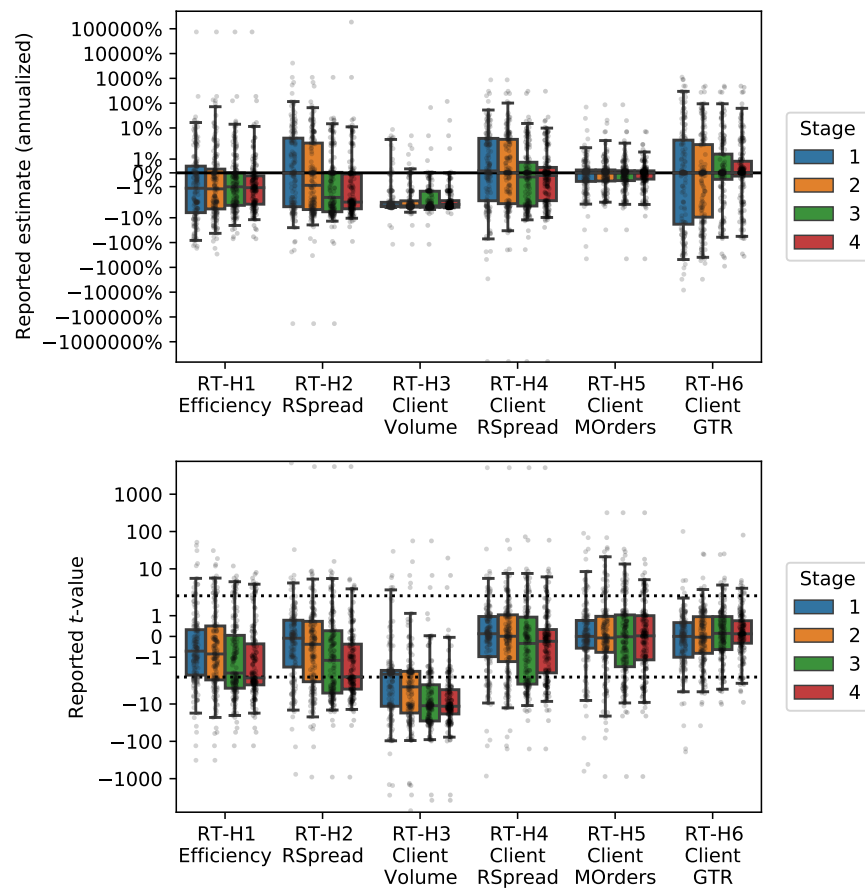


Figure 5: Dispersion research team beliefs

This plot illustrates the dispersion in beliefs across research teams for all six hypotheses. The teams were asked to report their belief on how large the standard deviation would be across all teams of the reported estimates and t -values. The whiskers span 95% of the reported beliefs, starting from the 2.5% quantile to the 97.5% quantile. The red dots represent the “truth,” the observed standard deviation across research teams. This figure is the only one that uses the unwinsorized sample as beliefs were solicited for the raw data.

