

29 This project is a mega-reproduction led by the Institute for Replication (I4R),
30 which evaluates the reproducibility and robustness of 110 published studies in eco-
31 nomics (79) and political science (31). Our focus is on studies published in 12
32 prestigious journals between 2022 and 2023. While each of these journals has a data
33 and code availability policy requiring authors to publicly share their materials upon
34 publication, most (though not all) also appoint a dedicated data editor. This edi-
35 tor is responsible for enforcing the journal’s data and code policy and conducting
36 internal computational reproducibility checks for accepted studies (see Supplementary
37 Materials 11.5).

38 Not all studies from our targeted journals were chosen for reproduction and
39 robustness, and our sample is thus not a random representative sample of studies
40 in economics and political science. Our approach leads to an over-representation of
41 studies using publicly available data ([18]). Another feature of our sample is that the
42 targeted journals have a data availability policy *and* enforce it. This is in contrast to
43 many top field journals in both economics and political science. Our sample should
44 thus be viewed as very selective both in terms of impact and high data and code avail-
45 ability rates, and might present an optimistic upper bound on reproducibility rates.
46 In fact, virtually all papers in our sample include replication packages with cleaned
47 data and code to reproduce the paper’s results, and about 30% also provide the raw
48 data and cleaning code used to generate the analytical data (Extended Data Figure
49 1, Levels 8, 9, and 10).

50 While this project relates to the broader reproducibility movement in psychology,
51 neuroscience, or biomedicine, it distinguishes itself from notable social science repli-
52 cation efforts along four key dimensions [24–26]. First, we are mostly reproducing
53 (non-experimental) studies using the same data as the original authors. Second, we
54 assess computational reproducibility and test the robustness of estimates to alternative
55 specification choices. Because of the unique nature of the underlying studies—largely
56 non-experimental work that uses observational data—we offer the first evidence about
57 the general robustness of economics and political science. Third, we concentrate on
58 recent studies for both economics and political science. Finally, this is an ongoing
59 initiative that aims to expand across disciplines, with the goal of mass reproducing
60 studies and reshaping research norms at scale. This paper reports findings from the
61 first 110 reproductions.

62 2 Definitions

63 We follow [27]’s nomenclature and define a claim **computationally reproducible** if
64 its results can be reproduced using the original study’s data and protocols. A claim
65 is **robust** if its results are robust to alternative reasonable analytical decisions on the
66 same data. Last, a claim is **replicable** if its results can be repeated using new data.

67 3 Teams and Choice of Study

68 The reproductions and replications in this project are generated in one of two streams.
69 First, I4R has a board of editors who recommend potential reproducers. Second, I4R
70 held 11 events called replication games (Games) ([28]). Games are one-day events

71 open to faculty, post-docs, graduate students and other researchers. Participants are
72 assigned to a small team of about 3–5 other researchers all working in the same subfield
73 (*e.g.*, development economics).

74 Participant teams are offered a short list of (average 5) studies in their subfield
75 of interest about three weeks before the games. They are asked to choose a paper as
76 a team, and familiarize themselves with the data and codes publicly posted by the
77 original authors (*i.e.*, replication package) prior to the games. After the Game, teams
78 submit a standardized reproduction report summarizing their results.

79 I4R emphasizes to reproducers that the goal is *not* to show that the results are
80 not reproducible. The goal is instead to test if the claims are reproducible and robust.
81 This is key as some reproducers might engage in reverse specification searching (*i.e.*,
82 selective reporting of insignificant results). I4R stresses the importance of reasonable
83 robustness checks and recoding [29]. Re-analyses are sensible tests of the research
84 question and expected to be statistically valid and theoretically informed.

85 We survey the reasons why teams selected their paper (Extended Data Figure 2).
86 While 13.6% of teams were assigned a study (*i.e.*, did not choose which study to work
87 on), about 45% of teams report “Methods used”, 36% of teams selected “because of
88 the journal of publication” and about 25% due to the “length of time to reproduce
89 results”.

90 If a large portion of reproducers select papers based on the assumption that their
91 findings are questionable, it could skew reproducibility rates downward, as such studies
92 might be more prone to revealing problematic outcomes. However, in this project,
93 only a minimal fraction of teams indicated that they chose their paper because of *ex*
94 *ante* beliefs that main results are (not) replicable (3.6%). We found that selecting a
95 paper due to the reproducers’ belief the paper is not robust is *inversely* correlated
96 with reproducer experience ($\rho = -0.19, p < 0.000$). A few teams (5%) indicated that
97 their choice was based on statistical power/sample size and trust of original authors.

98 4 Data and Computational Reproducibility

99 We find a computational reproducibility rate of 85%. That is, when provided with
100 the original data and code, independent researchers are able to reproduce the pub-
101 lished results in economics and political science studies 85% of the time using either:
102 (1) the raw and analytical data, or; (2) the analytical data when the raw data were
103 not provided. The remaining 15% of cases involved studies with only partial availabil-
104 ity of code or data, or instances where code failed to run or produced inconsistent
105 results (See Supplementary Materials 11.8 and Extended Data Figure 1). Fixing paths,
106 missing packages and software requirements were not considered failures of computa-
107 tional reproducibility. In those instances, we fixed paths, added missing packaged and
108 software requirements.

109 Our findings suggest high rates of computationally reproducible results, but far
110 from perfect for leading journals. Our results are in contrast with several studies
111 documenting low computational reproducibility rates in economics [19]; [13]; [22]. This
112 may in part reflect the effectiveness of editorial policies in journals that have introduced
113 data editors and mandatory sharing of replication packages.

114 To provide context to these findings, we mapped data and code availability in
115 all of our target journals between 2014 and 2023. As discussed in Supplementary
116 Materials 11.11, data and code sharing practices have dramatically improved during
117 this period. We found replication folders are attached to 59% of papers in 2014, while
118 replication folder provision increases to a seemingly stable value close to 90% in 2021–
119 2023 (Extended Data Figures 3, 4, 5 and 6). Additionally, for journals that introduced
120 data editors during this period, much of this improvement occurred during the first
121 year following this change.

122 5 Robustness

123 For robustness, we directly compare original point estimates to the revised point esti-
124 mates. This one-on-one comparison allows us to speak to the robustness of a specific
125 hypothesis test, in addition to the robustness of our entire sample. We are thus looking
126 at several claims within a study and conduct robustness reproducibility and robustness
127 for multiple claims.

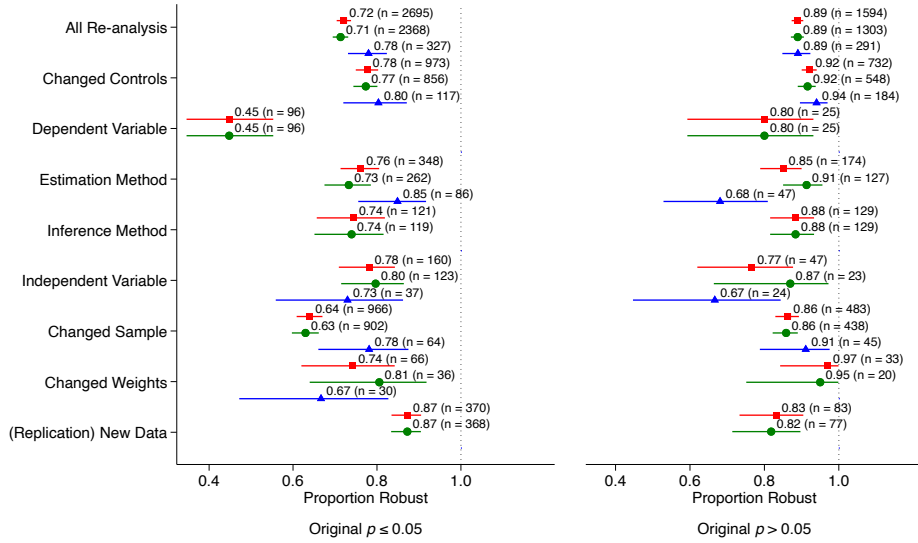
128 Reproducers are then free to conduct any robustness or recoding exercises. They
129 focus on the reproducibility of the claims and have access to the replication package,
130 allowing them to directly test the robustness of the main results. This is a crucial
131 advantage over the traditional review process as reproducers may uncover coding errors
132 and discrepancies between the paper and the codes. They may also uncover coding
133 decisions that were not discussed (or are hard to find) in the article.

134 However, this flexibility also brings some disadvantages. As with the journal review
135 process with reviewers, reproducers spend different amounts of time and effort on their
136 respective replication. Some reproducers are more experienced at coding, while others
137 are more familiar with methods, or simply unable to implement robustness checks due
138 to a lack of raw data (Extended Data Figure 7). This means that reproducibility efforts
139 and type of re-analysis vary across teams. Teams worked on average 13 active days
140 (std. dev. of 24) on the reproductions and robustness, and reports were on average 19
141 pages long (std. dev. of 14).

142 Figure 1 (top of left panel) shows a robustness rate of 72%. This result means that
143 when alternative analytical decisions were made on the same data, 72% of originally
144 statistically significant estimates ($p < 0.05$) remained statistically significant ($p <$
145 0.05) in the original direction.

146 **Figure 1.** Robustness Rate.

Fig. 1: Robustness Rate



Robustness rate for ... **Left panel:** ... originally statistically significant research **Second panel:** ... in economics **Third panel:** ... in political science **Right panel:** ... originally statistically insignificant research **All panels:** Squares, circles, and triangles represent proportions, with 95% Clopper-Pearson confidence intervals presented in whiskers. Red squares represent full sample. Green circles represent economics subsample. Blue triangles represent political science subsample. Each group of three estimates represent different types of re-analysis, non-mutually exclusive. The first 8 groups do not include re-analyses that use new data (replication), while the last one does. The first estimate group contains all types of re-analysis, then all types of re-analysis in economics, then all types of re-analysis in political science. The second represents re-analyses which changed the control variables, e.g., by adding or re-defining them. The third represents re-analyses which changed the dependent variable, e.g., by employing a different standardization or binarization. The fourth represents re-analyses which changed the estimation method, e.g., by adjusting a matching procedure. The fifth represents re-analyses which changed the inference method, e.g., changed the level on which standard errors are clustered. The sixth represents re-analyses which changed the main independent variable, e.g., by taking into account treatment intensity. The seventh represents re-analyses which changed the sample, e.g., by excluding outliers. The eighth represents re-analyses which changed the weights applied, or applied weights for the first time. The last represents replicability rates for re-analyses that introduced new data, e.g., comparable outcomes from more recent survey waves.

147 We find large differences by re-analysis type. The re-analysis type that has the
 148 highest robustness rate (78%) is changing the independent variable measure (exam-
 149 ples include log transformations, discretization, etc.). The re-analysis type that has
 150 the lowest robustness rate (45%) is any which included changing the dependent vari-
 151 able measure (e.g., categorizing the variable or log-transforming). When a replication
 152 (addition of new data, e.g., from more recent survey waves or an alternative source)
 153 is applied, the replication rate is 87%.

154 The average robustness rate is 71% and 78% for economics and political science,
 155 respectively, where the 6.7% difference is statistically significant (two-sample difference
 156 in proportions $z = -2.52$, $p = 0.012$, $n_1 = 2368$, $n_2 = 327$). The general pattern of the
 157 robustness rates is similar between economics and political science (with the exception
 158 of dependent variable and inference method, which were not applied by any of the

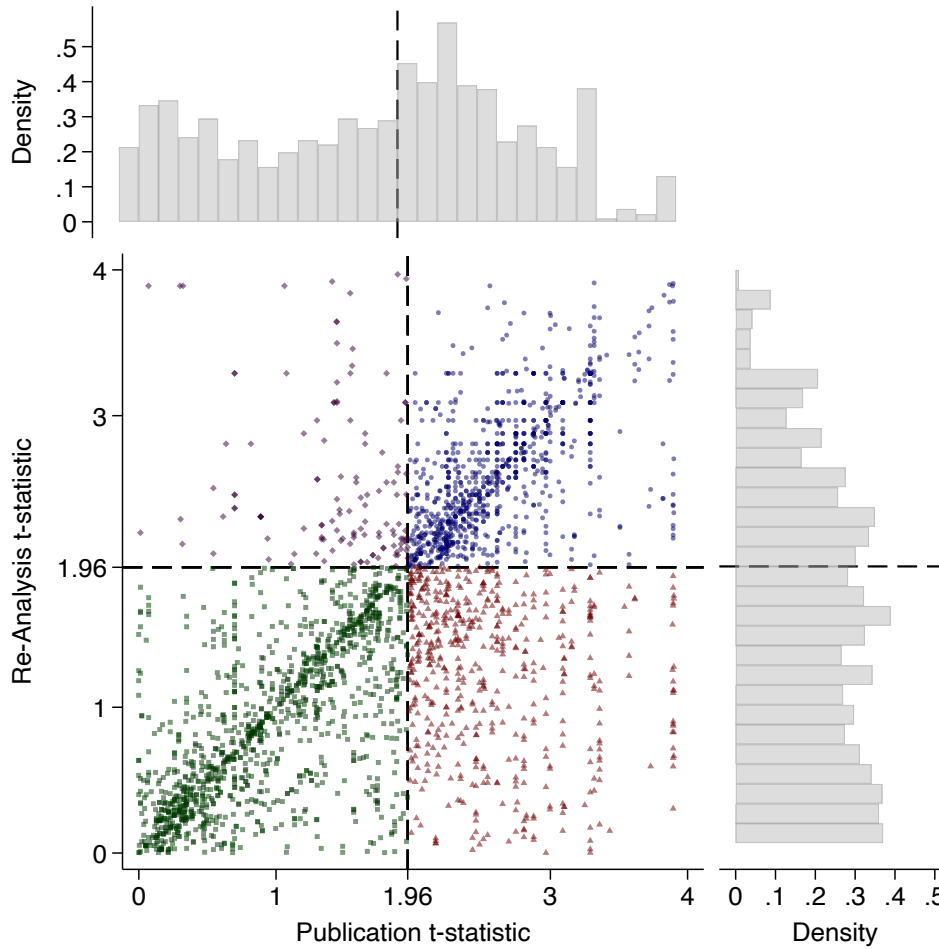
159 political science re-analyses). Focusing on robustness rates for originally statistically
160 *insignificant* findings, we find a robustness rate of 89%.

161 Supplementary Materials Appendix Table 1 shows shifts in statistical significance
162 between all significance regions. We find that 7.44% of re-analyses find an effect with
163 the opposite sign as the original result. In contrast, of the 62.84% of original analyses
164 that were statistically significant, most remained significant and of the same sign
165 (44.33%/62.84% = 70.5%). Of particular note is the 15.06% of re-analyses that find a
166 statistically *insignificant* result for originally statistically significant analyses.

167 We illustrate in Figure 2 the distribution of test statistics for the original point
168 estimates and the re-analyses. We find that 53% of the originally published test statis-
169 tics are statistically significant (to the right of the statistical significance threshold).
170 In contrast, 43% of re-analyses are statistically significant (above the statistical signif-
171 icance threshold in the vertical axis histogram). The simple difference in proportions
172 is statistically significant (difference of 10.4%, McNemar's $\chi^2 = 264.11$, $p < 0.001$,
173 $n = 4750$).

174 **Figure 2.** Statistical Significance of Publication and Re-analysis.

Fig. 2: Statistical Significance of Publication and Re-analysis



Top histogram: Distribution of publication tests of significance. T-statistics over 4 truncated for exposition. The histogram's bars are of width 0.14, with exactly 14 bars between 0 and the statistical threshold of $t = 1.96$ (corresponding to statistical significance at the 5% level). **Right histogram:** Distribution of re-analysis tests of significance. T-statistics over 4 truncated for exposition. **Scatterplot:** Each marker is a pair of test statistics, an originally published test statistic (horizontal value) and an associated re-analysis test statistic (vertical value). If the original and re-analysis test statistics were identical, this scatterplot would follow the 45 degree line. As either axis represents statistical significance, we have bifurcated each with a line at $t=1.96$, representing statistical significance at the 5% threshold. **Blue circles** indicate an originally statistically significant statistic that is also statistically significant under re-analysis. Represents 50% of sample. **Red triangles** indicate originally significant test statistics that are no longer statistically significant under re-analysis. Represents 14% of sample. **Green squares** indicate originally statistically insignificant test statistics that are the same under re-analysis. Represents 27% of sample. **Purple diamonds** indicate originally statistically insignificant test statistics that become statistically significant under re-analysis. Represents 3% of sample. **Not displayed** Not displayed are the 6% of test statistics that represent a sign reversal between the originally estimated effect and the effect estimated under re-analysis.

175 When expressed as t-statistics, the average originally published t-statistic is 1.797
176 whereas the average re-analysis t-statistic is 1.544. The difference between the pairs of
177 original study estimates and re-analysis estimates is statistically significant (Wilcoxon
178 signed-rank test $z = 15.477$, $p < 0.001$, $n = 3151$). Indeed, we reject the null hypoth-
179 esis of a two-sample Kolmogorov–Smirnov test that the two distributions come from
180 the same probability distribution ($p < 0.001$). Here, we also note the large increase in
181 test statistic density immediately after the statistical significance threshold (Extended
182 Data Figures 8 and 9), which offers strong evidence of publication bias in origi-
183 nally published research ([30, 31]). In contrast, this increase at the significance level
184 threshold is missing from the vertical axis histogram depicting the distribution of
185 re-analyses.

186 When expressed as p-values, the average originally published p-value is 0.167
187 whereas the average re-analysis p-value is 0.219; the difference is statistically significant
188 (Wilcoxon signed-rank test $z = -16.007$, $p < 0.001$, $n = 4063$).

189 In this project, we conduct multiple re-analyses per original study, and so it is
190 possible that much of the differences between original studies and their re-analyses
191 are driven or characterized by large changes in a small subset of studies rather than
192 indicative of more general shifts between original and re-analysis. In fact, we find
193 evidence of general shifts. The proportion of original studies that have at least one
194 statistically significant result is 95.3% whereas for the corresponding re-analyses this
195 is 92.9% (difference of 2.4%, McNemar’s $\chi^2 = 1.00$, $p < 0.625$, $n = 86$). Only 3.6%
196 of articles did not lose any statistical significance under replication, and the average
197 replication lost statistical significance for 29% of replication tests (median of 22%). In
198 only three original studies that reported statistically significant results, the reanalysis
199 found that all test statistics were not statistically significant.

200 6 Determinants of Robustness

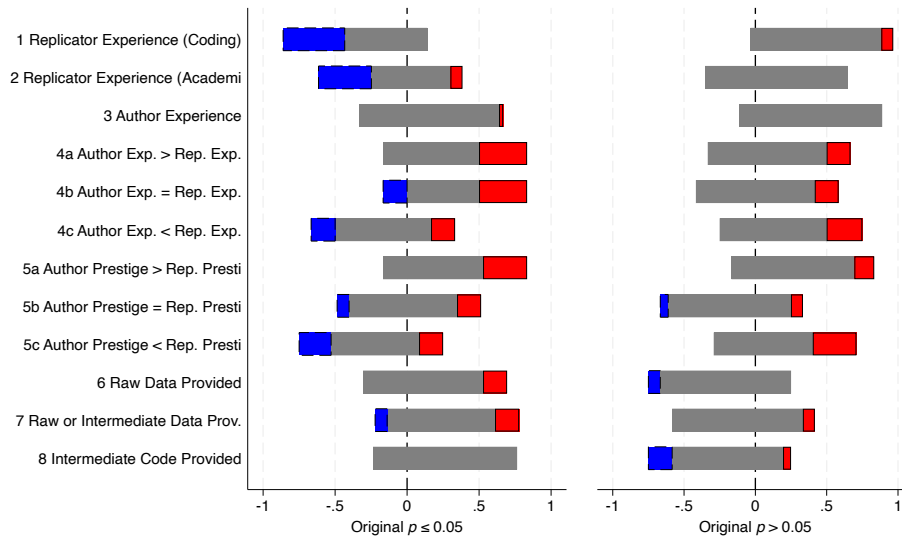
201 This section examines what, if any, characteristics of the authors, reproducers, or the
202 original articles are informative of the robustness rate.

203 While this analysis is merely exploratory, this project applied both a pre-
204 registration and many-analysts approach [32–36]. By pre-specifying which research
205 questions would be examined, and averaging the responses to those research questions
206 over multiple independent teams, the results here are guarded against specification
207 searching and confirmation bias.

208 About 110 co-authors were invited to participate regarding the proposed determi-
209 nants of robustness. We received answers from 10 individuals and ended up forming
210 six many-analysts teams. Each team answered several research questions. The results
211 are displayed in Figure 3.

212 **Figure 3.** Robustness Rate Determinants.

Fig. 3: Robustness Rate Determinants



Six independent teams answered twelve questions of the re-analysis database. Each bar represents a different question. **Left panel:** “Does reproducibility of an originally statistically significant result depend on...” **Right panel:** “Does reproducibility of an originally statistically *insignificant* result depend on...” **Both panels:** where the first bar represents “... the reproducers’ experience at coding.” **Blue, patterned outline** indicates the proportion of teams that indicated a negative and statistically significant relationship, in whichever manner the team defined so in their analysis. **Gray, no outline** indicates the proportion of teams that indicated a statistically insignificant relationship, where left of the zero line indicates negative and right of the zero line indicates positive. **Red, solid outline** indicates the proportion of teams that indicated a statistically significant and positive relationship. All teams equally weighted.

213 They began by analyzing originally statistically significant results and answer-
 214 ing the first question “Does reproducibility/replicability rate depend on reproducers’
 215 experience coding?” Specifically, most of the teams estimated a negative coefficient
 216 in a regression with reproducibility as the dependent variable and a measure of their
 217 choosing for reproducers’ experience as the primary independent variable, that is, the
 218 relationship is far more likely to be negative than positive. We interpret this result
 219 to mean that reproducers who are more experienced (broadly defined, as each of the
 220 many analysts defined experience independently) are better able to detect non-robust
 221 results in their chosen paper; likening the notion of the ‘trained eye’ of a detective
 222 finding subtle clues the untrained eye may miss at the scene. The remaining 11 pre-
 223 specified hypotheses that the analysts tested were whether reproducibility is associated
 224 with: (2) reproducers’ experience in academia, (3) the original authors’ experience in
 225 academia, whether authors have (4a) more, (4b) similar, or (4c) less experience than
 226 reproducers, (5a) more, (5b) similar, or (5c) less prestige (their institution, defined
 227 independently by the analysts) than reproducers, and whether (6) raw data was pro-
 228 vided (7) raw or intermediate data was provided, and (8) whether cleaning code was
 229 provided.

230 Among results that were originally statistically significant, the first hypothesis
231 yielded the clearest finding: the more experience a reproducer team had, the lower the
232 robustness rate they found. One plausible interpretation of our main results therefore
233 is that robustness in our full sample would likely have been lower if equally highly
234 qualified replicator teams had been assigned to each paper. However, according to
235 the results presented in the main text (Determinants of Robustness), the provision of
236 raw or intermediate data, or the necessary cleaning codes, does not seem to affect the
237 robustness of research.

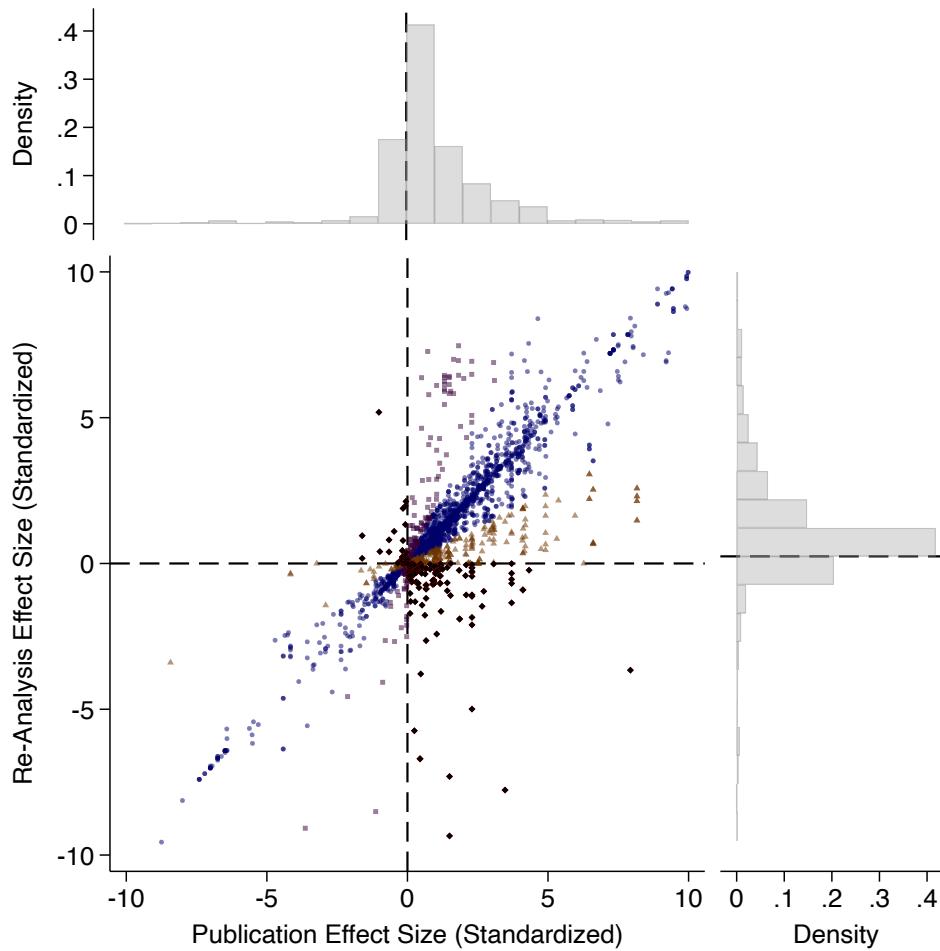
238 When analysts examined these same 12 hypotheses for originally statistically
239 *insignificant* results, the relationships are far more likely to be positive than negative,
240 but (as indicated by the proportion in gray) the relationships are often not statistically
241 significant.

242 7 Effect Size

243 Figure 4 displays publication and re-analysis effect sizes. In economics and political
244 science, effect sizes are largely reported as non-unit-less regression coefficients, whereas
245 in other sciences, effect sizes are often reported using more comparable measures such
246 as Cohen's-d. Because raw effect sizes vary widely between original studies, each of
247 the markers are standardized by the within-article average published effect size (e.g.,
248 estimated effects of 2, 4, and 6 are standardized within publication to be 0.5, 1.0, and
249 1.5).

250 **Figure 4.** Effect Size of Publication and Re-analysis.

Fig. 4: Effect Size of Publication and Re-analysis



Top histogram: Distribution of originally published effect size standardized by the average effect size within a published article. **Right histogram:** Distribution of re-analysis published effect size standardized by the average effect size within a published article. **Scatterplot:** Each marker is a pair of effect sizes, the originally published effect size (horizontal value) and an associated re-analysis effect size (vertical value). If an originally estimated and re-analysis effect size were of similar magnitude (and sign), the markers would gather tightly around the 45 degree line passing through the origin **Blue circles** indicate effect sizes which are similar (between 50% to 200% of original effect size) under re-analysis. Represents 69% of sample. **Red diamonds** indicate effect size estimates which switch sign under re-analysis. Represents 6% of sample. **Orange triangles** indicate effect size estimates which are 50% or less their original magnitude under re-analysis. Represents 9% of sample. **Purple squares** indicate effect size estimates which are double or larger than their original magnitude under re-analysis. Represents 16% of sample.

251 We find that, on average, the median effect size of a re-analysis is equivalent to
252 the published effect size (i.e., 99% the size of the published effect), while the mean
253 replicated effect is 9% larger than the original. Extended Data Figure 10 illustrates
254 the distribution effect sizes of re-analyses. This result is in stark contrast to previous
255 projects focused on replication with new data in psychology or social science exper-
256 iments uncovering replication rates ranging from 50 to 66% [24–26]. Three major
257 differences between our project and these replication efforts are that we focus on
258 robustness as opposed to replication with new data, our focus is on recent articles, and
259 that our sample is composed mostly of non-experimental studies using secondary data.

260 8 Coding Errors and Recoding

261 We investigate the prevalence of coding errors and discrepancies between the code
262 and article. Computational reproducibility pertains to the provided replication folder’s
263 ability to reproduce the exhibits and statistics displayed in the research (manuscripts,
264 appendices, *etc.*). Reproducers may be able to reproduce all exhibits exactly as they
265 appear (computationally reproducible), but the exhibits may have been constructed
266 with coding errors or discrepancies.

267 Except for minor inconveniences (*i.e.*, missing packages or broken pathways), we
268 identify coding errors in approximately 25% of the studies, with some studies contain-
269 ing multiple errors (Supplementary Materials 11.10). The prevalence of coding errors is
270 larger for economics (26%) than political science (16%). Types of errors include: defin-
271 ing the dependent variable, defining the main independent variable, defining control
272 variables, mis-specification of the estimation/model, inference or the sample. While not
273 all of these coding errors impacted the conclusions of the original studies, we uncover
274 several significant errors that warrant discussion. These major errors include instances
275 of duplicated observations on a large scale, incomplete interaction in a difference-in-
276 differences model, mislabeling the main treatment variable for a substantial number
277 (or all) of observations, and using different models, or estimators, than reported in
278 the article.

279 It is important to note that this 25% figure likely underestimates the true preva-
280 lence of coding errors. Reproducers may have missed some errors, and many replication
281 packages do not include raw data or data-cleaning code, limiting the ability to detect
282 additional issues.

283 A number of reproducers also recoded the analysis using a different statistical
284 software. Out of 23 recoding exercises, we find major differences for three studies and
285 minor differences for 10 studies. Two of the major differences were uncovered when
286 using a different software and looking at the authors’ code. Additionally, one team who
287 computationally reproduced the results using a different *version* of the software used
288 by the authors uncovered noteworthy differences in the magnitude and significance of
289 the estimates (Supplementary Materials 11.9).

290 9 Communication with Original Authors

291 I4R shares completed reproduction reports with original authors before public release
292 ([28]). Reports are reviewed typically by A.B. or another board member mainly for

293 tone and structure. I4R then disseminates the report and any author response simul-
294 taneously (see SI for the full list of reports). Reproducers may revise their reports
295 after receiving feedback from original authors.

296 About 95% of contacted authors responded (including one case where an author
297 was unreachable after leaving academia). Among respondents, 11% provided only brief
298 notes or indicated they could not respond, 59% offered informal feedback, and 30%
299 supplied a formal response. For comparison, [37] report that roughly 25% of authors
300 in their sample provided a formal response.

301 Roughly two-thirds of reproducers indicated that interactions with original authors
302 improved their reports, often by clarifying variables or procedures, supplying data or
303 data-access instructions, or helping adjust tone. In one case, original authors conducted
304 additional robustness checks in their non-public files at the reproducers' request.

305 Lastly, we assess agreement between authors and reproducers. Authors' final
306 responses were coded for whether disagreements remained after mediation; only
307 23% of articles showed any remaining disagreement. Further details appear in the
308 Supplementary Materials.

309 10 Discussion

310 A substantial information asymmetry exists between authors and the broader aca-
311 demic community, including reviewers and editors ([30]). Reviewers rarely see the
312 underlying data and code and may be unaware of crucial coding decisions, even as
313 journals routinely request multiple robustness checks. This limited visibility means
314 major errors or inconsistencies can go undetected.

315 Large-scale reproducibility initiatives offer a promising way to address these chal-
316 lenges in the social sciences and beyond. Our project provides a systematic, scalable
317 approach to evaluating reproducibility and robustness, with the goal of increasing
318 transparency and improving the credibility of published research. While stronger
319 incentives to conduct reproducibility and robustness remain necessary, we do not
320 attempt to evaluate which specific incentives would be most effective, as doing so
321 would require speculation beyond the scope of our data. Identifying the most effective
322 incentives is an important research question that we hope future work will address.

323 Given the low prevalence of diagnostic replication in published work [38], the
324 scale of this ongoing effort could shift research norms. By encouraging more rigorous
325 methodologies, deterring questionable research practices, and emphasizing collabo-
326 ration, it may help place greater weight on the reliability of results in publication
327 decisions.

328 Although our journal sample is selective, the findings are encouraging and suggest
329 a high level of computational reproducibility. These patterns—and the existence of a
330 large-scale, community-driven effort—may strengthen trust in published results.

331 We also asked reproducers about the quality of the replication packages they exam-
332 ined. Over 40% reported gaining a more optimistic view of the discipline, while only
333 about 5% developed a more negative opinion. This suggests that mass reproduction
334 of studies accompanied by replication packages can directly enhance researchers' trust
335 in scientific findings.

336 The initiative’s success and scalability have been driven by the intrinsic motivation
337 of participating researchers to support open science and improve their technical skills.
338 By late 2025, I4R had organized 80 replication games involving over 3,500 researchers,
339 with events held every other week. These efforts show that the skilled labor needed for
340 large-scale reproduction can emerge organically from an engaged research community.

341 The project also has the potential to advance science and improve equity. Publicly
342 posting data and code facilitates learning, speeds methodological diffusion, and enables
343 independent verification. Reproducing analyses in open-source software can also help
344 level the playing field for researchers who lack access to expensive licenses.

345 Our results have limitations. Only a small number of economics and political
346 science journals currently require data and code [17]; [18], and even fewer check
347 reproducibility [39]. Thus, our findings largely reflect leading journals with strong
348 data-sharing norms. Future research should assess reproducibility more broadly by
349 examining a random sample of papers from journals with and without data availability
350 policies.

351 Author list

352 Abel Brodeur, Derek Mikola, Nikolai Cook, Lenka Fiala, Thomas Brailey, Ryan Briggs,
353 Alexandra de Gendre, Yannick Dupraz, Jacopo Gabani, Romain Gauriot, Joanne
354 Haddad, Goncalo Lima, Jörg Ankel-Peters, Anna Dreber, Douglas Campbell, Lamis
355 Kattan, Diego Marino Fages, Fabian Mierisch, Pu Sun, Taylor Wright, Marie Connolly,
356 Fernando Hoces de la Guardia, Magnus Johannesson, Edward Miguel, Lars Villhuber,
357 Alejandro Abarca, Mahesh Acharya, Sossou Simplicie Adjisse, Ahwaz Akhtar, Eduardo
358 Alberto Ramirez Lizardi, Sabina Albrecht, Synøve Nygaard Andersen, Zubaria Andlib,
359 Falak Arrora, Thomas Ash, Etienne Bacher, Sebastian Bachler, Félix Bacon, Manuel
360 Bagues, Timea Balogh, Alisher Batmanov, Mara Barschkett, B. Kaan Basdil, Jaromír
361 Baxa, Sascha Becker, Monica Beeder, Louis-Philippe Beland, Abdel-Hamid Bello,
362 Daniel Benenson Markovits, Grant Benjamin, Thomas Bergeron, Moussa Blimpo,
363 Marco Binetti, Carl Bonander, Joseph Bonneau, Endre Borbáth, Nicolai Topstad Bor-
364 gen, Solveig Topstad Borgen, Jonathan Borowsky, Elisa Brini, Myriam Brown, Martin
365 Brun, Stephan Bruns, Nino Buliskeria, Andrea Calef, Alistair Cameron, Pamela
366 Campa, Santiago Campos-Rodríguez, Giulio Giacomo Cantone, Fenella Carpena,
367 Perry Carter, Paul Castañeda Dower, Ondrej Castek, Jill Caviglia-Harris, Gabriella
368 Chauca Strand, Shi Chen, Sya In Chzhen, Jong Chung, Jason Collins, Alexan-
369 der Coppock, Hugo Cordeau, Ben Couillard, Jonathan Crechet, Lorenzo Crippa,
370 Jeanne Cui, Christian Czymara, Haley Daarstad, Danh Chi Dao, Daniel Dao, Marco
371 David Schmandt, Astrid de Linde, Lucas De Melo, Lachlan Deer, Micole De Vera,
372 Velichka Dimitrova, Jan Fabian Dollbaum, Jan Matti Dollbaum, Michael Donnelly,
373 Luu Duc Toan Huynh, Tsvetomira Dumbalska, Jamie Duncan, Kiet Tuan Duong,
374 Thibaut Duprey, Christoph Dworschak, Sigmund Ellingsrud, Ali Elminejad, Yasmine
375 Eissa, Andrea Erhart, Giulian Etingin-Frati, Elaheh Fatemipour, Alexa Federice,
376 Jan Feld, Guidon Fenig, Mojtaba Firouzjaeiangalougah, Erlend Fleisje, Alexandre
377 FortiFriter-Chouinard, Julia Francesca Engel, Nadjim Fréchet, Reid Fortier, Tilman
378 Fries, Michael James Frith, Thomas Galipeau, Sebastián Gallegos, Areez Gangji,

379 Xiaoying Gao, Cloé Garnache, Attila Gáspár, Evelina Gavrilova, Arijit Ghosh, Gar-
380 reth Gibney, Grant Gibson, Geir Godager, Leonard Goff, Da Gong, Javier González,
381 Jeremy D. Gretton, Cristina Griffa, Idaliya Grigoryeva, Maja Grötting, Eric Gun-
382 termann, Jiaqi Guo, Alexi Gugushvili, Hooman Habibnia, Sonja Häffner, Jonathan
383 D. Hall, Olle Hammar, Amund Hanson Kordt, Barry Hashimoto, Jonathan S. Hart-
384 ley, Carina I. Hausladen, Tomáš Havránek, Harry He, Matthew Hepplewhite, Mario
385 Herrera-Rodriguez, Felix Heuer, Anthony Heyes, Anson T. Y. Ho, Jonathan Holmes,
386 Armando Holzknacht, Yu-Hsiang Dexter Hsu, Shiang-Hung Hu, Yu-Shiuan Huang,
387 Mathias Huebener, Christoph Huber, Kim P. Huynh, Zuzana Irsova, Ozan Isler,
388 Niklas Jakobsson, Raphaël Jananji, Tharaka A. Jayalath, Michael Jetter, Jenny John,
389 Rachel Joy Forshaw, Felipe Juan, Valon Kadriu, Sunny Karim, Edmund Kelly, Duy
390 Khanh Hoang Dang, Tazia Khushboo, Jin Kim, Gustav Kjellsson, Anders Kjelsrud,
391 Andreas Kotsadam, Jori Korpershoek, Lewis Krashinsky, Suranjana Kundu, Alexan-
392 der Kustov, Nurlan Lalayev, Audrée Langlois, Jill Laufer, Blake Lee-Whiting, Andreas
393 Leibing, Gabriel Lenz, Joel Levin, Peng Li, Tongzhe Li, Yuchen Lin, Ariel Listo,
394 Dan Liu, Xuewen Lu, Elvina Lukmanova, Alex Luscombe, Lester R. Lusher, Ke
395 Lyu, Hai Ma, Nicolas Mäder, Clifton Makate, Alice Malmberg, Adit Maitra, Marco
396 Mandas, Jan Marcus, Shushanik Margaryan, Lili Márk, Andres Martignano, Abi-
397 gail Marsh, Isabella Masetto, Anthony McCanny, Emma McManus, Ryan McWay,
398 Lennard Metson, Jonas Minet Kinge, Sumit Mishra, Myra Mohnen, Jakob Möller, Ros-
399 alie Montambeault, Sébastien Montpetit, Louis-Philippe Morin, Todd Morris, Scott
400 Moser, Fabio Motoki, Lucija Muehlenbachs, Andreea Musulan, Marco Musumeci,
401 Munirul Nabin, Karim Nchare, Florian Neubauer, Quan M. P. Nguyen, Tuan Nguyen,
402 Viet Nguyen-Tien, Ali Niazi, Giorgi Nikolaishvili, Ardyn Nordstrom, Patrick Nüß,
403 Angela Odermatt, Matt Olson, Henning Øien, Tim Ölkens, Miquel Oliver i Vert,
404 Emre Oral, Christian Oswald, Ali Ousman, Ömer Özak, Shubham Pandey, Alexan-
405 dre Pavlov, Martino Pelli, Romeo Penheiro, RyuGyung Park, Eva Pérez Martel,
406 Tereza Petrovičová, Linh Phan, Alexa Prettyman, Jakub Procházka, Aqila Putri,
407 Julian Quandt, Kangyu Qiu, Loan Quynh Thi Nguyen, Andaleeb Rahman, Carson
408 H. Rea, Adam Reiremo, Laëtitia Renée, Joseph Richardson, Nicholas Rivers, Bruno
409 Rodrigues, William Roelofs, Tobias Roemer, Ole Rogeberg, Julian Rose, Andrew
410 Roskos-Ewoldsen, Paul Rosmer, Barbara Sabada, Soodeh Saberian, Nicolas Sala-
411 manca, Georg Sator, Daniel Scates, Elmar Schlüter, Cameron Sells, Sharmi Sen, Ritika
412 Sethi, Anna Shcherbiak, Moyosore Sogaolu, Matt Soosalu, Erik Ø. Sørensen, Manali
413 Sovani, Noah Spencer, Stefan Staubli, Renske Stans, Anya Stewart, Felix Stips, Kieran
414 Stockley, Stephenson Strobel, Ethan Struby, John Tang, Idil Tanrisever, Thomas Tao
415 Yang, Ipek Tastan, Dejan Tatić, Benjamin Tatlow, Féraud Tchuisseu Seuyong, Rémi
416 Thériault, Vincent Thivierge, Wenjie Tian, Filip-Mihai Toma, Maddalena Totarelli,
417 Van Tran, Hung Truong, Nikita Tsoy, Kerem Tuzcuoglu, Diego Ubfal, Laura Villalo-
418 bos, Julian Walterskirchen, Joseph Tao-yi Wang, Vasudha Wattal, Matthew D. Webb,
419 Bryan Weber, Reinhard Weisser, Wei-Chien Weng, Christian Westheide, Kimberly
420 White, Jacob Winter, Timo Wochner, Matt Woerman, Jared Wong, Ritchie Woodard,
421 Marcin Wroński, Myra Yazbeck, Chung Yang, Luther Yap, Kareman Yassin, Hao Ye,
422 Jin Young Yoon, Chris Yurris, Tahreen Zahra, Mirela Zaneva, Aline Zayat, Jonathan
423 Zhang, Ziwei Zhao, Yaolang Zhong

424 Declarations

- 425 • Funding:
426 We acknowledge support from Coefficient Giving and the Social Sciences and
427 Humanities Research Council.
- 428 • Conflict of interest/Competing interests:
429 Any views expressed herein are the authors' personal opinions and not those of
430 Ontario Public Service. The work by Jeremy D. Gretton was not undertaken under
431 the auspices of Ontario Public Service as part of his employment responsibilities.
432 The views expressed in this paper are those of the authors. No responsibility for
433 them should be attributed to the Bank of Canada. The findings, interpretations,
434 and conclusions expressed in this work are entirely those of the authors and do not
435 necessarily reflect the views of the World Bank or its Board of Directors. The Center
436 for Crisis Early Warning (Kompetenzzentrum Krisenfrüherkennung) is funded by
437 the German Federal Ministry of Defense and the German Federal Foreign Office.
438 The views and opinions expressed in this article are those of the author(s) and do
439 not necessarily reflect the official policy or position of any agency of the German
440 government. The views expressed in this paper are those of the authors and do
441 not necessarily reflect the position of the Banco de España or the Eurosystem. All
442 remaining errors are the authors' responsibility.
- 443 • Data availability: The data are available on Zenodo (link: [https://zenodo.org/](https://zenodo.org/records/17792605)
444 [records/17792605](https://zenodo.org/records/17792605); DOI: 10.5281/zenodo.17792605) and OSF (link: [https://osf.io/](https://osf.io/8wsqx/)
445 [8wsqx/](https://osf.io/8wsqx/); DOI: 10.17605/OSF.IO/8WSQX). See OSF for our pre-analysis plan.
- 446 • Code availability: The codes are available on Zenodo ([https://zenodo.org/records/](https://zenodo.org/records/17792605)
447 [17792605](https://zenodo.org/records/17792605)) and OSF (<https://osf.io/8wsqx/>).
- 448 • Author contribution:
449 **Preparation of tables, figures, and manuscript:** Abel Brodeur (University of
450 Ottawa and Institute for Replication), Nikolai Cook (Wilfrid Laurier University),
451 Derek Mikola (University of Ottawa and Institute for Replication), Lenka Fiala
452 (University of Ottawa, Tilburg University and Institute for Replication)
453 **Conception or design of the work:** Jörg Ankel-Peters (RWI - Leibniz Institute
454 for Economic Research), Abel Brodeur (University of Ottawa and Institute for Repli-
455 cation), Marie Connolly (UQAM), Nikolai Cook (Wilfrid Laurier University), Anna
456 Dreber (Stockholm School of Economics), Fernando Hoces de la Guardia (Berkeley
457 Initiative for Transparency in the Social Sciences), Magnus Johannesson (Stockholm
458 School of Economics), Edward Miguel (UC Berkeley), Derek Mikola (University of
459 Ottawa and Institute for Replication), Lars Vilhuber (Cornell University)
460 **Analysis or interpretation of the reproducibility data:** Thomas Brailey
461 (University of Oxford), Ryan Briggs (University of Guelph), Abel Brodeur (Uni-
462 versity of Ottawa and Institute for Replication), Nikolai Cook (Wilfrid Laurier
463 University), Alexandra de Gendre (The University of Melbourne), Yannick Dupraz
464 (Paris Dauphine University, PSL University, LEDA, CNRS, IRD), Jacopo Gabani
465 (World Bank & Centre for health economics, university of York), Romain Gauriot
466 (Deakin University), Goncalo Lima (European University Institute and University
467 of Bologna), Derek Mikola (Institute for Replication)

468 **Analysis or interpretation of data And generating data And conception**
469 **of a reproduction:**

470 Douglas Campbell (Independent Researcher), Nikolai Cook (Wilfrid Laurier Univer-
471 sity), Joanne Haddad (Universitat Autònoma de Barcelona), Lamis Kattan (School
472 of Foreign Service, Georgetown University Qatar), Diego Marino Fages (Durham
473 University), Fabian Mierisch (Independent Researcher), Pu Sun (Dongbei University
474 of Finance and Economics), Taylor Wright (Brock University), Alejandro Abarca
475 (Texas Tech University), Mahesh Acharya (University of Calgary), Sossou Sim-
476 plice Adjisse (University of Wisconsin-Madison and African School of Economics),
477 Ahwaz Akhtar (George Washington University), Eduardo Alberto Ramirez Lizardi
478 (University of Oslo), Sabina Albrecht (University of Queensland), Synøve Nygaard
479 Andersen (University of Oslo), Zubaria Andlib (Lancaster University and Federal
480 Urdu University of Arts, Science and Technology), Falak Arrora (University of War-
481 wick), Thomas Ash (Anderson School of Management, UCLA), Etienne Bacher
482 (Luxembourg Institute of Socio-Economic Research), Sebastian Bachler (Univer-
483 sity of Innsbruck), Félix Bacon (Laval University), Manuel Bagues (University of
484 Warwick), Timea Balogh (UC Davis), Alisher Batmanov (UC San Diego), Mara
485 Barschkett (University of Bonn, IZA & DIW Berlin), B. Kaan Basdil (Risk Soft-
486 ware Technologies), Jaromír Baxa (Institute of Economic Studies, Faculty of Social
487 Sciences, Charles University, and Institute of Information Theory and Automa-
488 tion AS CR), Sascha Becker (University of Warwick and Monash University),
489 Monica Beeder (University of Southampton), Louis-Philippe Beland (Carleton Uni-
490 versity), Abdel-Hamid Bello (University of Montreal), Daniel Benenson Markovits
491 (Columbia University), Grant Benjamin (University of Toronto), Thomas Berg-
492 eron (Université de Montréal), Moussa P. Blimpo (University of Toronto), Marco
493 Binetti (Institute of Intercultural and International Studies, University of Bremen),
494 Carl Bonander (Karlstad Business School, Karlstad University), Joseph Bonneau
495 (UC Davis), Endre Borbáth (Ruprecht-Karls-Universität Heidelberg), Nicolai Top-
496 stad Borgen (Centre for Research on Equality in Education, University of Oslo),
497 Solveig Topstad Borgen (University of Oslo), Jonathan Borowsky (University of
498 Minnesota), Thomas Brailey (University of Oxford), Ryan Briggs (University of
499 Guelph), Elisa Brini (University of Florence), Myriam Brown (Laval University),
500 Martin Brun (Finnish Centre of Excellence in Tax Systems Research, Tampere
501 University), Stephan Bruns (Hasselt University, INCHER Kassel, METRICS Stan-
502 ford), Nino Buliskeria (Nazarbayev University), Andrea Calef (University College
503 London, School of Management), Alistair Cameron (Monash University), Pamela
504 Campa (Stockholm Institute of Transition Economics), Santiago Campos-Rodríguez
505 (University of California, Irvine), Giulio Giacomo Cantone (“Magna Graecia” Uni-
506 versity of Catanzaro), Fenella Carpena (Oslo Business School, Oslo Metropolitan
507 University), Perry Carter (NYU Abu Dhabi), Paul Castañeda Dower (University
508 of Wisconsin-Madison), Ondrej Cestek (Masaryk University), Jill Caviglia-Harris
509 (Salisbury University), Gabriella Chauca Strand (Institute of Medicine, Univer-
510 sity of Gothenburg), Shi Chen (School of Economics, Zhejiang University), Sya
511 In Chzhen (University of East Anglia), Jong Chung (Auburn University), Jason

512 Collins (University of Technology Sydney), Alexander Coppock (Northwestern Uni-
513 versity), Hugo Cordeau (University of Toronto), Ben Couillard (University of
514 Toronto), Jonathan Crechet (University of Ottawa), Lorenzo Crippa (University of
515 Strathclyde), Jeanne Cui (Beijing Normal University), Christian Czymara (Nether-
516 lands Interdisciplinary Demographic Institute), Haley Daarstad (UC Davis), Danh
517 Chi Dao (Queen’s University), Daniel Dao (University of Oxford), Marco David
518 Schmandt (TU Berlin), Astrid de Linde (University of Oslo), Lucas De Melo (Uni-
519 versity of Nottingham, NICEP), Lachlan Deer (University of Melbourne), Alexandra
520 de Gendre (The University of Melbourne), Micole De Vera (Banco de España),
521 Velichka Dimitrova (Social Research Institute, University College London), Jan
522 Fabian Dollbaum (University College Dublin), Jan Matti Dollbaum (University of
523 Fribourg and LMU Munich), Michael Donnelly (University of Toronto), Luu Duc
524 Toan Huynh (Queen Mary University of London), Tsvetomira Dumbalska (Univer-
525 sity of Oxford), Jamie Duncan (University of Toronto), Kiet Tuan Duong (University
526 of York), Yannick Dupraz (Paris Dauphine University, PSL University, LEDA,
527 CNRS, IRD), Thibaut Duprey (Bank of Canada), Christoph Dworschak (German
528 Institute for Development Evaluation & University of York), Sigmund Ellingsrud
529 (BI Norwegian Business School), Ali Elminejad (Nazarbayev University), Yasmine
530 Eissa (The American University in Cairo), Andrea Erhart (University of Innsbruck),
531 Giulian Etingin-Frati (ETH Zurich), Elaheh Fatemipour (University of Warwick),
532 Alexa Federice (UC Davis), Jan Feld (Victoria University of Wellington), Guidon
533 Fenig (University of Ottawa), Lenka Fiala (University of Ottawa, Tilburg University
534 and Institute for Replication), Mojtaba Firouzjaeiangalougah (Masaryk University),
535 Erlend Fleisje (Oslo Economics), Alexandre Fortier-Chouinard (Université Laval),
536 Julia Francesca Engel (Kiel University), Nadjim Fréchet (Concordia University),
537 Reid Fortier (VisualAIM), Tilman Fries (LMU Munich), Michael James Frith (Uni-
538 versity of Edinburgh), Jacopo Gabani (World Bank & Centre for health economics,
539 university of York), Thomas Galipeau (University of Toronto), Sebastián Galle-
540 gos (UAI Business School), Areez Gangji (Independent Researcher), Xiaoying Gao
541 (University of York), Cloé Garnache (Oslo Metropolitan University), Attila Gáspár
542 (ELTE Centre for Economic and Regional Studies; Central European University),
543 Romain Gauriot (Deakin University), Evelina Gavrilova (NHH Norwegian School
544 of Economics), Arijit Ghosh (RWI - Leibniz Institute for Economic Research), Gar-
545 reth Gibney (University of Galway), Grant Gibson (Canadian Research Data Centre
546 Network and McMaster University), Geir Godager (University of Oslo), Leonard
547 Goff (University of Calgary), Da Gong (State University of New York, Geneseo),
548 Javier González (Southern Methodist University), Jeremy D. Gretton (Ontario Pub-
549 lic Service’s Behavioural Insights Unit), Cristina Griffa (University of Chile), Idaliya
550 Grigoryeva (UC San Diego), Maja Grøtting (The Norwegian Institute of Public
551 Health), Eric Guntermann (UC Berkeley), Jiaqi Guo (University of Birmingham),
552 Alexi Gugushvili (University of Oslo), Hooman Habibnia (WU Vienna University of
553 Economics and Business), Sonja Häffner (Peace Research Institute Oslo), Jonathan
554 D. Hall (University of Alabama), Olle Hammar (Linnaeus University and Institute
555 for Futures Studies), Amund Hanson Kordt (University of Oslo), Barry Hashimoto
556 (Independent), Jonathan S. Hartley (Stanford University), Carina I. Hausladen

557 (University of Konstanz), Tomáš Havránek (Institute of Economic Studies, Fac-
558 ulty of Social Sciences, Charles University; and Faculty of International Relations,
559 Prague University of Economics and Business), Harry He (University of California,
560 San Diego), Matthew Hepplewhite (University of Oxford), Mario Herrera-Rodriguez
561 (CREST-Ecole Polytechnique; Programa Estado de la Nacion), Felix Heuer (RWI
562 – Leibniz Institute for Economic Research), Anthony Heyes (University of Birm-
563 ingham), Anson T. Y. Ho (Toronto Metropolitan University), Jonathan Holmes
564 (University of Ottawa), Armando Holzknicht (University of Innsbruck), Yu-Hsiang
565 Dexter Hsu (University of California, Davis), Shiang-Hung Hu (California Insti-
566 tute of Technology), Yu-Shiuan Huang (National Chengchi University), Mathias
567 Huebener (Federal Institute for Population Research (BiB)), Christoph Huber
568 (Aalto University), Kim P. Huynh (Indiana University, Department of Economics;
569 and Université d’Orléans and the Laboratoire d’Économie d’Orléans), Zuzana Irsova
570 (Institute of Economic Studies, Faculty of Social Sciences, Charles University, and
571 Anglo-American University, Prague), Ozan Isler (The University of Queensland),
572 Niklas Jakobsson (Karlstad University & FBK-IRVAPP), Raphaël Jananji (Univer-
573 sité de Montréal), Tharaka A. Jayalath (Global Water Security Center), Michael
574 Jetter (University of Western Australia), Jenny John (University of Ottawa), Rachel
575 Joy Forshaw (Heriot-Watt University), Felipe Juan (Howard University), Valon
576 Kadriu (University of Kassel and INCHER), Sunny Karim (Carleton University),
577 Edmund Kelly (University of Oxford), Duy Khanh Hoang Dang (University College
578 London), Tazia Khushboo (University of Calgary), Jin Kim (Chinese University of
579 Hong Kong), Gustav Kjellsson (Centre for Health Governance & HEPER, School of
580 Public Health & Community Medicine, University of Gothenburg), Anders Kjelsrud
581 (Oslo Metropolitan University), Jori Korpershoek (Erasmus University Rotterdam),
582 Andreas Kotsadam (Ragnar Frisch Centre for Economic Research), Lewis Krashin-
583 sky (Princeton University), Suranjana Kundu (World Inequality Lab, Paris School
584 of Economics), Alexander Kustov (University of Notre Dame), Nurlan Lalayev
585 (University of Warwick), Audrée Langlois (Université Laval), Jill Laufer (UC Center
586 Sacramento (UC Davis)), Blake Lee-Whiting (University of Western Ontario),
587 Andreas Leibing (Dresden University of Technology), Gabriel Lenz (UC Berkeley),
588 Joel Levin (UC San Diego), Peng Li (University of Bath), Tongzhe Li (University of
589 Guelph), Yuchen Lin (University of Warwick), Goncalo Lima (European University
590 Institute and University of Bologna), Ariel Listo (University of Maryland), Dan Liu
591 (Australian National University), Xuewen Lu (University of Calgary), Elvina Luk-
592 manova (New Economic School), Alex Luscombe (Government of Canada), Lester
593 R. Lusher (University of Pittsburgh), Ke Lyu (University of Nevada, Reno), Hai
594 Ma (McGill University), Nicolas Mäder (Knauss School of Business, University of
595 San Diego), Clifton Makate (Norwegian University of Life Sciences and Norwegian
596 Geotechnical Institute), Alice Malmberg (UC Davis), Adit Maitra (The University of
597 Melbourne), Marco Mandas (University of Cagliari), Jan Marcus (Freie Universität
598 Berlin), Shushanik Margaryan (University of Potsdam), Lili Márk (Central Euro-
599 pean University), Diego Marino Fages (Durham University), Andres Martignano
600 (University of Nottingham), Abigail Marsh (Finance Canada), Isabella Masetto

601 (London School of Economics and Political Science), Anthony McCanny (Univer-
602 sity of Toronto), Emma McManus (Health Organisation, Policy and Economics,
603 The University of Manchester), Ryan McWay (University of Minnesota), Lennard
604 Metson (London School of Economics and Political Science), Fabian Mierisch (Inde-
605 pendent Researcher), Jonas Minet Kinge (University of Oslo), Sumit Mishra (Krea
606 University), Myra Mohnen (University of Ottawa), Jakob Möller (WU Vienna
607 University of Economics and Business), Rosalie Montambeault (Université Laval),
608 Sébastien Montpetit (University of Warwick), Louis-Philippe Morin (University
609 of Ottawa), Todd Morris (University of Queensland), Scott Moser (University of
610 Nottingham, School of Politics and International Relations), Fabio Motoki (Uni-
611 versity of Texas Rio Grande Valley), Lucija Muehlenbachs (University of Calgary
612 and Resources for the Future), Andreea Musulan (University of Montreal, IVADO,
613 Mila), Marco Musumeci (University of Padova), Munirul Nabin (Deakin Univer-
614 sity), Karim Nchare (Vanderbilt University), Florian Neubauer (RWI - Leibniz
615 Institute for Economic Research), Quan M. P. Nguyen (University of Sussex), Tuan
616 Nguyen (Hasselt University), Viet Nguyen-Tien (London School of Economics), Ali
617 Niazi (University of Calgary), Giorgi Nikolaishvili (Wake Forest University), Ardyn
618 Nordstrom (Carleton University), Patrick Nüß (IWH Halle), Angela Odermatt (Uni-
619 versity of Oxford), Matt Olson (University of Pennsylvania Wharton), Henning Øien
620 (Department of Health Management and Health Economics, University of Oslo),
621 Tim Ölkens (Humboldt University zu Berlin), Miquel Oliver i Vert (Universitat de
622 Girona), Emre Oral (University of Mannheim), Christian Oswald (University of the
623 Bundeswehr Munich), Ali Ousman (McGill University), Ömer Özak (Department
624 of Economics, Southern Methodist University, IZA and GLO), Shubham Pandey
625 (Institute of Psychology, Osnabrück University), Alexandre Pavlov (Université de
626 Montréal), Martino Pelli (Asian Development Bank), Romeo Penheiro (University
627 of Houston), RyuGyung Park (Government Department at William & Mary), Eva
628 Pérez Martel (Universitat Autònoma de Barcelona), Jörg Ankel-Peters (RWI - Leib-
629 nitz Institute for Economic Research), Tereza Petrovičová (UCSD), Linh Phan (UC
630 Davis), Alexa Prettyman (Towson University), Jakub Procházka (Masaryk Univer-
631 sity), Aqila Putri (University of Maryland), Julian Quandt (WU Vienna University
632 of Economics and Business), Kangyu Qiu (McMaster University), Loan Quynh Thi
633 Nguyen (National Economics University), Andaleeb Rahman (Cornell University),
634 Carson H. Rea (Emory University), Adam Reiremo (Norwegian School of Eco-
635 nomics), Laëtitia Renée (Université de Montréal), Joseph Richardson (Lancaster
636 University), Nicholas Rivers (University of Ottawa), Bruno Rodrigues (Ministry
637 of Research and Higher Education, Luxembourg), William Roelofs (University of
638 Toronto), Tobias Roemer (University of Oxford), Ole Rogeberg (Ragnar Frisch Cen-
639 tre for Economic Research), Julian Rose (RWI - Leibniz Institute for Economic
640 Research), Andrew Roskos-Ewoldsen (UC Davis), Paul Rosmer (Humboldt Univer-
641 sity of Berlin & Berlin School of Economics), Barbara Sabada (Bank of Canada),
642 Soodeh Saberian (University of Manitoba), Nicolas Salamanca (The University of
643 Melbourne), Georg Sator (University of Nottingham & Institute for Advanced Stud-
644 ies Vienna), Daniel Scates (UC Davis), Elmar Schlüter (Justus Liebig University,
645 Giessen), Cameron Sells (Independent Researcher), Sharmi Sen (Monash University),

646 Ritika Sethi (University of Chicago), Anna Shcherbiak (WU Vienna University
647 of Economics and Business), Moyosore Sogaolu (GATE, Rotman, University of
648 Toronto), Matt Soosalu (Carleton University), Erik Ø. Sørensen (NHH Norwegian
649 School of Economics), Manali Sovani (Tufts University), Noah Spencer (University
650 of Toronto), Stefan Staubli (University of Calgary), Renske Stans (The Netherlands
651 Court of Audit), Anya Stewart (UC Davis), Felix Stips (Institute for Employment
652 Research (IAB)), Kieran Stockley (University of Nottingham), Stephenson Strobel
653 (McMaster University), Ethan Struby (Carleton College, Boston College, and Min-
654 nesota Supercomputing Institute), John Tang (Utrecht University), Idil Tannisever
655 (University of California, Irvine), Thomas Tao Yang (Australian National Univer-
656 sity), Ipek Tastan (University of Calgary), Dejan Tatić (WU Vienna University of
657 Economics and Business), Benjamin Tatlow (University of Nottingham), Féraud
658 Tehuisseu Seuyong (Université de Montréal), Rémi Thériault (New York Univer-
659 sity), Vincent Thivierge (University of Ottawa), Wenjie Tian (University of Ottawa),
660 Filip-Mihai Toma (Bucharest University of Economic Studies), Maddalena Totarelli
661 (Ifo Institute & Ludwig Maximilian University of Munich), Van-Anh Tran (Monash
662 University), Hung Truong (University of Ottawa), Nikita Tsoy (INSAIT, Sofia Uni-
663 versity), Kerem Tuzcuoglu (Amazon), Diego Ubfal (World Bank), Laura Villalobos
664 (Salisbury University), Julian Walterskirchen (University of Gothenburg), Joseph
665 Tao-yi Wang (Department of Economics and Taiwan Social Resilience Research
666 Center, National Taiwan University), Vasudha Wattal (The University of Manch-
667 ester), Matthew D. Webb (Carleton University), Bryan Weber (College of Staten
668 Island - CUNY), Reinhard Weisser (University of the West of England), Wei-Chien
669 Weng (University of California, Davis), Christian Westheide (Stockholm Business
670 School, Stockholm University & Leibniz Institute for Financial Research SAFE),
671 Kimberly White (Ludwig Maximilian University of Munich), Jacob Winter (Uni-
672 versity of Toronto), Timo Wochner (ETH Zurich & KOF Institute), Matt Woerman
673 (Colorado State University), Jared Wong (Yale University), Ritchie Woodard (Uni-
674 versity of East Anglia), Marcin Wroński (SGH Warsaw School of Economics),
675 Gustav Chung Yang (Harvard University), Myra Yazbeck (University of Ottawa),
676 Luther Yap (National University of Singapore), Kareman Yassin (Hitotsubashi Uni-
677 versity), Hao Ye (University of Pennsylvania / Community for Rigor), Jin Young
678 Yoon (Queen's University), Chris Yurris (McGill University), Tahreen Zahra (Car-
679 leton University), Mirela Zaneva (University of Oxford), Aline Zayat (University of
680 Ottawa), Jonathan Zhang (Duke University and Sanford School of Public Policy),
681 Ziwei Zhao (University of Lausanne and Swiss Finance Institute), Yaolang Zhong
682 (University of Warwick)

683 **Computational reproducibility:**

684 Abel Brodeur (University of Ottawa and Institute for Replication), Joanne Haddad
685 (Universitat Autònoma de Barcelona), Pu Sun (Dongbei University of Finance and
686 Economics)

687 **Local organizer Replication Games:**

688 Marie Connolly (UQAM), Romain Gauriot (Deakin University), Leonard Goff
689 (University of Calgary), Christoph Huber (Aalto University), Andreas Kotsadam

692 Figure Legends

693 Figure 1: Robustness Rate. Legend: Robustness rate for ... **Left panel:** ... originally
694 statistically significant research **Second panel:** ... in economics **Third panel:** ... in
695 political science **Right panel:** ... originally statistically *insignificant* research **All**
696 **panels:** Squares, circles, and triangles represent proportions, with 95% Clopper-
697 Pearson confidence intervals presented in whiskers. Red squares represent full sample.
698 Green circles represent economics subsample. Blue triangles represent political science
699 subsample. Each group of three estimates represent different types of re-analysis,
700 non-mutually exclusive. The first 8 groups do not include re-analyses that use new
701 data (replication), while the last one does. The first estimate group contains all
702 types of re-analysis, then all types of re-analysis in economics, then all types of
703 re-analysis in political science. The second represents re-analyses which changed the
704 control variables, e.g., by adding or re-defining them. The third represents re-analyses
705 which changed the dependent variable, e.g., by employing a different standardization
706 or binarization. The fourth represents re-analyses which changed the estimation
707 method, e.g., by adjusting a matching procedure. The fifth represents re-analyses
708 which changed the inference method, e.g., changed the level on which standard errors
709 are clustered. The sixth represents re-analyses which changed the main independent
710 variable, e.g., by taking into account treatment intensity. The seventh represents re-
711 analyses which changed the sample, e.g., by excluding outliers. The eighth represents
712 re-analyses which changed the weights applied, or applied weights for the first time.
713 The last represents replicability rates for re-analyses that introduced new data, e.g.,
714 comparable outcomes from more recent survey waves.

715
716 Figure 2. Statistical Significance of Publication and Re-analysis. Legend: **Top**
717 **histogram:** Distribution of publication tests of significance. T-statistics over 4
718 truncated for exposition. The histogram's bars are of width 0.14, with exactly 14
719 bars between 0 and the statistical threshold of $t = 1.96$ (corresponding to statistical
720 significance at the 5% level). **Right histogram:** Distribution of re-analysis tests of
721 significance. T-statistics over 4 truncated for exposition. **Scatterplot:** Each marker
722 is a pair of test statistics, an originally published test statistic (horizontal value) and
723 an associated re-analysis test statistic (vertical value). If the original and re-analysis
724 test statistics were identical, this scatterplot would follow the 45 degree line. As either
725 axis represents statistical significance, we have bifurcated each with a line at $t=1.96$,
726 representing statistical significance at the 5% threshold. **Blue circles** indicate an
727 originally statistically significant statistic that is also statistically significant under
728 re-analysis. Represents 50% of sample. **Red triangles** indicate originally significant
729 test statistics that are no longer statistically significant under re-analysis. Represents
730 14% of sample. **Green squares** indicate originally statistically insignificant test
731 statistics that are the same under re-analysis. Represents 27% of sample. **Purple**
732 **diamonds** indicate originally statistically insignificant test statistics that become

733 statistically significant under re-analysis. Represents 3% of sample. **Not displayed**
734 Not displayed are the 6% of test statistics that represent a sign reversal between the
735 originally estimated effect and the effect estimated under re-analysis.

736

737 Figure 3. Robustness Rate Determinants. Legend: Six independent teams
738 answered twelve questions of the re-analysis database. Each bar represents a different
739 question. **Left panel:** “Does reproducibility of an originally statistically significant
740 result depend on...” **Right panel:** “Does reproducibility of an originally statistically
741 *insignificant* result depend on...” **Both panels:** where the first bar represents “...
742 the reproducers’ experience at coding.” **Blue, patterned outline** indicates the pro-
743 portion of teams that indicated a negative and statistically significant relationship, in
744 whichever manner the team defined so in their analysis. **Gray, no outline** indicates
745 the proportion of teams that indicated a statistically insignificant relationship, where
746 left of the zero line indicates negative and right of the zero line indicates positive.
747 **Red, solid outline** indicates the proportion of teams that indicated a statistically
748 significant and positive relationship. All teams equally weighted.

749

750 Figure 4. Effect Size of Publication and Re-analysis. Legend: **Top histogram:**
751 Distribution of originally published effect size standardized by the average effect size
752 within a published article. **Right histogram:** Distribution of re-analysis published
753 effect size standardized by the average effect size within a published article. **Scatter-**
754 **plot:** Each marker is a pair of effect sizes, the originally published effect size (horizontal
755 value) and an associated re-analysis effect size (vertical value). If an originally esti-
756 mated and re-analysis effect size were of similar magnitude (and sign), the markers
757 would gather tightly around the 45 degree line passing through the origin **Blue cir-**
758 **cles** indicate effect sizes which are similar (between 50% to 200% of original effect
759 size) under re-analysis. Represents 69% of sample. **Red diamonds** indicate effect
760 size estimates which switch sign under re-analysis. Represents 6% of sample. **Orange**
761 **triangles** indicate effect size estimates which are 50% or less their original magni-
762 tude under re-analysis. Represents 9% of sample. **Purple squares** indicate effect size
763 estimates which are double or larger than their original magnitude under re-analysis.
764 Represents 16% of sample.

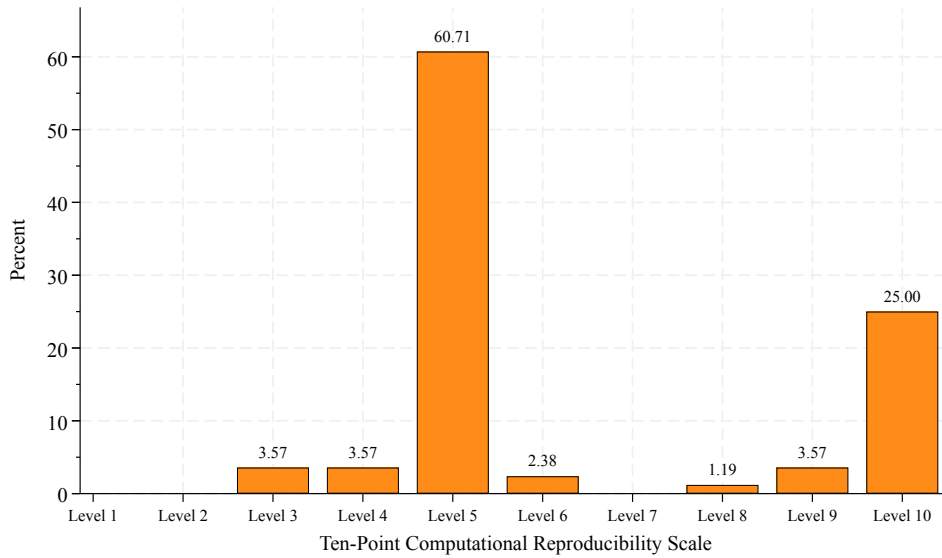
765 Extended Data Figure Legends

766 Extended Data Figure 1: 10-Point Computationally Reproducibility Score. Legend:
767 Each team assigned a reproducibility score on a scale of one to ten to the paper
768 reproduced. See Supplementary Materials for a description of each score. Level 10
769 (L10) means that all necessary materials are available and produce consistent results
770 with those presented in the paper, while level 5 (L5) means that analytic data sets
771 and analysis code are available and they produce the same results as presented in the
772 paper.

773

774 Extended Data Figure 2: Reasons Select Paper? (Select all which apply). Legend:
775 Data collected *via* survey of our reproducers after completing their reports. This

Fig. 5: 10-Point Computationally Reproducibility Score



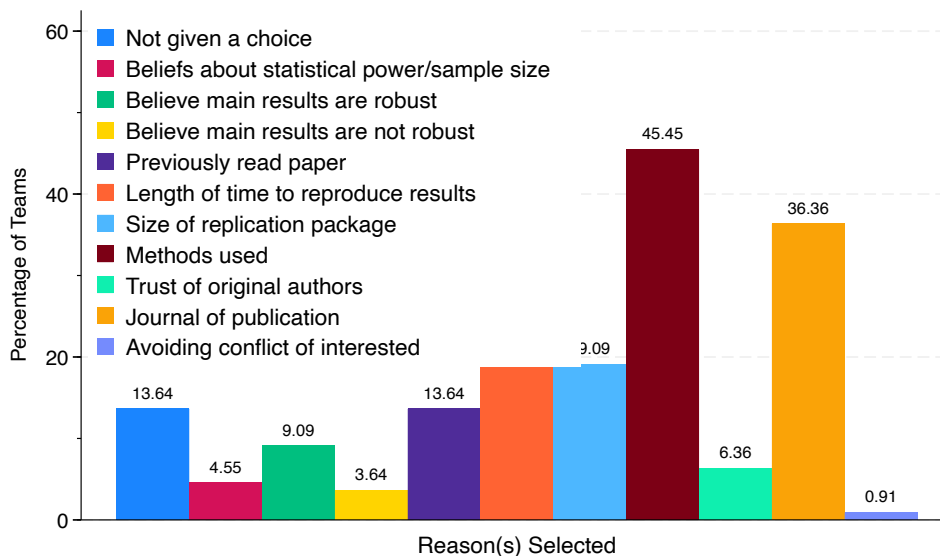
Notes: Each team assigned a reproducibility score on a scale of one to ten to the paper reproduced. See Supplementary Materials for a description of each score. Level 10 (L10) means that all necessary materials are available and produce consistent results with those presented in the paper, while level 5 (L5) means that analytic data sets and analysis code are available and they produce the same results as presented in the paper.

776 figure illustrates the responses to the question: “For what reasons did you select your
 777 specific paper to reproduce and/or replicate from the list of papers provided?
 778

779 Extended Data Figure 3: Percentage of Papers with a Replication Folder. Legend:
 780 The total sample is 1150 papers with 120 papers per year from 2019 to 2023 and 110
 781 papers per year from 2018 to 2014. Each journal has 10 papers per year except *Ameri-*
 782 *can Economic Review: Insights* which only formally became a journal in 2019 (and
 783 are omitted in earlier years). The journals sampled over correspond to those used in
 784 the manuscript’s main analysis, three from political science and nine from economics.
 785 The political science journals include: *American Journal of Political Science*, *Ameri-*
 786 *can Political Science Review*, and *Journal of Politics*. The economics journals include:
 787 *American Economic Review*, *Quarterly Journal of Economics*, *Review of Economic*
 788 *Studies*, *Journal of Political Economy*, *American Economic Journal: Macroeconomics*,
 789 *American Economic Journal: Applied Economics*, *American Economic Journal: Eco-*
 790 *nomic Policy*, *American Economic Review: Insights*, *Economic Journal*.
 791

792 Extended Data Figure 4: Percentage of Papers with a Replication Folder by
 793 Discipline. Legend: Panel (a) is for papers published in economics journals where

Fig. 6: For what reasons did you select your specific paper to reproduce and/or replicate from the list of papers provided? (Select all which apply)



Notes: Data collected *via* survey of our reproducers after completing their reports. This figure illustrates the responses to the question: “For what reasons did you select your specific paper to reproduce and/or replicate from the list of papers provided?”

794 Panel (b) is for papers published in political science. The total sample is the same
 795 as Extended Data Figure 3 is 1150 papers, where 850 papers are in the economics
 796 sample and 300 papers are in the political science sample.

797

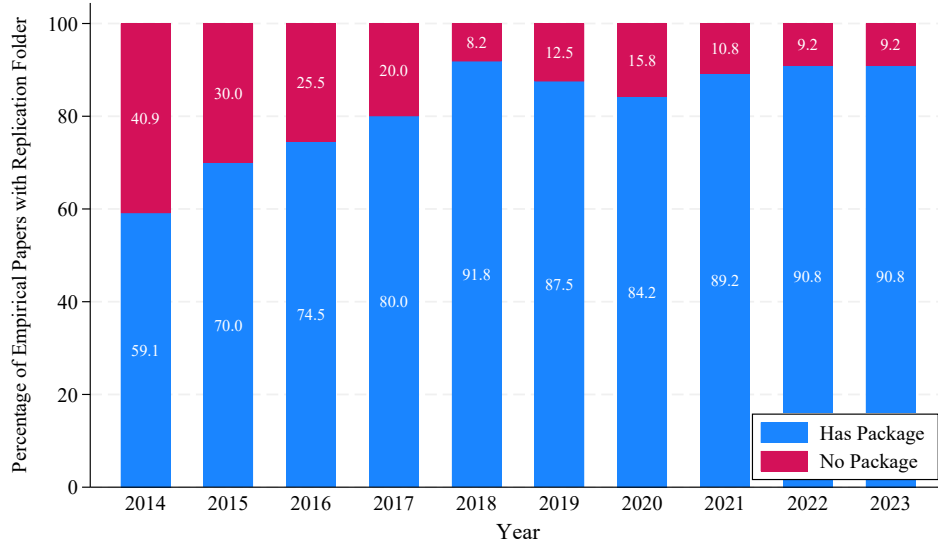
798 Extended Data Figure 5: Percentage Replication Folders’ with Contents Condi-
 799 tional on they Should Have a Replication Folder. Legend: Each subfigure represents
 800 the proportion of the replication folders which affirmatively (“Yes”) contained the
 801 variable (displayed as the title). The “Not Yes” in the legend corresponds to those
 802 replication folders which did not affirm (“No”) or had only “Some” of the required
 803 contents. Each sample is over those observations where categories are applicable (*i.e.*
 804 not all replication packages require the same contents).

805

806 Extended Data Figure 6: Percentage Replication Folders’ with Contents Condi-
 807 tional on they Should Have a Replication Folder. Legend: Each subfigure represents
 808 the proportion of the replication folders which affirmatively (“Yes”) contained the
 809 variable (displayed as the title). The “Not Yes” in the legend corresponds to those
 810 replication folders which did not affirm (“No”) or had only “Some” of the required
 811 contents. Each sample is over those observations where categories are applicable (*i.e.*
 812 not all replication packages require the same contents).

813

Fig. 7: Percentage of Papers with a Replication Folder



The total sample is 1150 papers with 120 papers per year from 2019 to 2023 and 110 papers per year from 2018 to 2014. Each journal has 10 papers per year except *American Economic Review: Insights* which only formally became a journal in 2019 (and are omitted in earlier years). The journals sampled over correspond to those used in the manuscript’s main analysis, three from political science and nine from economics. The political science journals include: *American Journal of Political Science*, *American Political Science Review*, and *Journal of Politics*. The economics journals include: *American Economic Review*, *Quarterly Journal of Economics*, *Review of Economic Studies*, *Journal of Political Economy*, *American Economic Journal: Macroeconomics*, *American Economic Journal: Applied Economics*, *American Economic Journal: Economic Policy*, *American Economic Review: Insights*, *Economic Journal*.

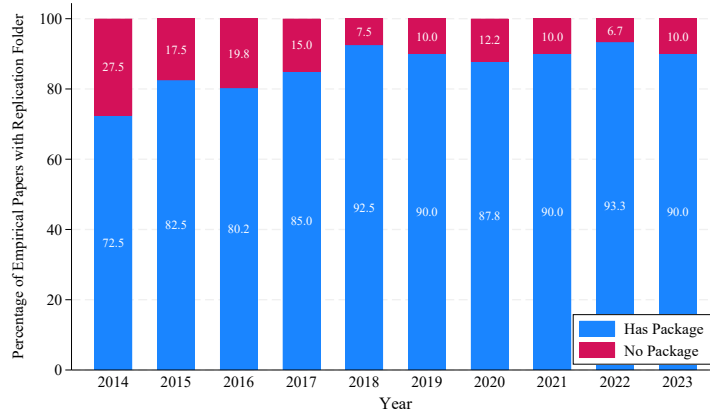
814 Extended Data Figure 7: Reasons Unable Conduct Robustness Checks. Legend:
 815 This Figure illustrates the share of teams who were unable to perform robustness
 816 checks (top-left), replications (top-right), key variable recodes (bottom-right) or
 817 extensions (bottom-left) for various reasons represented by the different coloured bars.

818
 819 Extended Data Figure 8: Distributions of t-Statistics for Original Studies and Re-
 820 Analyses. Legend: The top panels display a histogram of test statistics for $t \in [0, 5]$,
 821 with bins of width 0.1. The top left panel includes all original studies in our data set.
 822 The top right panel includes all re-analysis estimates in our data set. Vertical refer-
 823 ence lines are displayed at conventional two-tailed significance levels. We superimpose
 824 an Epanechnikov kernel (which includes renormalization at 0). The bottom figures
 825 display histograms of test statistics for p-values $\in [0.0025, 0.1500]$, with bins of width
 826 0.0025, among original studies and those from re-analyses, respectively.

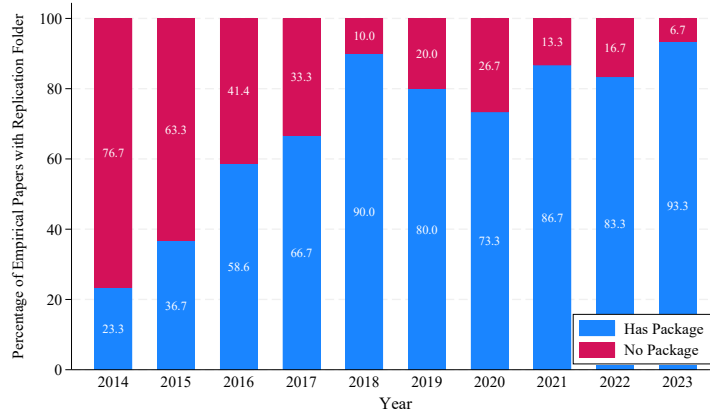
827

Fig. 8: Percentage of Papers with a Replication Folder by Discipline

(a) Economics



(b) Political Science

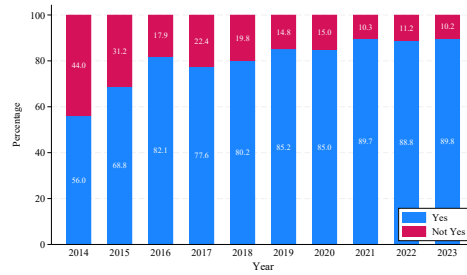


Panel (a) is for papers published in economics journals where Panel (b) is for papers published in political science. The total sample is the same as Figure 7 is 1150 papers, where 850 papers are in the economics sample and 300 papers are in the political science sample.

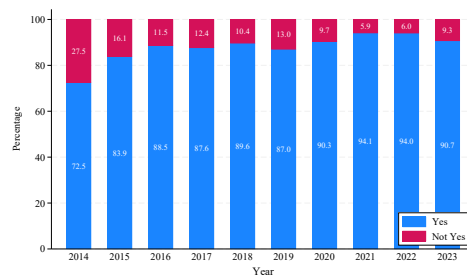
828 Extended Data Figure 9: Distributions of t -Statistics and p -values by Field.
829 Legend: We restrict the sample to articles published in the indicated field, journals.
830 Top panels display histograms of test statistics for $t \in [0, 5]$, with bins of width 0.1
831 respectively. Vertical reference lines are displayed at conventional two-tailed significance
832 levels. We superimpose an Epanechnikov kernel density curve (which includes
833 renormalization at 0). Bottom panels display histograms of test statistics for p -values
834 $\in [0.0025, 0.1500]$, with bins of width 0.0025.
835

Fig. 9: Percentage Replication Folders' with Contents Conditional on they Should Have a Replication Folder

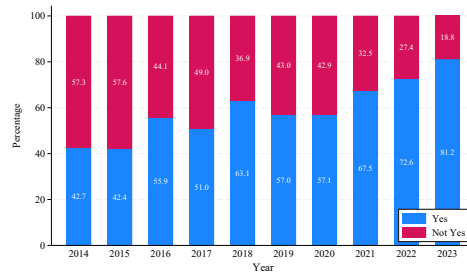
(a) README



(b) Analysis Code



(c) Cleaning Code

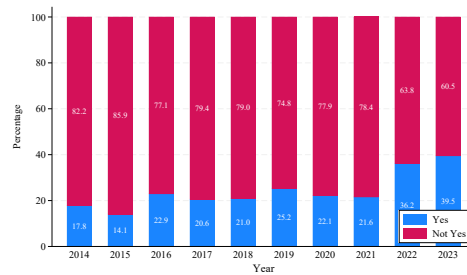


Each subfigure represents the proportion of the replication folders which affirmatively (“Yes”) contained the variable (displayed as the title). The “Not Yes” in the legend corresponds to those replication folders which did not affirm (“No”) or had only “Some” of the required contents. Each sample is over those observations where categories are applicable (*i.e.* not all replication packages require the same contents).

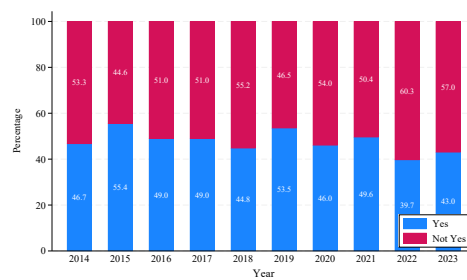
836 Extended Data Figure 10: Relative Reproduced Effect Size. Legend: 48% of rel-
 837 ative effect sizes are exactly equal to or greater than 1. This figure illustrates the
 838 ratio of re-analysis estimates and original estimates. The standardized effect sizes
 839 are normalized so that 1 equals the original effect size. A positive value indicates that the

Fig. 10: Percentage Replication Folders' with Contents Conditional on they Should Have a Replication Folder

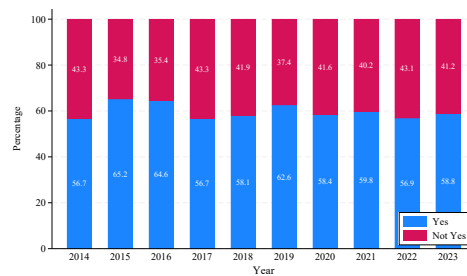
(a) Raw Data



(b) Final Data



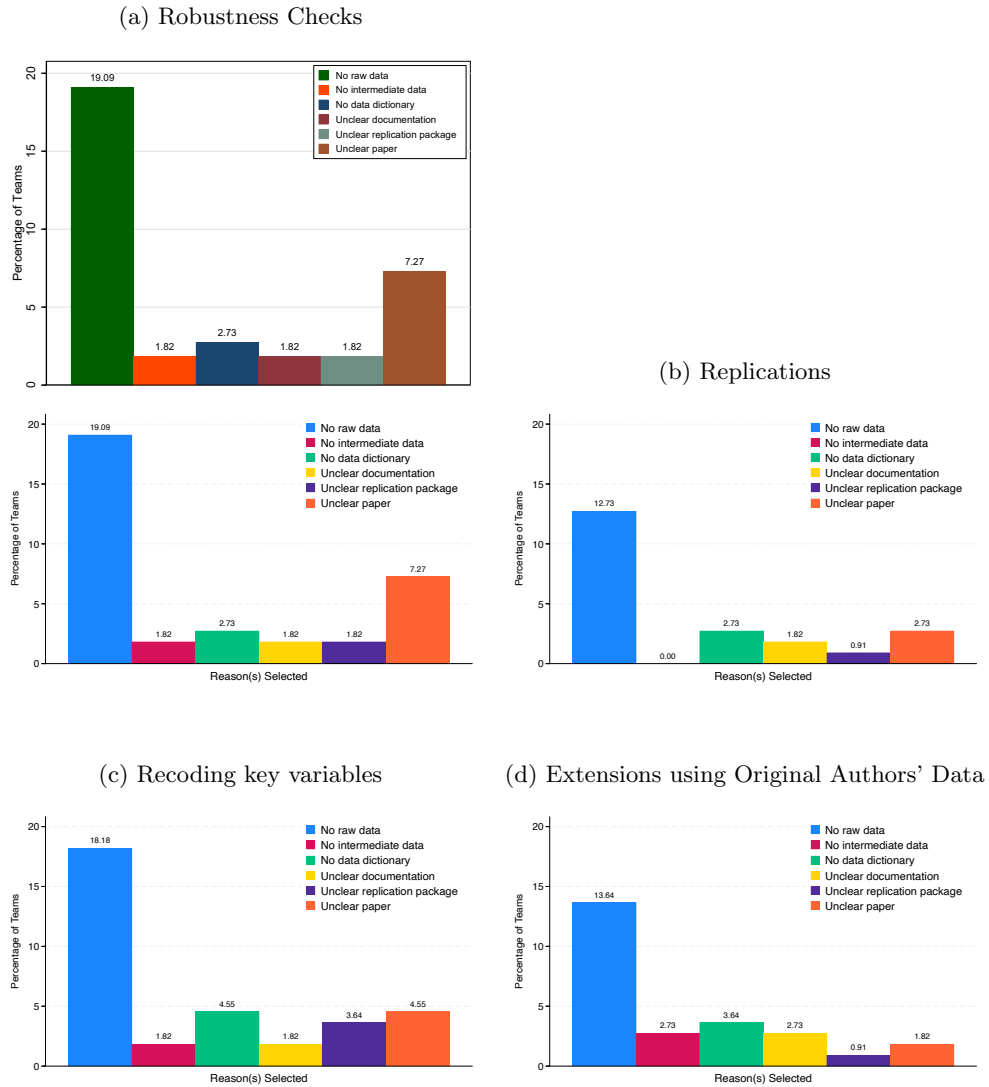
(c) Final Data + Raw or Intermediate Data with Cleaning Code



Each subfigure represents the proportion of the replication folders which affirmatively (“Yes”) contained the variable (displayed as the title). The “Not Yes” in the legend corresponds to those replication folders which did not affirm (“No”) or had only “Some” of the required contents. Each sample is over those observations where categories are applicable (*i.e.* not all replication packages require the same contents).

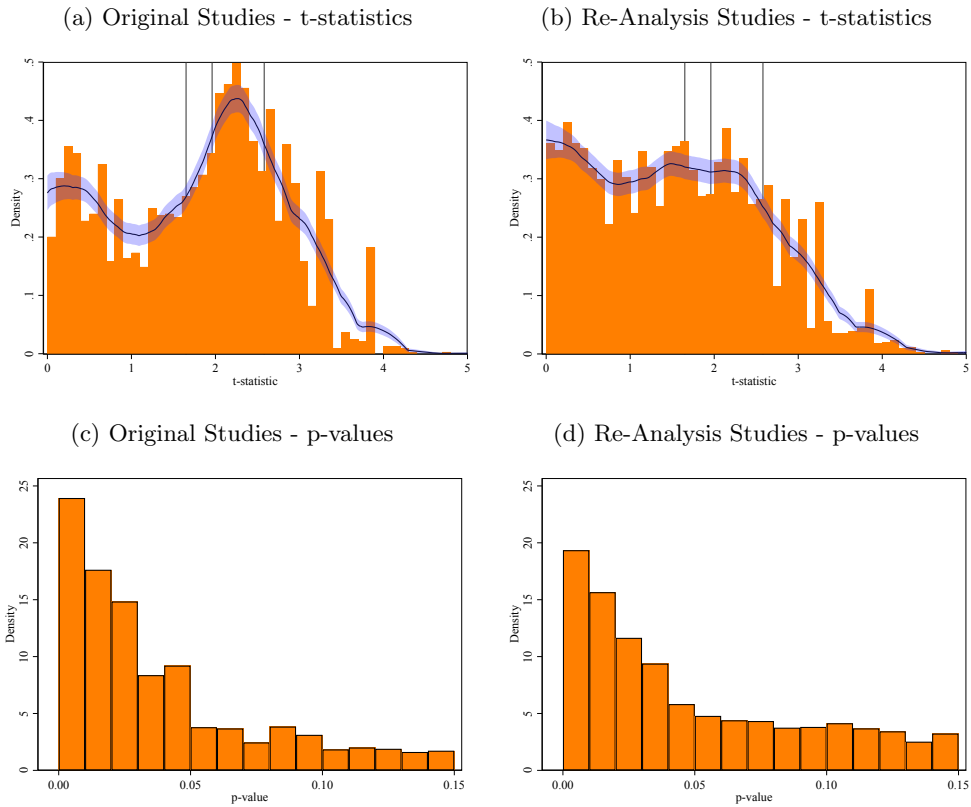
840 re-analysis estimate is in the same direction as in the original study. A negative value
 841 indicates that the re-analysis estimate is not in the same direction as in the original

Fig. 11: For which of the following reasons were you unable to conduct robustness checks, recoding exercises, extensions, or a replication using new data, prior to communications with the original authors? (Select all which apply)



Notes: This Figure illustrates the share of teams who were unable to perform robustness checks (top-left), replications (top-right), key variable recodes (bottom-right) or extensions (bottom-left) for various reasons represented by the different coloured bars.

Fig. 12: Distributions of t-Statistics for Original Studies and Re-Analyses



Notes: The top panels display a histogram of test statistics for $t \in [0, 5]$, with bins of width 0.1. The top left panel includes all original studies in our data set. The top right panel includes all re-analysis estimates in our data set. Vertical reference lines are displayed at conventional two-tailed significance levels. We superimpose an Epanechnikov kernel (which includes renormalization at 0). The bottom figures display histograms of test statistics for p-values $\in [0.0025, 0.1500]$, with bins of width 0.0025, among original studies and those from re-analyses, respectively.

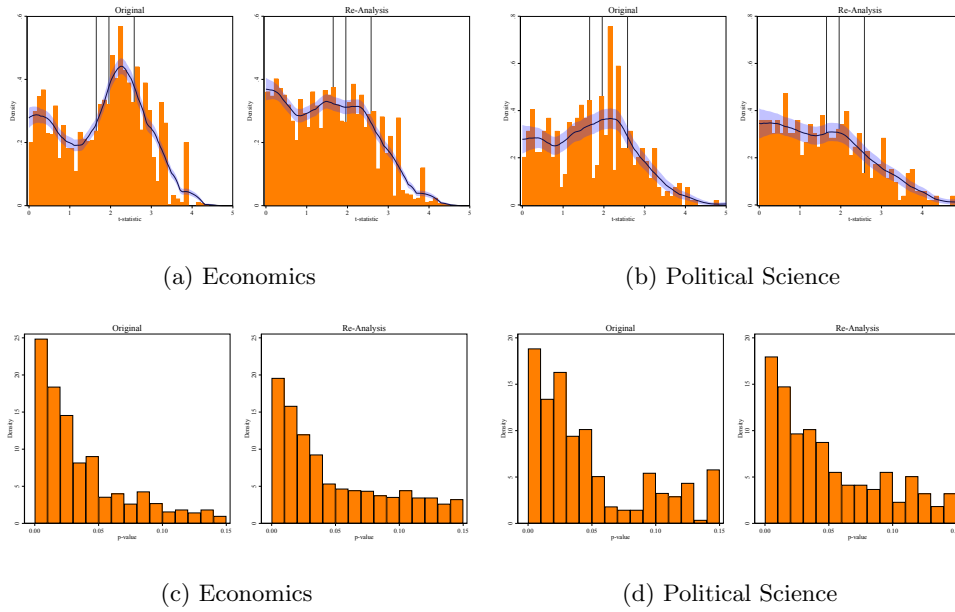
842 study. Outliers (3%) are excluded for visibility.

843

844 References

- 845 1. Vazire, S. Quality Uncertainty Erodes Trust in Science. *Collabra: Psychology* **3**,
846 1 (2017).

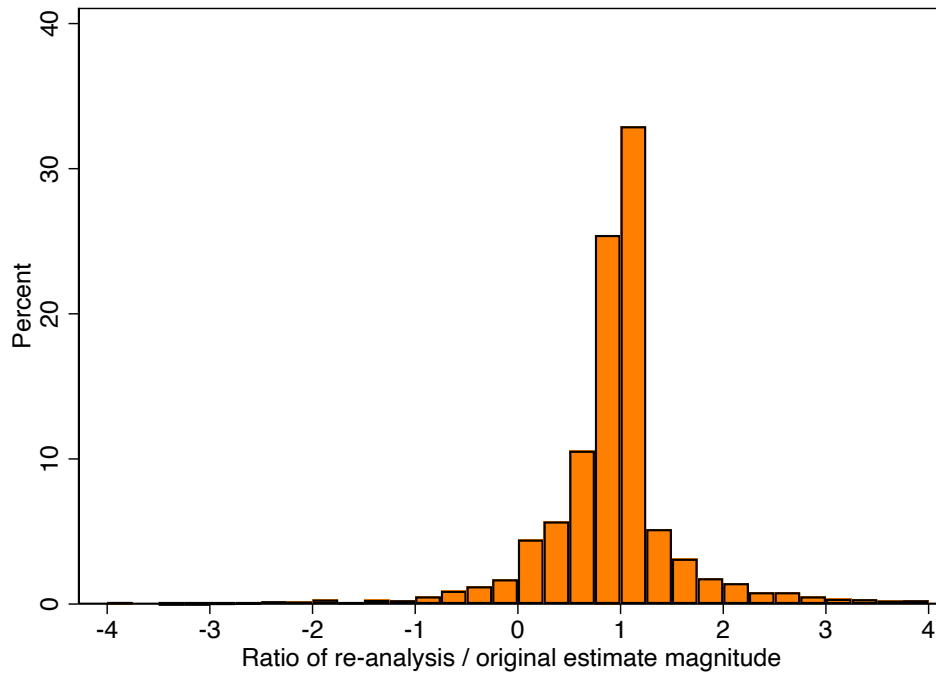
Fig. 13: Distributions of t -Statistics and p -values by Field



Notes: We restrict the sample to articles published in the indicated field. journals. Top panels display histograms of test statistics for $t \in [0, 5]$, with bins of width 0.1 respectively. Vertical reference lines are displayed at conventional two-tailed significance levels. We superimpose an Epanechnikov kernel density curve (which includes renormalization at 0). Bottom panels display histograms of test statistics for p -values $\in [0.0025, 0.1500]$, with bins of width 0.0025.

847 2. Donoho, D. L., Maleki, A., Rahman, I. U., Shahram, M. & Stodden, V. Reproducible
848 research in computational harmonic analysis. *Computing in Science & Engineering* **11**, 8–18 (2008).
849
850 3. King, G. Replication, Replication. *PS: Political Science & Politics* **28**, 444–452
851 (1995).
852 4. Goodman, S. N., Fanelli, D. & Ioannidis, J. P. What does research reproducibility
853 mean? *Science Translational Medicine* **8**, 341ps12–341ps12 (2016).
854 5. Marcoci, A. *et al.* Predicting the replicability of social and behavioural science
855 claims in COVID-19 preprints. *Nature human behaviour* **9**, 287–304 (2025).
856 6. Milkowski, M., Hensel, W. M. & Hohol, M. Replicability or reproducibility? On
857 the replication crisis in computational neuroscience and sharing only relevant
858 detail. *Journal of Computational Neuroscience* **45**, 163–172 (2018).
859 7. Moonesinghe, R., Houry, M. J. & Janssens, A. C. J. W. Most Published
860 Research Findings Are False—but a Little Replication Goes a Long Way. *PLoS*
861 *Medicine* **4**, e28 (2007).
862 8. National Academies of Sciences, Engineering, and Medicine. *Reproducibility and*
863 *Replicability in Science* <https://www.nap.edu/catalog/25303> (National
864 Academies Press, 2019).

Fig. 14: Relative Reproduced Effect Size



Notes: 48% of relative effect sizes are exactly equal to or greater than 1. This figure illustrates the ratio of re-analysis estimates and original estimates. The standardized effect sizes are normalized so that 1 equals the original effect size. A positive value indicates that the re-analysis estimate is in the same direction as in the original study. A negative value indicates that the re-analysis estimate is not in the same direction as in the original study. Outliers (3%) are excluded for visibility.

- 865 9. Peterson, D. & Panofsky, A. Self-Correction in Science: The Diagnostic and
866 Integrative Motives for Replication. *Social Studies of Science* **51**, 583–605 (2021).
- 867 10. Pérignon, C., Gadouche, K., Hurlin, C., Silberman, R. & Debonnel, E. Certify
868 reproducibility with confidential data. *Science* **365**, 127–128 (2019).
- 869 11. Brandon, A. & List, J. A. Markets for Replication. *Proceedings of the National
870 Academy of Sciences* **112**, 15267–15268 (2015).
- 871 12. Freese, J. & Peterson, D. Replication in Social Science. *Annual Review of
872 Sociology* **43**, 147–165 (2017).
- 873 13. Gertler, P., Galiani, S. & Romero, M. How to Make Replication the Norm. *Nature
874* **554**, 417–9 (2018).
- 875 14. Maniadis, Z. & Tufano, F. The Research Reproducibility Crisis and Economics
876 of Science. *Economic Journal* **127** (2017).
- 877 15. Munafò, M. R. *et al.* A Manifesto for Reproducible Science. *Nature Human
878 Behaviour* **1**, 1–9 (2017).

- 879 16. Nosek, B. A. *et al.* Replicability, Robustness, and Reproducibility in Psychological
880 Science. *Annual Review of Psychology* **73**, 719–748 (2022).
- 881 17. Askarov, Z., Doucouliagos, A., Doucouliagos, H. & Stanley, T. The Significance
882 of Data-sharing Policy. *Journal of the European Economic Association* **21**, 1191–
883 1226 (2023).
- 884 18. Brodeur, A., Cook, N. & Neisser, C. P-Hacking, Data Type and Data-Sharing
885 Policy. *Economic Journal* **134**, 985–1018 (2024).
- 886 19. Chang, A. C. & Li, P. Is Economics Research Replicable? Sixty Published Papers
887 From Thirteen Journals Say "Often Not". *Critical Finance Review* **11**, 185–206
888 (2022).
- 889 20. Christensen, G. & Miguel, E. Transparency, Reproducibility, and the Credibility
890 of Economics Research. *Journal of Economic Literature* **56**, 920–80 (2018).
- 891 21. Dafoe, A. Science Deserves Better: the Imperative to Share Complete Replication
892 Files. *PS: Political Science & Politics* **47**, 60–66 (2014).
- 893 22. McCullough, B., McGeary, K. A. & Harrison, T. D. Do Economics Journal
894 Archives Promote Replicable Research? *Canadian Journal of Economics* **41**,
895 1406–1420 (Nov. 2008).
- 896 23. Pérignon, C. *et al.* *Computational Reproducibility in Finance: Evidence from*
897 *1,000 Tests* HEC Paris Paper. 2023.
- 898 24. Camerer, C. F. *et al.* Evaluating Replicability of Laboratory Experiments in
899 Economics. *Science* **351**, 1433–1436 (2016).
- 900 25. Camerer, C. F. *et al.* Evaluating the Replicability of Social Science Experiments
901 in Nature and Science Between 2010 and 2015. *Nature Human Behaviour* **2**, 637–
902 644 (2018).
- 903 26. Open Science Collaboration. Estimating the Reproducibility of Psychological
904 Science. *Science* **349**, aac4716 (2015).
- 905 27. Dreber, A. & Johannesson, M. A Framework for Evaluating Reproducibility and
906 Replicability in Economics. *Economic Inquiry* (2023).
- 907 28. Brodeur, A., Dreber, A., Hoces de la Guardia, F. & Miguel, E. Replication
908 Games: How to Make Reproducibility Research More Systematic. *Nature* **621**,
909 684–686 (2023).
- 910 29. Simonsohn, U., Simmons, J. P. & Nelson, L. D. Specification Curve Analysis.
911 *Nature Human Behaviour* **4**, 1208–1214 (2020).
- 912 30. Brodeur, A., Lé, M., Sangnier, M. & Zylberberg, Y. Star Wars: The Empir-
913 ics Strike Back. *American Economic Journal: Applied Economics* **8**, 1–32 (Jan.
914 2016).
- 915 31. Brodeur, A., Cook, N. & Heyes, A. Methods Matter: P-Hacking and Publication
916 Bias in Causal Analysis in Economics. *American Economic Review* **110**, 3634–
917 3660 (2020).
- 918 32. Botvinik-Nezer, R. *et al.* Variability in the Analysis of a Single Neuroimaging
919 Dataset by Many Teams. *Nature* **582**, 84–88 (2020).
- 920 33. Breznau, N. *et al.* Observing Many Researchers Using the Same Data and
921 Hypothesis Reveals a Hidden Universe of Uncertainty. *Proceedings of the National*
922 *Academy of Sciences* **119**, e2203150119 (2022).

- 923 34. Huntington-Klein, N. *et al.* The Influence of Hidden Researcher Decisions in
 924 Applied Microeconomics. *Economic Inquiry* **59**, 944–960 (2021).
- 925 35. Menkveld, A. J. *et al.* Non-Standard Errors. *Journal of Finance* (Forthcoming).
- 926 36. Silberzahn, R. *et al.* Many Analysts, One Data Set: Making Transparent How
 927 Variations in Analytic Choices Affect Results. *Advances in Methods and Practices*
 928 *in Psychological Science* **1**, 337–356 (2018).
- 929 37. Fišar, M. *et al.* Reproducibility in Management Science. *Management Science*
 930 (2023).
- 931 38. Ankel-Peters, J., Fiala, N. & Neubauer, F. Do Economists Replicate? *Journal of*
 932 *Economic Behavior & Organization* **212**, 219–232 (2023).
- 933 39. Vilhuber, L., Turrilo, J. & Welch, K. Report by the AEA Data Editor. *AEA*
 934 *Papers and Proceedings* **110**, 764–75. [https://www.aeaweb.org/articles?id=10.](https://www.aeaweb.org/articles?id=10.1257/pandp.110.764)
 935 [1257/pandp.110.764](https://www.aeaweb.org/articles?id=10.1257/pandp.110.764) (May 2020).

936 11 Methods

937 Our focus is on 12 journals. The journals are the following for economics: *Amer-*
 938 *ican Economic Review*, *American Economic Review: Insights*, *American Economic*
 939 *Journal: Applied Economics*, *American Economic Journal: Economic Policy*, *Ameri-*
 940 *can Economic Journal: Macroeconomics*, *The Economic Journal*, *Journal of Political*
 941 *Economy*, *Quarterly Journal of Economics*, *Review of Economic Studies*. For political
 942 science, the journals are: *American Journal of Political Science*, *American Political*
 943 *Science Review*, and *Journal of Politics*.

944 We have two streams to generate reproductions.

945 ***I4R’s Board.***

946 First, I4R has a board of editors who recommend potential reproducers. All board
 947 members are nominated by the lead author, A.B. He then reaches out to the board
 948 for suggestions of reproducers who could be a good fit for the studies in the targeted
 949 journals.

950 ***Replication Games.***

951 Our second stream to generate reproductions and replications is the replication games
 952 (Games). Games are one-day meet-ups open to faculty, post-docs, graduate students
 953 and other researchers. Participants join a small team of about 3–5 researchers all
 954 working in the same subfield (*e.g.*, development economics).

955 11.1 Types of Re-Analyses

956 We group re-analyses into eight groups: (i) alternative control variables, (ii) change
 957 the sample, (iii) change (coding of) the dependent variable, (iv) change (coding of)
 958 the main independent variable, (v) change estimation method, (vi) change inference
 959 method, (vii) change weighting scheme and (viii) replication using new data.

960 11.2 Robustness for Figures

961 While the bulk of our analysis compares coefficients and statistical significance from
962 the original study and the work of reproducers, many results in papers are also dis-
963 played in figures. For those which are plots of coefficients (i.e., event studies) we
964 encouraged reproducers to give the underlying statistics used to create the graph. This
965 was often at the discretion of the reproducers: it could be taxing to write new code
966 to compare and extract those values. In one example, the underlying programs which
967 were written by the original authors were too complicated to modify with robustness
968 checks. Excepting anecdotal examples, many teams found it feasible to reproduce a
969 figure as part of a robustness check or direct replication. In those circumstances, we
970 (A.B. and D.M.) tried to subjectively describe if we believed the results were the same.
971 This was usually taken with the discussion of the reproducers and reading the original
972 paper. We find that 189 out of 263 figures—71.9 percent—we believe to have display
973 the same result as the original paper and can be reasonably compared.

974 11.3 Non Comparable Re-Analyses

975 As mentioned earlier, a direct comparison is not possible between the original analysis
976 and the reproducers' analysis for about 15% of re-analyses. In applied microeconomics
977 and politics papers, this may be due to a change in the estimator or a change in
978 the scale of the dependent or main independent variable. There are also scenarios
979 where the original paper uses methods where coefficient estimates and p-values are
980 not the objective of the analysis. This is apparent in a few empirical macroeconomics
981 papers teams looked at. A common "robustness check" would be to adjust parameters
982 which enter a model, possibly using accepted values in the field or estimated from an
983 alternative dataset.

984 82 articles have at least one non-comparable estimate. Only a small proportion
985 (10 re-analyses) were not directly comparable for all reported re-analysis estimates.
986 For not directly comparable re-analyses, we report the proportion that reproducers
987 indicated were of the same statistical significance as the original and same sign. For our
988 four definitions of reproducibility and replication rates these are: When the original
989 estimate is statistically significant at the 5% level, 85% of those we considered not
990 directly comparable indicated their re-analysis was of the same significance (93%
991 for the 10% level). When the original estimate was not statistically significant at the 5%
992 level, 88% of those we considered not directly comparable indicated their re-analysis
993 was of the same (non)significance (92% for the 10% level).

994 11.4 Study Selection

995 Not all studies from our targeted journals have been reproduced or replicated. Our
996 approach leads to an over-representation of studies using publicly available data.
997 Another feature of our sample is that the targeted journals have a data availability
998 policy *and* enforce it. This is in contrast to many top field journals in both economics
999 and political science. Our sample should thus be viewed as very selected both in terms
1000 of impact and high data and code availability rates. In fact, approximately 45% of

1001 replication packages in our sample included raw data and complete cleaning code. An
1002 additional 13.5% provided partial cleaning code.

1003 11.5 Journal Policy

1004 The *American Journal of Political Science* does not have a data editor. Instead, the
1005 computational reproducibility is carried out by the staff at the Odum Institute for
1006 Research in Social Science, at the University of North Carolina, Chapel Hill. The jour-
1007 nals which do not conduct reproducibility checks are the *American Political Science*
1008 *Review*, the *Journal of Political Economy* and the *Quarterly Journal of Economics*.
1009 The other journals conduct computational reproducibility internally.

1010 Data editors make sure that the replication packages include the data and codes,
1011 and that the documentation (e.g., Readme) is complete. In the event that the authors
1012 cannot share some or all the data, they request that information is provided on how
1013 other researchers could obtain the data set(s). Their teams also run the codes and
1014 make sure that the output is similar to what is reported in the article. They do not
1015 look for coding errors nor run robustness checks.

1016 11.6 Many-Analysts Approach

1017 Our approach and research questions, which we detail below, were pre-registered. Our
1018 pre-analysis plan was pre-registered here: <https://osf.io/8wsqx/>. The pre-analysis plan
1019 was pre-registered prior to sharing the Meta Database with analysts. See the SI for
1020 more information on the Meta Database.

1021 The six analyst teams tackled the following eight questions:

- 1022 1. “Does reproducibility/replicability rate depend on replicators’ experience coding?”
- 1023 2. “Does reproducibility/replicability rate depend on replicators’ academic experi-
1024 ence?”
- 1025 3. “Does reproducibility/replicability rate depend on the authors’ experience?”
- 1026 4. “Does reproducibility/replicability rate depend on the interaction of the authors’
1027 experience and replicators’ experience?” In particular:
 - 1028 (a) Are reproducibility/replicability rate higher when authors’ experience is
1029 high, and replicators’ experience is low (in comparison to similar levels)?
 - 1030 (b) Are reproducibility/replicability rate higher when authors’ experience
1031 and replicators’ experience is similar (in comparison to dissimilar
1032 levels)?
 - 1033 (c) Are reproducibility/replicability rate higher when authors’ experience is
1034 low, and replicators’ experience is high (in comparison to similar levels)?
- 1035 5. “Does reproducibility/replicability rate depend on the interaction of the authors’
1036 prestige and replicators’ prestige?” In particular:
 - 1037 (a) Are reproducibility/replicability rate higher when authors’ have high
1038 prestige, and replicators’ experience have low prestige (in comparison
1039 to similar levels)?
 - 1040 (b) Are reproducibility/replicability rate higher when authors’ and replica-
1041 tors’ prestige is similar (in comparison to dissimilar levels)?

- 1042 (c) Are reproducibility/replicability rate higher when authors' have low
1043 prestige, and replicators' experience have high prestige (in comparison
1044 to similar levels)?
- 1045 6. "Does reproducibility/replicability rate depend on the original authors providing
1046 raw data?"
- 1047 7. "Does reproducibility/replicability rate depend on the original authors providing
1048 raw or intermediate data?"
- 1049 8. "Does reproducibility/replicability rate depend on the original authors providing
1050 cleaning code?"

1051 **11.6.1 Data for Analysts**

1052 Analysts were not given access to raw data (database, team leader surveys, individual
1053 surveys). Rather, they were given access to intermediate/analytical data which was
1054 cleaned and merged in a manner which would be consistent for their analysis. Giving
1055 researchers a downstream dataset allowed A.B. and D.M. to make restrictions on
1056 what the analysts could do. The clearest example of this would be defining dependent
1057 variables which were not allowed to be changed - providing a consistent definition
1058 between analysts. Asking certain research questions also restricted the data given to
1059 the analysts. These restrictions were done in ways so that any analysis done would be
1060 more comparable.

1061 The backbone of the data provided to analysts was the Meta Database, of which
1062 questions from the team leader surveys and individual surveys were added. Much of
1063 the information from the individual surveys were aggregated to the report level.

1064 The data given to the analysts changed as reproduction reports, team leader and
1065 individual surveys were completed. In total, we provided 13 updated databases for
1066 analysts between November 6th, 2023 and February 12th, 2024. We did this to give
1067 analysts time to create scripts which would work with partial datasets as we worked
1068 to gather reports and surveys. This allowed analysts to expedite their analysis once
1069 the full dataset was constructed.

1070 The goal was to have each team answer each research question independently.
1071 Each team received the same instructions and data. We allowed full flexibility to all
1072 teams. Teams were allowed to use any statistics package, statistical model, inference,
1073 weighting scheme, *etc.* Teams were free to choose the independent variables and how
1074 to code them. Teams were also free to construct their own derived variables from the
1075 dataset given to them.

1076 We provided the four dependent variables and the database to all teams. They
1077 were allowed to use any of the provided variables and new data. The only restriction
1078 imposed on teams is that they needed to use our four main dependent variables.

1079 **11.6.2 Team Construction**

1080 We asked a subset of coauthors on this paper (reproducers) if they would like to help
1081 analyze our database. We informed them that we would "have different teams inde-
1082 pendently working together at answering the same research questions (e.g., what is
1083 the reproducibility/replicability rate for each specific type of robustness checks/re-
1084 coding)." The subset of coauthors who received an invitation to volunteer were: (1)

1085 contacted between September 21st and October 8th 2023 *and* (2) had completed,
1086 or were near completion of, their reproduction report. We sent invitations (a simple
1087 sign-up form) in an email which also asked the reproducers to respond to individ-
1088 ual and team leader surveys which formed parts of our previous analysis. About
1089 110 co-authors were invited between September 21st and October 8th. 10 individuals
1090 ultimately signed-up as “many-analysts.”

1091 In our request for volunteers, we asked volunteers if they: (1) had a team who
1092 wanted to do research on the project; (2) wanted to be added to a team; (3) wanted to
1093 work on the analysis alone. No one joined as teams, most people wanted to be added
1094 to a team, and the remainder wanted to work alone. For those that wanted to work
1095 together, we assembled teams as best we could so they were close enough in timezones.
1096 We had two teams of three, one team of two, and two individuals. A.B. and D.M.
1097 also acted as a team of two, yielding six teams in total. No members of any teams left
1098 during the analysts phase.

1099 Although the PI ultimately provided each volunteer with a payment of \$3,000
1100 CAD, this compensation was not disclosed or anticipated at the time they agreed to
1101 participate.

1102 11.7 Database: Sample Composition

1103 The database described above provides 6,583 re-analyzed test statistics from 103 repro-
1104 duction reports. (Seven reports did not include robustness checks.) The other test
1105 statistics are estimates obtained by re-coding the analysis.

1106 Supplementary Materials Appendix Table 11 provides summary statistics for the
1107 full sample and by journal. In total, 83 reproduction reports were completed through
1108 Games in comparison to 27 through the editorial board stream. 79 reproduction reports
1109 are for the field of economics against 31 for political science.

1110 There is no universally agreed upon criterion for reproduction. As a first criterion,
1111 we follow much of the literature and define reproducibility as obtaining a statistically
1112 significant effect in the same direction (positive or negative) as the original study.
1113 Throughout, we rely on four main dependent variables:

1114 **First Dependent Variable:** dummy variable indicating whether the re-analysis
1115 is statistically significant at 5% level and same sign. For this dependent variable,
1116 we only keep original estimates statistically significant at the 5% level.

1117 **Second Dependent Variable:** dummy variable indicating whether the re-analysis
1118 is statistically significant at 10% level and same sign. For this dependent variable,
1119 we only keep original estimates statistically significant at the 10% level.

1120 **Third Dependent Variable:** dummy variable indicating whether the re-analysis
1121 remains not statistically significant at 5% level. For this dependent variable, we
1122 only keep original estimates statistically insignificant at the 5% level.

1123 **Fourth Dependent Variable:** dummy variable indicating whether the re-analysis
1124 remains not statistically significant at 10% level. For this dependent variable, we
1125 only keep original estimates statistically insignificant at the 10% level.

1126 The average number of re-analyzed test statistics per article is about 60. The stan-
1127 dard deviation is very high (73), with a maximum of 421. This is unsurprising given

1128 that some teams, for instance, focused most of their attention to (blindly) recoding
1129 using the raw data (either provided by the authors or re-downloaded by the repro-
1130 ducers), while other teams have focused solely on conducting robustness checks for
1131 multiple central hypotheses. As an illustrative example, imagine that an original arti-
1132 cle has three main outcome variables and relies on two main specifications. If the
1133 reproducers conduct five different robustness checks for each outcome variable and
1134 specification, then this would lead to 30 re-analyzed test statistics.

1135 As a robustness check, we deal with this issue by adjusting the weight of each test
1136 statistics by the inverse number of such statistics in the reproduction report such that
1137 each reproduction report has the same weight.

1138 Supplementary Materials Appendix Table 2 provides descriptive statistics. The
1139 articles in our sample are all recently published with a relatively small number of
1140 Google Scholar citations (44 on average) as of the completion of a reproduction report.
1141 The original authors are more experienced than reproducers with 11 years of experi-
1142 ence (*i.e.*, years since completing their Ph.D.) against 3. Original authors have on
1143 average 4,269 Google Scholar citations in comparison to 478 for reproducers. Those dif-
1144 ferences are mostly driven by the larger share of graduate students among reproducers
1145 than for original authors (49% against 6%). There are about 2.6 original authors per
1146 article in comparison to 3.2 for reproducers. About 15% of reproducers have recently
1147 published in a Top 5 or one of the three leading political science journals in our sample.
1148 Approximately 30% have published in those journals or in one of the other journals
1149 in our sample.

1150 While reproducers have less academic experience than original authors on average,
1151 their level of expertise as programmers is quite advanced. About 10%, 48% and 33%
1152 of reproducers report that their level of expertise is “Expert”, “Proficient” and “Com-
1153 petent,” respectively. Moreover, about 55% of reproducers had already produced a
1154 replication package for their own work or journal publication.

1155 11.8 Computational Reproducibility

1156 We rely on the Social Science Reproduction Platform (SSRP)’s 10-point scale to docu-
1157 ment computational reproducibility. This scale is useful as it is standardized and offers
1158 more details than a simple indicator for whether the results are computationally repro-
1159 ducible (Visit <https://bitss.github.io/ACRE/assessment.html#score> for more details
1160 on SSRP and this scale). On this scale, a rating of 1 signifies the incapacity to repro-
1161 duce results due to the absence of data or code, while a rating of 10 indicates the
1162 capability to faithfully reproduce results from the raw data (unaltered files obtained
1163 by the authors from the sources cited in the paper) to the final numerical results as
1164 published in the paper.

1165 The following is a direct reproduction from the Guide for Accelerating Computa-
1166 tional Reproducibility in the Social Sciences.

1167 **Level 1 (L1):** No data or code are available. Possible improvements include adding:
1168 raw data, analysis data, cleaning code, and analysis code.

1169 **Level 2 (L2):** Code scripts are available (partial or complete), but no data are
1170 available. Possible improvements include adding: raw data and analysis data.

1171 **Level 3 (L3):** Analytic data and code are partially available, but raw data and
1172 cleaning code are missing. Possible improvements include: completing analysis data
1173 and/or code, adding raw data, and adding analysis code.

1174 **Level 4 (L4):** All analytic data sets and analysis code are available, but the code
1175 fails to run or produces results inconsistent with the paper (not CRA). Possible
1176 improvements include: debugging the analysis code or obtaining raw data.

1177 **Level 5 (L5):** Analytic data sets and analysis code are available and they produce
1178 the same results as presented in the paper (CRA). The reproducibility package may
1179 be improved by obtaining the original raw data.

1180 Note: This is the highest level that most published research papers can attain
1181 currently. Computational reproducibility from raw data is required for papers that
1182 are reproducible at Level 6 and above.

1183 **Level 6 (L6):** Cleaning code scripts are available (partial or complete), but raw
1184 data is missing. Possible improvements include: adding raw data.

1185 **Level 7 (L7):** Cleaning code is available and complete, and raw data is partially
1186 available. Possible improvements: adding raw data.

1187 **Level 8 (L8):** All the materials (raw data, analytic data, cleaning code, and analy-
1188 sis code) are available. However, the cleaning code fails to run or produces different
1189 results from those presented in the paper (not CRR) or the analysis code fails to run
1190 or produces results inconsistent with the paper (not CRA). Possible improvements:
1191 debugging the cleaning or analysis code.

1192 **Level 9 (L9):** All the materials (raw data, analytic data, cleaning code, and anal-
1193 ysis code) are available. The analysis code produces the same output as presented
1194 in the paper (CRA). However, the cleaning code fails to run or produces differ-
1195 ent results from those presented in the paper (not CRR). Possible improvements:
1196 debugging the cleaning code.

1197 **Level 10 (L10):** All necessary materials are available and produce consistent
1198 results with those presented in the paper. The reproduction involves minimal effort
1199 and can be conducted starting from the analytic data (CRA) and the raw data
1200 (CRR). Note that Level 10 is aspirational and may be unattainable for most
1201 research published today.

1202 Each team was asked to assign a reproducibility score on a scale of one to ten to
1203 the paper reproduced. This involved documenting the completeness of the data and
1204 code, and whether the materials produce results consistent with those in the article.
1205 Their focus for computational reproducibility is only for the claims that they have
1206 investigated rather than all exhibits in the article.

1207 The results are presented in Extended Data Figure 1. This figure shows the varia-
1208 tion across papers, with the highest concentration of scores concentrated at levels 10
1209 and 5. Indeed, over 85% (Levels 5 and 10) of results examined in our sample were fully
1210 reproducible using either: (1) the raw and analytical data, or; (2) the analytical data
1211 when the raw data were not provided. Level 10 (L10) means that all necessary mate-
1212 rials are available and produce consistent results with those presented in the paper.
1213 Level 5 (L5) means that analytic data sets and analysis code are available, and they
1214 produce the same results as presented in the paper. In other words, L5 indicates that
1215 the reproducers successfully (computationally) reproduced the numerical results using

1216 the analytical data, but the raw data were not provided, while L10 indicates that the
1217 reproducers successfully (computationally) reproduced the numerical results using the
1218 raw data and cleaning and analytical codes.

1219 The remaining 15% includes studies for which analytic code and data are partially
1220 available and studies for which some of the codes (cleaning or analytic) fail to run or
1221 produce results inconsistent with the paper. These findings suggest very high rates of
1222 computationally reproducible results.

1223 Our results are in stark contrast with several studies documenting low compu-
1224 tational reproducibility rates ([13, 19, 40]). This is perhaps unsurprising given that
1225 most of the articles in our sample were already computationally reproduced by data
1226 editors. This highlights the open science movement has improved computational repro-
1227 ducibility of research findings in leading economics and political science journals. Our
1228 approach is also different as we are targeting newer studies and only articles for which
1229 (at least) analytical data were available to the teams of reproducers. A more compa-
1230 rable (and recent) study is [37] which assess the reproducibility of nearly 500 articles
1231 published in the journal *Management Science*. They find that more than 95% of arti-
1232 cles could be reproduced if data accessibility and software requirements were not an
1233 obstacle for reproducers.

1234 11.9 Recoding

1235 We now turn to recoding exercises conducted by a subset of teams. Those teams either
1236 recoded using a different software language or used the same software without looking
1237 at the original authors' code. In total, 19 teams of reproducers engaged in computa-
1238 tionally reproducing and checking for coding errors using a different statistical software
1239 than the original authors. This may be due to reproducers being more comfortable in
1240 another software language or the availability of specific commands (to run a robust-
1241 ness check). Five teams also recoded the empirical analysis without looking at the
1242 authors' code/programs.

1243 Recoding in a different software opens up the ability for others to benefit and
1244 understand the empirical foundations of published articles in ways that the original
1245 authors may not have been able to convey. For instance, verifying reproducibility
1246 by translating it into R or Python makes the study itself accessible to many more
1247 researchers.

1248 Recoding also helps to assess the importance of differing assumptions embedded
1249 within programming languages (e.g., different types of Random Number Generations,
1250 rounding rules and numerical precision). We categorized recoding exercises done by
1251 reproducers into three categories: (i) identical numerical results, (ii) minor differences,
1252 and (iii) major differences. Minor differences involve small numerical discrepancies
1253 between the authors' estimates and those obtained by the reproducers. Those dif-
1254 ferences do not lead to important changes in significance or magnitude. In contrast,
1255 major differences lead to major differences in one or multiple claims.

1256 11.10 Coding Errors and Discrepancies

1257 We now turn to documenting the prevalence of coding errors and discrepancies between
1258 the code and the published article. Of note, a paper might be fully reproducible,
1259 but the programs may contain coding errors. Similarly, there might be important
1260 discrepancies between what the article states and what the programs compute, while
1261 remaining computationally reproducible.

1262 We do not document trivial coding errors such as versioning issues and missing
1263 packages/paths. Those coding errors are typically easily fixed by the reproducers.
1264 We instead focus on coding errors which could have had an impact on claims and
1265 conclusions of articles.

1266 We uncover minor or major coding errors in 26 of the 110 studies in our sample,
1267 with some studies containing multiple errors. The errors can be broadly categorized
1268 into errors of the dependent variable (4 articles), main independent variable (5), control
1269 variables (10), estimation (2), inference (2), sample/observations (8) and other (5).
1270 While not all coding errors lead to changes in the conclusions of the original study,
1271 we uncovered several major coding errors worth discussing. Some examples of major
1272 errors include: a very large number of duplicated observations, failing to fully interact
1273 a difference-in-differences regression specification, miscoding the treatment variable
1274 for a large number of (or all) observations, and clear model misspecification.

1275 The prevalence of coding errors is larger for economics (26%) than political science
1276 (16%). A plausible explanation is that replication packages from economic articles
1277 have more lines of code than those in political science, mechanically increasing the
1278 likelihood of at least one coding error.

1279 We also uncovered transcription issues for 13 studies, typically involving small
1280 numerical differences or rounding errors not impacting the claims or conclusions of
1281 the article.

1282 11.11 Time Trends in Data and Code Availability

1283 To document time trends in data and code availability in economics and political
1284 science between 2014 and 2023, we randomly sampled 10 empirical articles per year
1285 for each of our 12 target journals. We define an article as empirical if it relies on real
1286 or simulated data at any point in the text. Thus, a theoretical article that is motivated
1287 with a descriptive analysis of labor market trends, or an econometric paper showing
1288 properties of an estimator on synthetic data would both be classified as empirical for
1289 the purposes of our study.

1290 To randomly select papers, we proceeded as follows: First, we noted the number
1291 of issues per journal per year. Second, we drew ten issues (with replacement) for each
1292 year. Third, for each issue, we generated a random permutation of numbers between
1293 1 and 35, giving us the order in which papers from a given issue should be considered.
1294 So, for example, if the first issue drawn was 4, and the first number in our permutation
1295 sequence was 10, we would consider the tenth article in the fourth issue for coding.
1296 We skipped an article and proceeded with the next number in the permutation if the
1297 article in question a) was not empirical, b) was not a standard article (we excluded
1298 comments, replies and corrections, retraction notices, and editor notes, even if they

1299 were empirical in nature), c) was a duplicate that had already been considered (e.g.,
1300 issue number one, article number five, was drawn twice in a row), or d) did not exist
1301 (our chosen journals typically publish around ten articles per issue, so higher numbers
1302 in the permutation often went unused).

1303 In our coding we considered whether the journal website or the article pdf contain a
1304 link to a replication package, whether this package is accessible, and what the contents
1305 of the package are. We tracked the availability of a Readme file, cleaning and analytical
1306 code, and raw, intermediate, and final data. Note that our coding of code availability
1307 is optimistic in the sense that we only note whether a particular type of code exists;
1308 we did not verify its completeness or correctness. However, when authors explicitly
1309 indicated that a code was incomplete, we noted this information.

1310 Of note, the *American Economic Review: Insights* only formally became a journal
1311 in 2019. For the five years earlier, we did not collect for this journal, leading to 10
1312 fewer papers per year.

1313 **Methods references**

- 1314 13. Gertler, P., Galiani, S. & Romero, M. How to Make Replication the Norm. *Nature*
1315 **554**, 417–9 (2018).
- 1316 19. Chang, A. C. & Li, P. Is Economics Research Replicable? Sixty Published Papers
1317 From Thirteen Journals Say "Often Not". *Critical Finance Review* **11**, 185–206
1318 (2022).
- 1319 37. Fišar, M. *et al.* Reproducibility in Management Science. *Management Science*
1320 (2023).
- 1321 40. Wood, B. D., Müller, R. & Brown, A. N. Push Button Replication: Is Impact
1322 Evaluation Evidence for International Development Verifiable? *PloS one* **13**,
1323 e0209416 (2018).

SI: Reproducibility and Robustness of Economics and Political Science Research

Abel Brodeur et al. (Author list and contributions are provided in the SI)^{1*}

^{1*}Department of Economics and Institute for Replication, University of
Ottawa, 75 Laurier Avenue East, Ottawa, K1N 6N5, Ontario, Canada.

Corresponding author(s). E-mail(s): abrodeur@uottawa.ca;

Methods and Additional Results

Replication Games

So far, teams have been as small as one individual or as large as seven. The locations of Games are chosen based on (i) local interests, (ii) geography, (iii) possibility to have the Games as part of a major conference, and (iv) EDI considerations.

I4R groups graduate students with faculty members and senior researchers, ensuring a mix of junior and more senior economists in each team. A virtual meeting with the organizers before the Games allows each team to ask questions and discuss a game plan. During the Games, A.B., D.M. or one of I4R's co-directors, provide live assistance to the teams.

Participants are offered a short list of (about 5) studies in their field of interest about three weeks before the Games. They are asked to choose a paper as a team, read it and familiarize themselves with the replication package prior to the Games.

Teams are asked to develop a game plan for the Games; each team member should know what they are supposed to do during the Games. A virtual meeting with the organizers before the Games allows each team to ask questions and discuss a game plan. During the games, A.B., D.M. or one of I4R's co-directors, provide live assistance to each team. Teams then have to write a (templated - <https://osf.io/8dkxc/>) report summarizing their work and results in the following months. Of note, virtually all teams kept working on their reproduction after the Games and some even started the re-analysis prior to the Games.

Participants are offered the possibility to virtually attend Games. In our sample of completed reports, about 68% of participants attended the games in-person, while 32% virtually attended the events. Most teams are fully virtual or in-person, with only

31 a small share of teams having a mix of virtual and in-person participants. Mixed teams
32 are typically due to a variety of reasons (*e.g.*, canceled flight for one participant), or
33 late registrations.

34 We asked a subset of games participants the following question: “Why did you
35 choose to participate in the Replication Games?” We offered seven potential options,
36 with an empty box to provide additional reasons. We find that a majority of respon-
37 dents chose the responses “Learn about academic replications and reproductions”,
38 “Expand your network”, and “Contribute to Open Science”. Other popular responses
39 include “Improve your ability to program and code” and “Improve your ability to
40 conduct research”.

41 Teams have on average worked 13 active days on their reproduction (std. dev. of
42 24). Supplementary Materials Appendix Figure 1 shows the distribution of days across
43 reports, trimmed at over 100 days.

44 About half the teams worked from 5 to 20 days on their reproduction report. Most
45 of the remaining teams worked between 25 to 85 active days. A very small fraction
46 worked less than 5 days. This is due to the reproducers not being able to conduct
47 robustness checks. In contrast, about 8% of teams worked more than 100 days. This
48 is typically due to uncovering major coding errors or issues with the original study
49 and having to engage in multiple rounds of back and forth with the original authors.
50 There is also the potential for people to have spent many days on their paper even if
51 the number of hours were low. Reports are on average 19 pages long, with a standard
52 deviation of 14.

53 In terms of retention for the Replication Games, over 90% of registered partici-
54 pants ended up participating in the event. Furthermore, within one year of completing
55 the first two replication Games (October and November 2022), 85% of teams had
56 completed a report.

57 The goal for all reproducers is clearly stated; testing whether the main claims are
58 reproducible and robust. I4R emphasizes to reproducers that the goal is NOT to show
59 that the results are not reproducible. The goal is instead to test if the results are
60 reproducible to recoding and/or robustness checks. This is key as some reproducers
61 might engage in reverse specification searching (*i.e.*, selective reporting of insignificant
62 results). Moreover, we ask reproducers from I4R’s Board stream to provide a pre-re-
63 analysis plan. The game plan acts as a pre-re-analysis plan for the second stream.

64 In practice, some teams in both streams did not write a pre-re-analysis plan and
65 virtually all teams that did write one ended up deviating from it. The latter is because
66 it is very unclear from only reading the original paper what is the range of re-analyses
67 that is feasible. Reproducers had to carefully look at the replication package provided
68 by the authors to gauge whether specific robustness checks were implementable given
69 data availability. Our re-analyses should thus all be considered as not pre-registered.

70 Supplementary Materials Appendix Table 2 provides summary statistics.

71 **Types of Re-Analyses**

72 One of our main objectives is to document the relative importance of several robustness
73 checks and re-analyses in impacting the magnitude and significance of the original

74 point estimates. We group the robustness checks and coding exercises conducted by
75 the reproducers into eight groups.

76 **Alternative control variables:** Removing, adding or changing control variables.
77 In our sample, there are 1,939 new re-analyses involving alternative controls.

78 **Change the sample:** Decreasing or increasing the sample size. In our sample,
79 there are 1,774 new re-analyses involving changing the sample size. Reproducers
80 may change the sample by adding/removing years, geographical units or individu-
81 als. For instance, a team could check if the results are robust to adding/removing
82 a state to/from the analytical sample.

83 **Change (coding of) the dependent variable:** The reproducers may change
84 the coding of the dependent variable. In our sample, there are 285 new re-analyses
85 involving changing the dependent variable. Examples include using an alternative
86 standardization of the outcome variable, alternative calculation of factor loadings
87 for an index in a different software, using an indicator variable for an outcome
88 above a certain threshold, and using a composite index of several indicators as the
89 dependent variable.

90 **Change (coding of) the main independent variable:** The reproducers may
91 change the coding of the main independent variable. In our sample, there are
92 264 new re-analyses involving changing the main independent variable. Examples
93 include using a continuous variable instead of a dummy or a factor variable for
94 treatment, broadening the definition of treated based on physical proximity, using
95 a different type of TV program in television exposure, and allowing for non-linear
96 effects in institutional exposure.

97 **Change estimation method:** This category involves any changes to the estima-
98 tion method. In our sample, there are 605 new re-analyses involving changing the
99 estimation method. Examples include using non-linear models and changing the
100 variables used for matching.

101 **Change inference method:** This category involves changing the inference
102 method. In our sample, there are 542 new re-analyses involving changing the infer-
103 ence method. Examples include bootstrapping the standard errors and clustering
104 at a different level.

105 **Change weighting scheme:** This category involves changing the weighting
106 scheme. In our sample, there are 126 new re-analyses involving changing the
107 weighting scheme. Examples include removing a weighting scheme used by the
108 authors.

109 **Replication using new data:** Replication using new data involve both collecting
110 new data or using data from another data source. In our sample, there are 469
111 new re-analyses involving using new data. Replicators have used new data for the
112 dependent, independent or control variables.

113 In practice, many reproducer teams performed multiple robustness checks *simul-*
114 *taneously* in a single robustness exercise, or, combined two independent robustness
115 checks into a new, third robustness check. We tracked all the changes reproducers made
116 when comparing to an original estimate and coded accordingly. In our sample, about
117 809 re-analyses fall into at least two categories of simultaneous robustness checks.

118 Supplementary Materials Appendix Table 3 provides a decomposition of reports
119 and test statistics by type of re-analyses. The most popular re-analyses involve using
120 alternative control variables and changing the sample. In contrast, only 14 reports had
121 any robustness check which changed the weighting scheme and only 15 reproduction
122 reports had any robustness checks which used new data.

123 The types of re-analyses are quite similar for economics and political science.
124 Using alternative control variables, changing the sample and changing the estimation
125 method/model are among the most popular re-analyses for both fields. One notice-
126 able difference is that reproducers are more likely to change the method of inference
127 for economic articles than in political science.

128 Database

129 In what follows, we describe our database. The database is mainly built from three
130 sources of raw data: (1) reports; (2) surveys of individual reproducers; and (3) sur-
131 veys of teams of reproducers. We also collected information from publicly available
132 *curricula vitae* of all original authors and reproducers.

133
134 **Reports:** Two of the lead authors (A.B. and D.M.) and research assistants read
135 reports and copied test statistics into an Excel file. We also coded and grouped robust-
136 ness and replicability exercises, and information on computational reproducibility and
137 coding errors. The work being entered by RAs was checked by A.B. or D.M. for com-
138 pleteness and accuracy. If any part of any entry was unclear, they were checked again
139 and discussed.

140 Only a subset of results was considered suitable for our research. We follow the
141 following criteria. We exclude extensions of the original authors' research, effects by
142 heterogeneity, or mediation analyses. These analyses correspond to situations where
143 there are no "original" estimates to which we can reasonably compare the reproducers'
144 estimates. Most often, reproducers included tables and figures which were the output
145 of a computational reproduction using the original authors' replication package. These
146 are always left out for the re-analyses. After being checked, reproducers would then
147 be contacted with their subset of the database and asked to confirm our transcribing
148 of their reports into the database.

149 Coding errors and discrepancies are also excluded from the re-analyses. We discuss
150 coding errors and discrepancies between original authors' values in their published
151 paper compared to what their replication package produces in the Methods Section.

152 We report some additional information in our database. We collect information
153 on the journal, year of publication, number of Google Scholar citations at the time
154 of entry into the database, the research field, the position of the test in the original
155 article and the number of original authors and reproducers. We also collect informa-
156 tion from *curricula vitae* of all the original authors and reproducers. We obtained
157 information on their academic affiliation at the time of publication, their position
158 at the main institution and year the PhD was earned. In addition, we gather for
159 each author and reproducer the following information (at the time of completing
160 the replication): the total number of Google Scholar citations and whether they had
161 published in a Top-5 economic journal, a leading political science journal, and/or one

162 of the other economic journals we are reproducing/replicating. The Top-5 economic
163 journals are the *American Economic Review*, *Econometrica*, the *Journal of Political*
164 *Economy*, the *Quarterly Journal of Economics* and the *Review of Economic Studies*.
165 The leading political science journals considered here are the *American Journal of*
166 *Political Science*, *American Political Science Review* and *Journal of Politics*.

167

168 **Surveys:** We asked all reproducers to fill out an individual survey. We also asked
169 one author per reproduction report to fill out a team survey. Both surveys gave addi-
170 tional information on the academic and programming experience of reproducers, how
171 long their report took to create and the completeness of the original authors' replica-
172 tion package, and whether they improved it. Teams were invited to answer the surveys
173 following the completion of transcribing their report.

174 The team survey provides additional information on data availability, computa-
175 tional reproducibility, the reasons the paper to be reproduced/replicated was chosen,
176 how long it took to run the code provided in the replication package, reasons they were
177 unable to conduct specific robustness exercises, *etc.* We also asked whether there was
178 any communication with the original authors for clarifications and how it improved
179 the quality of the report.

180 The individual survey also provides us information about whether the reproducers
181 participated in the Games, whether they virtually attended, why they participated
182 in the Games, and their general experience, and how it improved their networking
183 and coding skills. We conclude the individual survey with subjective questions such
184 as "How does the quality of the replication package affect your view of the discipline
185 as a whole?"

186 **Communication with Original Authors**

187 Once a report is completed, A.B. reviews it if it falls within his expertise. Otherwise,
188 someone else on I4R's board reviews the report. This review involves checking the tone
189 and structure of the report. A.B. then shares the report with the original authors. A.B.
190 emailed all the original authors unless there were more than 5 authors. A reminder was
191 sent a few months later if the original authors did not respond to the initial email. If
192 the authors did not respond to the reminder, the report was released after 6 months.

193 Reproducers may change their report after receiving the original authors' response,
194 allowing them to include their feedback. This is especially important if a re-analysis
195 was judged unreasonable. I4R then allows the original authors to change their response
196 as well. Of note, the reproducers may remain anonymous. In practice, about 11% of
197 reproducers have decided to remain anonymous.

198 Original authors have been incredibly fast at providing a response, perhaps since
199 papers being reproduced have just been published. See [1] (pages 133-244) for a link
200 to authors' responses.

201 In some instances, original authors requested to see the reproducers' replication
202 package, which we provided. See Supplementary Materials Appendix Table 4 for a
203 breakdown by discipline.

204 How often do reproducers and original authors agree? This is a key question as
205 reproducers have freedom to conduct any recoding or sensitivity analysis. This free-
206 dom might lead to disagreement on the validity of some re-analyses. We document
207 (dis)agreements in multiple ways. First, authors’ final responses (i.e., post-mediation)
208 were coded as whether there remained disagreements between authors and reproduc-
209 ers. The coding was done by A.B. and three ambiguous cases were discussed at length
210 with D.M.

211 Overall, we find that there are remaining disagreements for only 23% of articles in
212 our sample. This percentage goes up to over 75% if we restrict the sample to articles for
213 which the original authors wrote a formal response, suggesting that the majority of formal
214 responses we obtained include some sources of disagreements. Disagreements are
215 mostly due to the validity of the re-analyses. There were no remaining disagreements
216 on the presence of coding errors, but authors and reproducers sometimes disagreed
217 on their importance. Disagreements on the scope of the re-analyses and definition of
218 reproducibility were quite rare, and there were also disagreements involving the tone
219 or interpretation of the re-analyses/errors.

220 We observed a general lack of adversariality between original authors and repro-
221 ducers [2]. The broad lack of adversariality is potentially due to the high rates of
222 reproducibility and replicability, but also perhaps on the institutionalization of replica-
223 tions and the fact that discussion between original authors and reproducers is mediated
224 by the I4R. Moreover, original authors may feel less targeted by our reproducers as
225 our aim is to mass-reproduce and replicate studies published in leading economic and
226 political science outlets.

227 We asked reproducers whether their team or I4R contacted, or attempted to con-
228 tact, the original authors for clarifications. About 40% responded “yes”. About 10%
229 reached out because the replication package was unclear, while 17% needed help to
230 computationally reproduce the original authors’ results. Another 17% were unable
231 to access the original authors’ data. Other reasons include verifying coding errors,
232 clarifications about the design model parameters or other coding decisions.

233 Study Selection

234 As a benchmark, A.B. investigated whether studies published at the *Journal of Devel-*
235 *opment Economics* (JDE) using publicly available data complied with the journal’s
236 mandatory data sharing policy. He manually checked the presence of a replication
237 package on JDE’s website for all articles published in four volumes in 2022. Out of 75
238 studies, 47 did not provide a replication package or mentioned that data and codes
239 will be made available upon request. The remaining 28 studies can be categorized as
240 follows: 13 report relying on confidential data; 14 provided a link to a replication pack-
241 age; and one provided only Stata codes and information on how to obtain the data.
242 He then contacted (through I4R’s email) all authors who did not provide a replication
243 package. Seven ended up providing a package. Some authors mentioned that they did
244 not know that the policy existed. A few mentioned that they shared the replication
245 materials with JDE and were surprised that it was not posted.

246 We explore the reasons why teams selected their paper. All teams answered the
247 following question: “For what reasons did you select your specific paper to reproduce

248 and/or replicate from the list of papers provided?” 12 options were offered, including
249 *Other (please specify)*. Options were not mutually exclusive, so any one team could
250 provide multiple reasons for why they selected their paper. Extended Data Figure 2
251 summarizes the percentage of teams who selected each category. Of note 13.6% of
252 teams were assigned a study (*i.e.*, did not choose which study to work on), so they did
253 not answer this question. About 45% of teams report “Methods used”, 36% of teams
254 selected because of the journal of publication, about 25% due to the “Length of time
255 to reproduce results” and about 19% due to the “Size of replication package”. This is
256 in line with our provided guidelines for choosing a study.

257 If a large portion of reproducers select papers based on the assumption that their
258 findings are questionable, it could skew reproducibility rates downward, as there’s a
259 tendency to pick studies more prone to revealing problematic outcomes. However, in
260 this project, only a minimal fraction of teams indicated that they chose their paper
261 because of *ex ante* beliefs that main results are (not) robust/replicable (3.6%). Few
262 teams also selected papers based on statistical power/sample size and trust of original
263 authors.

264 Supplementary Materials Appendix Table 5 explores if our sample is representa-
265 tive of all subfields within economics. We compare JEL Codes of economic papers
266 that we reproduced relative to those of a random sample of representative journal
267 articles published in the top 100 journals in Economics (as ranked by IDEAS/RePec).
268 This comparison benchmark comes from [3]. A comparison of the two samples sug-
269 gest that some subfields are under-represented. Our sample under-represents, among
270 other fields, C-Mathematical and Quantitative Methods, G-Financial Economics and
271 F-International Economics.

272 **Many Analysts: Additional Results**

273 Each row in Supplementary Materials Appendix Table 6 represents one of the
274 eight research questions. The four columns represent four broad categories regarding
275 research teams’ coefficient estimate(s) to the research question: (1) negative and sta-
276 tistically significant, (2) negative and not-statistically significant, (3) positive and not
277 statistically significant and (4) positive and statistically significant. The left-to-right
278 order of the column categories corresponds to where the associated analyst t-statistic
279 would fall on the real number line. While the dependent variable (which does not
280 change in this table) is the same for each team, each team chooses their own primary
281 independent variable. Each cell represents the proportion of analyst-estimated rela-
282 tionships by category. The cells are team-weighted so that if a many-analyst team
283 presents three estimates and another team presents a single estimate, the first team’s
284 estimates enter the proportion as 1/3 each.

285 The cell in the first row and first column tells us that 42.8% of results from the
286 many-analysts find a negative and statistically significant relationship between the
287 coding experience of a reproducer and the reproducibility rate for estimates that were
288 originally statistically significant at the 5% level (*i.e.*, lower reproducibility rate for
289 more experienced reproducers).

290 From the second column, it becomes clear that, if there is a relationship between
291 reproducers experience coding and the reproducibility rate, it seems to be almost

292 definitively negative with a combined proportion of 86% of results returned as negative
293 and statistically significant or negative and not statistically significant at the 5% level.
294 Only 14% of estimates find a positive relationship between the reproducers' experience
295 coding and the reproducibility rate - of which none of the estimated positive relation-
296 ships estimated were statistically significant. (The associated row in Supplementary
297 Materials Appendix Table 7 , which looks at the replication for the 10% threshold finds
298 the same pattern.) This result potentially suggests that reproducers with more experi-
299 ence coding are better suited to detecting and correcting less-than-robust estimations
300 - possibly because of having greater expertise with the methods used.

301 Supplementary Materials Appendix Table 6 presents results where the dependent
302 variable takes a value of one if an originally 5% statistically significant result was
303 reproduced also at the 5% level. Supplementary Materials Appendix Table 7 has the
304 same structure, but uses the 10% threshold. Supplementary Materials Appendix Table
305 8 then examines whether an originally *not* 5% statistically significant result was repro-
306 duced, while Supplementary Materials Appendix Table 9 continues this with the 10%
307 threshold.

308 For the second research question - whether the reproduction robustness rate
309 depends on the reproducers' academic experience, a somewhat similar albeit less
310 starkly negative result is found with some proportion moving into the positive and sta-
311 tistically significant category. That said, the ratio of negative-and-significant results
312 to positive-and-significant results remains above 4 to 1. The associated row in Supple-
313 mentary Materials Appendix Table 7, which looks at robustness for the 10% threshold
314 finds the same pattern, although with 75% of many-analysts results being negatively
315 signed.

316 For the third research question - whether the replication rate depends on the
317 author's experience seems to be centered on the null. Combined, the negative and not
318 statistically significant and the positive and not statistically significant cells contain
319 97.2% of results. The null hypothesis dominates in Supplementary Materials Appendix
320 Tables 7,8, and 9 (which examine reproducibility rates for originally statistically signif-
321 icant at the 10% level, not statistically significant at the 5% level, and not statistically
322 significant at the 10% level, respectively) as well.

323 For the fourth research question, (which has three sub-questions depending on the
324 relative hierarchy of reproducer and original author experience) there seems to be a
325 positive relationship when authors have more or the same level of experience as the
326 reproducer (research question 4a and 4b). This relationship, however, weakens to a
327 likely null when authors have comparatively less experience than their reproducers.
328 Supplementary Materials Appendix Tables 7, 8 and 9, find similar patterns.

329 For the fifth research question, which has the same comparative structure as the
330 fourth while focusing now on the relative prestige of the authors and reproducers, the
331 same (albeit weaker) pattern is found. When authors have more prestige than their
332 reproducers, there is a very positive relationship with replication rate. When origi-
333 nal authors and reproducers have similar prestige levels, this relationship becomes
334 much more likely to be a null (since the middle two columns so outsize the outer
335 two columns). When the authors have less prestige than the reproducers, then the

336 relationship seems to be negative: 22% finding a negative and statistically signifi-
337 cant relationship. In Appendix Table 7, we see the same pattern. When examining
338 replication rate of originally statistically insignificant results, the null hypothesis
339 dominates.

340 The null hypothesis seems to dominate for the final three research questions, with
341 statistical significance not being achieved in either direction for more than one-sixth
342 of the teams’ analyses. This means that the replication rate does not seem to have a
343 relationship with whether the authors provided raw data (research question 6), both
344 raw and intermediate data (research question 7) or cleaning codes (research question
345 8). See Supplementary Materials Appendix Tables 7, 8, and 9 as well.

346 In Supplementary Materials Appendix Table 10, we reproduce the analyses in
347 Supplementary Materials Appendix Tables 6, 7, 8 and 9 while only including estimates
348 if the analyst team indicated that, in their opinion, the estimated effect size was
349 economically meaningful. Results are broadly consistent with those described above
350 without the restriction.

351 The results may reflect that our focus is on journals with a data and code availabil-
352 ity policy. The provision (or not) of raw data, intermediate data, or cleaning codes,
353 may thus be due to data type rather than selective data/code provision by original
354 authors. Our results are consistent with [4] who document no relationship between
355 the presence of a data and code availability policy and the incidence of p-hacking,
356 including for research leveraging harder-to-access (e.g., administrative) data. They
357 also document a statistically insignificant relationship between voluntary provision of
358 data by authors on their homepages and selective reporting.

359 Replication Packages and Expectations

360 In an assessment of reproducers’ expectations regarding the quality of replication pack-
361 ages, we ask reproducers the following question in the individual survey: “Which of the
362 following best describes how the replication package aligned with your expectations”.
363 We find that more than half of reproducers report that the replication package aligned
364 reasonably with expectations, and an additional 26% of reproducers indicated that
365 the replication packages exceeded their initial expectations. Fewer than 10% report
366 that the replication package was worse than expected, possibly indicating that for this
367 small proportion of reproducers, the provided materials did not meet the anticipated
368 quality standards or may have lacked certain elements critical for an effective repro-
369 duction process. Overall, we find it encouraging that most reproducers found that the
370 provided materials exceeded or aligned well with their initial expectations.

371 Re-Analyses, P-Hacking, and Publication Bias

372 **t- and p-Curves:** We present both t-statistics and p-curves in Supplementary Materi-
373 als Appendix Figure 2. The top left panel provides the distribution of t-statistics from
374 the *originally* published estimates. We restrict the visualization to $t \in [0, 5]$, present
375 bins of width 0.1, and present an Epanechnikov kernel (with standard errors in blue,
376 along with renormalization at zero) which softens valleys and peaks. We provide refer-
377 ence lines at the conventional two-tailed significance levels. Roughly 60%, 50%, and

378 25% of test statistics are significant at the 10%, 5% and 1% levels, respectively. We
379 note especially that the distribution exhibits a peak (global maximum) just above the
380 two-star statistical significance threshold of $t = 1.96$ and a valley before the one-star
381 statistical significance threshold between $t = 1.0$ and $t = 1.65$. We take this as our
382 first piece of evidence that the original studies in our sample suffer from (marginal)
383 p-hacking and publication bias. The bottom left panel provides the equivalent p-curve
384 for p-values $\in [0.0025, 0.1500]$, with bins of width 0.0025. We have removed $p < 0.0025$
385 (for a two-tailed test this is roughly $t = 3$) for illustrative purposes only, as inclusion
386 of that mass in the left-most bar of the p-curve leads the resolution of the remaining
387 bars to be quite low. We note that, much like the peak after $t = 1.96$ and the valley
388 just before, the p-curve exhibits a too-tall bar just to the left of the $p = 0.05$ thresh-
389 old. Whether interpreted through the t-curve or p-curve, we consider this to be our
390 first piece of evidence that the sample of original studies suffers from some form of
391 p-hacking and publication bias.

392 We present t-and-p-curves using data from [5] in the right panels to serve as a
393 benchmark with which to compare the original studies. The top right panel presents
394 the distribution of t-statistics associated with hypothesis tests from articles published
395 in 25 leading economics journals in 2015 and 2018. These articles rely on one of four
396 popular identification methods (i.e., difference-in-differences, instrumental variable,
397 randomized controlled trials, and regression discontinuity design). Overall, the distri-
398 bution from our original studies sample is similar to that in [5], although with visually
399 markedly more bunching around the 5% significance threshold.

400 This could be due to at least three reasons. First, the extent of p-hacking and
401 publication bias might be larger in our sample. Second, reproducers might focus on the
402 most central claim(s) in original studies, while [5] focus on all claims. Arguably, the
403 central claim(s) could be more p-hacked or suffer from more publication bias. Third,
404 reproducers might choose to reproduce studies finding an effect or focus on replicating
405 claims that reject the null hypothesis.

406 Extended Data Figure 8 directly compares the distribution of test statistics for origi-
407 nal studies and our re-analyses. Just as in Supplementary Materials Appendix Figure
408 2, the top panels present t-distributions while the bottom panels present p-curves, and
409 the left panels present the original studies while the right panels now present statisti-
410 cal significance for the re-analyses. (See Supplementary Materials Appendix Figure
411 3 for the weighted distributions. For the re-analyses, we use the inverse of the num-
412 ber of test statistics presented in the reproduction report to weigh observations.) We
413 use this visual analysis to test whether re-analyses are less likely to reject the null
414 hypothesis than their original counterparts. If they are, we would expect to see less of
415 a peak (global maximum) just beyond the 5% statistical significance threshold and a
416 shift in the mass of test statistics leftward to the statistical insignificance region, i.e.,
417 if re-analyses ‘re-distribute’ the mass of test statistics without (or with less of) the
418 distorting effects of publication bias or p-hacking.

419 Our findings are striking. Moving from left to right in the top panels - from the
420 original to the re-analysis test statistics - there is a large shift in the mass of test statis-
421 tics from the *just* statistically significant at the 5% level region to the statistically

422 insignificant and 10% significance regions ($[0.10 > p > 0.05]$). We note this follow-
423 ing the global maximum has shifted in mass into where the valley was, and noting
424 also the much greater mass where $t = 0$. This visual result suggests that re-analyses
425 decrease the statistical significance of many originally published test statistics. This
426 is confirmed by a Kolmogorov–Smirnov test which rejects the null of equality of dis-
427 tributions ($p < 0.000$). A similar result emerges from visual inspection of the bottom
428 panels which display the same statistical significance distributions using p-values. An
429 over-abundance of just statistically significant results here is reflected in a particu-
430 larly large bar just to the left of $p = 0.05$. Under the assumption of no p-hacking and
431 publication bias, the p-curve should be non-increasing - this particularly large bar is
432 too large. We note that, in the same manner as the t-statistics no longer displaying a
433 marked peak once they have been re-analyzed, the p-curve resulting from re-analysis is
434 much better characterized as non-increasing (particularly at the statistical significance
435 thresholds).

436 The top panels of Extended Data Figure 9 reproduce the top panel of Extended
437 Data Figure 8 for economics and political science while the bottom panels of Extended
438 Data Figure 9 reproduce the bottom panel of Extended Data Figure 8. A reduction
439 in the peak of t-statistics or a reduction of the p-value bar just to left of $p = 0.05$ can
440 be seen for both economics and political science.

441 Supplementary Materials Appendix Figure 4 extends the visual analysis by offer-
442 ing a direct comparison of the statistical significance of an original estimate and its
443 corresponding re-analysis. Depicted is a histogram of $(p_{\text{replication}} - p_{\text{original}})$ with bars
444 of width 0.05. Interpretation of this difference-statistic is as follows. If the original
445 estimate and its re-analysis have very similar p-values, then the difference-statistic
446 will be close to zero. If the re-analysis p-value is high (indicating statistical insignifi-
447 cance) while the original p-value is low (indicating statistical significance), then this
448 difference-statistic will add to the right tail of the distribution. Notably, this is what
449 we see—a large proportion of re-analyses find similar p-values as the original (repre-
450 sented by both tall bars just above and just below zero), while we also see that the
451 right tail (which indicates re-analyses finding a lower statistical significance on aver-
452 age) being much thicker than the left tail (which indicates an original study finding a
453 lower statistical significance than its re-analysis). This trend is robust to weights and
454 is present in economics as well as in political science (second through fourth panels of
455 Supplementary Materials Appendix Figure 4.

456 So far, we have not distinguished between re-analyses that find an effect in the
457 same versus opposite direction as the original estimate. This is potentially problematic
458 if a large fraction of re-analyses finds a significant effect in the opposite direction. In
459 Supplementary Materials Appendix Figure 5 we make this distinction. Whenever the
460 re-analysis estimates an effect that is in the opposite direction, we assign the t-statistic
461 (top panels) or p-value (bottom panels) a negative value. We see that the statistical
462 significance of an original estimate with a re-analysis with an oppositely-signed effect
463 are often statistically significant. There is also still positive t-statistics, highlighting
464 the mass peak’s disappearance when moving from original to re-analysis.

465 Overall, our graphical analysis suggests that re-analyses can lead to both increases
466 and decreases in statistical significance, although the average effect is a reduction. In

467 all cases, there appears to be a downward shift of an over-abundance of just marginally
468 significant test statistics at the 5% level to the less and not statistically significant
469 regions.

470 Supplementary Materials Appendix Table 1 explicitly presents the change in
471 statistical significance from the original to a re-analysis at the test-statistic level. See
472 the main text for a description of this table.

473
474 **Formal Tests for P-Hacking and Publication Bias:** We next formally
475 document how re-analyses display a markedly different presence of p-hacking and pub-
476 lication bias. We first rely on caliper tests ([6]) which analyze test statistics within
477 a narrow range slightly above and below a statistical significance threshold. The
478 rationale behind this approach is rooted in the assumption that in the absence of
479 manipulation, be it due to publication bias or p-hacking, we would anticipate a com-
480 parable frequency of test statistics falling just below a significance threshold and those
481 falling just above it.

482 We estimate probit models where the dependent variable is a dummy variable that
483 takes the value one if a test statistic is statistically significant at the 5%-level, and
484 zero otherwise:

$$Pr(\textit{Significant}_{pr} = 1) = \Phi(\alpha + \lambda \textit{Reanalysis}_r) \quad (1)$$

485 where $\textit{Significant}_{pr}$ is a dummy variable for whether p-value p in report r is statis-
486 tically significant at the 10%, 5% or 1%-level. We rely on probit models throughout
487 and present the average marginal effects and associated standard errors clustered at
488 the report-level. The variable of interest is $\textit{Reanalysis}_r$, which represents a dummy
489 variable that takes a value of one if the p-value is associated with a re-analysis, and
490 zero if it is associated with the original publication.

491 The estimates are reported in Supplementary Materials Appendix Table 13 for the
492 5% significance threshold. In column 1, we restrict the sample to $[0.05 \pm 0.04]$. The
493 other columns repeat the specification in column 1 but with narrower bandwidths.
494 We find that re-analysis test statistics are about 10 percentage points less likely to be
495 statistically significant than an originally published test statistic. See Supplementary
496 Materials Appendix Table 14 for the 10% threshold. The point estimates for the 10%
497 level are similar, albeit less statistically significant.

498 We then rely on an application of [7]. The results are presented in Supplemen-
499 tary Materials Appendix Table 15. The columns μ , τ , and df represent the model's
500 estimated parameters (using an underlying t -distribution and symmetric sign prob-
501 abilities). The fourth column $[0, 1.645]$ presents the relative publication probability
502 for a t -statistic in the $[0, 1.645]$ interval compared to one in the reference interval of
503 $(2.576, \infty)$.

504 We find that a not statistically significant 'original analysis' test statistic is 17.16%
505 as likely as a very statistically significant test statistic to be observed (published).
506 Similarly, for the $(1.645, 1.96]$ interval, the original analyses offer only a 38.29% rela-
507 tive publication probability. These findings suggest that original articles in our sample
508 suffer from severe publication bias. As a comparison, we estimate that the same rela-
509 tive 'publication' probability for our re-analyses. This comparison serves only as a

510 benchmark since re-analyses are not submitted for publication and thus do not suffer
511 from publication bias. Nonetheless, we see this comparison as insightful. We find that
512 the relative ‘publication’ probability for a re-analysis jumps to 27.31% from 17.16%.
513 This trend continues for the $(1.645, 1.96]$ interval, where we observe a 64.30% relative
514 publication probability in a re-analysis versus 38.29%. For the relative publication
515 probability of test statistics significant at the 5% level, the original analyses offer an
516 almost equal probability of 107.40%, whereas the re-analysis is now slightly lower than
517 the original at 89.94%.

518 The second and third panels offer a similar analysis for the economics and political
519 science subsamples, respectively. The economics subsample behaves similarly to that of
520 the full sample. The political science subsample behaves similarly, with the exception
521 of the not statistically significant interval where the original analysis is more likely to
522 have not statistically significant result published.

523 We adopt diverse methodologies introduced by [5] and [8] as our foundation. Our
524 initial focus is on randomization tests, as designed by [5] to affirm the visually apparent
525 discontinuities near conventional statistical thresholds. We assess whether the con-
526 centration of test statistics just above versus just below these thresholds significantly
527 differs between the original studies and the re-analyses.

528 We operate under the assumption that the underlying distribution of p-values
529 (for any research method) is continuous and infinitely differentiable. Any observed
530 discontinuity in p-values is inferred to result from p-hacking or publication bias.

531 It is pertinent to note that publication bias is likely to operate predominantly in
532 a single direction (towards significance), as an excess of successes is more indicative
533 of bias than a scarcity. Hence, one-sided p-values are considered for our tests. The
534 outcomes are detailed in Supplementary Materials Appendix Table 16 for the 5%
535 threshold. In the first panel we use observations where $(0.01 < p < 0.09)$. The lower
536 panels use smaller windows. In the first panel, 78.3% of the original analysis p-values
537 within this window are significant. A test for whether this proportion is statistically
538 greater than 0.50 yields a p-value of 0.000. Similarly, we obtain very small p-values
539 for the smaller windows, confirming the presence of p-hacking or publication bias in
540 the sample of original studies.

541 We further test for the presence of p-hacking and publication bias by employing
542 the methodology and code by [8], and conducting six distinct tests to assess p-hacking
543 and publication bias: Binomial, Fisher’s, Discontinuity, CS1, CS2B, and LCM. The
544 outcomes are detailed in Supplementary Materials Appendix Figure 6. This figure
545 presents p-curves and test statistics for the battery of p-hacking tests for the full
546 sample in the first panel, for the economics subsample in the second, and the political
547 science subsample in the third.

548 In the absence of p-hacking and publication bias, the p-curve should be non-
549 increasing; a spike just to the left of the 0.05 threshold is indicative of p-hacking. This
550 spike is present in the full sample, though larger in the political science subsample
551 than the economics subsample.

552 Tests based on non-increasingness include the Binomial Test and Fisher’s test.
553 Only for the political science subsample is there sufficient evidence to reject the null
554 that the density (PDF) of p-values is non-increasing. In the absence of p-hacking, the

555 PDF is continuous. Again, only for the political science subsample is there sufficient
556 evidence to reject the null that the density (PDF) of p-values is continuous.

557 Under general assumptions, p-curves are completely monotonic (the CS1 test) and
558 are upper bounded in PDF and its derivatives (CS2B test). Here the trend reverses,
559 in that only the full sample and the economics subsample offer sufficient evidence to
560 reject the null of monotonicity and violations of the upper bound and derivatives of
561 the PDF.

562 Last, a consequence of hypothesizing the non-increasingness of the PDF is that the
563 PDF is also concave. The LCM test (Least Concave Majorant) assesses concavity of
564 the CDF of p-values. Again, only the full sample and the economics subsample offer
565 sufficient evidence to reject the null of concavity.

566 Overall, we take this mixed evidence to indicate the presence of p-hacking in both
567 the economics and political science subsamples, as well as the full sample.

568 **Recoding: Additional Results**

569 Supplementary Materials Appendix Table 12 shows our results. Out of 23 recoding
570 exercises, we find major differences for three studies and minor differences for 10
571 studies. Two of the major differences were uncovered when using a different software
572 and looking at the authors' code.

573 Additionally, one team that computationally reproduced the results using a differ-
574 ent *version* of the software used by the authors uncovered noteworthy differences in
575 the magnitude and significance of the estimates. About half the main claims were no
576 longer reproducible (i.e., same sign and statistically insignificant or different sign) due
577 to a change in the defaults used by base R when generating random numbers start-
578 ing in version 3.6.0. This is the only instance where using a different version of the
579 software led to major differences in the size and significance of the estimates.

580 These results suggest that most teams who recoded using a different software
581 language or without looking at the authors' code could obtain similar or very similar
582 results.

583 **Time Trends in Data and Code Availability: Additional Results**

584 **Replication Folder Availability:** Extended Data Figure 3 displays the percentage
585 of papers which have a replication package over the sample of 1150 papers which *should*
586 have a replication package according to the [AEA 2020 Definition](#). We see a general
587 increase in the trend of replication folders being provided between 2014 and 2020. We
588 found replication folders are attached to 59.1% and 70.0% of papers in 2014 and 2015,
589 respectively. Replication folder provision then increases to a seemingly stable value
590 close to 90% in 2021, 2022 and 2023.

591 Extended Data Figure 4 breaks down the previous figure's sample into those jour-
592 nals sampled from economics and political science. While the increasing trend within
593 both samples exists also in the subsamples, political science starts from a much lower
594 inclusion of replication packages. Political science papers have percentages of papers
595 with replication folders equal to 23.3%, 36.7%, 58.6% and 66.7% in years 2014, 2015,

596 2016, and 2017, respectively. In comparison, economics papers have percentages equal
597 to 72.5%, 82.5%, 80.2%, and 85.0% in years 2014, 2015, 2016 and 2017, respectively.

598 **Contents of Replication Folders:** The data presented in Extended Figures 5
599 through 6 are subsamples of varying sizes, reflecting the variation in what is required
600 for each paper’s replication folder. Each figure represents the percentage which possess
601 the associated variable (README, cleaning code, raw data, *etc.*). We display the
602 percentage of a binary variable equal to “Yes” if the package contained “All” of the
603 field, and “Not Yes” if the variable had only “Some” or “None” of that variable. We
604 are generally finding that replication folders are improving in their contents over time,
605 especially regarding the inclusion of READMEs and cleaning code. In earlier years,
606 replication folders were more likely to include analysis code and final data. Part of the
607 reason is that providing raw data and cleaning code makes redundant the inclusion
608 of final data (since one can generate the final data from the package). Extended Data
609 Figure 6 uses an alternative measure for “final data” which adds replication folders
610 that have complete (“Yes”) raw/intermediate data *and* have complete (“Yes”) cleaning
611 code to those replication folders which have “final data” explicitly included. This
612 alternative definition yields “final” data inclusion to be around 60% for the whole
613 sample with a range 56.7% as a minimum in 2014 and 2017 and a maximum of 65.2%
614 in 2015.

615 **Comparison Between Journals with a Data Editor and Journals With-**
616 **out a Data Editor in 2023:** One question raised is the importance of data editors
617 in improving the quality of replication folders. We split the sample presented earlier
618 in this section into those journals which did not have a data editor in 2023 and those
619 which did. Recall the journals which had a data editor in 2023 include: *American*
620 *Journal of Political Science*, *Journal of Politics*, *American Economic Review*, *Review*
621 *of Economic Studies*, *American Economic Journal: Macroeconomics*, *American Eco-*
622 *nomic Journal: Applied Economics*, *American Economic Journal: Economic Policy*,
623 *American Economic Review: Insights*, and *Economic Journal*. Journals that did not
624 have a data editor in 2023 include: *American Political Science Review*, *Journal of*
625 *Political Economy*, and *Quarterly Journal of Economics*.

626 We continue to use the binary variables presented in the previous section and
627 calculate a simple t-test to understand the difference in means which we present in
628 Supplementary Materials Appendix Table 18. In general, journals with data editors
629 are more likely to have replication folders (difference about 19%), are more likely to
630 contain READMEs (difference about 16%), and contain complete code (differences in
631 cleaning and analysis code being about 26% and 17%, respectively). The provision
632 of data is more nuanced. Journals with data editors were more likely to provide raw
633 data than those without data editors (difference about 19%) but less likely to provide
634 final or intermediate data (differences equal to -20% and -17%, respectively. Again, we
635 believe this is likely due to raw data being sufficient for producing final data. When
636 using our alternative definition of “final” data inclusion that takes into account raw
637 (or intermediate data) with complete cleaning code, the difference between journals
638 with data editors and those without data editors reduces to about -6%.

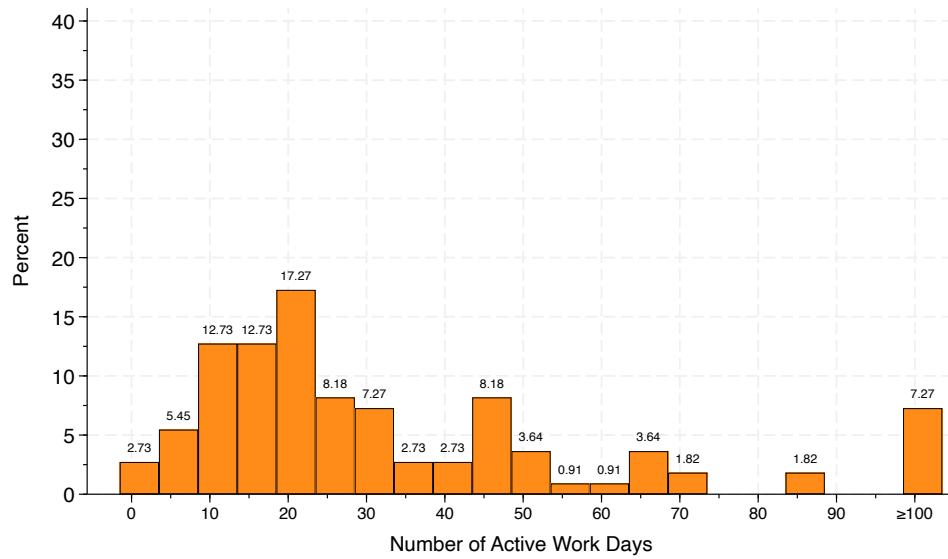
639 **Additional Discussion**

640 We aim for high-quality reproduction reports and believe our process contributes posi-
641 tively to the scientific community for at least four reasons. First, original authors are
642 allowed to respond and may point out flaws in the reproducers' work. In practice,
643 original authors and reproducers do not disagree on the completeness of the repli-
644 cation package (e.g., whether raw data is provided) nor on the presence of major
645 coding errors. Disagreements are almost always about the validity of robustness and
646 replicability. Second, A.B. or a co-director at I4R checks the tone of both the origi-
647 nal authors' response and reproducers' report. Third, while reproducers may make
648 mistakes, so do reviewers and editors. Our reproducers have the advantage of having
649 access to the replication package. They may identify coding errors and uncover coding
650 decisions which may not be discussed in the main body of the article. For example,
651 multiple studies in our sample do not mention the use of a weighting scheme for their
652 main analysis. This coding decision is obvious to a reproducer, but not to an editor or
653 reviewer. Relatedly, our teams of reproducers spent on average 13 active days work-
654 ing on their reproducibility and replicability. This may compare favorably to a typical
655 referee report, which is not prepared with peers and may involve subjectivity about
656 the contribution of the paper to the literature. Fourth, reproducers learn throughout
657 the process and benefit from this experience. This, in itself, is a positive contribution.

658 **Barriers to Reproducibility and Robustness:** We ask the following question
659 in the team survey: "For which of the following reasons were you unable to conduct
660 robustness checks, recoding exercises, extensions, or a replication using new data, prior
661 to communications with the original authors? (Select all which apply)". Extensive
662 Data Figure 6 provides a summary of the responses for these four categories. Out of
663 110 teams, 64 did not respond to the question. This suggests that the majority of
664 teams felt their replication packages contained enough to create a reproduction report
665 for I4R. That said, the lack of raw data restricted most what reproducers could do
666 when analyzing a paper across all four categories. Raw data inhibited 19% of teams
667 when trying to do robustness checks and 18% of teams wanting to recode key variables.
668 12% of teams also believed the lack of raw data inhibited their ability to perform a
669 reproduction and 13% of teams believed it inhibited their ability to perform extensions.
670 ([9] also provide evidence that non-reproducibility for the journal *Management Science*
671 is due to non-availability/accessibility of data.) The remaining reasons for potential
672 hurdles reproducers could have faced (like no intermediate data, no data dictionary,
673 unclear documentation, and/or unclear replication package) did not affect most teams.
674 About 7% of teams felt the original paper was unclear to the point of not being able
675 to perform robustness checks. We thus see a lack of raw data provided in a replication
676 package as a significant barrier to reproducibility and replicability, even in our selected
677 sample of journals which have data and code availability policies.

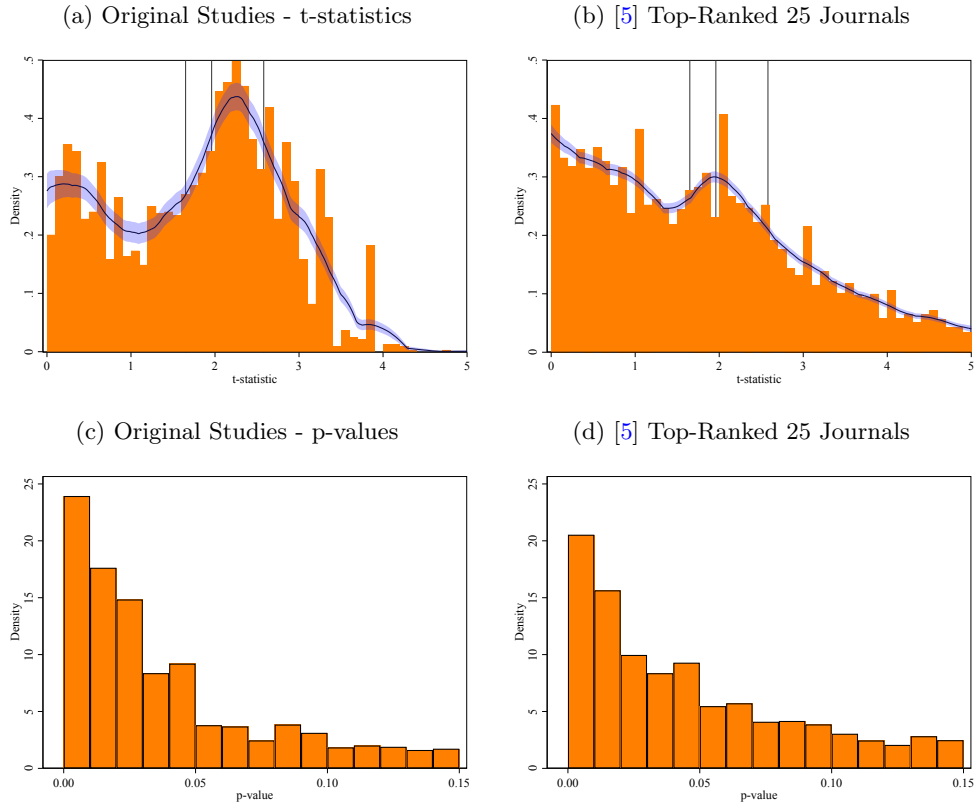
678 **On the Benefits for Reproducers:** We document several benefits of conducting
679 reproductions and replications. We ask the following question in the individual survey:
680 "Please indicate the degree to which your experience with I4R has contributed to
681 your improvement in the following areas." We offer six choices: (i) Networking, (ii)
682 coding skills, (iii) capacity to write a good replication package, (iv) learning difference
683 between reproduction and replication, (v) further ability as a researcher and (vi)

684 communicate issues with a paper to others. Supplementary Materials Appendix Table
685 17 provides a breakdown of the responses. We find that about 70% of reproducers
686 responded that their experience with I4R contributed either a lot or moderately to
687 their: (1) capacity to write a good replication package and (2) learning the difference
688 between reproduction and replication. Reproducers further said their experience with
689 I4R contributed at least moderately to furthering their ability as a researcher (about
690 53%) and their ability to communicate issues with a paper to others (about 60%).

Fig. 1: Histogram of Number of Active Work Days

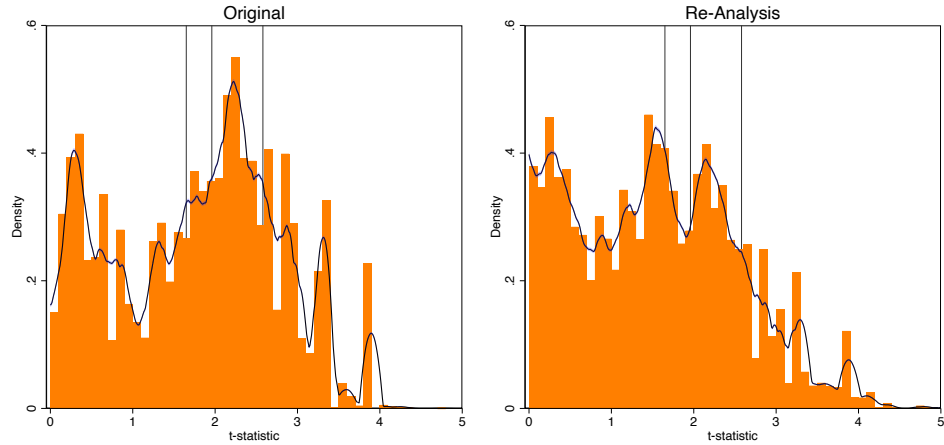
Legend: Data collected *via* survey of our reproducers after completing their reports. This figure illustrates the number of active days each team worked on their report.

Fig. 2: Distributions of t-Statistics and p-Values for Original Studies and [5]

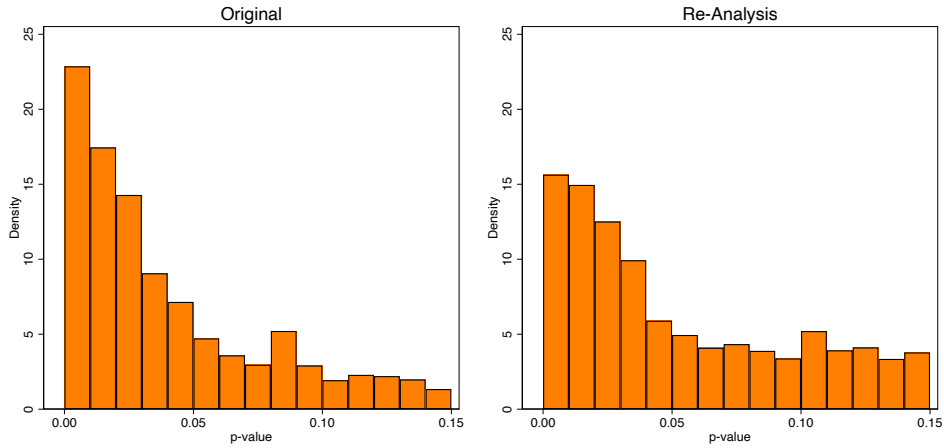


Legend: The top figures display a histogram of test statistics for $t \in [0, 5]$, with bins of width 0.1. The top left figure includes all original studies in our data set. As a comparison, the top right figure plots the corresponding histogram of z-statistics from the top-ranked 25 economics journals published in 2015 and 2018 (from [5]). Vertical reference lines are displayed at conventional two-tailed significance levels. We superimpose an Epanechnikov kernel density curve (which includes renormalization at 0). The bottom figures display histograms of test statistics for p-values $\in [0.0025, 0.1500]$, with bins of width 0.0025, among original studies and those from [5], respectively.

Fig. 3: Weighted Distributions of Statistics for Original Studies and Re-Analyses



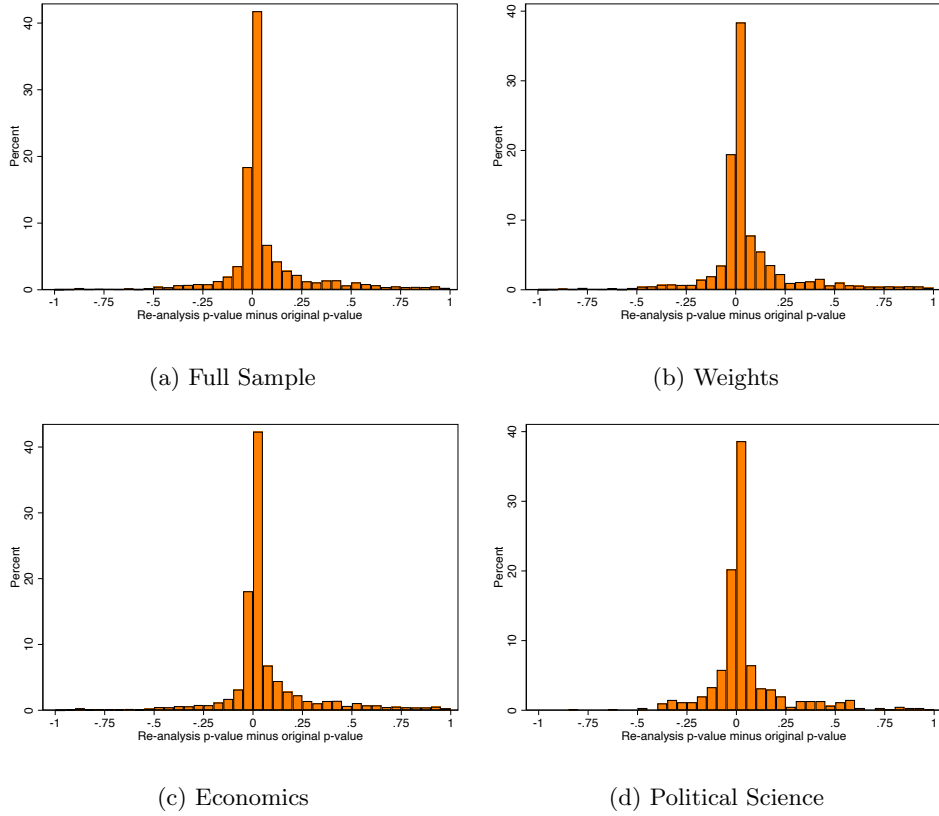
(a) t-statistic



(b) p-value

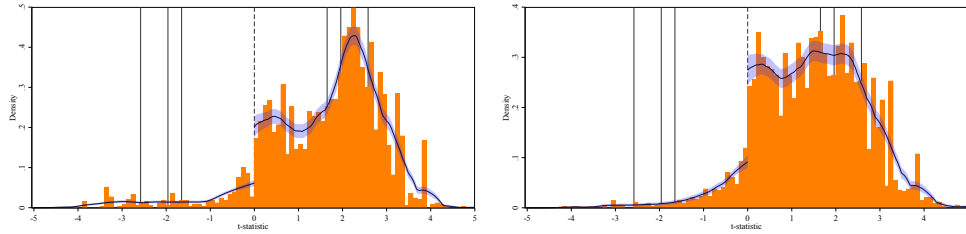
Legend: Top figures display histograms of test statistics for $t \in [0, 5]$, with bins of width 0.1, among original studies and re-analyses, respectively. Vertical reference lines are displayed at conventional two-tailed significance levels. We superimpose an Epanechnikov kernel density curve (which includes renormalization at 0). We use the inverse of the number of tests presented in the same article to weight observations. Bottom figures display histograms of test statistics for p-values $\in [0.0025, 0.1500]$, with bins of width 0.0025, among original studies and re-analyses, respectively. We use the inverse of the number of tests presented in the same article to weight observations.

Fig. 4: Distribution of $p_{\text{re-analysis}} - p_{\text{original}}$ by Weights and Fields



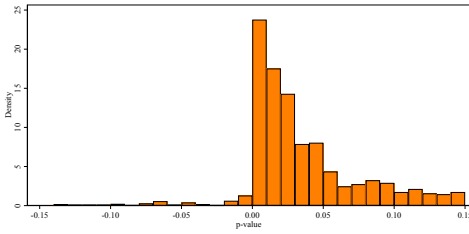
Legend: Second panel: We use the inverse of the number of test statistics in each reproduction report to weight observations. Third and fourth panel: The sample is restricted to original articles published in the indicated field. All panels: This figure presents the distribution of $(p_{\text{re-analysis}} - p_{\text{original}})$

Fig. 5: t and p -curves where negative represents a sign change from original to reproducer

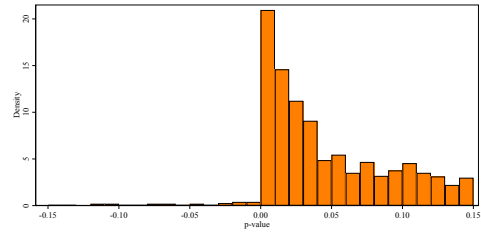


(a) Original

(b) Re-analysis



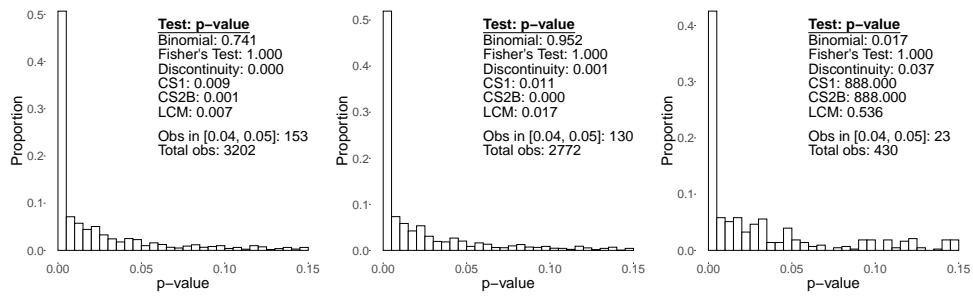
(c) Original



(d) Re-analysis

Legend: Top panels display a histogram of test statistics for $t \in [0, 5]$, with bins of width 0.1. We have added a dashed reference line at $t = 0$, demarcating the areas where the reproducers' and original estimates agree in sign. For both sides of the zero line, vertical reference lines are displayed at conventional two-tailed significance levels. We superimpose an Epanechnikov kernel density curve (which includes renormalization at 0), separately estimated for the positive and negative masses. Bottom panels display a histogram of test statistics for $p \in [0.00, 0.15]$, with bins of width 0.01. The left panels display statistics associated with originally published estimates. The right panels display statistics associated with reproducers' estimates. If the reproducer's estimated effect was of the opposite sign than the originally published estimate, we set the sign of the associated statistic to be negative.

Fig. 6: Applying [8]'s Tests



Notes: *Legend:* This figure present p-curves and results for the battery of p-hacking tests proposed in [8] for the full sample in the first panel, for the economics subsample in the second, and the political science subsample in the third. An error code of “888.00” represents an inability for that test to be calculated.

Table 1: Shifts in Statistical Significance Regions

Original Significance Level	Re-Analysis Significance Level			Total
	Sign Change	Not Sig.	Sig. at 5%	
Not Significant	4.99	28.47	3.71	37.16
Significant 5%	2.45	15.06	45.33	62.84
Total	7.44	43.53	49.03	100.00

Legend: This table illustrates shifts across significance and insignificance regions. Each row focuses on an initial level of statistical significance. Each column reports the share of re-analyses that ended up in each statistical significance region.

Table 2: Summary Statistics: Original Authors and Reproducers

	Mean (1)	Standard Deviation (2)	Minimum (3)	Maximum (4)
Test Statistics per Report	59.84	72.67	0	421
Year	2022.13	0.33	2022	2023
Economic Articles	0.72	0.45	0	1
Proportion of Economics Papers in Top 5	0.43	0.50	0	1
GS Citations (<i>As of Report Completed</i>)	43.98	71.39	0	573
Original Authors				
Number Original Authors	2.63	1.23	1	6
Share Graduate Student	0.06	0.18	0	1
Avg. Experience (<i>Years since PhD</i>)	11.21	6.34	0	31.50
Avg. GS Citations	4269.05	8882.00	31	55633.5
Replicators				
Number Replicators	3.25	1.22	1	7
Share Published Top 5 Econ/Targeted Poli Sci	0.15	0.36	0	1
Share Pub. Targeted Journals	0.30	0.46	0	1
Share Pub. Top 5/Targeted Poli Sci (<i>Past 5 Years</i>)	0.14	0.34	0	1
Share Pub. Targeted Journals (<i>Past 5 Years</i>)	0.26	0.44	0	1
Share Team Graduate Student	0.49	0.34	0	1
Avg. Experience (<i>Years since PhD</i>)	3.12	3.10	0	13.50
Avg. GS Citations	478.49	1016.67	0	6095.33
Comfortable programming in Stata	0.74	0.44	0	1
Comfortable programming in R	0.64	0.48	0	1
Comfortable programming in MATLAB	0.14	0.34	0	1

Legend: Each observation is an article. We do not weight test statistics. The Top 5 journals in economics are the American Economic Review, Econometrica, Journal of Political Economy, Quarterly Journal of Economics, and Review of Economic Studies. The 3 leading political science journals in our sample are the American Journal of Political Science, American Political Science Review and Journal of Politics. Panels two and three focus on the original authors and reproducers, respectively. Average experience is the mean of years since PhD. GS citations in the top panel refers to the number of Google Scholar citations for the original article as of the completion of the reproduction report. Average GS citations in the bottom panels refers to the number of Google Scholar citations at the time the report is completed.

Table 3: Summary Statistics by Types of Re-Analyses

	# Articles (1)	# Articles RGs (2)	# Articles Editor (3)	# Tests (4)
All Re-Analyses	103	81	22	6583
All Simultaneous Robustness Checks	51	41	10	809
Full Sample				
By Re-Analyses: Change in				
Control variables	58	45	13	1939
Sample	75	57	18	1774
Dependent Variable	23	18	5	285
Main Independent Variable	20	19	1	264
Estimation Method	33	28	5	605
Inference Method	23	19	4	542
Weighting Scheme	14	10	4	126
Use New Data	15	13	2	469
Economics				
By Re-Analyses: Change in				
Control variables	45	36	9	1612
Sample	55	47	8	1647
Dependent Variable	19	17	2	279
Main Independent Variable	15	15	0	195
Estimation Method	22	21	1	433
Inference Method	19	15	4	507
Weighting Scheme	9	8	1	80
Use New Data	13	11	2	461
Political Science				
By Re-Analyses: Change in				
Control variables	13	9	4	327
Sample	20	10	10	127
Dependent Variable	4	1	3	6
Main Independent Variable	5	4	1	69
Estimation Method	11	7	4	172
Inference Method	4	4	0	35
Weighting Scheme	5	2	3	46
Use New Data	2	2	0	8

Legend: This table shows the number of articles and test statistics for all re-analyses (top panel), by types of re-analyses (2nd panel), by types of re-analyses for economic articles (3rd panel) and by types of re-analyses for political science articles (bottom panel), respectively. The second and third columns show the number of reports created *via* replication games and editor stream, respectively.

Table 4: Communication with Original Authors

	# Authors Contacted (1)	% Responded (2)	% Short Note (3)	% Feedback (4)	% Formal Response (5)
Economics	75	93%	11%	61%	28%
Political Science	31	97%	14%	53%	33%
Total	106	94%	11%	59%	30%

Legend: This table provides information about original authors' responses. The second column shows that 94% of original authors that A.B. reached out to responded to his email. The remaining columns restrict the sample to those that responded.

Table 5: JEL Codes in our Sample

Top 10 JEL Codes in our Sample	Our Sample (All)		Representative Sample	
	Rank	%	Rank	%
D: Microeconomics	1	54.4	1	15.2
J: Labor and Demographic Economics	2	33.8	5	8.4
O: Economic Dev., Innov., Tech. Change, and Growth	3	33.8	6	7.9
I: Health, Education, and Welfare	4	29.4	10	6.3
H: Public Economics	5	17.6	9	6.3
N: Economic History	6	17.6	15	1.4
C: Mathematical and Quantitative Methods	7	16.2	2	15.1
E: Macroeconomics and Monetary Economics	8	13.2	4	10.7
L: Industrial Organization	9	13.2	11	5.6
G: Financial Economics	10	5.8	3	13.9
Q: Ag. and NR Econ & Envr. and Ecological Econ	11	7.4	7	7.7
P: Pol. Econ. and Comp. Economic Systems	12	5.8	17	0.8
Z: Other Special Topics	13	8.3	16	1
M: Bus. Admin and Bus. Econ & Mktg & Accg & Personnel Econ	14	3.3	13	1.8
R: Urban, Rural, Regional, Real Estate, and Trans. Economics	15	5.8	12	2.9
F: International Economics	16	2.5	8	7.6
K: Law and Economics	17	8.3	14	1.4
A: Gen. Econ & Teaching	18	NA	18	0.4
B: History of Econ Thought, Methodol., Heterodox Approaches	19	NA	19	0.4
Y: Miscellaneous Categories	20	NA	20	0.2

Legend: This table compares the JEL Codes in our sample and in a representative sample of economics papers ([3]). The JEL Codes are only available for some of the economic journals.

Table 6: Many-Analysts’ Robustness Rate and Reproducer Characteristics For Published Results Originally Statistically Significant at the 5% Level

RQ	Category				Total
	Neg. & Sig.	Neg. & Not Sig.	Pos. & Not Sig.	Pos. & Sig.	
1	42.78	43.33	13.89	0.00	100.00
2	36.75	24.79	30.13	8.33	100.00
3	0.00	33.33	63.89	2.78	100.00
4a	0.00	16.67	50.00	33.33	100.00
4b	16.67	0.00	50.00	33.33	100.00
4c	16.67	50.00	16.67	16.67	100.00
5a	0.00	16.67	52.78	30.56	100.00
5b	8.33	40.28	34.72	16.67	100.00
5c	22.22	52.78	8.33	16.67	100.00
6	0.00	30.56	52.78	16.67	100.00
7	8.33	13.89	61.11	16.67	100.00
8	0.00	23.61	76.39	0.00	100.00

Legend: Six many-analyst teams independently answered eight pre-registered research questions concerning the possible relationship between reproduction robustness rate and selected author/reproducer characteristics. This table restricts the analysis to originally published estimates that were statistically significant at the 5% level. The **columns** represent one of four categories that a many-analyst could classify their analysis; if a many-analyst found a relationship that was statistically significant (at the 5% level) and positive, it was included in the column ‘Pos. & Sig.’ The **rows** represent eight pre-registered research questions (two have three sub-questions). They are: 1- Does reproducibility/replicability rate depend on replicators’ experience coding? 2- Does reproducibility/replicability rate depend on replicators’ academic experience? 3- Does reproducibility/replicability rate depend on the authors’ experience? 4a- Does reproducibility/replicability rate depend on the interaction of the authors’ experience and replicators’ experience? In particular, (i) Are reproducibility/replicability rate higher when authors’ experience is high, and replicators’ experience is low (in comparison to similar levels)? 4b- Does reproducibility/replicability rate depend on the interaction of the authors’ experience and replicators’ experience? In particular, (ii) Are reproducibility/replicability rate higher when authors’ experience and replicators’ experience is similar (in comparison to dissimilar levels)? 4c- Does reproducibility/replicability rate depend on the interaction of the authors’ experience and replicators’ experience? In particular, (iii) Are reproducibility/replicability rate higher when authors’ experience is low, and replicators’ experience is high (in comparison to similar levels)? 5a- Does reproducibility/replicability rate depend on the interaction of the authors’ prestige and replicators’ prestige? In particular, (i) Are reproducibility/replicability rate higher when authors’ have high prestige, and replicators’ experience have low prestige (in comparison to similar levels)? 5b- Does reproducibility/replicability rate depend on the interaction of the authors’ prestige and replicators’ prestige? In particular, (ii) Are reproducibility/replicability rate higher when authors’ and replicators’ prestige is similar (in comparison to dissimilar levels)? 5c- Does reproducibility/replicability rate depend on the interaction of the authors’ prestige and replicators’ prestige? In particular, (iii) Are reproducibility/replicability rate higher when authors’ have low prestige, and replicators’ experience have high prestige (in comparison to similar levels)? 6- Does reproducibility/replicability rate depend on the original authors providing raw data? 7- Does reproducibility/replicability rate depend on the original authors providing raw or intermediate data? 8- Does reproducibility/replicability rate depend on the original authors providing cleaning code? For example, the **top row can be interpreted** as no many-analysts find a positive and statistically significant relationship between replicators’ experience coding and replication rate. 13.89% of many-analyst teams find a positive but not statistically significant relationship. 42.78% find a negative and statistically significant relationship, and 43.33% of many-analyst teams find a negative and not statistically significant relationship. Most many-analysts provide more than one estimate per research question; this table weights many-analysts equally.

Table 7: Many-Analysts’ Robustness Rate and Reproducer Characteristics For Published Results Originally Statistically Significant at the 10% Level

RQ	Category				Total
	Neg. & Sig.	Neg. & Not Sig.	Pos. & Not Sig.	Pos. & Sig.	
1	28.33	68.89	2.78	0.00	100.00
2	37.96	37.04	16.67	8.33	100.00
3	0.00	47.22	50.00	2.78	100.00
4a	0.00	8.33	33.33	58.33	100.00
4b	16.67	8.33	41.67	33.33	100.00
4c	8.33	58.33	0.00	33.33	100.00
5a	5.56	19.44	25.00	50.00	100.00
5b	16.67	36.11	30.56	16.67	100.00
5c	13.89	69.44	0.00	16.67	100.00
6	0.00	16.67	66.67	16.67	100.00
7	8.33	0.00	55.56	36.11	100.00
8	0.00	16.67	75.00	8.33	100.00

Legend: Six many-analyst teams independently answered eight pre-registered research questions concerning the possible relationship between reproduction robustness rate and selected author/reproducer characteristics. This table restricts the analysis to originally published estimates that were statistically significant at the 10% level. The **columns** represent one of four categories that a many-analyst could classify their analysis; if a many-analyst found a relationship that was statistically significant (at the 5% level) and positive, it was included in the column ‘Pos. & Sig.’ The **rows** represent eight pre-registered research questions (two have three sub-questions). They are: 1- Does reproducibility/replicability rate depend on replicators’ experience coding? 2- Does reproducibility/replicability rate depend on replicators’ academic experience? 3- Does reproducibility/replicability rate depend on the authors’ experience? 4a- Does reproducibility/replicability rate depend on the interaction of the authors’ experience and replicators’ experience? In particular, (i) Are reproducibility/replicability rate higher when authors’ experience is high, and replicators’ experience is low (in comparison to similar levels)? 4b- Does reproducibility/replicability rate depend on the interaction of the authors’ experience and replicators’ experience? In particular, (ii) Are reproducibility/replicability rate higher when authors’ experience and replicators’ experience is similar (in comparison to dissimilar levels)? 4c- Does reproducibility/replicability rate depend on the interaction of the authors’ experience and replicators’ experience? In particular, (iii) Are reproducibility/replicability rate higher when authors’ experience is low, and replicators’ experience is high (in comparison to similar levels)? 5a- Does reproducibility/replicability rate depend on the interaction of the authors’ prestige and replicators’ prestige? In particular, (i) Are reproducibility/replicability rate higher when authors’ have high prestige, and replicators’ experience have low prestige (in comparison to similar levels)? 5b- Does reproducibility/replicability rate depend on the interaction of the authors’ prestige and replicators’ prestige? In particular, (ii) Are reproducibility/replicability rate higher when authors’ and replicators’ prestige is similar (in comparison to dissimilar levels)? 5c- Does reproducibility/replicability rate depend on the interaction of the authors’ prestige and replicators’ prestige? In particular, (iii) Are reproducibility/replicability rate higher when authors’ have low prestige, and replicators’ experience have high prestige (in comparison to similar levels)? 6- Does reproducibility/replicability rate depend on the original authors providing raw data? 7- Does reproducibility/replicability rate depend on the original authors providing raw or intermediate data? 8- Does reproducibility/replicability rate depend on the original authors providing cleaning code? Most many-analysts provide more than one estimate per research question; this table weights many-analysts equally.

Table 8: Many-Analysts’ Robustness Rate and Reproducer Characteristics For Published Results Originally **Not** Statistically Significant at the 5% Level

RQ	Category				Total
	Neg. & Sig.	Neg. & Not Sig.	Pos. & Not Sig.	Pos. & Sig.	
1	0.00	3.33	88.33	8.33	100.00
2	0.00	35.19	64.81	0.00	100.00
3	0.00	11.11	88.89	0.00	100.00
4a	0.00	33.33	50.00	16.67	100.00
4b	0.00	41.67	41.67	16.67	100.00
4c	0.00	25.00	50.00	25.00	100.00
5a	0.00	16.67	69.44	13.89	100.00
5b	5.56	61.11	25.00	8.33	100.00
5c	0.00	29.17	40.28	30.56	100.00
6	8.33	66.67	25.00	0.00	100.00
7	0.00	58.33	33.33	8.33	100.00
8	16.67	58.33	19.44	5.56	100.00

Legend: Six many-analyst teams independently answered eight pre-registered research questions concerning the possible relationship between reproduction robustness rate and selected author/reproducer characteristics. This table restricts the analysis to originally published estimates that were not statistically significant at the 5% level. The **columns** represent one of four categories that a many-analyst could classify their analysis; if a many-analyst found a relationship that was statistically significant (at the 5% level) and positive, it was included in the column ‘Pos. & Sig.’ The **rows** represent eight pre-registered research questions (two have three sub-questions). They are: 1- Does reproducibility/replicability rate depend on replicators’ experience coding? 2- Does reproducibility/replicability rate depend on replicators’ academic experience? 3- Does reproducibility/replicability rate depend on the authors’ experience? 4a- Does reproducibility/replicability rate depend on the interaction of the authors’ experience and replicators’ experience? In particular, (i) Are reproducibility/replicability rate higher when authors’ experience is high, and replicators’ experience is low (in comparison to similar levels)? 4b- Does reproducibility/replicability rate depend on the interaction of the authors’ experience and replicators’ experience? In particular, (ii) Are reproducibility/replicability rate higher when authors’ experience and replicators’ experience is similar (in comparison to dissimilar levels)? 4c- Does reproducibility/replicability rate depend on the interaction of the authors’ experience and replicators’ experience? In particular, (iii) Are reproducibility/replicability rate higher when authors’ experience is low, and replicators’ experience is high (in comparison to similar levels)? 5a- Does reproducibility/replicability rate depend on the interaction of the authors’ prestige and replicators’ prestige? In particular, (i) Are reproducibility/replicability rate higher when authors’ have high prestige, and replicators’ experience have low prestige (in comparison to similar levels)? 5b- Does reproducibility/replicability rate depend on the interaction of the authors’ prestige and replicators’ prestige? In particular, (ii) Are reproducibility/replicability rate higher when authors’ and replicators’ prestige is similar (in comparison to dissimilar levels)? 5c- Does reproducibility/replicability rate depend on the interaction of the authors’ prestige and replicators’ prestige? In particular, (iii) Are reproducibility/replicability rate higher when authors’ have low prestige, and replicators’ experience have high prestige (in comparison to similar levels)? 6- Does reproducibility/replicability rate depend on the original authors providing raw data? 7- Does reproducibility/replicability rate depend on the original authors providing raw or intermediate data? 8- Does reproducibility/replicability rate depend on the original authors providing cleaning code? Most many-analysts provide more than one estimate per research question; this table weights many-analysts equally.

Table 9: Many-Analysts’ Robustness Rate and Reproducer Characteristics For Published Results Originally **Not** Statistically Significant at the 10% Level

RQ	Category				Total
	Neg. & Sig.	Neg. & Not Sig.	Pos. & Not Sig.	Pos. & Sig.	
1	0.00	11.67	71.67	16.67	100.00
2	0.00	35.19	64.81	0.00	100.00
3	0.00	36.11	63.89	0.00	100.00
4a	0.00	16.67	75.00	8.33	100.00
4b	0.00	38.89	52.78	8.33	100.00
4c	0.00	16.67	66.67	16.67	100.00
5a	0.00	45.83	29.17	25.00	100.00
5b	0.00	66.67	25.00	8.33	100.00
5c	0.00	37.50	37.50	25.00	100.00
6	0.00	83.33	16.67	0.00	100.00
7	0.00	61.11	30.56	8.33	100.00
8	16.67	58.33	16.67	8.33	100.00

Legend: Six many-analyst teams independently answered eight pre-registered research questions concerning the possible relationship between reproduction robustness rate and selected author/reproducer characteristics. This table restricts the analysis to originally published estimates that were not statistically significant at the 10% level. The **columns** represent one of four categories that a many-analyst could classify their analysis; if a many-analyst found a relationship that was statistically significant (at the 5% level) and positive, it was included in the column ‘Pos. & Sig.’ The **rows** represent eight pre-registered research questions (two have three sub-questions). They are: 1- Does reproducibility/replicability rate depend on replicators’ experience coding? 2- Does reproducibility/replicability rate depend on replicators’ academic experience? 3- Does reproducibility/replicability rate depend on the authors’ experience? 4a- Does reproducibility/replicability rate depend on the interaction of the authors’ experience and replicators’ experience? In particular, (i) Are reproducibility/replicability rate higher when authors’ experience is high, and replicators’ experience is low (in comparison to similar levels)? 4b- Does reproducibility/replicability rate depend on the interaction of the authors’ experience and replicators’ experience? In particular, (ii) Are reproducibility/replicability rate higher when authors’ experience and replicators’ experience is similar (in comparison to dissimilar levels)? 4c- Does reproducibility/replicability rate depend on the interaction of the authors’ experience and replicators’ experience? In particular, (iii) Are reproducibility/replicability rate higher when authors’ experience is low, and replicators’ experience is high (in comparison to similar levels)? 5a- Does reproducibility/replicability rate depend on the interaction of the authors’ prestige and replicators’ prestige? In particular, (i) Are reproducibility/replicability rate higher when authors’ have high prestige, and replicators’ experience have low prestige (in comparison to similar levels)? 5b- Does reproducibility/replicability rate depend on the interaction of the authors’ prestige and replicators’ prestige? In particular, (ii) Are reproducibility/replicability rate higher when authors’ and replicators’ prestige is similar (in comparison to dissimilar levels)? 5c- Does reproducibility/replicability rate depend on the interaction of the authors’ prestige and replicators’ prestige? In particular, (iii) Are reproducibility/replicability rate higher when authors’ have low prestige, and replicators’ experience have high prestige (in comparison to similar levels)? 6- Does reproducibility/replicability rate depend on the original authors providing raw data? 7- Does reproducibility/replicability rate depend on the original authors providing raw or intermediate data? 8- Does reproducibility/replicability rate depend on the original authors providing cleaning code? Most many-analysts provide more than one estimate per research question; this table weights many-analysts equally.

Table 10: Many-Analysts' Robustness Rate and Reproducer Characteristics - Only if Analyst Indicated the Effect Size was Meaningful

RQ	Category				Total
	Neg. & Sig.	Neg. & Not Sig.	Pos. & Not Sig.	Pos. & Sig.	
1	54.17	45.83	0.00	0.00	100.00
2	47.33	28.67	14.00	10.00	100.00
3	0.00	27.78	38.89	33.33	100.00
4a	0.00	0.00	50.00	50.00	100.00
4b	20.00	0.00	40.00	40.00	100.00
4c	16.67	50.00	16.67	16.67	100.00
5a	0.00	16.67	52.78	30.56	100.00
5b	12.50	25.00	37.50	25.00	100.00
5c	33.33	41.67	0.00	25.00	100.00
6	0.00	30.00	50.00	20.00	100.00
7	20.00	6.67	53.33	20.00	100.00
8	0.00	34.00	66.00	0.00	100.00

RQ	Category				Total
	Neg. & Sig.	Neg. & Not Sig.	Pos. & Not Sig.	Pos. & Sig.	
1	50.00	50.00	0.00	0.00	100.00
2	55.00	25.00	10.00	10.00	100.00
3	0.00	41.67	25.00	33.33	100.00
4a	0.00	0.00	12.50	87.50	100.00
4b	25.00	0.00	25.00	50.00	100.00
4c	8.33	58.33	0.00	33.33	100.00
5a	6.67	13.33	20.00	60.00	100.00
5b	25.00	25.00	25.00	25.00	100.00
5c	16.67	63.33	0.00	20.00	100.00
6	0.00	20.00	60.00	20.00	100.00
7	20.00	0.00	26.67	53.33	100.00
8	0.00	37.50	50.00	12.50	100.00

RQ	Category				Total
	Neg. & Sig.	Neg. & Not Sig.	Pos. & Not Sig.	Pos. & Sig.	
1	0.00	0.00	83.33	16.67	100.00
2	0.00	100.00	0.00	0.00	100.00
3	0.00	41.67	58.33	0.00	100.00
4a	0.00	33.33	33.33	33.33	100.00
4b	0.00	33.33	33.33	33.33	100.00
4c	0.00	33.33	33.33	33.33	100.00
5a	0.00	0.00	72.22	27.78	100.00
5b	11.11	72.22	0.00	16.67	100.00
5c	0.00	37.50	16.67	45.83	100.00
6	12.50	75.00	12.50	0.00	100.00
7	0.00	50.00	33.33	16.67	100.00
8	50.00	33.33	5.56	11.11	100.00

RQ	Category				Total
	Neg. & Sig.	Neg. & Not Sig.	Pos. & Not Sig.	Pos. & Sig.	
1	0.00	0.00	83.33	16.67	100.00
2	0.00	100.00	0.00	0.00	100.00
3	0.00	75.00	25.00	0.00	100.00
4a	0.00	0.00	83.33	16.67	100.00
4b	0.00	50.00	33.33	16.67	100.00
4c	0.00	12.50	75.00	12.50	100.00
5a	0.00	12.50	50.00	37.50	100.00
5b	0.00	83.33	0.00	16.67	100.00
5c	0.00	37.50	25.00	37.50	100.00
6	0.00	87.50	12.50	0.00	100.00
7	0.00	38.89	44.44	16.67	100.00
8	33.33	50.00	0.00	16.67	100.00

Legend: This table presents the same analysis as in Tables 6, 7, 8, and 9 while only including analyst results that were indicated by the analysis that "in your opinion, is the estimated effect size economically meaningful?" The first panel corresponds to Table 6. The second panel corresponds to Table 7. The third panel corresponds to Table 8. The fourth panel corresponds to Table 9. The rows correspond to the same research questions, and the columns represent the same effect sign and statistical significance categories. The cells remain weighted in the same manner.

Table 11: Summary Statistics by Journal

Discipline and Journal	# Articles Total (1)	# Articles RGs (2)	# Articles Editor (3)	# Tests (4)	Data Editor (5)
Economics	79	67	12	5,494	
American Economic Review	17	12	5	1,392	Yes
American Economic Review: Insights	2	0	2	149	Yes
American Economic J.: Applied Economics	9	6	3	260	Yes
American Economic J.: Economic Policy	11	11	0	811	Yes
American Economic J.: Macroeconomics	3	3	0	25	Yes
Economic Journal	20	18	2	1,262	Yes
Journal of Political Economy	8	8	0	1,283	No
Quarterly Journal of Economics	4	4	0	101	No
Review of Economic Studies	5	5	0	211	Yes
Political Science	31	16	15	1,089	
American Journal of Political Science	13	6	7	539	External
American Political Science Review	6	3	3	214	No
Journal of Politics	12	7	5	336	Yes
Total	110	83	27	6,583	

Legend: This table provides an overview of test statistics and articles reproduced and/or replicated by journal. Columns 1 and 4 indicate the number of article and test statistics per journal, respectively. Columns 3 and 4 report the number of articles per stream, where RGs is an acronym for Replication Games. Column 5 indicates if the journal has a data editor.

Table 12: Recoding Using Same or Different Softwares

	Identical (1)	Minor Differences (2)	Major Differences (3)	Total (4)
Same Software (Without Looking)	2	2	1	5
Different Software (Without Looking)	1	1	0	2
Different Software (Looking)	8	7	2	17
Total	10	10	3	23

Legend: This table illustrates the number of reports recoding the analysis (i) in the same software without looking at the authors' code/programs, (ii) using a different software language without looking at the authors' code/programs or (iii) using a different software language looking at the authors' code/programs.

Table 13: Caliper Tests, Significance at 5% Level

	Significant at 5% Level			
	(1)	(2)	(3)	(4)
Re-Analysis=1	-0.080*** (0.029)	-0.094** (0.045)	-0.073 (0.051)	-0.136** (0.068)
Observations	2,027	1,353	801	420
Threshold	0.05	0.05	0.05	0.05
Window	0.04	0.03	0.02	0.01

Legend: The dependent variable takes a value of one if $p \leq 0.05$. The variable Re-Analysis takes a value of one if the p -value is associated with a re-analysis, and zero if it is associated with the original publication. For example, in column 1 a Re-Analysis p -value is 8.9% less likely to be statistically significant than an original publication p -value at the 5% level in the small window of $0.01 \leq p \leq 0.09$. Standard errors clustered at the paper level. Estimated using probit, with marginal effects presented.

Table 14: Caliper Tests, Significance at 10% Level

	Significant at 10% Level			
	(1)	(2)	(3)	(4)
Re-Analysis=1	-0.085 (0.053)	-0.097 (0.063)	-0.134* (0.072)	-0.169* (0.091)
Observations	812	634	445	212
Threshold	0.10	0.10	0.10	0.10
Window	0.04	0.03	0.02	0.01

Legend: The dependent variable takes a value of one if $p \leq 0.10$. The variable Re-Analysis takes a value of one if the p -value is associated with a re-analysis, and zero if it is associated with the original publication. Standard errors clustered at the paper level. Estimated using probit, with marginal effects presented.

Table 15: Applying [7]

	μ	τ	df	[0, 1.645]	(1.645, 1.96]	(1.96, 2.576]
Original Analysis	0.0006	0.0024	1.2705	0.1716	0.3829	1.0740
Re-Analysis	0.0001	0.0000	1.2836	0.2731	0.6430	0.8994
Original Economics	0.0002	0.0011	1.1969	0.1522	0.3910	1.0556
Re-Analysis Economics	0.0000	0.0000	1.1942	0.2705	0.6107	0.9020
Original Political Science	0.0155	0.0254	2.1907	0.3078	0.3496	1.1846
Re-Analysis Political Science	0.0069	0.0155	2.4069	0.2653	0.6693	0.7916

Legend: An application of [7]. The columns μ , τ , and df represent the model's estimated parameters (using an underlying t -distribution and symmetric sign probabilities). The fourth column [0, 1.645] presents the relative publication probability for a t -statistic in the [0, 1.645] interval compared to one in the reference interval of (2.576, ∞).

Table 16: Randomization Tests, Significance at 5% Level

	Original Analysis
Proportion Significant in $.05 \pm .04$	0.783
One-Sided p-value against 0.50	0.000
Number of Tests in $.05 \pm .04$	2027.000
Proportion Significant in $.05 \pm .03$	0.751
One-Sided p-value against 0.50	0.000
Number of Tests in $.05 \pm .03$	1353.000
Proportion Significant in $.05 \pm .02$	0.689
One-Sided p-value against 0.50	0.000
Number of Tests in $.05 \pm .02$	801.000
Proportion Significant in $.05 \pm .01$	0.690
One-Sided p-value against 0.50	0.000
Number of Tests in $.05 \pm .01$	420.000

Legend: Following [5], in this table we present the results of binomial proportion tests where a success is defined as a statistically significant observation at the 5% level. In the first panel we use observations where $(0.01 < p < 0.09)$. The lower panels use smaller windows. We test if the proportion is statistically greater than 0.50. The associated p-values are then reported. We also include the number of observations in the third row. We do not weight articles.

Table 17: Please indicate the degree to which your experience with I4R has contributed to your improvement in the following areas (select all which apply):

	Nothing	A Little	Moderately	A Lot	Don't Know	Not Applicable
Networking	10.40	46.82	27.17	10.69	2.89	2.02
Coding Skills	19.08	40.17	26.88	10.98	1.73	1.16
Capacity to write a good replication package	5.19	21.90	46.97	23.63	1.15	1.15
Learning difference between reproduction and replication	6.65	19.36	36.71	33.53	3.47	0.29
Further ability as a researcher	5.20	39.02	38.15	17.05	0.29	0.29
Communicate issues with a paper to others	3.75	28.82	41.50	23.05	0.58	2.31

Legend: This table provides information on reproducers' feelings about how I4R contributed to their improvement in various areas. Each row represents a different category. Values are percentages and all rows in a category sum to 100. All values are unweighted.

Table 18: Mean Differences in Replication Package Contents in 2023 by those with and without a Data Editor

	Mean		Difference	P-value
	Has Data Editor	No Data Editor		
Is link to replication folder on website?	0.956	0.767	0.189	0.002
Does replication package contain a README?	0.922	0.767	0.156	0.021
Does replication package contain cleaning code?	0.822	0.567	0.256	0.004
Does replication package contain analysis code?	0.933	0.767	0.167	0.011
Does replication package contain raw data?	0.422	0.233	0.189	0.065
Does replication package contain intermediate data?	0.200	0.400	-0.200	0.029
Does replication package contain final data?	0.367	0.533	-0.167	0.110
(i) Final Data or (ii) cleaning code + raw/intermediate data?	0.544	0.600	-0.056	0.599

Legend: Columns (1) and (2) display the mean of the binary statistic labelled for each row. Column (3) shows the difference between column (1) and column (2). Column (4) shows the p-value of the applicable t-test. The binary variables are equal to one if the contents of the package were believed to be there and equal to zero if none or only some of the contents were included in the replication package. Samples vary across all rows as the statistics omits observations where it did not apply for whatever reason (*e.g.* simulations may have no data and no cleaning code). Journals which *had* a data editor in 2023 include: *American Journal of Political Science*, *Journal of Politics*, *American Economic Review*, *Review of Economic Studies*, *American Economic Journal: Macroeconomics*, *American Economic Journal: Applied Economics*, *American Economic Journal: Economic Policy*, *American Economic Review: Insights*, and *Economic Journal*. Journals that did *not have* a data editor in 2023 include: *American Political Science Review*, *Journal of Political Economy*, and *Quarterly Journal of Economics*.

693 **List of Articles and Reproduction**

694 **Reproduction Report**

695 **Title Original Study:** Antinormative Messaging, Group Cues, and the Nuclear
696 Ban Treaty

697 **doi:** <https://doi.org/10.1086/714924>, *Journal of Politics*

698 **Abstract:** Herzog, Baron, and Gibbons (2022) explore the effects of exposure to
699 official elite rhetoric and group cues on public support against the international
700 nuclear weapons prohibition norm. The authors find that elite cues, in particular
701 security and institutional cues, increase individuals' opposition to the Treaty on
702 the Prohibition of Nuclear Weapons (TPNW). However, elite cues do not seem to
703 have an effect on changing individuals' broader attitudes towards nuclear weapons,
704 as measured by individuals' existing opposition to nuclear arms. We replicate and
705 expand the authors' methods and results to test the robustness of the effects found
706 in the study. First, we reproduce the main finding using the authors' original data
707 and method. We do not find any coding errors that undermine the authors' analysis

708 or conclusions. Second, we test the robustness of the results by (1) using a dif-
709 ferent operationalization of party identity, and (2) calculating additional subgroup
710 analysis for gender. We find no significant differences between our replicated and
711 the original results, however females' support for the TPNW is more responsive to
712 security cues, while males' support is more responsive to institutions cues.

713 **Link to Full Report:** <https://econpapers.repec.org/paper/zbwi4rdps/97.htm>

714 **Replication Package:** <https://osf.io/xbvzg/>

715 **Link to Original Authors' Response:** <https://econpapers.repec.org/paper/zbwi4rdps/98.htm>

716
717 **Original Authors' Package:** [https://dataverse.harvard.edu/dataset.xhtml;
718 jsessionid=5bc4beb0bd5aef9d7a5ba5284fc6?persistentId=doi%3A10.7910%
719 2FDVN%2FGLT4FX&version=&q=&fileTypeGroupFacet=%22Text%22&
720 fileAccess=&fileSortField=size](https://dataverse.harvard.edu/dataset.xhtml;jsessionid=5bc4beb0bd5aef9d7a5ba5284fc6?persistentId=doi%3A10.7910%2FDVN%2FGLT4FX&version=&q=&fileTypeGroupFacet=%22Text%22&fileAccess=&fileSortField=size)

721 **Reproduction Report**

722 **Title Original Study:** Ascriptive Characteristics and Perceptions of Impropriety
723 in the Rule of Law: Race, Gender, and Public Assessments of Whether Judges Can
724 Be Impartial

725 **doi:** <https://doi.org/10.1111/ajps.12599>, American Journal of Political Science

726 **Abstract:** Ono & Zilis (2022) investigated the effects of ascriptive characteristics
727 of US American judges, such as race, gender, and ethnicity, on citizens' perceptions
728 of the judges' professional impropriety and bias in their rulings. They conducted
729 two studies, comparing citizens' perceptions of different ascriptive characteristics
730 and judgments about the judges' biases and the need for recusal from cases. They
731 found that political and ideological predispositions shape perceptions of judicial
732 impropriety. In this comment, we recode the analysis using a different software and
733 conduct robustness checks. We were able to reproduce the main results.

734 **Link to Full Report:** <https://osf.io/yf48r/>

735 **Replication Package:** <https://osf.io/yf48r/>

736 **Link to Original Authors' Response:** No response.

737 **Original Authors' Package:** [https://dataverse.harvard.edu/dataset.xhtml?](https://dataverse.harvard.edu/dataset.xhtml?persistentId=doi:10.7910/DVN/ZHOL6Y)
738 [persistentId=doi:10.7910/DVN/ZHOL6Y](https://dataverse.harvard.edu/dataset.xhtml?persistentId=doi:10.7910/DVN/ZHOL6Y)

739 **Reproduction Report**

740 **Title Original Study:** Assortative Matching at the Top of the Distribution:
741 Evidence from the World’s Most Exclusive Marriage Market

742 **doi:** <https://doi.org/10.1257/app.20180463>, American Economic Journal: Applied
743 Economics

744 **Report’s Abstract:** Goni (2022) relies on a novel data on peerage marriages in
745 Britain to examine the impact of matching technology on marital sorting. He relies
746 on the London Season interruption (1861 – 1863) as a natural experiment that
747 raised search costs and reduced market segregation. In his preferred specification, he
748 exploits exogenous variation in womens’ probability to marry during the interrup-
749 tion for their age in 1861 and finds that the interruption increased the probability
750 of marrying a commoner; reduced the probability of marrying an heir, increased
751 the difference in spouses’ family landholdings (in absolute value); decreased the
752 difference in spouses’ family landholdings (husband – wife); and increased the like-
753 likelihood of never getting married (See Table 2, columns 1 to 6, respectively). First,
754 we reproduce the papers’ main findings and find no coding errors. Second, we test
755 the robustness of the results to (1) the use of additional fixed effects and (2) sample
756 restrictions. Finally, we examine the heterogeneous effects of this interruption by
757 age and year. We find that original estimates are robust and are not significantly
758 affected using these alternative specifications.

759 **Link to Full Report:** <https://econpapers.repec.org/paper/zbwi4rdps/47.htm>

760 **Link to Replicators’ Package:** <https://osf.io/pqsem/>

761 **Original Author’s Response:** “I now reviewed the report carefully and with
762 interest, and I am glad to see that the authors succeeded in replicating all results
763 and found no coding errors. I hope the replication package was clear and easy to
764 work with. I am also happy to see that they performed several additional robustness
765 checks and heterogeneity analysis, and that these show that the original estimates
766 “are robust and are not significantly affected using these alternative specifications”
767 (p. 1). Given this, and the replicators’ conclusion that “the study’s main findings
768 demonstrate robustness and reliability” (p. 7), I think that there is nothing sub-
769 stantial for me to write in a response in the form of a discussion paper. This is
770 because both the replication exercise and the additional analysis found no major
771 issues in the original work to respond to. I would also like to thank the authors
772 for the fairness and professionalism of their report, and also for the time and effort
773 they put in producing it, from which I ultimately benefit — as it adds to the cred-
774 ibility of my original paper — as well as the profession as a whole benefits — as
775 making replication exercises more common is important for economics.

776 Please let me know if I can be of any further assistance regarding this Repro-
777 duction Report in the future. I am at your or the authors’ disposal, in case I can
778 be of help in clarifying anything in the replication package or in the analysis of
779 the original paper. As I stated above, I believe that increasing replication rates is
780 important for our field, as it is making original datasets publicly available — even
781 when, as in the case of my paper, the data collection is an important part of my
782 contribution, and in this situation, many do not grant public access to the original
783 data.”

784 **Original Authors' Package:** [https://www.openicpsr.org/openicpsr/project/](https://www.openicpsr.org/openicpsr/project/140921/version/V1/view)
785 [140921/version/V1/view](https://www.openicpsr.org/openicpsr/project/140921/version/V1/view)

786 **Reproduction Report**

787 **Title Original Study:** Black Workers in White Places: Daytime Racial Diversity
788 and White Public Opinion

789 **doi:** <https://doi.org/10.1086/716289>, Journal of Politics

790 **Report’s Abstract:** In this replication study, we revisit the main empirical claims
791 of Hamel and Wilcox-Archuleta’s (HW) 2022 study on the impact of daytime racial
792 diversity on White Americans’ voting behavior and racial attitudes. HW introduce
793 a novel zip code level measure of racial diversity that accounts for the influx of
794 Black workers during daytime, showing that conventional purely residential based
795 measures often underestimate the true degree of experienced racial diversity. Using
796 survey data from the CCES, their findings suggest a negative correlation between
797 racial flux and White Americans’ Democratic voting tendencies and a positive
798 correlation with racial resentment and opposition to affirmative action, all while
799 controlling for the residential share of Blacks in the zip code. We assess the repli-
800 cability of these findings by: (1) replicating the main results using the provided
801 replication code, (2) reconstructing the racial flux measure and survey from raw
802 data, (3) conducting multiverse analyses, and (4) replicating the analysis using an
803 alternative data source. Our replication validates the robustness and accuracy of
804 HW’s initial conclusions, emphasizing the role of daytime racial diversity in shaping
805 White Americans’ political and racial attitudes.

806 **Link to Full Report:** <https://econpapers.repec.org/paper/zbwi4rdps/61.htm>

807 **Link to Replicators’ Package:** <https://osf.io/ue4pm/>

808 **Original Authors’ Response:** “We enjoyed reading the replication, and don’t
809 see a need to write a response.

810 Thank you for doing this important work.”

811 **Original Authors’ Package:** [https://dataverse.harvard.edu/dataset.xhtml;
812 jsessionid=da88db7b3367419a2d1a87a9e687?persistentId=doi%3A10.7910%
813 2FDVN%2FFMOR6K&version=&q=&fileTypeGroupFacet=&fileAccess=
814 &fileSortField=type](https://dataverse.harvard.edu/dataset.xhtml;jsessionid=da88db7b3367419a2d1a87a9e687?persistentId=doi%3A10.7910%2FDVN%2FFMOR6K&version=&q=&fileTypeGroupFacet=&fileAccess=&fileSortField=type)

815 **Reproduction Report**

816 **Title Original Study:** Brahmin Left Versus Merchant Right: Changing Political
817 Cleavages in 21 Western Democracies, 1948–2020

818 **doi:** <https://doi.org/10.1093/qje/qjab036>, Quarterly Journal of Economics

819 **Report’s Abstract:** Gethin, Martínez-Toledano and Piketty (2022) analyze the
820 long-run evolution of political cleavages using a new database on socioeconomic
821 determinants of voting from approximately 300 elections in 21 Western democracies
822 between 1948 and 2020. They find that, in the 1950s and 1960s, voting for the
823 ”left” was associated with lower-educated and low-income voters. After that, voting
824 for the ”left” has gradually become associated with higher-educated voters, while
825 high income voters have continued to vote for the ”right”. In the 2010s, there is
826 a disconnection between the effects of income and education on voting. In this
827 replication, we first conduct a computational reproduction, using the replication
828 package provided by the authors. Second, we do a robustness replication testing to
829 what extent the original results are robust to i) restricting the sample to ”core” left
830 and right parties, ii) analyzing the top 80% versus bottom 20%, iii) weighting by
831 population, iv) dropping control variables, and v) using country fixed effects. The
832 main results of the paper are found to be largely replicable and robust.

833 **Link to Full Report:** <https://econpapers.repec.org/paper/zbwi4rdps/19.htm>

834 **Link to Replicators’ Package:** <https://osf.io/2hpeq/>

835 **Original Authors’ Response:** “Thank you for your mail and for your interesting
836 report! We are happy to see that you were able to easily replicate our results and
837 that our main conclusions were found to be largely robust. In this context, we do
838 not think that an answer from our side would be particularly useful: we are happy
839 with the report as it is.

840 Thank you for the very valuable work that your institute is producing in testing
841 the replicability and robustness of published studies!”

842 **Original Authors’ Package:** [https://dataverse.harvard.edu/dataset.xhtml?](https://dataverse.harvard.edu/dataset.xhtml?persistentId=doi:10.7910/DVN/XUSWG6)
843 [persistentId=doi:10.7910/DVN/XUSWG6](https://dataverse.harvard.edu/dataset.xhtml?persistentId=doi:10.7910/DVN/XUSWG6)

844 **Reproduction Report**

845 **Title Original Study:** Bubbles, Crashes, and Economic Growth: Theory and
846 Evidence

847 **doi:** <https://doi.org/10.1257/mac.20220015>, American Economic Journal: Macroeconomics
848

849 **Report’s Abstract:** Guerron-Quintana, Hirano, and Jinnai (2023) explore the
850 short-, medium-, and long-run effects of financial bubbles on economic growth by
851 way of a macroeconomic general equilibrium framework. In their model, a key
852 theoretical result is that, in net terms, the “crowding in” of capital investment
853 during a bubble ushers the economy onto a higher balanced growth path post-
854 bubble than it was on pre-bubble (Figure 10), thus (seemingly) suggesting that
855 economic bubbles are growth-enhancing. In turn, the main result of the paper is
856 that this positive view of bubbles is a fallacy so long as the latter are recurrent,
857 namely because a counterfactual economy in which bubbles never occur in the first
858 place grows at a significantly faster pace (Figure 10). The reason for this is that
859 the expectation of future bubbles stifles capital investment and, as such, reduces
860 economic growth in the long run.

861 We successfully reproduce the paper’s main figures using the original code pro-
862 vided in the replication package. Given the hard-coded nature of all empirical data
863 used in the paper, most of our efforts are devoted to reproducing the employed
864 empirical data itself and, in turn, conducting a direct replication with our own mea-
865 sures. Using various specifications of the HP filter, we are successful in qualitatively,
866 but not quantitatively reproducing the paper’s main time series (stock-market-to-
867 GDP ratio). Nevertheless, even without updating the model’s parameterization,
868 the paper’s main empirical findings (i.e. Figures 8-10) are largely robust to our own
869 measure. In turn, we are successful in quantitatively reproducing the second key
870 time series (credit-to-GDP ratio), albeit only with a highly unusual specification
871 of the HP filter’s smoothing parameter (10^{10} instead of 1600 for quarterly data).
872 We find that, unlike in the case of the stock-to-GDP ratio, the paper’s (auxiliary)
873 findings are not robust to our own credit-to-GDP series

874 **Link to Full Report:** <https://osf.io/d76tn/>

875 **Link to Replicators’ Package:** <https://osf.io/d76tn/>

876 **Original Authors’ Response:** Provided a short response and answered a
877 question. Did not provide a final response as of November 2025.

878 **Original Authors’ Package:** [https://www.openicpsr.org/openicpsr/project/
879 173441/version/V1/view](https://www.openicpsr.org/openicpsr/project/173441/version/V1/view)

880 **Reproduction Report**

881 **Title Original Study:** Campaign Contributions and Roll-Call Voting in the U.S.
882 House of Representatives: The Case of the Sugar Industry

883 **doi:** <https://doi.org/10.1017/S0003055422000466>, American Political Science
884 Review

885 **Report's Abstract:** In their study, Grier et al. (2023) explore the causal relation-
886 ship between campaign contributions and roll-call voting. Their analysis focuses on
887 the influence of campaign contributions on two specific anti-sugar votes conducted
888 in 2013 and 2018. The authors identify a substantial increase in inflationadjusted
889 sugar contributions from the sugar industry to incumbent politicians between these
890 two voting events. The aim of our research is to replicate and validate the authors'
891 main models. In addition to cross-platform replication, we conduct several robust-
892 ness checks to further examine the reliability of their findings. These include (1)
893 clustering the standard errors, (2) utilizing an Ordinary Least Squares (OLS) model
894 instead of the authors' logistic regression, and (3) altering the dependent variable
895 to represent the change in the vote from 2013 to 2018. Our results largely confirm
896 the authors' findings and reveal additional insights regarding the money buys vote
897 hypothesis.

898 **Link to Full Report:** <https://econpapers.repec.org/paper/zbwi4rdps/57.htm>

899 **Link to Replicators' Package:** <https://osf.io/4hjb9/>

900 **Original Authors' Final Response:** "We thank the Institute for Replication for
901 their diligent work replicating and performing some extensions to our 2023 APSR
902 paper. Replication is an important and often undervalued work in the scientific
903 process. Of course we are quite pleased to see that our results do replicate and that
904 the extensions performed largely support the results and ideas we advanced in our
905 paper. Keep up the good work!"

906 **Original Authors' Package:** [https://dataverse.harvard.edu/dataset.xhtml?](https://dataverse.harvard.edu/dataset.xhtml?persistentId=doi:10.7910/DVN/2IFZR9)
907 [persistentId=doi:10.7910/DVN/2IFZR9](https://dataverse.harvard.edu/dataset.xhtml?persistentId=doi:10.7910/DVN/2IFZR9)

908 **Reproduction Report**

909 **Title Original Study:** Can Information Reduce Ethnic Discrimination? Evidence
910 from Airbnb

911 **doi:** <https://doi.org/10.1257/app.20190188>, American Economic Journal: Applied
912 Economics

913 **Report's Abstract:** Laouénan & Rathelot (2022) investigate the mechanism
914 underlying ethnic discrimination using self-collected panel data from Airbnb
915 between 2014 and 2017. They find that hosts from minority groups charge 3.2%
916 less than those from the majority group within the same neighbourhood. Using a
917 theoretical framework, they estimate that the ethnic price gap vanishes as more
918 information (reviews) become available conditional on observables. The point esti-
919 mates for their main results are statistically significant at the 1% level. This finding
920 suggests that ethnic discrimination is due to statistical discrimination rather than
921 taste-based discrimination. First, we reproduce the original article's main findings
922 using R, whereby the authors of the original article use STATA. We can repro-
923 duce the main findings in R except for a few marginal discrepancies at the second
924 or third decimal place. Second, we extend two robustness analyses reported in the
925 original article. These robustness analyses impose restrictions on the sample and
926 these restrictions are not justified in the article. Once these restrictions are not
927 imposed, the picture becomes more complex and the robustness analysis warrants
928 more discussion. However, only a small fraction of the observations causes some
929 ambiguity and there might be good reasons to impose restrictions. Transparently
930 presenting the robustness analyses with and without restrictions, motivating the
931 restrictions and discussing its implications for the main findings would have been
932 desirable. Generally, the original article does a great job with regard to repro-
933 ducibility by providing data, code and documentation that ease the reproduction
934 of a complex analysis. We conclude that our reproduction and replication support
935 the main findings of the original article.

936 **Link to Full Report:** <https://osf.io/zn98a/>

937 **Link to Replicators' Package:** [https://github.com/TuanNguyen04/Replication-
938 Airbnb](https://github.com/TuanNguyen04/Replication-Airbnb)

939 **Original Authors' Response:** The authors provided initial feedback which the
940 replicators took it into account.

941 **Original Authors' Package:** [https://www.openicpsr.org/openicpsr/project/
942 120078/version/V1/view](https://www.openicpsr.org/openicpsr/project/120078/version/V1/view)

943 **Reproduction Report**

944 **Title Original Study:** Can Technology Solve the Principal-Agent Problem?
945 Evidence from China's War on Air Pollution

946 **doi:** <https://doi.org/10.1257/aeri.20200373>, American Economic Review: Insights

947 **Report's Abstract:** Greenstone et al. examine the effect of the introduction of
948 automatic air pollution monitoring on the reporting of local air pollution in China.
949 Using 654 regression discontinuity designs (RDDs) based on city-level variation
950 in the day that monitoring was automated, they find an immediate and lasting
951 increase of 35 percent in reported PM10 concentrations post-automation. More-
952 over, they find that automation's introduction increases online searches for face
953 masks and air filters by 200 percent and 28 percent, respectively, using an RDD.
954 Results are consistent when using an event study design. First, we were able to
955 computationally replicate the results. Second, we find that results are robust to
956 more flexible specifications of the weather variables, to re-constructed weather vari-
957 ables using the same matching procedure as the authors (i.e., closest station) and
958 meteorological data with additional weather stations, to alternative construction of
959 the weather variables using an inverse distance weighted approach of the surround-
960 ing weather stations, and to more flexible choices of fixed effects (up to the city
961 level). Finally, we find limited evidence of discontinuity in objective measures of
962 ground pollution (i.e., AOD) for a sub-sample using alternative weather variables.
963 The estimate, however, is economically insignificant. Moreover, no discontinuity is
964 observed in the full sample. Therefore, we believe this result does not invalidate
965 the original study's findings.

966 **Link to Full Report:** <https://osf.io/b7dn2/>

967 **Link to Replicators' Package:** [https://osf.io/m8hfr/?view_only=](https://osf.io/m8hfr/?view_only=9f6632ec96c0451daf0f8889b9ad2b25)
968 [9f6632ec96c0451daf0f8889b9ad2b25](https://osf.io/m8hfr/?view_only=9f6632ec96c0451daf0f8889b9ad2b25)

969 **Original Authors' Response:** <https://osf.io/b7dn2/>

970 **Original Authors' Package:** [https://www.openicpsr.org/openicpsr/project/](https://www.openicpsr.org/openicpsr/project/125321/version/V1/view)
971 [125321/version/V1/view](https://www.openicpsr.org/openicpsr/project/125321/version/V1/view)

972 **Reproduction Report**

973 **Title Original Study:** Can't We All Just Get Along? How Women MPs Can
974 Ameliorate Affective Polarization in Western Publics

975 **doi:** <https://doi.org/10.1017/S0003055422000491>, American Political Science
976 Review

977 **Report's Abstract:** We present a replication and extension of Adams et al.
978 (2023), examining the influence of women Members of Parliament (MPs) on affective
979 polarization. Conducted during the 2023 Montreal Replication Games, our
980 analysis reaffirms the original findings through the authors' base R code and a tidy-
981 verse simplification. Our results highlight that the mitigating effect on polarization
982 is predominantly observed among left-wing respondents, with null effects noted
983 for centrist and right-wing parties. This discrepancy is attributed to left-wing parties'
984 explicit commitment to gender equality. Further analysis reveals the study's
985 robustness across different countries and years (1996-2007) while addressing data
986 structure and imputation methods to ensure reliability. Our findings underscore
987 the nuanced role of women MPs in political dynamics, particularly among left-wing
988 voters, against democratic backsliding concerns.

989 **Link to Full Report:** <https://osf.io/69px3/>

990 **Link to Replicators' Package:** <https://osf.io/69px3/>

991 **Original Authors' Response:** Thank you for replicating our paper Can't We
992 All Just Get Along? How Women MPs Can Ameliorate Affective Polarization in
993 Western Publics (APSR 2023) as part of the Montreal Replication Games. We
994 appreciate the attention to detail and rigor applied to the replication project. We are
995 pleased that our initial results replicate well. We appreciate your robust approach
996 to testing the stability of our findings using a country and year 'leave-one-out' cross-
997 validation strategy. We also thank you for catching the coding error which dropped
998 a handful of cases from the original analysis; we are glad that the results remain
999 substantively the same when this error is corrected. We also are interested in the
1000 results from the extension you undertook, finding that our results are primarily
1001 driven by left-wing parties' supporters, in particular parties from the green, radical
1002 left and social Democratic parties. On the other hand, the point estimates are
1003 positive for all parties excepting the conservative and radical right parties, which
1004 can be expected to have the most conservative views on gender roles. We note that
1005 the authors' interpretation, that "the portion of women MPs affects the attitudes
1006 of left-wing voters and not the attitudes of the voters most likely to undermine
1007 democracy" is true, but that the results also suggest that far-right parties, who
1008 most aggressively challenge liberal democratic norms, may be able to "soften" their
1009 image among left-wing voters by running female candidates. This is consistent
1010 with the argument made by Catalano Week et al (2023), that radical right parties
1011 strategically run women to broaden their appeal. Again, we deeply appreciate your
1012 replication and insightful extension of our research.

1013 **Original Authors' Package:** [https://dataverse.harvard.edu/dataset.xhtml?](https://dataverse.harvard.edu/dataset.xhtml?persistentId=doi:10.7910/DVN/AHQVRV)
1014 [persistentId=doi:10.7910/DVN/AHQVRV](https://dataverse.harvard.edu/dataset.xhtml?persistentId=doi:10.7910/DVN/AHQVRV)

1015 **Reproduction Report**

1016 **Title Original Study:** Changing Hearts and Minds? Why Media Messages
1017 Designed to Foster Empathy Often Fail

1018 **doi:** <https://doi.org/10.1086/719416>, Journal of Politics

1019 **Report's Abstract:** This paper focuses on computational reproducibility and
1020 robustness replicability of Gubler et al.'s(2022) studies which examine the effect of
1021 media messages on empathic concern, dissonance, and out-group policy attitudes.
1022 The original paper tests four hypotheses using two online experiments with large
1023 samples from one US state ($N1 = 5,800$; $N2 = 2,200$). Regarding the first experi-
1024 ment, we successfully reproduced the effect that initial antipathy weakens the effect
1025 of humanizing treatment on empathic concern (H1). However, we show that the
1026 moderating effect is negligible and has little practical significance. Moreover, the
1027 individual effect estimates in our analyses slightly differed from the original paper
1028 due to different procedure of data cleaning and minor coding errors in the original
1029 paper. The most relevant difference was the opposite effect of gender than reported
1030 in the original paper. We also show that empathic concern might mediate the effect
1031 of humanizing treatment on attitudes toward immigrants (H3). The original study
1032 rejected the mediation hypothesis due to not finding a total effect of humanizing
1033 treatment on attitudes. In contrast, we found that humanization treatment has a
1034 positive indirect effect on attitudes through empathic concern. At the same time, it
1035 also has a direct negative effect on attitudes. For the second experiment (H1, H2a,
1036 H2b, H3), we attempted to reproduce the results using a different software. We
1037 partially succeeded once receiving support from the authors of the original study.
1038 We note throughout the report issues we have encountered.

1039 **Link to Full Report:** <https://osf.io/zes6g/>

1040 **Link to Replicators' Package:** See Report's Online Appendix for the codes.

1041 **Original Authors' Response:** <https://osf.io/zes6g/>

1042 **Original Authors' Package:** [https://dataverse.harvard.edu/dataset.xhtml?
1043 persistentId=doi:10.7910/DVN/FUCDTT](https://dataverse.harvard.edu/dataset.xhtml?persistentId=doi:10.7910/DVN/FUCDTT)

1044 **Reproduction Report**

1045 **Title Original Study:** Changing Tides: Public Attitudes on Climate Migration

1046 **doi:** <https://doi.org/10.1086/715163>, Journal of Politics

1047 **Report's Abstract:** See entry below.

1048 **Link to Full Report:** <https://www.socialsciencereproduction.org/reproductions/791/published/index>

1049 **Link to Replicators' Package:** <https://github.com/alexkustov/Replication-of-Arias-and-Blair-2021>

1050 **Original Authors' Response:** "Thank you very much for reaching out! We are
1051 very pleased to hear that the results of our study were replicated, and do not need
1052 to provide an answer."

1053 **Original Authors' Package:** [https://dataverse.harvard.edu/dataset.xhtml;
1054 jsessionid=e19065118cc12d43f1b412109d41?persistentId=doi%3A10.7910%
1055 2FDVN%2FFDML2N&version=&q=&fileAccess=&fileTag=&fileSortField=
1056 name&fileSortOrder=desc](https://dataverse.harvard.edu/dataset.xhtml;jsessionid=e19065118cc12d43f1b412109d41?persistentId=doi%3A10.7910%2FDVN%2FFDML2N&version=&q=&fileAccess=&fileTag=&fileSortField=name&fileSortOrder=desc)
1057
1058

1059 **Reproduction Report**

1060 **Title Original Study:** Checking and Sharing Alt-Facts

1061 **doi:** <https://doi.org/10.1257/pol.20210037>, American Economic Journal: Economic
1062 Policy

1063 **Report's Abstract:** Henry, Zhuravskaya, and Guriev (2022) examine whether
1064 people are willing to share "alternative facts" espoused by right-wing populist par-
1065 ties before the 2019 European elections in France and how this interacted with
1066 the availability of fact-checking information. They find that both imposed and
1067 voluntary fact-checking reduce the likelihood of sharing false statements by approx-
1068 imately 45%, and that imposed and voluntary fact-checking have similar effect sizes.
1069 We reproduce these findings and introduce several alternative estimates to assess
1070 the robustness of the original results, including resolving an inconsistency in the
1071 handling of pre-treatment controls. Overall, our results align with the results of the
1072 original paper. The differences we find are small in absolute magnitude but, since
1073 many effects were small, not always trivial in terms of relative differences. This
1074 replication supports the conclusions of the original paper.

1075 **Link to Full Report:** <https://econpapers.repec.org/paper/zbwi4rdps/34.htm>

1076 **Link to Replicators' Package:** <https://doi.org/10.5281/zenodo.7858829>

1077 **Link to Original Authors' Response:** "Many thanks! No, we won't be writing
1078 a response."

1079 **Original Authors' Package:** [https://www.openicpsr.org/openicpsr/project/
1080 140161/version/V1/view](https://www.openicpsr.org/openicpsr/project/140161/version/V1/view)

1081 **Reproduction Report**

1082 **Title Original Study:** Child Marriage Bans and Female Schooling and Labor
1083 Market Outcomes: Evidence from Natural Experiments in 17 Low- and Middle-
1084 Income Countries

1085 **doi:** <https://doi.org/10.1257/pol.20200008>, American Economic Journal: Economic
1086 Policy

1087 **Report's Abstract:** By studying child marriage bans in 17 developing countries,
1088 Wilson (2022) finds that raising the minimum legal age of marriage to 18 success-
1089 fully increased the age at first marriage, the age at first birth, and the likelihood of
1090 employment. Additionally, the bans reduced child marriage and increased educa-
1091 tional attainment in urban areas. We replicate these findings by collecting the raw
1092 data from the same sources as the paper and analysing the data following the pro-
1093 cedures described in the paper, without referring to the data and codes provided
1094 by the author. Our findings are consistent with the results of the paper in terms of
1095 the statistical significance of point estimates and differ in magnitude by a negligible
1096 amount.

1097 **Link to Full Report:** <https://osf.io/5yhxc/>

1098 **Link to Replicators' Package:** <https://osf.io/5yhxc/>

1099 **Original Authors' Response:** We could not reach out the author.

1100 **Original Authors' Package:** [https://www.openicpsr.org/openicpsr/project/
1101 130784/version/V1/view](https://www.openicpsr.org/openicpsr/project/130784/version/V1/view)

1102 **Reproduction Report**

1103 **Title Original Study:** Concentration Bias in Intertemporal Choice

1104 **doi:** <https://doi.org/10.1093/restud/rdab043>, Review of Economic Studies

1105 **Report's Abstract:** Dertwinkel-Kalt et al. (2022) examine the effect of concen-
1106 tration bias - the tendency to overweight advantages that are concentrated in time
1107 relative to costs that are spread over multiple time periods - on intertemporal choice
1108 in a laboratory experiment. In their preferred empirical specification, the authors
1109 report that concentration bias leads to a 22.4% higher willingness to work than
1110 explained by a standard model of intertemporal discounting. We conduct a compu-
1111 tational replication of the main results of the paper using the same procedures and
1112 original data. Our results confirm the sign, magnitude and statistical significance
1113 of the author's reported estimates across each of their five main findings.

1114 **Link to Full Report:** <https://econpapers.repec.org/paper/zbwi4rdps/42.htm>

1115 **Link to Replicators' Package:** <https://osf.io/d42xr/>

1116 **Original Authors' Response:** "We thank Deer, Ellingsrud, Heuer, and Kordt
1117 (2023) for conducting the Reproduction Report and appreciate that their "results
1118 confirm the sign, magnitude and statistical significance of [our] reported estimates
1119 across each of [our] five main findings" (p. 1). We don't have anything substantive
1120 to add to this. "

1121 **Original Authors' Package:** <https://zenodo.org/records/5091975>

1122 **Reproduction Report**

1123 **Title Original Study:** Cooperative Property Rights and Development: Evidence
1124 from Land Reform in El Salvador

1125 **doi:** <https://doi.org/10.1086/717042>, Journal of Political Economy

1126 **Report's Abstract:** Montero (2022) explores a discontinuity in a land reform in
1127 El Salvador and reports two main findings. First, relative to outside-owned hacien-
1128 das operated by contract workers, the productivity of worker-owned cooperatives is
1129 higher for staple crops and lower for cash-crop. Second, cooperative property rights
1130 increase workers' incomes and compress wage distributions. In this comment, we
1131 show that the latter result rests on two mistakes: three-quarters of the observations
1132 are duplicates and income inequality is calculated over too few workers to be mean-
1133 ingful. When corrected, the data sources and research design provide no credible
1134 evidence regarding the causal effects of ownership structure on income levels and
1135 inequality.

1136 **Link to Full Report:** <https://econpapers.repec.org/paper/zbwi4rdps/20.htm>

1137 **Link to Replicators' Package:** <https://doi.org/10.7910/DVN/AMD3NO>

1138 **Link to Original Authors' Response:** [https://www.journals.uchicago.edu/doi/](https://www.journals.uchicago.edu/doi/10.1086/725234)
1139 [10.1086/725234](https://www.journals.uchicago.edu/doi/10.1086/725234)

1140 **Original Authors' Package:** [https://www.journals.uchicago.edu/doi/suppl/10.](https://www.journals.uchicago.edu/doi/suppl/10.1086/717042/suppl_file/20190161data.zip)
1141 [1086/717042/suppl_file/20190161data.zip](https://www.journals.uchicago.edu/doi/suppl/10.1086/717042/suppl_file/20190161data.zip)

1142 **Reproduction Report**

1143 **Title Original Study:** Decentralization Can Increase Cooperation among Public
1144 Officials

1145 **doi:** <https://doi.org/10.1111/ajps.12606>, American Journal of Political Science

1146 **Report's Abstract:** Molina-Garzón, Grillos, Zarychta, and Andersson (2022)
1147 examine how health sector decentralization affects cooperation between public offi-
1148 cials. Using a public goods game conducted in Honduras, they find that officials
1149 who work under decentralized regimes contributed 0.8 more lempiras per round to
1150 a group solidarity fund, compared to officials who work under centralized regimes.
1151 They also find that most of this increase in investment under decentralized regimes
1152 occurred during rounds of the game in which the participants were able to commu-
1153 nicate with each other. Finally, they find that decentralization was associated with
1154 a 14 percentage point increase in the proportion of potential cross-level network
1155 ties between participants that were realized. In this paper, I examine whether these
1156 results are robust to (1) the omission of some individual-level controls that may
1157 have been affected by the decentralization treatment, and (2) the use of a linear
1158 regression model instead of a Poisson regression model for the network analysis. I
1159 find that omitting the individual-level controls leads to similar conclusions about
1160 the effect of decentralization on individual contributions in the public goods game,
1161 but the interaction effect between decentralization and communication becomes
1162 statistically insignificant at the 0.05 level. For the network analysis, I find that using
1163 a linear regression instead of a Poisson regression has little bearing on the magni-
1164 tude of the effect of decentralization on the proportion of ties realized, though the
1165 effect of decentralization becomes statistically insignificant for one version of the
1166 network model.

1167 **Link to Full Report:** <https://osf.io/q3dpt/>

1168 **Link to Replicators' Package:** <https://osf.io/q3dpt/>

1169 **Link to Original Authors' Response:** <https://osf.io/q3dpt/>

1170 **Original Authors' Package:** [https://dataverse.harvard.edu/dataset.xhtml?](https://dataverse.harvard.edu/dataset.xhtml?persistentId=doi:10.7910/DVN/ZLHYSZ)
1171 [persistentId=doi:10.7910/DVN/ZLHYSZ](https://dataverse.harvard.edu/dataset.xhtml?persistentId=doi:10.7910/DVN/ZLHYSZ)

1172 **Reproduction Report**

1173 **Title Original Study:** Declining Worker Turnover: The Role of Short-Duration
1174 Employment Spells

1175 **doi:** <https://doi.org/10.1257/mac.20190230>, American Economic Journal: Macroeconomics
1176

1177 **Report's Abstract:** Using a Diamond-Mortensen-Pissarides (DMP) model with
1178 noisy signals on worker-firm match quality calibrated on data from 30 US states
1179 for 1999 and 2017, Pries and Rogerson argue that improved screening may explain
1180 the decrease in short-term employment spells observed in the US labor market.
1181 Using a decomposition exercise in a "reduced form" model, the authors show that
1182 changes in short-term employment spells (and) are almost entirely accounted for
1183 by changes in the rate of learning on match quality and in the probability of a good
1184 match . Then, using a decomposition exercise in a "structural" model, they show in
1185 their main calibration strategy that changes in and are mainly driven by changes
1186 in and , parameters pertaining to learning about match quality. First, we reproduce
1187 the authors' codes in R and Python, two popular free open source programming
1188 languages. We find identical results to the paper. Second, we test the robustness
1189 of results to (1) using an earlier starting year, (2) adding additional states in the
1190 analysis, and (3) increasing the value of the 1999 mean vacancy duration parameter.
1191 The direction and relative size of the effect of each parameter on and is preserved
1192 in all robustness tests, corroborating the authors' argument.

1193 **Link to Full Report:** <https://econpapers.repec.org/paper/zbwi4rdps/93.htm>

1194 **Link to Replicators' Package:** [https://github.com/AlexandrePavlov/](https://github.com/AlexandrePavlov/PriesRogerson2022Replication)
1195 [PriesRogerson2022Replication](https://github.com/AlexandrePavlov/PriesRogerson2022Replication)

1196 **Original Authors' Response:** Declined to respond.

1197 **Original Authors' Package:** [https://www.openicpsr.org/openicpsr/project/](https://www.openicpsr.org/openicpsr/project/120568/version/V1/view)
1198 [120568/version/V1/view](https://www.openicpsr.org/openicpsr/project/120568/version/V1/view)

1199 **Reproduction Report**

1200 **Title Original Study:** Digital Addiction

1201 **doi:** <https://doi.org/10.1257/aer.20210867>, American Economic Review

1202 **Report's Abstract:** Using an original economic model of digital addiction and a
1203 randomized experiment, Hunt Allcott, Matthew Gentzkow, and Lena Song (2022)
1204 isolate the effect of habit formation and self-control problems on how people use
1205 their smartphones. They find a persistent effect of temporary incentives on reduc-
1206 ing social media usage. With the model-free results, the study shows that (after the
1207 incentive was in effect), participants in the bonus group reduced use by 56, 19 and
1208 12 minutes in periods 3, 4 and 5, respectively, suggesting a persistent effect. But
1209 before the incentive was in effect in period 2, social media use reduced use by 5.1
1210 minutes per day. Participants who used the limit functionality reduced FITSBY use
1211 by over 20 minutes per day, suggesting an impact of self-control problems on social
1212 media use. All these estimates are statistically significant. We perform a direct
1213 replication of the paper. Upon re-calculating the core dependent variable (FITSBY
1214 use by period), we find a small but concerning discrepancy: For a small number
1215 of observations, the aggregated dependent variable does not equal the sum of the
1216 disaggregated categories. Thankfully, this discrepancy does not have a major effect
1217 on the results. Using the provided data, we re-coded the core figures from scratch
1218 and found that we could replicate them all. We also compare the pre-analysis plan
1219 (PAP) with the main study to identify gaps and perform computational repro-
1220 duction/replication of the structural model and model-free analysis. We only find
1221 minor differences between the PAP and the main paper, almost all of which are
1222 acknowledged in the paper.

1223 **Link to Full Report:** <https://osf.io/8kvdf/>

1224 **Link to Replicators' Package:** <https://osf.io/8kvdf/>

1225 **Link to Original Authors' Response:** <https://osf.io/8kvdf/>

1226 **Original Authors' Package:** [https://www.openicpsr.org/openicpsr/project/
1227 163822/version/V2/view](https://www.openicpsr.org/openicpsr/project/163822/version/V2/view)

1228 **Reproduction Report**

1229 **Title Original Study:** Do Thank-You Calls Increase Charitable Giving? Expert
1230 Forecasts and Field Experimental Evidence

1231 **doi:** <https://doi.org/10.1257/app.20210068>, American Economic Journal: Applied
1232 Economics

1233 **Report's Abstract:** Samek and Longfield estimate the effect of 'thank you calls'
1234 on the extensive and intensive margins of subsequent donations. Based on a series
1235 of experimental interventions, the authors find no statistically discernable effect
1236 of thank-you calls on either the likelihood of donating again, or on the size of
1237 any subsequent donations made within the period of the study. In a companion
1238 exercise the researchers quantify the ability of experts in charitable fundraising and
1239 non-experts (using the Understanding America Survey) to predict the behaviours
1240 elicited by the experiment. Experts and non-experts (incorrectly) make the same
1241 predictions of an increase to the extensive margin of donation behaviour induced
1242 by the thank you call, and while both groups overestimate the intensive margin,
1243 the non-experts overestimated by a smaller magnitude. We were able to reproduce
1244 the papers findings completely, discovering only one difference in an appendix table
1245 related to the average gift amount — treatment for experiment 1 where only the
1246 constant term of the regression was affected. Upon careful examination of the code
1247 we found a few small errors that did not affect the results (one of the errors in the
1248 code did not seem to be carried through and used anywhere). Finally, we conducted
1249 several extensions of the original analysis which demonstrated that the findings
1250 are robust to heterogeneity of treatment effect by initial donation size, as well as
1251 different specifications of the regression analysis.

1252 **Link to Full Report:** <https://osf.io/fe2tr/>

1253 **Link to Replicators' Package:** https://gitlab.com/c3754/replication-games/-/tree/main/replication%20games%20MTL%202023%20charity?ref_type=heads

1254 **Link to Original Authors' Response:** Waiting for the authors' response.

1255 **Original Authors' Package:** <https://www.openicpsr.org/openicpsr/project/149481/version/V1/view>
1256
1257

1258 **Reproduction Report**

1259 **Title Original Study:** Do Transitional Justice Museums Persuade Visitors?
1260 Evidence from a Field Experiment

1261 **doi:** <https://doi.org/10.1086/714765>, Journal of Politics

1262 **Report's Abstract:** Balcells et al. (2022) explore the effect of transitional justice
1263 museums through a field experiment in Santiago, Chile, and attendance at the
1264 government's remembrance museum, the Museum of Memory and Human Rights
1265 which looks at the time of Pinochet's dictatorship. The authors want to understand
1266 how such experiences shape an individual's perceptions of trust in government
1267 institutions, and transitional justice policies, and how they are affected emotionally.
1268 Additionally, they seek to measure how long they last over time. They do this by
1269 creating treatment (museum attendance) and control (non-attendance) groups and
1270 administering pre-and post-treatment surveys and estimating the 'complier average
1271 causal effect' (CACE). They find that satisfaction with the current government
1272 significantly increases for the treatment group, looking over the entire population
1273 ($= 0.15, p = .04$) as measured with a 4-point Likert scale and support for a military
1274 government significantly drops by 11% ($= 0.11, p = .002$) across ideological stances.
1275 We first reproduce their results and find no major coding errors. Second, we test
1276 the robustness of the effects by 1) testing for heterogeneous effects by gender, 2) we
1277 combine the emotion variables into two indices, a mobilization and demobilization
1278 index, and 3) conduct a causal mediation analysis to see how confidence in the
1279 church may mediate effects found in the study.

1280 **Link to Full Report:** <https://osf.io/m3hwg/>

1281 **Link to Replicators' Package:** <https://osf.io/m3hwg/>

1282 **Original Authors' Response:** "We thank all involved for their interest in our
1283 work. We are happy to hear that the results from our paper successfully replicated.
1284 We are intrigued by the additional analyses performed by the replicators. We hope
1285 their insights and results can inform future theorizing and empirical studies of the
1286 impact of Transitional Justice."

1287 **Original Authors' Package:** [https://dataverse.harvard.edu/dataset.xhtml;
1288 jsessionid=a651e0dbc9a5c140152aa84be2e0?persistentId=doi%3A10.7910%
1289 2FDVN%2FTNFDDX&version=&q=&fileTypeGroupFacet=&fileAccess=
1290 Public&fileSortField=type](https://dataverse.harvard.edu/dataset.xhtml;jsessionid=a651e0dbc9a5c140152aa84be2e0?persistentId=doi%3A10.7910%2FDVN%2FTNFDDX&version=&q=&fileTypeGroupFacet=&fileAccess=Public&fileSortField=type)

1291 **Reproduction Report**

1292 **Title Original Study:** Does Competence Make Citizens Tolerate Undemocratic
1293 Behavior?

1294 **doi:** <https://doi.org/10.1017/S0003055422000119>, American Political Science
1295 Review

1296 **Report's Abstract:** We replicate the analysis conducted by Frederiksen, 2022a.
1297 We focus on assessing the computational and robustness replicability of their work.
1298 We find that their main exhibits and supplementary analysis are replicable, both
1299 when running their original Stata replication package, and when we attempt to
1300 replicate their findings from scratch in R. We also conduct additional robustness
1301 checks by estimating additional specifications and by subsetting the dataset by the
1302 time taken by the respondent to complete the survey. We again find that their work
1303 is robust to our battery of alternative specifications.

1304 **Link to Full Report:** <https://econpapers.repec.org/paper/zbwi4rdps/28.htm>

1305 **Link to Replicators' Package:** [https://github.com/tjbrailey/nottingham_](https://github.com/tjbrailey/nottingham_replication_2023)
1306 [replication_2023](https://github.com/tjbrailey/nottingham_replication_2023)

1307 **Link to Original Authors' Final Response:** "Thanks a lot for this initiative
1308 and not least for replicating my results."

1309 **Original Authors' Package:** [https://dataverse.harvard.edu/dataset.xhtml?](https://dataverse.harvard.edu/dataset.xhtml?persistentId=doi:10.7910/DVN/NGFLRO)
1310 [persistentId=doi:10.7910/DVN/NGFLRO](https://dataverse.harvard.edu/dataset.xhtml?persistentId=doi:10.7910/DVN/NGFLRO)

1311 **Reproduction Report**

1312 **Title Original Study:** Does Patient Demand Contribute to the Overuse of
1313 Prescription Drugs?

1314 **doi:** <https://doi.org/10.1257/app.20190722>, American Economic Journal: Applied
1315 Economics

1316 **Report's Abstract:** We replicate Lopez et al.'s (2022) study on gatekeeping costs
1317 and the potential evidence for patient-driven and doctor-driven demand. Using
1318 their publicly available source materials, we first re-run their analysis "as is" to see
1319 if their results can be exactly replicated. We then expand the analysis to include
1320 patients previously excluded for not being acutely ill, offering a broader perspective
1321 on medication demand among all patient types. The findings confirm Lopez et al.'s
1322 results.

1323 **Link to Full Report:** <https://osf.io/x7g9z/>

1324 **Link to Replicators' Package:** <https://osf.io/x7g9z/>

1325 **Link to Original Authors' Response:** Provided feedback to an initial report.
1326 Final response: "Thank you very much for sharing the updated report. We appreciate
1327 that the authors of the replication reworked the paper and have no further
1328 response or comments."

1329 **Original Authors' Package:** [https://www.openicpsr.org/openicpsr/project/
1330 126722/version/V1/view](https://www.openicpsr.org/openicpsr/project/126722/version/V1/view)

1331 **Reproduction Report**

1332 **Title Original Study:** Does Public Opinion Affect the Preferences of Foreign
1333 Policy Leaders? Experimental Evidence from the UK Parliament

1334 **doi:** <https://doi.org/10.1086/719007>, Journal of Politics

1335 **Report’s Abstract:** The study by Chu and Recchia (2022) tests the hypothesis
1336 that providing public opinion information can shift policymakers’ opinions in the
1337 direction of what the public favors. They surveyed 101 British Members of Par-
1338 liament (MPs) about their views regarding the United Kingdom’s presence in the
1339 South China Sea. Their results demonstrated that MPs who received information
1340 about the public opinion poll expressed viewpoints closer to that of public opinion.
1341 The authors reported an effect that is “substantively meaningful and statistically
1342 significant at the .10 level.” Our computational replication of the original study
1343 found that the paper is fully computationally reproducible. We successfully repli-
1344 cated the authors’ results but found that the main findings are no longer significant
1345 when analyzed using unweighted data (see Table 1). We also conducted several
1346 robustness checks on sub-samples of the data to examine the key analyses both
1347 with and without weights. Here, we found that the results are once again robust
1348 and significant when weights are used, but no longer significant when weights are
1349 not used. As a further robustness check, we found no moderating effect of gender.
1350 Overall, our replication efforts suggest that the main finding of the original study
1351 may be sensitive to the use of survey weights.

1352 **Link to Full Report:** <https://osf.io/bqz6w/>

1353 **Link to Replicators’ Package:** [https://osf.io/vwt2n/?view_only=](https://osf.io/vwt2n/?view_only=84e52a7c684942a4880410b3c89ff4c6)
1354 [84e52a7c684942a4880410b3c89ff4c6](https://osf.io/vwt2n/?view_only=84e52a7c684942a4880410b3c89ff4c6)

1355 **Original Authors’ Response:** “ Thank you for your note and engaging with our
1356 work. We don’t have a formal reply, though this is an honest question: isn’t it stan-
1357 dard practice to use weights when using YouGov’s data, since making valid claims
1358 about representativeness depends on using their weights? YouGov’s MP panels
1359 operate similarly to their public opinion poll, in that their claims to representa-
1360 tiveness rely on using weights, provided by YouGov. I [Chu] think your write-up
1361 mentioned that there’s a debate about using weights, and cited MTurk data, but I
1362 think that MTurk is quite different, and yes, I agree I do not use weights for MTurk
1363 data except unless requested by a reviewer for robustness checks, etc.. But I don’t
1364 think MTurk and the MP representative poll we used is a good comparison in the
1365 context of evaluating the validity of weighting. In any case, happy to adapt if there
1366 is a clear consensus on this. Thanks again.”

1367 **Original Authors’ Package:** [https://dataverse.harvard.edu/dataset.xhtml?](https://dataverse.harvard.edu/dataset.xhtml?persistentId=doi:10.7910/DVN/BNINNL)
1368 [persistentId=doi:10.7910/DVN/BNINNL](https://dataverse.harvard.edu/dataset.xhtml?persistentId=doi:10.7910/DVN/BNINNL)

1369 **Reproduction Report**

1370 **Title Original Study:** Effective for Whom? Ethnic Identity and Nonviolent
1371 Resistance

1372 **doi:** <https://doi.org/10.1017/S0003055421000940>, American Political Science
1373 Review

1374 **Report's Abstract:** Manekin and Mitts (2022) investigate the success chances of
1375 minority ethnic groups when engaging in non-violent protests demanding political
1376 change. First, using observational data, the authors find that the success rate for
1377 nonviolent campaign tactics is lower for excluded/minority ethnic groups than for
1378 non-excluded/majority ethnic groups. Second, the authors use two original survey
1379 experiments to show that non-violent protest by ethnic minorities is perceived as
1380 more violent and requiring more policing than identical protest by majorities. This
1381 report reproduces the paper computationally and conducts several sensitivity anal-
1382 yses for both the observational and the experimental parts of the paper. We can
1383 confirm the general direction of the postulated effects, but evidence becomes less
1384 consistent (effect magnitudes and significance levels are not robust to some of the
1385 changes).

1386 **Link to Full Report:** <https://econpapers.repec.org/paper/zbwi4rdps/86.htm>

1387 **Link to Replicators' Package:** <https://zenodo.org/records/10193470>

1388 **Original Authors' Response:** Cannot provide a response.

1389 **Original Authors' Package:** [https://dataverse.harvard.edu/dataset.xhtml?](https://dataverse.harvard.edu/dataset.xhtml?persistentId=doi:10.7910/DVN/SHHVCA)
1390 [persistentId=doi:10.7910/DVN/SHHVCA](https://dataverse.harvard.edu/dataset.xhtml?persistentId=doi:10.7910/DVN/SHHVCA)

1391 **Reproduction Report**

1392 **Title Original Study:** Enabling or Limiting Cognitive Flexibility? Evidence of
1393 Demand for Moral Commitment

1394 **doi:** <https://doi.org/10.1257/aer.20201333>, American Economic Review

1395 **Report’s Abstract:** We computationally reproduce Saccardo and Serra-Garcia
1396 (2023) where subjects exploit cognitive flexibility by viewing their incentives first
1397 and partially ignoring product quality information, and hence, recommend the
1398 incentivized product. We find one major coding error for the variable Selfishness.
1399 Additionally, two of the “moral cost” questions more likely capture spitefulness.
1400 After correcting the erroneous coding or dropping the two questions, we find
1401 stronger support for the authors’ main conclusion regarding Selfishness driving
1402 incentive information avoidance with double effect size. Finally, we find weak evi-
1403 dence that subjects update their posterior beliefs differently depending on the
1404 product they are incentivized to recommend.

1405 **Link to Full Report:** <https://osf.io/nwds7>

1406 **Link to Replicators’ Package:** <https://osf.io/yfdet/>

1407 **Link to Original Authors’ Response:** [https://www.aeaweb.org/doi/10.1257/](https://www.aeaweb.org/doi/10.1257/aer.20201333.appx)
1408 [aer.20201333.appx](https://www.aeaweb.org/doi/10.1257/aer.20201333.appx)

1409 **Original Authors’ Package:** [https://www.openicpsr.org/openicpsr/project/](https://www.openicpsr.org/openicpsr/project/180741/version/V1/view)
1410 [180741/version/V1/view](https://www.openicpsr.org/openicpsr/project/180741/version/V1/view)

1411 **Reproduction Report**

1412 **Title Original Study:** Entertaining Beliefs in Economic Mobility
1413 **doi:** <https://doi.org/10.1111/ajps.12702>, American Journal of Political Science
1414 **Report’s Abstract:** In Entertaining Beliefs in Economic Mobility (AJPS 2023)
1415 Kim finds that watching “rags-to-riches” style reality TV programs strengthens
1416 Americans’ belief in the American dream. Through thoughtful and clever exper-
1417 imental and observational analysis, she demonstrates that exposure to television
1418 programs containing everyday people working hard to earn large prizes increases
1419 Americans’ belief that success can be internally attributed and that economic mobil-
1420 ity is possible. We computationally replicate Kim’s results, finding no major errors
1421 in her coding or statistical procedure. We also include several robustness checks.
1422 First, we merge her two experimental samples, which increases the precision of her
1423 main quantity of interest such that it attains conventional levels of statistical signif-
1424 icance. Second, we recreate tables and visualizations for alternative specifications
1425 of her main observational results. The original results are robust to these alterna-
1426 tive models, but we do find that if sports programming is operationalized in the
1427 same manner as “rags-to-riches” programming, the sign, magnitude, and signifi-
1428 cance of watching either programming type are similar. We also uncover a partisan
1429 interaction effect, as only Democrats change their beliefs in economic mobility with
1430 increased TV viewing.
1431 **Link to Full Report:** <https://osf.io/xf5w2/>
1432 **Link to Replicators’ Package:** [https://github.com/jacobawinter/rep_games_](https://github.com/jacobawinter/rep_games_2023)
1433 [2023](https://github.com/jacobawinter/rep_games_2023)
1434 **Original Author’s Response:** “Thanks for this! I have no particular response
1435 per se. I’m grateful for your collective efforts to make social science much more
1436 transparent.”
1437 **Original Authors’ Package:** [https://dataverse.harvard.edu/dataset.xhtml?](https://dataverse.harvard.edu/dataset.xhtml?persistentId=doi:10.7910/DVN/FVRZYU)
1438 [persistentId=doi:10.7910/DVN/FVRZYU](https://dataverse.harvard.edu/dataset.xhtml?persistentId=doi:10.7910/DVN/FVRZYU)

1439 **Reproduction Report**

1440 **Title Original Study:** Evaluating Deliberative Competence: A Simple Method
1441 with an Application to Financial Choice

1442 **doi:** <https://doi.org/10.1257/aer.20210290>, American Economic Review

1443 **Report's Abstract:** Ambuehl et al. (2022) explore ways to evaluate interven-
1444 tions designed to enhance decision-making quality when individuals misjudge the
1445 outcomes of their choices. The authors propose a novel outcome metric that can
1446 distinguish between interventions better than conventional metrics such as financial
1447 literacy and directional behavioral responses. The proposed metric, which trans-
1448 forms price-metric bias into interpretable welfare loss measures, can be applied
1449 to evaluate various training programs on financial products. Table 4 of the paper
1450 reports the authors' significant main point estimates at the 1% level. In this repli-
1451 cation exercise, we first replicate the main findings of the original paper. Then,
1452 we modify the clustering method by using k-means with demographic variables as
1453 inputs, then we re-calculate standard errors with jackknife estimators. Finally, we
1454 include subjects who were excluded by the authors due to multiple switching in
1455 the multiple price lists. We find that all of these replications result in robust find-
1456 ings. Additionally, we successfully replicate Figure 4 from the paper. Notably, this
1457 replication demonstrates the insensitivity of the results to the choice of distance
1458 metric.

1459 **Link to Full Report:** <https://osf.io/scgbt/>

1460 **Link to Replicators' Package:** <https://osf.io/scgbt/>

1461 **Link to Original Authors' Response:** Authors provided feedback.

1462 **Original Authors' Package:** [https://www.openicpsr.org/openicpsr/project/
1463 171681/version/V1/view](https://www.openicpsr.org/openicpsr/project/171681/version/V1/view)

1464 **Reproduction Report**

1465 **Title Original Study:** Exposure and Preferences: Evidence from Indian Slums

1466 **doi:** <https://doi.org/10.1111/ajps.12570>, American Journal of Political Science

1467 **Report's Abstract:** Successful computational reproducibility. The replicators
1468 could not conduct the robustness checks without the help of the author.

1469 **Link to Full Report:** No report.

1470 **Original Author's Response:** "Thanks for your email. I am not interested in
1471 participating."

1472 **Original Authors' Package:** [https://dataverse.harvard.edu/dataset.xhtml?](https://dataverse.harvard.edu/dataset.xhtml?persistentId=doi:10.7910/DVN/AV8PLT)
1473 [persistentId=doi:10.7910/DVN/AV8PLT](https://dataverse.harvard.edu/dataset.xhtml?persistentId=doi:10.7910/DVN/AV8PLT)

1474 **Reproduction Report**

1475 **Title Original Study:** Finance and Green Growth

1476 **doi:** <https://doi.org/10.1093/ej/ueac081>, Economic Journal

1477 **Report's Abstract:** De Haas and Popov (2023) estimate the effect of country-level
1478 financial sector size and structure on decarbonization to show that countries with
1479 relatively more equity versus debt financing have more emission-efficient economies.

1480 We uncover multiple coding errors that change the magnitude and the precision
1481 of the coefficients of interest. These coding errors include misreporting of standard
1482 errors, and misspecifying generalized method of moments (GMM) estimators. We
1483 further provide robustness tests of the results to (1) restricting the sample to con-
1484 sistent sets of countries across the country and country-by industry samples, and
1485 (2) using a limited information maximum likelihood (LIML) estimator to address a
1486 weak-instrument problem. We find that the results from the robustness checks are
1487 qualitatively different from the original results but similar to the corrected results.

1488 **Link to Full Report:** <https://econpapers.repec.org/paper/zbwi4rdps/95.htm>

1489 **Link to Replicators' Package:** <https://osf.io/h8ct2/>

1490 **Link to Original Authors' Response:** <https://econpapers.repec.org/paper/zbwi4rdps/96.htm>

1492 **Original Authors' Package:** <https://zenodo.org/records/7220094>

1493 **Reproduction Report**

1494 **Title Original Study:** Flight to Safety: COVID-Induced Changes in the Intensity
1495 of Status Quo Preference and Voting Behavior

1496 **doi:** <https://doi.org/10.1017/S0003055421000691>, American Political Science
1497 Review

1498 **Report’s Abstract:** Bisbee and Honig (2022) examine the effect of the COVID-
1499 19 pandemic on voting for Bernie Sanders in the 2020 Democratic Party primary
1500 using a difference-in-differences design, finding evidence that exposure to COVID-
1501 19 resulted in a 7-15 percentage point increase in voting for Biden. The study also
1502 uses a regression design with district-level fixed effects to estimate the effect of
1503 the COVID-19 pandemic on voting for anti-establishment candidates during the
1504 US 2020 House primaries. It finds evidence that an increase in COVID cases was
1505 associated with a decline in voting for anti-establishment candidates in general,
1506 and for those endorsed by the Tea Party. We re-run the code for all tests in this
1507 paper, successfully reproducing its results in a preliminary replication. We then
1508 use the De Chaisemartin and D’Haultfoeuille difference-in-differences estimator to
1509 replicate their main results, finding that though the coefficient remains negative,
1510 the results are not statistically significant. We also replicate their tests regard-
1511 ing US House primary candidates using a different measure of anti-establishment
1512 candidates. Here, we find that the interaction term between anti-establishment can-
1513 didates and COVID-19 remain statistically significant, with the same sign. Finally,
1514 we employ an expanded dataset that includes Congressional primary candidates
1515 that were omitted in the initial dataset, as well as a re-coded extremism variable
1516 that also includes candidates endorsed by Donald Trump. These updated find-
1517 ings corroborate the paper’s initial results. However, due to a restrictive number
1518 of observations that interfered with our application of the De Chaisemartin and
1519 D’Haultfoeuille estimator, we believe that the expanded U.S. House primary results
1520 constitute the more robust half of our replication.

1521 **Link to Full Report:** <https://econpapers.repec.org/paper/zbwi4rdps/36.htm>

1522 **Link to Replicators’ Package:** [https://github.com/Dmscates/
1523 Bisbee-and-Honig-2022-Flight-to-Safety-Replication](https://github.com/Dmscates/Bisbee-and-Honig-2022-Flight-to-Safety-Replication)

1524 **Original Authors’ Response:** “You guys are amazing. Thank you for doing this!
1525 We are impressed by your rigor and grateful for the introduction to DCD’H DiD
1526 estimator that we’ll have to add to the repertoire. We were working on the condi-
1527 tional accept when the flurry of generalized DiD work (Goodman-Bacon, Callaway,
1528 etc.) was blowing up [...] We also appreciate the manner in which you communi-
1529 cated with us during the course of your re-analysis, and the thoughtfulness of your
1530 report. [...] Although I’m sure it is a logistical nightmare and likely would add even
1531 more delays to the publication pipeline, it would be very pro-science if this type
1532 of replication were part of a journal’s own pre-publication process. (This is what I
1533 naively thought replication meant back when I got my first publication, and have
1534 been disappointed in the process ever since.)”

1535 **Original Authors’ Package:** [https://dataverse.harvard.edu/dataset.xhtml?
1536 persistentId=doi:10.7910/DVN/S5YMS7](https://dataverse.harvard.edu/dataset.xhtml?persistentId=doi:10.7910/DVN/S5YMS7)

1537 **Reproduction Report**

1538 **Title Original Study:** Gender Differences in Cooperative Environments? Evi-
1539 dence from The U.S. Congress

1540 **doi:** <https://doi.org/10.1093/ej/ueab069>, Economic Journal

1541 **Report's Abstract:** Gagliarducci and Paserman (2022) study gender differences
1542 in cooperative behavior among politicians using information from the U.S. House
1543 of Representatives between 1988 and 2010 on (i) the number of co-sponsors on bills
1544 and (ii) the share of co-sponsors from the rival party. Through different empirical
1545 strategies, they show that women-sponsored bills tend to have more co-sponsors,
1546 but the gap is only statistically significant among Republicans. Moreover, Repub-
1547 lican women recruit a significantly larger share of co-sponsors from the rival party
1548 than Republican men, whereas the opposite is true among Democrats. GP argue
1549 that the observed pattern is consistent with a commonality of interest driving coop-
1550 eration, rather than gender per se, since during this period Republican women
1551 were ideologically closer to the rival party than their male colleagues, while female
1552 Democrats were further away. We examine the robustness of these findings to (i)
1553 the correction of some errors in two control variables of the dataset used by GP and
1554 (ii) clustering the standard errors at the individual level, instead of individual-term.
1555 These changes have a relatively minor impact on results: most coefficients are still
1556 statistically significant and the main conclusions from the analysis are confirmed.
1557 Furthermore, we extend the analysis to the 2011-2020 period. The analysis of gen-
1558 der differences in bipartisan cooperation confirms GP's hypothesis that ideological
1559 distance plays an important role. However, results are slightly different when we
1560 analyze overall cooperation. The gender gap in favor of women is larger in magni-
1561 tude than in GP and it is statistically significant in several specifications, providing
1562 support for the hypothesis that gender also matters for cooperation.

1563 **Link to Full Report:** <https://econpapers.repec.org/paper/zbwi4rdps/75.htm>

1564 **Link to Replicators' Package:** [https://www.dropbox.com/scl/fo/
1565 dmvgx5wlgql3tz98dx47u/h?rlkey=af83xiqrkw70rbjcdunoach9&dl=0](https://www.dropbox.com/scl/fo/dmvgx5wlgql3tz98dx47u/h?rlkey=af83xiqrkw70rbjcdunoach9&dl=0)

1566 **Link to Original Authors' Response:** <https://osf.io/w48tf/>

1567 **Original Authors' Package:** <https://zenodo.org/records/5111360>

1568 **Reproduction Report**

1569 **Title Original Study:** Good Reverberations? Teacher Influence in Music Com-
1570 position since 1450

1571 **doi:** <https://doi.org/10.1086/718370>, Journal of Political Economy

1572 **Report's Abstract:** Borowiecki (2022) studies the influence of teachers on the
1573 style of their students in the domain of musical composition. The author finds that
1574 realized student-teacher pairs are on average 0.2-0.3 standard deviations more sim-
1575 ilar to unrealized, but possible, student-teacher pairs. In this report we provide the
1576 results of our replication of Borowiecki (2022). We direct our attention to the fol-
1577 lowing tasks: 1) Replicating the outcome variables used in the paper, starting from
1578 the raw data, and generating alternative measures of similarity between students
1579 and teachers 2) Testing the validity of the random teacher-student pairing, a key
1580 assumption for the validity of the estimation strategy employed in the paper. We
1581 can replicate most of the outcome variables, but not all of them, due to incom-
1582 plete raw data. Our alternative measures of similarity confirm the robustness of
1583 the original results. We find significantly different characteristics between paired
1584 and unpaired students, suggesting that matching between students and teachers
1585 does not occur randomly. However, controlling for these characteristics in the main
1586 regressions leads to quantitatively similar results to the ones reported in the original
1587 paper.

1588 **Link to Full Report:** <https://econpapers.repec.org/paper/zbwi4rdps/27.htm>

1589 **Link to Replicators' Package:** <https://www.dropbox.com/scl/fo/6hecmgjsq3mjo5ekkv8lm/h?rlkey=ftuoe4mf5f9jon0hiabb4brtn&dl=0>

1590 **Link to Original Authors' Response:** <https://osf.io/79e2z/>

1591 **Original Authors' Package:** https://www.journals.uchicago.edu/doi/suppl/10.1086/718370/suppl_file/20210405data.zip
1592
1593

1594 **Reproduction Report**

1595 **Title Original Study:** Hate Crimes and Gender Imbalances: Fears over Mate
1596 Competition and Violence against Refugees

1597 **doi:** <https://doi.org/10.1111/ajps.12595>, American Journal of Political Science

1598 **Report's Abstract:** Dancygier et al. (2022) ascribe anti-refugee hate crime in
1599 Germany from 2015 to 2017 to the fear of mate competition felt by native Ger-
1600 man men, amplified by growing refugee populations and existing gender gaps. In
1601 a replication of this article, we discovered that the substantively and statistically
1602 significant relationship between perceptions of mate competition and support for
1603 anti-refugee violence found in a 2016–17 survey of adults in Germany were robust
1604 when analyzed with ensembles of regression trees permitting arbitrary interactions
1605 in a large design matrix. However, statistically significant pairwise comparisons
1606 between survey respondents' perceptions of mate competition across strata of the
1607 municipality-level gender gap as recorded by German censuses were not robust to
1608 controlling the family-wise Type I error rate. Moreover, statistically significant rela-
1609 tionships between the gender gap and the incidence of hate crime in Germany in the
1610 authors' panel regressions vanished in a wide range of specifications with munic-
1611 ipality fixed effects—in certain cases, being replaced with statistically significant
1612 estimates of the opposite sign.

1613 **Link to Full Report (and Initial Version of the Report):** [https://osf.io/](https://osf.io/5n3ds/)
1614 [5n3ds/](https://osf.io/5n3ds/)

1615 **Link to Replicators' Package:** <https://osf.io/5n3ds/>

1616 **Link to Original Authors' Response:** <https://osf.io/5n3ds/>

1617 **Original Authors' Package:** [https://dataverse.harvard.edu/dataset.xhtml?](https://dataverse.harvard.edu/dataset.xhtml?persistentId=doi:10.7910/DVN/QXJDJ5)
1618 [persistentId=doi:10.7910/DVN/QXJDJ5](https://dataverse.harvard.edu/dataset.xhtml?persistentId=doi:10.7910/DVN/QXJDJ5)

1619 **Reproduction Report**

1620 **Title Original Study:** Historical Lynchings and the Contemporary Voting
1621 Behavior of Blacks

1622 **doi:** <https://doi.org/10.1257/app.20190549>, American Economic Journal: Applied
1623 Economics

1624 **Report's Abstract:** Williams (2022) ties the political participation of Blacks to
1625 historical lynchings that occurred in the United States. Her findings document lower
1626 Black voter registration rates in southern counties with greater number of historical
1627 lynchings. We show that this effect is driven by four outlier counties with relatively
1628 high Black lynching rates. Excluding these counties from the analysis yields a point
1629 estimate that is no longer statistically significant. Dropping the ninety-fifth per-
1630 centile lynching rates and correcting the errors in voter registration rates rule out
1631 the effect size reported by Williams (2022), which now becomes close to zero and
1632 statistically insignificant. We also show that the main results are highly sensitive
1633 to the way lynching and voter registration rates are measured.

1634 **Link to Full Report:** <https://econpapers.repec.org/paper/zbwi4rdps/32.htm>

1635 **Link to Replicators' Package:** <https://osf.io/hv7wp/>

1636 **Original Author's Response:** No response to our emails.

1637 **Original Author's Package:** [https://www.openicpsr.org/openicpsr/project/
1638 136741/version/V1/view](https://www.openicpsr.org/openicpsr/project/136741/version/V1/view)

1639 **Reproduction Report**

1640 **Title Original Study:** Hobo Economicus

1641 **doi:** <https://doi.org/10.1093/ej/ueab103>, Economic Journal

1642 **Report's Abstract:** Peter Leeson, August Hardy and Paola Suarez (2022) test
1643 maximizing behaviour of panhandlers at several Metrorail stations in Washington,
1644 D.C. Their main findings are that "stations with more panhandling opportunities
1645 attract more panhandlers" (the first statement) and that "cross-station differences
1646 in hourly panhandling receipts are statistically indistinguishable from zero" (the
1647 second statement). We test computational reproducibility and robustness replica-
1648 bility of their results. We can reproduce both statements, in Stata and R. Our
1649 robustness replications for the first statement confirm the authors' results in the
1650 vast majority of cases (replication was successful in 91% of the cases). Our robust-
1651 ness replications for the second statement might raise doubts on this finding. We
1652 run weighted ANOVA tests, we change the bounds in minutes used by authors by
1653 5 minutes in their robustness checks, we run Bartlett's tests of equality of vari-
1654 ances of means, and run pair-wise tests of equality of means. In three out of four
1655 cases we cannot replicate the results, and the differences (of either means, medi-
1656 ans or variances of donations) across Metrorail stations are statistically different
1657 from zero. We hypothesize that panhandlers have a general idea about which sta-
1658 tions have more passers-by, and will rationally go more often there. However, they
1659 are unlikely to have information about smaller variations in the number of passers-
1660 by (e.g., variations in passers-by at the same station over time due to non-public
1661 events), and therefore might find it difficult to perfectly maximize donations.

1662 **Link to Full Report:** <https://econpapers.repec.org/paper/zbwi4rdps/55.htm>

1663 **Link to Replicators' Package:** <https://osf.io/s4bca/>

1664 **Original Authors' Response:** Declined to respond.

1665 **Original Authors' Package:** <https://zenodo.org/records/5719541>

1666 **Reproduction Report**

1667 **Title Original Study:** How Do Beliefs about the Gender Wage Gap Affect the
1668 Demand for Public Policy?

1669 **doi:** <https://doi.org/10.1257/pol.20200559>, American Economic Journal: Economic
1670 Policy

1671 **Report's Abstract:** We conduct a replication of Settele (2022), a online survey
1672 experiment designed to find out how individual's beliefs about the gender wage
1673 gap affect their policy preferences. We reproduce Results 1 and 2 of the study: how
1674 prior beliefs around the wage gap are distributed among individuals and how a
1675 information treatment causally affects the policy demand. Our re-coded replication
1676 shows that the reported results are robust.

1677 **Link to Full Report:** <https://econpapers.repec.org/paper/zbwi4rdps/12.htm>

1678 **Link to Replicators' Package:** <https://osf.io/j2ubt/>

1679 **Original Authors' Response:** "I very much appreciate the effort of you and your
1680 team to replicate not just my paper, but many others too. It is quite impressive to
1681 see the scope of your project and I am curious about your future plans with this
1682 initiative.

1683 I just read the Reproduction Report for my paper and I think it is great. In
1684 particular, Figures 2 and 3 are really cool. (They are actually new, and offer a really
1685 insightful way of looking at the data. I should have come up with them myself!)

1686 Regarding your question, I don't think the Reproduction Report requires a
1687 formal response from my side. I fully agree with the authors' interpretation of the
1688 results, and just want to say thank you for their great work. Please go ahead and
1689 publish the report whenever it suits you."

1690 **Original Authors' Package:** [https://www.openicpsr.org/openicpsr/project/
1691 134041/version/V1/view](https://www.openicpsr.org/openicpsr/project/134041/version/V1/view)

1692 **Reproduction Report**

1693 **Title Original Study:** How Effective Are Monetary Incentives to Vote? Evidence
1694 from a Nationwide Policy?

1695 **doi:** <https://doi.org/10.1257/app.20200482>, American Economic Journal: Applied
1696 Economics

1697 **Report's Abstract:** Successful computational reproducibility. No re-analyses
1698 conducted.

1699 **Link to Original Authors' Response:** Not contacted.

1700 **Original Authors' Package:** [https://www.journals.uchicago.edu/doi/suppl/10.](https://www.journals.uchicago.edu/doi/suppl/10.1086/720458/suppl_file/20190733data.zip)
1701 [1086/720458/suppl_file/20190733data.zip](https://www.journals.uchicago.edu/doi/suppl/10.1086/720458/suppl_file/20190733data.zip)

1702 **Reproduction Report**

1703 **Title Original Study:** How Much Should We Trust the Dictator’s GDP Growth
1704 Estimates?

1705 **doi:** <https://doi.org/10.1086/720458>, Journal of Political Economy

1706 **Report’s Abstract:** In this brief commentary, we have conducted a robustness
1707 reproducibility and replicability of Martinez’s 2022 paper entitled “How much
1708 should we trust the dictator’s GDP growth estimates?” by selecting different clus-
1709 ters and omitting fixed effect terms. Concurrently, we conduct sub-sample analyses
1710 and employ alternative measurements for the sake of robustness and direct replica-
1711 bility. Our results are generally robust, yet they also raise some intriguing questions.
1712 First, we attempt to remove the year fixed effect in the model specifications, but
1713 the elimination of the year fixed effect from the baseline equation did not account
1714 for unobserved variables across year, suggesting the variable bias by Oster (2019).
1715 Second, the entirety of the baseline results is influenced by the periods 2007-2013
1716 (for a five-year interval) and 2010-2013 (for a three-year interval). Third, when uti-
1717 lizing a more varied dataset for the autocracy measurement, the effect vanished for
1718 countries that are partially unfree.

1719 **Link to Full Report:** <https://osf.io/4sk52/>

1720 **Link to Replicators’ Package:** <https://osf.io/4sk52/>

1721 **Original Author’s Final Response:** “As before, please extend my gratitude to
1722 the replicators for their thoughtful work and for taking into consideration my pre-
1723 vious comments. I am reassured by the fact that they were able to replicate all the
1724 original results in the paper. I also find it reassuring that the results prove robust
1725 to additional robustness tests concerning alternative clustering structures for the
1726 standard errors or alternative data sources on political regimes (albeit with some
1727 loss of precision). The heterogeneous effects by subperiod are also quite intrigu-
1728 ing and potentially reflect changes in the geopolitical incentives to overstate GDP
1729 growth in non-democracies. My paper is certainly not the final word on this topic
1730 and these results could be the first step towards new and exciting research.”

1731 **Original Authors’ Package:** [https://www.journals.uchicago.edu/doi/suppl/10.
1732 1086/720458/suppl_file/20190733data.zip](https://www.journals.uchicago.edu/doi/suppl/10.1086/720458/suppl_file/20190733data.zip)

1733 **Reproduction Report**

1734 **Title Original Study:** Ideological Asymmetries and the Determinants of Politi-
1735 cally Motivated Reasoning

1736 **doi:** <https://doi.org/10.1111/ajps.12624>, American Journal of Political Science

1737 **Report's Abstract:** Guay and Johnston (2022) examine asymmetric politically
1738 motivated reasoning on the part of liberals and conservaites. In our replication
1739 of the paper we examine four potential issues with the analysis: confounding in
1740 the numeracy task, heterogeneity across ideological constraints, the use of control
1741 variables, and heterogeneity in the moderator index items. None of these potential
1742 issues are in fact issues. The results are quite robust. We found only one minor
1743 issue with the codebook, which does not affect the results.

1744 **Link to Full Report:** <https://econpapers.repec.org/paper/zbwi4rdps/79.htm>

1745 **Link to Replicators' Package:** <https://osf.io/mh5sk/>

1746 **Link to Original Authors' Response:** "Thank you again for examining our
1747 paper so closely [...] we changed the codebook and appreciate this replication effort."

1748 **Original Authors' Package:** [https://dataverse.harvard.edu/dataset.xhtml?](https://dataverse.harvard.edu/dataset.xhtml?persistentId=doi:10.7910/DVN/CGHTPZ)
1749 [persistentId=doi:10.7910/DVN/CGHTPZ](https://dataverse.harvard.edu/dataset.xhtml?persistentId=doi:10.7910/DVN/CGHTPZ)

1750 **Reproduction Report**

1751 **Title Original Study:** Immigration and Redistribution

1752 **doi:** <https://doi.org/10.1093/restud/rdac011>, Review of Economic Studies

1753 **Report's Abstract:** Alesina et al. (2023) examine how people perceive the num-
1754 ber and characteristics of migrants and how those perceptions affect their support
1755 for redistribution. They find that respondents from the United States, United
1756 Kingdom, Sweden, Italy, Germany and France markedly overestimate the share of
1757 immigrants in each country, with the average respondent in all countries except
1758 Sweden overestimating by more than a factor of two. We reproduce these results
1759 using the original code and data and test the robustness by (i) including participants
1760 excluded for time to complete the survey, (ii) extending the analysis of mispercep-
1761 tions to all survey respondents, and (iii) using alternative authoritative estimates
1762 of the proportion of immigrants. We find that these checks marginally change the
1763 estimates of the size of the misperception but do not change the conclusions to be
1764 drawn from the analysis. Alesina et al. (2023) also test the effect on support for
1765 redistribution of showing videos on immigrant characteristics. We computationally
1766 reproduced the treatment effects on support for redistribution.

1767 **Link to Full Report:** <https://econpapers.repec.org/paper/zbwi4rdps/40.htm>

1768 **Link to Replicators' Package:** <https://osf.io/ajm9g/>

1769 **Original Authors' Response:** “Dear Institute for Replication team,

1770 Thank you very much for taking the time to replicate our paper. We appreciate
1771 the important work you do. We are happy to see that our results replicated well
1772 and that our robustness checks were confirmed.

1773 With best wishes, Armando Miano and Stefanie Stantcheva”

1774 **Original Authors' Package:** <https://zenodo.org/records/5997521>

1775 **Reproduction Report**

1776 **Title Original Study:** Indecent Disclosures: Anticorruption Reforms and Political
1777 Selection

1778 **doi:** <https://doi.org/10.1111/ajps.12646>, American Journal of Political Science

1779 **Report’s Abstract:** This short report summarises a replication exercise performed
1780 on data from Szakonyi (2021). The original work applies a difference-in-differences
1781 design to the case of an anti-corruption reform implemented in Russia for local
1782 election candidates, mandating financial disclosures. The author applies this design
1783 by comparing the electoral outcomes of municipalities that happened to hold elec-
1784 tions right after the reform with those that held elections right before the reform.
1785 For both groups, the design uses information from the previous electoral cycle as a
1786 pre-treatment benchmark. Using only data provided by the author in the original
1787 dataset, I first verified that results are reproducible when using alternative soft-
1788 ware. Second, I performed two simple placebo tests to obtain evidence on violations
1789 of the design’s identifying assumptions. These placebo tests return null results,
1790 reassuring on the reproducibility of the original findings.

1791 **Link to Full Report:** <https://osf.io/gx4d6/>

1792 **Link to Replicators’ Package:** <https://osf.io/gx4d6/>

1793 **Original Author’s Response:** “Many thanks to the replicator for taking the time
1794 to replicate and extend the paper. The placebo tests are very helpful in illuminating
1795 whether the identifying assumptions hold. I will make sure to run versions of these
1796 in future analyses.”

1797 **Original Authors’ Package:** [https://dataverse.harvard.edu/dataset.xhtml?](https://dataverse.harvard.edu/dataset.xhtml?persistentId=doi:10.7910/DVN/KDUMRM)
1798 [persistentId=doi:10.7910/DVN/KDUMRM](https://dataverse.harvard.edu/dataset.xhtml?persistentId=doi:10.7910/DVN/KDUMRM)

1799 **Reproduction Report**

1800 **Title Original Study:** Inflammatory Political Campaigns and Racial Bias in
1801 Policing

1802 **doi:** <https://doi.org/10.1093/qje/qjac037>, Quarterly Journal of Economics

1803 **Report's Abstract:** Grosjean et al. (2023) (GMY2023) estimate the causal effect
1804 of a Trump rally on the number of black drivers stopped by police officers, using a
1805 difference-in-difference approach. In their preferred specification, the authors find
1806 that after a Trump rally, the probability that a stopped driver is black increases by
1807 5.74%. This effect is significant at the 1% level. In this report we focus on repro-
1808 ducing the main claim of the paper. First, we reproduce the paper's main findings
1809 and uncover an issue with counties that experience multiple Trump rally treat-
1810 ments, given the original modelling choices taken in GMY2023. When we remove
1811 counties that experience multiple rallies, the estimated effect size drops to 2.46%
1812 and loses statistical significance. Second, we attempt to conduct a direct replica-
1813 bility check, by employing a new data set as a source for the dependent variable.
1814 We use data from the National Incident Based Reporting System (NIBRS). We
1815 observe no effect of Trump rallies both on the original data, covered by NIBRS
1816 and on the NIBRS data. Third, we conduct a robustness replicability exercise by
1817 coding an event-study difference-in-difference design at the day level. We estimate
1818 the event-study in a $[7; +7]$ window. We do not discover any systematic effect of
1819 Trump rallies on the dependent variable from GMY2023.

1820 **Link to Full Report:** <https://osf.io/xadb6/>

1821 **Link to Replicators' Package:** <https://osf.io/c7j58/>

1822 **Original Authors' Response:** <https://osf.io/xadb6/>

1823 **Original Authors' Package:** [https://dataverse.harvard.edu/dataset.xhtml?](https://dataverse.harvard.edu/dataset.xhtml?persistentId=doi:10.7910/DVN/A3B9HE)
1824 [persistentId=doi:10.7910/DVN/A3B9HE](https://dataverse.harvard.edu/dataset.xhtml?persistentId=doi:10.7910/DVN/A3B9HE)

1825 **Reproduction Report**

1826 **Title Original Study:** Interaction, Stereotypes, and Performance: Evidence from
1827 South Africa

1828 **doi:** <https://doi.org/10.1257/aer.20181805>, American Economic Review

1829 **Report's Abstract:** Corno, La Ferrara and Burns (2022) exploit the random allo-
1830 cation of freshman roommates in a large South African university to gauge the
1831 impact of a roommate's race on racial attitudes as measured by an implicit associ-
1832 ation test, and on school performance. They notably find that (a) white students
1833 randomly assigned to black roommates have less negative racial stereotypes, and
1834 (b) black students randomly assigned to live with white students have higher GPAs.
1835 We first reproduce all regression tables in Corno et al. (Corno et al. (2022)), and
1836 then test for robustness by varying the controls and conducting influential analysis.
1837 Overall, we find the results for finding (a) and (b) and robust in 15% and 40% of
1838 the robustness checks we ran, and the t/z scores are on average 78% and 85% as
1839 large as the original study.

1840 **Link to Full Report:** <https://osf.io/w7vpu/files>

1841 **Link to Replicators' Package:** <https://osf.io/w7vpu/files>

1842 **Link to Original Authors' Response:** Did not receive formal response as of
1843 November 2025.

1844 **Original Authors' Package:** [https://www.openicpsr.org/openicpsr/project/
1845 174501/version/V1/view](https://www.openicpsr.org/openicpsr/project/174501/version/V1/view)

1846 **Reproduction Report**

1847 **Title Original Study:** Interventions and Cognitive Spillovers

1848 **doi:** <https://doi.org/10.1093/restud/rdab087>, Review of Economic Studies

1849 **Report's Abstract:** In the paper of, Altmann et al. (2022) the authors investigate
1850 whether positive effects which are due to behavioral policy interventions in policy-
1851 targeted domains come along with negative effects in policy non-targeted domains.
1852 Using lab and online experiments where subjects have to solve one policy-focused
1853 decision task and one non-focused background task, the authors show that increas-
1854 ing incentives or steering attention to the former led to higher attention spans,
1855 lower default adherence rates, and a higher choice quality in the decision task.
1856 However, because of steering participants focus to the decision task, lower choice
1857 quality and lower attention spans in the background task emerged as a consequence,
1858 which was particularly pronounced among individuals with lower cognitive capa-
1859 bilities and complex decision tasks. Essentially, the authors also describe that the
1860 negative effects in the background tasks offset the positive effects in the decision
1861 task, ultimately yielding a net-zero effect overall. Therefore, the authors emphasize
1862 policymakers to also consider the potential negative cognitive spillovers in order to
1863 not overestimate the benefits of behavioral policy interventions. All the results the
1864 authors in the main text report are significant on 5% and 1% significance levels.
1865 All findings presented in the main text of the paper can be replicated using the
1866 original Stata code and verified thoroughly using R. Additionally, we performed
1867 two robustness tests to ensure the reliability of the paper's main results, and they
1868 remained consistent. Hence, the reported findings in the paper appear to be robust.

1869 **Link to Full Report:** [https://www.econstor.eu/bitstream/10419/272845/1/
1870 I4R-DP043.pdf](https://www.econstor.eu/bitstream/10419/272845/1/I4R-DP043.pdf)

1871 **Link to Replicators' Package:** <https://osf.io/kugbs/>

1872 **Original Authors' Response:** "We do not have any comments regarding the
1873 replication. We just want to briefly express our gratitude for the thorough and
1874 excellent work of the authors of the replication study."

1875 **Original Authors' Package:** <https://zenodo.org/records/5652808>

1876 **Reproduction Report**

1877 **Title Original Study:** Jumping the Gun: How Dictators Got Ahead of Their
1878 Subjects

1879 **doi:** <https://doi.org/10.1093/ej/ueac073>, Economic Journal

1880 **Report's Abstract:** Hariri and Wingender add new nuance to the traditional
1881 wisdom that economic modernisation is a path to democracy. They show that the
1882 diffusion of repressive, military technologies, causes a decline in the number of
1883 democratisations in the following years, and argue that this is because of a greater
1884 ability to forcefully oppress popular dissent. We conduct a robustness replication
1885 exercise, focussed on three tests: i) Are findings robust to alternative weightings of
1886 individual technologies in the instrument for country-aggregate military technol-
1887 ogy? ii) Is high leverage in individual countries, regions or time periods driving the
1888 global findings? iii) Are the strength of the IV and its independence of important
1889 macroeconomic indicators a chance occurrence? The main findings of the paper are
1890 largely robust to these tests.

1891 **Link to Full Report:** <https://econpapers.repec.org/paper/zbwi4rdps/50.htm>

1892 **Link to Replicators' Package:** <https://osf.io/4cx86/>

1893 **Original Authors' Response:** “we do not have any comments to the Reproduc-
1894 tion Report, so I'm just sending you this email to applaud the initiative. You and
1895 the Institute for Replication is doing a great service to the profession.”

1896 **Original Authors' Package:** <https://zenodo.org/records/7077694>

1897 **Reproduction Report**

1898 **Title Original Study:** Liquidity Constraints in the U.S. Housing Market

1899 **doi:** <https://doi.org/10.1093/restud/rdab063>, Review of Economic Studies

1900 **Report's Abstract:** Successful computational reproducibility. No re-analyses
1901 conducted.

1902 **Link to Original Authors' Response:** Authors provided feedback and sugges-
1903 tions.

1904 **Original Authors' Package:** <http://doi.org/10.5281/zenodo.5112964>

1905 **Reproduction Report**

1906 **Title Original Study:** Local Elites as State Capacity: How City Chiefs Use Local
1907 Information to Increase Tax Compliance in the Democratic Republic of the Congo
1908 **doi:** <https://doi.org/10.1257/aer.20201159>, American Economic Review

1909 **Report’s Abstract:** Balán et al. (2022) evaluate the impact of “local elites”
1910 involvement in local tax collection in a large city in the Democratic Republic of
1911 Congo. Using a randomized controlled trial to vary the identities of tax collectors,
1912 they find that local elites’ involvement raises tax compliance and total revenue by 50
1913 and 44 percent, respectively. The paper argues that the primary mechanism behind
1914 the results is better targeting made possible by local elites’ superior information
1915 about property holders’ willingness and ability to pay. In this replication comment,
1916 we first reproduce the paper’s main results. Then, we assess the robustness of the
1917 results by (1) employing randomization inference for statistical tests; (2) control-
1918 ling for baseline characteristics that are not balanced; and (3) using an alternative
1919 method to examine the claims supporting the preferred mechanism of better tar-
1920 geting. We find robust estimates in (1). However, the results are less robust both
1921 in terms of statistical significance and magnitude for (2) and (3). We conclude that
1922 the average treatment effect is robust, while the main claim about mechanisms,
1923 the information channel, is less robust to alternative estimation approaches. We
1924 contextualize and discuss the significance of these results, including the negligible
1925 revenue potential even under full compliance.

1926 **Link to Full Report:** <https://ideas.repec.org/p/zbw/i4rdps/191.html>

1927 **Link to Replicators’ Package:** [https://github.com/SossouAdjisse/](https://github.com/SossouAdjisse/LocalTaxReplicationProject.git)
1928 [LocalTaxReplicationProject.git](https://github.com/SossouAdjisse/LocalTaxReplicationProject.git)

1929 **Link to Original Authors’ Response:** [https://ideas.repec.org/p/zbw/i4rdps/](https://ideas.repec.org/p/zbw/i4rdps/192.html)
1930 [192.html](https://ideas.repec.org/p/zbw/i4rdps/192.html)

1931 **Original Authors’ Package:** [https://www.openicpsr.org/openicpsr/project/](https://www.openicpsr.org/openicpsr/project/147561/version/V1/view)
1932 [147561/version/V1/view](https://www.openicpsr.org/openicpsr/project/147561/version/V1/view)

1933 **Reproduction Report**

1934 **Title Original Study:** Major Reforms in Electricity Pricing: Evidence from a
1935 Quasi-Experiment

1936 **doi:** <https://doi.org/10.1093/ej/ueab076>, Economic Journal

1937 **Report's Abstract:** Labandeira et al. (2022) examine the effect of a policy in
1938 Spain that modified the electricity bill structure for all Spanish households. The pol-
1939 icy simultaneously increased fixed costs and decreased marginal costs on household
1940 electricity bills. Using fixed effects and instrumental variables (IV) specifications,
1941 the main causal finding in the paper is that the reform reduced house- hold electric-
1942 ity consumption for Spanish households by 15%. Their point estimate is statistically
1943 significant at the 1% level. In a similar specification, they find the reform reduced
1944 household expenditures on electricity by 9.8%, statistically significant at the 1%
1945 level. The code provided by the authors is computationally reproducible. We found
1946 two coding errors in different IV specifications, which had served as robustness
1947 checks to their main results. Correcting the errors removes statistical significance in
1948 two of four IV results, but increases the point estimates and statistical significance
1949 in the other two IV results. We also perform robustness checks. The IV estimates
1950 lose statistical significance in two of four robustness checks (with point estimates
1951 changing 1.1% to -39%). However, the OLS regressions are robust to changing
1952 covariates (sign and significance remained for 12 of 14 tests of the OLS specification,
1953 with changes in the estimates ranging from -157% to 64%, but averaging -3.3%).

1954 **Link to Full Report:** <https://osf.io/bysa7/>

1955 **Link to Replicators' Package:** <https://osf.io/bysa7/>

1956 **Original Authors' Response:** Original authors provided feedback. Multiple
1957 rounds of back and forth with replicators.

1958 **Original Authors' Package:** <https://zenodo.org/records/5423782>

1959 **Reproduction Report**

1960 **Title Original Study:** Market Access and Quality Upgrading: Evidence from
1961 Four Field Experiments

1962 **doi:** <https://doi.org/10.1257/aer.20210122>, American Economic Review

1963 **Report’s Abstract:** Bold et al. (2022b) investigate the effect of providing access
1964 to a market (i.e. a buyer) which rewards quality with a premium on farm productiv-
1965 ity and farming incomes from smallholder maize farmers in western Uganda, using
1966 a series of randomized experiments and a difference-in-differences approach. We
1967 successfully reproduce the results of this study using the publicly provided replica-
1968 tion packet. Then test the robustness of these results by re-defining treatment and
1969 outcome variables, testing for model misspecification and the leverage of outliers,
1970 and testing for non-random selection in the Fisher-permutation process. Our results
1971 show that the findings in Bold et al. (2022b) are robust to a variety of decisions in
1972 the research process. This evokes confidence in the internal validity of the findings.

1973 **Link to Full Report:** <https://econpapers.repec.org/paper/zbwi4rdps/72.htm>

1974 **Link to Replicators’ Package:** [https://journaldata.zbw.eu/dataset/
1975 bold-et-al-american-economic-review-2022](https://journaldata.zbw.eu/dataset/bold-et-al-american-economic-review-2022)

1976 **Original Authors’ Response:** “Thank you very much for sharing the report (and
1977 taking the time to replicate the study). We have no comments.”

1978 **Original Authors’ Package:** [https://www.openicpsr.org/openicpsr/project/
1979 158401/version/V1/view](https://www.openicpsr.org/openicpsr/project/158401/version/V1/view)

1980 **Reproduction Report**

1981 **Title Original Study:** Market-Based Monetary Policy Uncertainty

1982 **doi:** <https://doi.org/10.1093/ej/ueab086>, Economic Journal

1983 **Report's Abstract:** Bauer et al. (2022) derive market-based monetary policy
1984 uncertainty and uncover an 'FOMC uncertainty cycle' characterized by a fall of
1985 uncertainty after FOMC announcements and its subsequent built-up. Then, the
1986 authors show that the financial markets' response to monetary policy announce-
1987 ments depends on the level of short-rate uncertainty on the day before the FOMC
1988 announcement. First, we reproduced the paper's findings, though with Matlab
1989 version-specific issues. Second, we tested the robustness of the two main results of
1990 the paper. We show that the uncertainty cycle in the monetary policy uncertainty
1991 is confirmed when the crisis period is included in the sample or when the median
1992 instead of the average of changes in the monetary policy uncertainty is considered.
1993 However, the FOMC uncertainty cycle does not appear when the monetary pol-
1994 icy uncertainty index (Husted et al. 2020) or the daily economic policy uncertainty
1995 index (Baker et al. 2016) are used as uncertainty proxies.

1996 **Link to Full Report:** <https://econpapers.repec.org/paper/zbwi4rdps/77.htm>

1997 **Link to Replicators' Package:** <https://osf.io/qx8aw/>

1998 **Original Authors' Response:** "Thank you, glad to see that this work found our
1999 results to be rock solid!

2000 We won't write a response. Do let us know if you have any other questions
2001 about our work."

2002 **Original Authors' Package:** <https://zenodo.org/records/5566246>

2003 **Reproduction Report**

2004 **Title Original Study:** Market-Based Monetary Policy Uncertainty

2005 **doi:** <https://doi.org/10.1093/ej/ueab086>, Economic Journal

2006 **Report's Abstract:** This report replicates and examines Bauer et al.'s (2021)
2007 paper on monetary policy transmission to financial markets. The paper introduces
2008 novel measures of monetary policy uncertainty and analyses its drivers. It also
2009 investigates the impact of uncertainty changes on interest rates and financial asset
2010 prices. We assess reproducibility, consolidate market uncertainty measures using
2011 PCA and Factor Analysis, and rigorously test the reduction of uncertainty after
2012 Federal Market Open Committee (FOMC) announcements. Our findings support
2013 the paper's claim of reduced uncertainty on meeting days. Additionally, we explore
2014 the implications of the uncertainty channel on various financial assets, such as Gold,
2015 the Swiss Franc, European stock indexes, and Bitcoin.

2016 **Link to Full Report:** <https://econpapers.repec.org/paper/zbwi4rdps/76.htm>

2017 **Link to Replicators' Package:** [https://github.com/YaolangZhong/Nottingham_](https://github.com/YaolangZhong/Nottingham_Replication_Game/tree/main/replication_code)
2018 [Replication_Game/tree/main/replication_code](https://github.com/YaolangZhong/Nottingham_Replication_Game/tree/main/replication_code)

2019 **Original Authors' Response:** "Thank you, glad to see that this work found our
2020 results to be rock solid!

2021 We won't write a response. Do let us know if you have any other questions
2022 about our work."

2023 **Original Authors' Package:** <https://zenodo.org/records/5566246>

2024 **Reproduction Report**

2025 **Title Original Study:** Measuring the Welfare Effects of Shame and Pride

2026 **doi:** <https://doi.org/10.1257/aer.20190433>, American Economic Review

2027 **Report's Abstract:** This Reproduction Report examines and extends the research
2028 conducted by Butera, Metcalfe, Morrison, and Taubinsky (2022) on "The Welfare
2029 Effects of Pride and Shame." The original paper explores the welfare implications of
2030 public recognition as a motivator for desirable behavior and introduces an empirical
2031 methodology to measure Public Recognition Utility (PRU), which quantifies the
2032 utility individuals experience when their actions are publicly recognized. This report
2033 focuses on the real effort experiment reported in the paper that was conducted
2034 using a classroom sample, a lab sample, and an online sample. I computationally
2035 reproduce the original results and verify their robustness. While reproducing the
2036 results, I found two minor coding errors in the replication package. Correcting
2037 these errors slightly changes some estimates reported in the paper but does not
2038 turn over any results. The main treatment effect findings are further robust to
2039 using different sets of controls and sample selection criteria. Moreover, I conduct a
2040 heterogeneity analysis which reveals significant variations in how participants value
2041 public recognition. Overall, the replication study confirms the original conclusions
2042 while providing additional insights into the heterogeneity of PRU shapes on an
2043 individual level.

2044 **Link to Full Report:** <https://econpapers.repec.org/paper/zbwi4rdps/64.htm>

2045 **Link to Replicators' Package:** [https://github.com/tilmanfries/
2046 welfare-shame-pride-replication-report](https://github.com/tilmanfries/welfare-shame-pride-replication-report)

2047 **Original Authors' Final Response:** "Thanks again for all your hard work on
2048 this."

2049 **Original Authors' Package:** [https://www.openicpsr.org/openicpsr/project/
2050 145141/version/V1/view](https://www.openicpsr.org/openicpsr/project/145141/version/V1/view)

2051 **Reproduction Report**

2052 **Title Original Study:** Mental Health Costs of Lockdowns: Evidence from Age-
2053 Specific Curfews in Turkey

2054 **doi:** <https://doi.org/10.1257/app.20200811>, American Economic Journal: Applied
2055 Economics

2056 **Report's Abstract:** This report presents a replication of Altindag et al. (2022)
2057 performed at the Oslo Replication Games in 2022. Altindag et al. (2022) estimate
2058 the effects of an age-specific lockdown on mental health outcomes and mobility
2059 among adults aged 65 and older in Turkey, using a regression discontinuity design.
2060 The authors find a decline in mobility with a one-day decrease in the number of
2061 days being outside and an increase in the probability of never going out by 30 per-
2062 centage points. These point estimates are statistically significant at the 1% level.
2063 The mobility restrictions lead to a worsening in mental health outcomes of approx-
2064 imately 0.2 standard deviations, statistically significant at the 10% level in their
2065 preferred specification. In this paper we accomplish two things. First, we success-
2066 fully reproduce Altindag et al.'s main findings. Second, we test the ro-bustness
2067 of the results to a small number of changes to their preferred estimations by (1)
2068 not clustering the standard errors on the running variable, (2) not including con-
2069 trol variables, and (3) calculating the optimal bandwidth using another technique.
2070 Point estimates for mobility outcomes are stable to all three manipulations, and
2071 standard errors only change marginally. Point estimates and standard errors for the
2072 mental health outcomes are somewhat more sensitive, especially to changing the
2073 optimal bandwidth selection method. However, the observed changes are reason-
2074 ably expected when applying data-driven model selection methods to noisy data
2075 (to avoid over-fitting, it is likely preferable to apply a less data-driven approach
2076 like the original authors did). Our general impression is that the original analyses
2077 and results are both theoretically plausible and credible, despite some defensible
2078 model dependencies.

2079 **Link to Full Report:** <https://econpapers.repec.org/paper/zbwi4rdps/16.htm>

2080 **Link to Replicators' Package:** <https://osf.io/25u7b/>

2081 **Link to Original Authors' Response:** [https://econpapers.repec.org/paper/](https://econpapers.repec.org/paper/zbwi4rdps/17.htm)
2082 [zbwi4rdps/17.htm](https://econpapers.repec.org/paper/zbwi4rdps/17.htm)

2083 **Original Authors' Package:** [https://www.openicpsr.org/openicpsr/project/](https://www.openicpsr.org/openicpsr/project/131981/version/V1/view)
2084 [131981/version/V1/view](https://www.openicpsr.org/openicpsr/project/131981/version/V1/view)

2085 **Reproduction Report**

2086 **Title Original Study:** Mortality, Temperature, and Public Health Provision:
2087 Evidence from Mexico

2088 **doi:** <https://doi.org/10.1257/pol.20180594>, American Economic Journal: Economic
2089 Policy

2090 **Report's Abstract:** Cohen and Dechezleprêtre (2022) investigate the hetero-
2091 geneous impact of temperature on mortality across Mexico, and how affordable
2092 healthcare services that target the low-income population attenuate the mortal-
2093 ity effects of weather events. They find that while extreme temperatures are more
2094 dangerous than less extreme temperatures, the increased frequency of non-extreme
2095 temperatures mean these temperatures cause more deaths. First, we reproduce the
2096 paper's main findings, uncovering a minor coding error that has a trivial effect
2097 on the main results. Second, we test the robustness of the results to clustering at
2098 the state level, omitting precipitation, and using a different weighting scheme. The
2099 original results are robust to all of these changes.

2100 **Link to Full Report:** <https://econpapers.repec.org/paper/zbwi4rdps/90.htm>

2101 **Link to Replicators' Package:** <https://osf.io/q52e4/>

2102 **Original Authors' Response:** Cohen: "We thank The Institute for Replication.
2103 Next time, I will make sure I do not forget Feb. 29th in the code!"

2104 **Original Authors' Package:** [https://www.openicpsr.org/openicpsr/project/
2105 125201/version/V1/view](https://www.openicpsr.org/openicpsr/project/125201/version/V1/view)

2106 **Reproduction Report**

2107 **Title Original Study:** Motivated Beliefs and Anticipation of Uncertainty Reso-
2108 lution

2109 **doi:** <https://doi.org/10.1257/aeri.20200829>, American Economic Review: Insights

2110 **Report’s Abstract:** Drobner (2022) examines the effect of manipulating experi-
2111 mental subjects’ expectations about uncertainty resolution in learning about their
2112 performance on their belief updating patterns in an ego-relevant domain. In their
2113 preferred empirical specification, the author finds that individuals update their
2114 beliefs optimistically as they exhibit a higher belief adjustment in response to good
2115 compared to bad news only when they do not expect resolution of underlying uncer-
2116 tainty about their performance in an IQ test and neutrally when they know they will
2117 find out their relative performance at the end of the experiment. First, we reproduce
2118 the all of the paper’s findings without identifying any coding errors. Second, we test
2119 the robustness of the results to (1) adding individual covariates and (2) excluding
2120 subjects who exhibit a fundamental error in their belief updating from the analysis.
2121 We find no substantial changes in the main coefficients of interest with the inclu-
2122 sion of demographic variables in the analysis, consistent with demonstrated balance
2123 in covariates between the two experimental groups. Yet, several of the main esti-
2124 mates lose statistical significance and change from conservatism (under-updating)
2125 to over-inference (over-updating) in some conditions on the subset of participants
2126 excluding those who exhibit fundamental errors in belief updating.

2127 **Link to Full Report:** <https://econpapers.repec.org/paper/zbwi4rdps/65.htm>

2128 **Link to Replicators’ Package:** <https://osf.io/evt3a/>

2129 **Original Authors’ Response:** “Thanks for sharing the report. I think it’s a great
2130 initiative and feel free to publish this report on your webpage. I will not be able to
2131 provide an “answer”.”

2132 **Original Authors’ Package:** [https://www.openicpsr.org/openicpsr/project/
2133 139262/version/V1/view](https://www.openicpsr.org/openicpsr/project/139262/version/V1/view)

2134 **Reproduction Report**

2135 **Title Original Study:** Multinationals' Sales and Profit Shifting in Tax Havens
2136 **doi:** <https://doi.org/10.1257/pol.20200203>, American Economic Journal: Economic
2137 Policy

2138 **Report's Abstract:** We perform a robustness replication analysis of Laffitte and
2139 Toubal (2022), which considers how multinational corporations shift profit to "tax
2140 havens", jurisdictions where they face lower tax burdens. We find that the main
2141 results of Laffitte and Toubal (2022), are fairly robust to alternative versions of
2142 three important researcher choices: i) the definition of tax havens; ii) the use of
2143 a continuous measure of tax-friendliness rather than a binary classification of tax
2144 havens; and iii) a sample that omits two small but "extreme" tax havens: Bermuda
2145 and Barbados. In all cases, results remain of the same sign and retain statistical
2146 significance, though the magnitudes are somewhat attenuated in our robustness
2147 exercises.

2148 **Link to Full Report:** <https://econpapers.repec.org/paper/zbwi4rdps/37.htm>

2149 **Link to Replicators' Package:** <https://osf.io/3sbmr/>

2150 **Original Authors' Response:** "Thanks for your email and for replicating our
2151 exercise. Your work is useful. We recognize that the results remain consistent
2152 even when considering different interpretations of the haven concept and a smaller
2153 sample of observations.

2154 We are also pleased to hear that the replication file we shared with the AEJ:
2155 Policy has proven helpful."

2156 **Original Authors' Package:** [https://www.openicpsr.org/openicpsr/project/
2157 148301/version/V1/view](https://www.openicpsr.org/openicpsr/project/148301/version/V1/view)

2158 **Reproduction Report**

2159 **Title Original Study:** Multiracial identity and political preferences

2160 **doi:** <https://doi.org/10.1086/714760>, Journal of Politics

2161 **Report's Abstract:** The growing concern regarding reproducibility and replica-
2162 bility of social science re- sults has powered the adoption of open data and code
2163 requirements at journals and norms among researchers. However, even when these
2164 norms and requirements are fol- lowed, changes to the software used in data cleaning
2165 and analysis can render papers non-reproducible. This paper details the challenges
2166 of reproducibility in the face of software updates. We present a case study of a pub-
2167 lished article whose results are no longer reproducible due to changes in the software
2168 used. We then discuss the tools and techniques researchers can use to ensure that
2169 their research remains reproducible despite changes in the software used.

2170 **Link to Full Report:** <https://osf.io/ecymu/>

2171 **Link to Replicators' Package:** [https://github.com/taylorjwright/r_and_p_](https://github.com/taylorjwright/r_and_p_versioning)
2172 [versioning](https://github.com/taylorjwright/r_and_p_versioning)

2173 **Original Authors' Response:** Back and forth between authors and replicators,
2174 but did not obtain a final response as of November 2025.

2175 **Original Authors' Package:** [https://dataverse.harvard.edu/dataset.xhtml?](https://dataverse.harvard.edu/dataset.xhtml?persistentId=doi:10.7910/DVN/BLVJJH)
2176 [persistentId=doi:10.7910/DVN/BLVJJH](https://dataverse.harvard.edu/dataset.xhtml?persistentId=doi:10.7910/DVN/BLVJJH)

2177 **Reproduction Report**

2178 **Title Original Study:** News Shocks Under Financial Frictions
2179 **doi:** <https://doi.org/10.1257/mac.20170066>, American Economic Journal: Macroeconomics
2180

2181 **Report's Abstract:** Görtz et al. (2022) estimate the effects of innovations to
2182 future total factor productivity (TFP) on financial markets. In a Bayesian vector
2183 autoregression, they identify a TFP news shock as one that explains the largest
2184 share of 40- quarter ahead forecast error variance (FEV) of TFP. Their estimated
2185 impulse responses functions show that a positive news shock significantly decreases
2186 credit market spreads and increases credit market supply. They also find that a
2187 shock that explains the maximum of the FEV of the "excess bond premium" (EBP)
2188 (Gilchrist and Zakrajsek 2012) causes similar responses. These results are consistent
2189 with an estimated DSGE model with financial frictions. We estimate the main IRFs
2190 of the study using the original data and a frequentist estimation approach. We
2191 obtain similar point estimates for the dynamic responses to TFP news and EBP
2192 max-share shocks. We also update their macroeconomic and financial time series,
2193 as some of the data has been revised substantially since their original estimate.
2194 We use the updated data to re-estimate the above-mentioned IRFs, and we find
2195 that the results are robust to this change in the data. Finally, we investigate the
2196 computational reproducibility of their DSGE results, and find that their provided
2197 code (consistent with warnings in their README file) does not execute in the most
2198 recent version of Dynare or Matlab. Using the version indicated in their replication
2199 files, we encounter issues estimating the posterior mode.

2200 **Link to Full Report:** <https://econpapers.repec.org/paper/zbwi4rdps/51.htm>

2201 **Link to Replicators' Package:** [https://github.com/gionikola/](https://github.com/gionikola/replication-game-ucsd)
2202 [replication-game-ucsd](https://github.com/gionikola/replication-game-ucsd)

2203 **Original Authors' Final Response:** "Thank you for the update and considering
2204 our work for replication."

2205 **Original Authors' Package:** [https://www.openicpsr.org/openicpsr/project/](https://www.openicpsr.org/openicpsr/project/130141/version/V1/view)
2206 [130141/version/V1/view](https://www.openicpsr.org/openicpsr/project/130141/version/V1/view)

2207 **Reproduction Report**

2208 **Title Original Study:** Non-Linearities, State-Dependent Prices and the Trans-
2209 mission Mechanism of Monetary Policy

2210 **doi:** <https://doi.org/10.1093/ej/ueab049>, Economic Journal

2211 **Report's Abstract:** Ascari and Haber (2022) fill the gaps in the literature by
2212 showing evidence in favor of the state-dependent sticky price model's predictions
2213 using the macro-aggregates. We report a replication and robustness check of the
2214 study. We employ several additional macroeconomic control variables and different
2215 alternative measurements for monetary policy shocks and find that the original
2216 results remain qualitatively robust. Our analysis further shows that the turbulent
2217 periods of inflation in the 1970s and 1980s have an important role in claiming the
2218 robustness of the original results.

2219 **Link to Full Report:** <https://osf.io/kbwap/>

2220 **Link to Replicators' Package:** <https://osf.io/kbwap/>

2221 **Original Authors' Response:** No response.

2222 **Original Authors' Package:** [https://oup.silverchair-cdn.com/](https://oup.silverchair-cdn.com/oup/backfile/Content_public/Journal/ej/132/641/10.1093_ej_ueab049/)
2223 [oup/backfile/Content_public/Journal/ej/132/641/10.1093_ej_ueab049/](https://oup.silverchair-cdn.com/oup/backfile/Content_public/Journal/ej/132/641/10.1093_ej_ueab049/)

2224 [1/ueab049_replication_package.zip?Expires=1765387359&Signature=](https://oup.silverchair-cdn.com/oup/backfile/Content_public/Journal/ej/132/641/10.1093_ej_ueab049/)

2225 [Ysj-YfzVPuVPtgGm1ZBFtFNI1APv5x1Rgajkm2orCIsJXt2dZG3Amp92XbuA6m0iP-4LwOFv~PFKA4t1qXm](https://oup.silverchair-cdn.com/oup/backfile/Content_public/Journal/ej/132/641/10.1093_ej_ueab049/)

2226 [_&Key-Pair-Id=APKAIE5G5CRDK6RD3PGA](https://oup.silverchair-cdn.com/oup/backfile/Content_public/Journal/ej/132/641/10.1093_ej_ueab049/)

2227 **Reproduction Report**

2228 **Title Original Study:** Not All Elections Are Created Equal: Election Quality
2229 and Civil Conflict

2230 **doi:** <https://doi.org/10.1086/714778>, Journal of Politics

2231 **Report's Abstract:** Utilizing a time-series cross-sectional dataset on developing
2232 countries, Donno et al. (2022) examine how variation in election quality shapes
2233 opportunities and incentives for civil conflict. Across a number of models in their
2234 analysis, they find that civil conflict is more likely when elections are not free and
2235 fair. They also find that for countries with low integrity elections, the probability of
2236 conflict occurring is higher if it has experienced conflict before. We begin by repro-
2237 ducing Donno et al.'s (2022) main models and findings, which yielded no coding
2238 errors or differences in effect estimates. Afterwards, for replication purposes we run
2239 a series of robustness and conceptual replication tests. For our first replication, we
2240 examine the heterogeneous effect between electoral integrity and ethnic fractional-
2241 ization on conflict. Our second test examines whether a subsample of authoritarian
2242 regimes should have been included in the authors' original analysis.

2243 **Link to Full Report:** <https://osf.io/unhkr/>

2244 **Link to Replicators' Package:** [https://drive.google.com/drive/folders/1Vlwfr3_](https://drive.google.com/drive/folders/1Vlwfr3_Q7c56XFPsQuKUNSnEU_lMFUQz?usp=sharing)
2245 [Q7c56XFPsQuKUNSnEU_lMFUQz?usp=sharing](https://drive.google.com/drive/folders/1Vlwfr3_Q7c56XFPsQuKUNSnEU_lMFUQz?usp=sharing)

2246 **Link to Original Authors' Response:** <https://osf.io/unhkr/>

2247 **Original Authors' Package:** [https://dataverse.harvard.edu/dataset.xhtml?](https://dataverse.harvard.edu/dataset.xhtml?persistentId=doi:10.7910/DVN/8B31FG)
2248 [persistentId=doi:10.7910/DVN/8B31FG](https://dataverse.harvard.edu/dataset.xhtml?persistentId=doi:10.7910/DVN/8B31FG)

2249 **Reproduction Report**

2250 **Title Original Study:** Parties as Disciplinarians: Charisma and Commitment
2251 Problems in Programmatic Campaigning

2252 **doi:** <https://doi.org/10.1111/ajps.12638>, American Journal of Political Science

2253 **Report’s Abstract:** Hollyer, Klašnja, and Titiumik (2022) analyse the trade-off
2254 that political parties face between running programmatic campaigns and fielding
2255 charismatic candidates, whose electoral appeal may come at the cost of undermin-
2256 ing the party brand. They argue that higher electoral volatility prompts parties
2257 to rely on charismatic candidates, even though they might not be as loyal to
2258 the party’s programmatic stance. They substantiate their argument with a cross-
2259 national dataset and a quantitative case study in Brazil. We computationally
2260 reproduced and conducted further robustness tests for their cross-national study by
2261 translating the Stata code to R. Next, we conducted a computational reproduction
2262 and some additional robustness tests for the quantitative case study. We find that
2263 their cross-national analysis is reproducible, albeit with some minor discrepancies.
2264 The quantitative case study is also largely reproducible and both are robust in sev-
2265 eral ways. We conclude by making some suggestions about data dissemination and
2266 robustness checks for authors of regression discontinuity designs.

2267 **Link to Full Report:** <https://econpapers.repec.org/paper/zbwi4rdps/54.htm>

2268 **Link to Replicators’ Package:** [https://osf.io/93gfx/?view_only=](https://osf.io/93gfx/?view_only=7063353244d646ffaf7bfd53013e3143)
2269 [7063353244d646ffaf7bfd53013e3143](https://osf.io/93gfx/?view_only=7063353244d646ffaf7bfd53013e3143)

2270 **Original Authors’ Response:** “Thanks for your note and for all the work of
2271 Kelly, Odermatt, and Metson in replicating our paper. [...] Our read of the Repro-
2272 duction Reports that the findings in our paper hold in the Kelly et al replication.
2273 [...] Our sense is that the discrepancies between the replication and original paper
2274 are sufficiently small, and the task of comparing the replication R code to the orig-
2275 inal Stata code is likely to be sufficiently demanding of time, that the opportunity
2276 cost of a thorough response is high. So, I think we’ll forgo the opportunity to draft
2277 a response, and just let the replication stand without reply.

2278 We’ll leave it to any sufficiently interested parties with expertise in both Stata
2279 and R to iron out the discrepancies between the replication and original paper.”

2280 **Original Authors’ Package:** [https://dataverse.harvard.edu/dataset.xhtml?](https://dataverse.harvard.edu/dataset.xhtml?persistentId=doi:10.7910/DVN/AWSQTW)
2281 [persistentId=doi:10.7910/DVN/AWSQTW](https://dataverse.harvard.edu/dataset.xhtml?persistentId=doi:10.7910/DVN/AWSQTW)

2282 **Reproduction Report**

2283 **Title Original Study:** Patience, Risk-Taking, and Human Capital Investment
2284 Across Countries

2285 **doi:** <https://doi.org/10.1093/ej/ueab105>, Economical Journal

2286 **Report's Abstract:** Hanushek et al. (2021) test how country-level measures of
2287 patience and risk-taking from the Global Preference Survey predict student per-
2288 formance on the Programme for International Student Assessment (PISA) math
2289 test. They find that country-level patience positively predicts math test scores and
2290 country-level risk-taking negatively predicts math test scores. They find similar
2291 results when holding country of residence characteristics constant and focusing on
2292 the preferences of the country of origin of migrants. We have checked the com-
2293 putational reproducibility and find that the data and analysis script provided by
2294 the authors allowed us to exactly reproduce the main tables in the paper. We
2295 also checked the robustness replicability by testing how robust the results are to
2296 decisions about imputation, weighting, operationalization of dependent variables,
2297 choice of control variables, and the inclusion of highly leveraged observations. We
2298 see that results are generally robust, though statistical significance of the risk-
2299 taking coefficient in the migrant analysis hinges on whether a control for OECD
2300 country of residence is included. Finally, we check the conceptual replicability of
2301 the results by using data from the Trends in International Mathematics and Science
2302 Study (TIMSS) instead of PISA - a different dataset with a different standardized
2303 test. This exercise shows that their results are robust to expanding the analysis to
2304 countries participating in both PISA and TIMSS.

2305 **Link to Full Report:** <https://econpapers.repec.org/paper/zbwi4rdps/48.htm>

2306 **Link to Replicators' Package:** <https://osf.io/kgt8z/>

2307 **Link to Original Authors' Response:** "We would like to thank the replicators
2308 and compliment them for their thoughtful replication and extension of our paper.
2309 We are particularly impressed by the extension to the TIMSS data, which is actually
2310 great support for the underlying idea. We do not see a reason to formulate a formal
2311 response for your website.

2312 Thank you all for your valuable work!"

2313 **Original Authors' Package:** [https://www.openicpsr.org/openicpsr/project/
2314 153101/version/V2/view](https://www.openicpsr.org/openicpsr/project/153101/version/V2/view)

2315 **Reproduction Report**

2316 **Title Original Study:** Peer Effects in Academic Research: Senders and Receives
2317 **doi:** <https://doi.org/10.1093/ej/ueac031>, Economical Journal

2318 **Report's Abstract:** In this report, we provide an overview from reproducing and
2319 replicating Bosquet et al. (2022). As a first step, we successfully reproduce all the
2320 results in the paper, as well as figure A1. All results were fully reproducible and
2321 match the published version of the paper. Next, we carry out three sensitivity
2322 analysis. We examine how the main results change from the weights used, additional
2323 controls, and author-university pairs. The main results are robust to these checks.

2324 **Link to Full Report:** <https://osf.io/czkgw/>

2325 **Link to Replicators' Package:** <https://osf.io/czkgw/>

2326 **Link to Original Authors' Response:** The authors responded to the replicators'
2327 questions. Bosquet then responded to the final report: "I would simply thank the
2328 team of replicators and I am happy to see that the tested results are robust to the
2329 tested alternatives. As written in my previous email, I think those kinds of efforts
2330 are very useful for the community and the credibility of published results so thanks
2331 as well for that."

2332 **Original Authors' Package:** <https://zenodo.org/records/6457037>

2333 **Reproduction Report**

2334 **Title Original Study:** Playing Politics with Environmental Protection: The
2335 Political Economy of Designating Protected Areas

2336 **doi:** <https://doi.org/10.1086/718978>, Journal of Politics

2337 **Report's Abstract:** Mangonnet et al. (2022) examine whether political alignment
2338 at the national and sub-national levels explain the spatial designation of Protected
2339 Areas (PAs) in Brazil. Their identification relies on spatial discontinuities in polit-
2340 ical alignment across municipalities. They find that a president-mayor coalition
2341 alignment reduces the incidence of PAs by about one percentage point, whereas
2342 they find no party alignment effects. We were able to reproduce the paper's findings
2343 using the same code and software. Alternative software routines reproduce their
2344 results with small and inconsequential numerical differences. Moreover, robustness
2345 replications find consistent results for one out the two treatments. Finally, we find
2346 no evidence of fabrication of data.

2347 **Link to Full Report:** <https://econpapers.repec.org/paper/zbwi4rdps/73.htm>

2348 **Link to Replicators' Package:** <https://osf.io/t76jd/>

2349 **Original Authors' Response:** "We are grateful to Laura Villalobos, Jill Caviglia-
2350 Harris, Tharaka Jayalath, and the team at the Institute for Replication for
2351 generously replicating our work. We encourage readers to follow their alternative
2352 software routines for faster estimations."

2353 **Original Authors' Package:** [https://dataverse.harvard.edu/dataset.xhtml?](https://dataverse.harvard.edu/dataset.xhtml?persistentId=doi:10.7910/DVN/N6LIMH)
2354 [persistentId=doi:10.7910/DVN/N6LIMH](https://dataverse.harvard.edu/dataset.xhtml?persistentId=doi:10.7910/DVN/N6LIMH)

2355 **Reproduction Report**

2356 **Title Original Study:** Policy Deliberation and Voter Persuasion: Experimental
2357 Evidence from an Election in the Philippines

2358 **doi:** <https://doi.org/10.1111/ajps.12566>, American Journal of Political Science

2359 **Report's Abstract:** I would characterize my robustness replication as almost
2360 entirely successful. The design checks I report all support a straightforward under-
2361 standing of the design. My effect and uncertainty estimates barely differ from the
2362 original estimates (when compared with like estimation procedures), with any dis-
2363 crepancies attributable to simulation error. One small area of difference was the
2364 weighting scheme employed by the authors to correct for “over-representation” of
2365 meeting attendees in the treatment group. As discussed below, I do not understand
2366 the design reason for this choice and when I simulate its properties, I can obtain
2367 small amounts of bias. The net consequence of their approach was usually to make
2368 coefficient estimates smaller, so we don't have a major difference in conclusion
2369 except perhaps in a secondary analysis of mechanisms.

2370 **Link to Full Report:** <https://osf.io/y8ubt/>

2371 **Link to Replicators' Package:** <https://osf.io/y8ubt/>

2372 **Link to Original Authors' Response:** <https://osf.io/y8ubt/>

2373 **Original Authors' Package:** [https://dataverse.harvard.edu/dataset.xhtml?](https://dataverse.harvard.edu/dataset.xhtml?persistentId=doi:10.7910/DVN/S3HACJ)
2374 [persistentId=doi:10.7910/DVN/S3HACJ](https://dataverse.harvard.edu/dataset.xhtml?persistentId=doi:10.7910/DVN/S3HACJ)

2375 **Reproduction Report**

2376 **Title Original Study:** Political Turnover, Bureaucratic Turnover, and the Quality
2377 of Public Services

2378 **doi:** <https://doi.org/10.1257/aer.20171867>, American Economic Review

2379 **Report's Abstract:** The politically motivated replacement in local governments
2380 is a pervasive fact in our modern democracies. Whether it has causal effects on
2381 the quality of public services, such as education, is a critical question and yet
2382 understudied. This paper uses a regression discontinuity design (RDD) for close
2383 elections to replicate Akthari, Moreira and Trucco (2022) who find negative effects
2384 on the quality of public education in Brazil (.05-.08 standard deviations of lower test
2385 scores). I first reproduce these main results, finding minor computational differences
2386 that have no effect on the conclusions. I also show that the estimates for Brazil
2387 are in general robust to different specifications following Brodeur, Cook and Heyes
2388 (2020). Finally, I implement the same RDD framework now applied to Chilean
2389 administrative records to find null effects on test scores. Taken together, these
2390 results suggest that political turnover has weakly negative effects on service quality.

2391 **Link to Full Report:** <https://econpapers.repec.org/paper/zbwi4rdps/39.htm>

2392 **Link to Replicators' Package:** <https://osf.io/q43vz/>

2393 **Link to Original Authors' Response:** <https://osf.io/kv4pj/>

2394 **Original Authors' Package:** [https://www.openicpsr.org/openicpsr/project/
2395 150323/version/V1/view](https://www.openicpsr.org/openicpsr/project/150323/version/V1/view)

2396 **Reproduction Report**

2397 **Title Original Study:** Pre-Colonial Warfare and Long-Run Development in India
2398 **doi:** <https://doi.org/10.1093/ej/ueab089>, Economic Journal

2399 **Report's Abstract:** We test the reproducibility and replicability of Dincecco et
2400 al. (2022), which reports a positive relationship between pre-colonial interstate
2401 warfare and long-run development patterns across India. Overall, we confirm that
2402 all of the study's estimates are computationally reproducible by using both the
2403 provided replication package in Stata and code written by the present authors in
2404 R. We test for and find no evidence of data manipulation in the final datasets.
2405 Concerning direct replicability, we consider different ways of measuring distance to
2406 conflicts and also alternative proxies for both the dependent variable and variables
2407 which capture channels by which the main effects operate. We are able to replicate
2408 the magnitude and significance of the estimated coefficient on conflict exposure in
2409 most of the tests, noting that while most estimates are substantively in line with
2410 the original study, some alternative measures of distance to conflict imply different
2411 magnitudes for estimates, and proxy estimates are sensitive to both the time period
2412 and type of conflict considered.

2413 **Link to Full Report:** <https://econpapers.repec.org/paper/zbwi4rdps/35.htm>

2414 **Link to Replicators' Package:** <https://osf.io/af6m2/>

2415 **Link to Original Authors' Response:** <https://osf.io/af6m2/>

2416 **Original Authors' Package:** <https://zenodo.org/records/5583263>

2417 **Reproduction Report**

2418 **Title Original Study:** Public Infrastructure and Economic Development: Evi-
2419 dence from Postal Systems

2420 **doi:** <https://doi.org/10.1111/ajps.12594>, American Journal of Political Science

2421 **Report's Abstract:** Rogowski et al. (2022) use secondary data to study the impact
2422 of historic postal infrastructure on economic development, both cross-country and
2423 within the US. Their results suggest a large positive effect of post offices on economic
2424 development that is robust across various sensitivity checks. We successfully com-
2425 putationally reproduce all results. In a robustness assessment, we find the results
2426 to be robust to simple changes in the analysis but observe some sensitivity to
2427 accounting for spatial trends in the cross-country analysis. Additionally, we correct
2428 a coding inconsistency, showing that in the corrected version, one main robustness
2429 check for the US-analysis is no longer supporting the result. Despite this, we find
2430 the results to be overall robust given the numerous analyses and robustness checks
2431 in the original paper.

2432 **Link to Full Report:** <https://econpapers.repec.org/paper/zbwi4rdps/92.htm>

2433 **Link to Replicators' Package:** [https://osf.io/j3ydr/?view_only=](https://osf.io/j3ydr/?view_only=ad14a07cb3a741ca9bbfab391ad7c6fe)
2434 [ad14a07cb3a741ca9bbfab391ad7c6fe](https://osf.io/j3ydr/?view_only=ad14a07cb3a741ca9bbfab391ad7c6fe)

2435 **Original Authors' Response:** "Thanks so much for reproducing the findings in
2436 our paper and exploring extensions of our results. We also appreciate your sharing
2437 the report with us. [...] I [Rogowski] confirm that we are comfortable letting your
2438 report stand and that we will not write a response to it. "

2439 **Original Authors' Package:** [https://dataverse.harvard.edu/dataset.xhtml?](https://dataverse.harvard.edu/dataset.xhtml?persistentId=doi:10.7910/DVN/33K3EF)
2440 [persistentId=doi:10.7910/DVN/33K3EF](https://dataverse.harvard.edu/dataset.xhtml?persistentId=doi:10.7910/DVN/33K3EF)

2441 **Reproduction Report**

2442 **Title Original Study:** Re-Assessing Elite-Public Gaps in Political Behavior
2443 **doi:** <https://doi.org/10.1111/ajps.12583>, American Journal of Political Science
2444 **Report’s Abstract:** Kertzer (2022) conducts a meta-analysis of parallel experi-
2445 ments on samples of political elites and ordinary citizens. He examines whether the
2446 average treatment effect for elites is significantly different from the average treat-
2447 ment effect for citizens, finding that only 19 of 162 (11.7%) difference-in-difference
2448 estimates are statistically significant after adjusting for the false discovery rate. He
2449 also finds that elites and masses hold similar foreign policy attitudes after control-
2450 ling for their demographic characteristics. In this Reproduction Report, we begin
2451 by running robustness and heterogeneity tests for the first claim. We find that the
2452 results survive many robustness tests. We also find, however, that only a small num-
2453 ber of the these treatments significantly affected masses (N=28) or elites (N=30).
2454 This low rate suggests the possibility that almost all of these experiments failed to
2455 successfully manipulate either masses or elites. If so, we may not be able to con-
2456 clude that masses and elites respond similarly to experiments with confidence until
2457 political scientists produce more experiments with actual treatment effects or with
2458 successful manipulation checks in cases of null effects. In the second part of this
2459 Reproduction Report, we conceptually replicate the second Kertzer analysis, find-
2460 ing a strong correlation between elite and mass political decisions and attitudes,
2461 thus confirming Kertzer’s analysis.
2462 **Link to Full Report:** [https://www.econstor.eu/bitstream/10419/266385/1/](https://www.econstor.eu/bitstream/10419/266385/1/I4R-DP010.pdf)
2463 [I4R-DP010.pdf](https://www.econstor.eu/bitstream/10419/266385/1/I4R-DP010.pdf)
2464 **Link to Replicators’ Package:** <https://osf.io/93urk/>
2465 **Original Authors’ Response:** “Thank you for your email and for the invitation.
2466 [...] please send my appreciation to the authors for their interest in the manuscript;
2467 I find their analysis very interesting.”
2468 **Original Authors’ Package:** [https://dataverse.harvard.edu/dataset.xhtml?](https://dataverse.harvard.edu/dataset.xhtml?persistentId=doi:10.7910/DVN/LHOTOK)
2469 [persistentId=doi:10.7910/DVN/LHOTOK](https://dataverse.harvard.edu/dataset.xhtml?persistentId=doi:10.7910/DVN/LHOTOK)

2470 **Reproduction Report**

2471 **Title Original Study:** Rebel on the Canal: Disrupted Trade Access and Social
2472 Conflict in China, 1650–1911

2473 **doi:** <https://doi.org/10.1257/aer.20201283>, American Economic Review

2474 **Report’s Abstract:** Cao and Chen (2022a) study the role of disruption of trade
2475 on social conflict in China in the 19th century. Identification builds on the closure
2476 of China’s Grand Canal in 1826 in a difference-in-differences framework. In their
2477 preferred analytical specification, the authors find that counties along the canal
2478 experienced a 117 percent increase in rebelliousness after the initial closure of the
2479 canal in 1826 relative to their non-canal counterparts. First, we reproduce the
2480 paper’s main findings using the official replication package. Second, we examine
2481 whether a sub-sample of counties/prefectures/provinces drives the result. Third,
2482 we test the robustness of the results to alternative treatment periods.

2483 **Link to Full Report:** <https://osf.io/dhn6e/>

2484 **Link to Replicators’ Package:** <https://osf.io/dhn6e/>

2485 **Link to Original Authors’ Response:** <https://osf.io/dhn6e/>

2486 **Original Authors’ Package:** [https://www.openicpsr.org/openicpsr/project/
2487 157781/version/V1/view](https://www.openicpsr.org/openicpsr/project/157781/version/V1/view)

2488 **Reproduction Report**

2489 **Title Original Study:** Recessions, Mortality, and Migration Bias: Evidence from
2490 the Lancashire Cotton Famine

2491 **doi:** <https://doi.org/10.1257/app.20190131>, American Economic Journal: Applied
2492 Economics

2493 **Report's Abstract:** Vellore Arthi, Brian Beach and W. Walker Hanlon (2022)
2494 investigate the effect of the Lancashire Cotton Famine on mortality, accounting
2495 for the migration response to the downturn. They use difference-in-differences to
2496 estimate the effect of the cotton famine on mortality. To account for the migration
2497 response to the cotton famine, they construct a linked dataset giving mortality
2498 rates by district of residence during the cotton famine, rather than by district of
2499 residence at the time of death. They find that the cotton famine increased mortality
2500 in cotton-textile producing districts, and that accounting for migration matters,
2501 in the sense that their estimates would have been markedly different had they
2502 not accounted for it. I check that ABH results are fully reproducible using their
2503 data and code, and that their claims are robust to (1) decreasing the age window
2504 for building the linked dataset, (2) modifying the specification and (3) computing
2505 different standard errors. The only significant discrepancy in results is that I find
2506 stronger effects of the cotton famine when I decrease the age window for building
2507 the linked dataset, likely because this reduces measurement errors.

2508 **Link to Full Report:** <https://econpapers.repec.org/paper/zbwi4rdps/25.htm>

2509 **Link to Replicators' Package:** [https://www.openicpsr.org/openicpsr/project/
2510 192272/version/V1/view](https://www.openicpsr.org/openicpsr/project/192272/version/V1/view)

2511 **Original Authors' Response:** "Thanks for the interest in our work. We've had
2512 a chance to review the report and it looks like everything replicated. Since there
2513 are no outstanding queries, we are happy to sign off on this."

2514 **Original Authors' Package:** [https://www.openicpsr.org/openicpsr/project/
2515 128521/version/V1/view](https://www.openicpsr.org/openicpsr/project/128521/version/V1/view)

2516 **Reproduction Report**

2517 **Title Original Study:** Reshaping Adolescents' Gender Attitudes: Evidence from
2518 a School-Based Experiment in India

2519 **doi:** <https://doi.org/10.1257/aer.20201112>, American Economic Review

2520 **Report's Abstract:** Dhar et al. (2022) examine the effect of a gender attitude
2521 change program in secondary schools in India. In their preferred specification, the
2522 authors show that the program made the students report more gender-egalitarian
2523 attitudes by 0.18 of a standard deviation, and shifted self-reported behaviors to
2524 be more aligned with gender-progressive norms by 0.20 standard deviations (both
2525 significant at 1% level). In contrast, they found no effect on girls' aspirations,
2526 as these were already high before the intervention. The effects did not attenuate
2527 between the first end-line (right after the programme was completed) and the second
2528 (two years later). To put the paper's results in perspective, we first comment on
2529 the authors' deviations from their pre-registration and pre-analysis plans, provide
2530 detailed power calculations, and add multiple-hypothesis-testing-adjusted standard
2531 errors. Second, we show that the paper's results are perfectly reproducible. Third,
2532 we show that the results are robust to excluding control variables, and alternative
2533 ways of constructing indices and dealing with non-response.

2534 **Link to Full Report:** <https://econpapers.repec.org/paper/zbwi4rdps/24.htm>

2535 **Link to Replicators' Package:** <https://osf.io/r5jfe/>

2536 **Final Original Authors' Response:** "Thanks. the revision looks good. I actually
2537 don't think we need to have a formal response any more. [...] Thus, I don't think
2538 there is anything substantive for us to include in a discussion paper/response. That
2539 reflects the fact that the Reproduction Reports fair and there is nothing major to
2540 respond to, so it's good news, from both the perspective of the integrity of our
2541 original paper and the professionalism of the replication."

2542 **Original Authors' Package:** [https://www.openicpsr.org/openicpsr/project/
2543 149882/version/V1/view](https://www.openicpsr.org/openicpsr/project/149882/version/V1/view)

2544 **Reproduction Report**

2545 **Title Original Study:** Run-off Elections in the Laboratory

2546 **doi:** <https://doi.org/10.1093/ej/ueab051>, Economic Journal

2547 **Report's Abstract:** Bouton et al. (2022) make a causal claim by manipulating
2548 the voting system under which participants vote (runoff or plurality) and exam-
2549 ining whether this manipulation affects the proportion of strategic voting. They
2550 estimate the effect of the voting system on the proportion of strategic voting for
2551 the participant population, using random effect regression where standard errors
2552 are clustered on group level. Regarding replication results, we reproduced the orig-
2553 inal study's main findings. Our analysis confirms that there are minor and mostly
2554 non-significant disparities in electoral outcomes and voters' welfare between the
2555 two voting systems, consistent with the original study's conclusions. Specifically,
2556 we conducted tests to assess the study's computational reproducibility and direct
2557 replicability. While the authors provided the raw data, they did not include a script
2558 for cleaning it or a codebook describing its contents. Consequently, we developed a
2559 data cleaning script to ensure accurate and consistent data processing.

2560 **Link to Full Report:** <https://osf.io/a8cev/>

2561 **Link to Replicators' Package:** [https://github.com/carinahausladen/](https://github.com/carinahausladen/runoff-elections)
2562 [runoff-elections](https://github.com/carinahausladen/runoff-elections)

2563 **Original Authors' Response:** The authors provided feedback which was taken
2564 into account.

2565 **Original Authors' Package:** <https://doi.org/10.1093/ej/ueab051>

2566 **Reproduction Report**

2567 **Title Original Study:** School Spending and Student Outcomes: Evidence from
2568 Revenue Limit Elections in Wisconsin

2569 **doi:** <https://doi.org/10.1257/pol.20200226>, American Economic Journal: Economic
2570 Policy

2571 **Report’s Abstract:** Baron (2022) explores the independent effects of operational
2572 expenditure and capital expenditure on student outcomes in school districts across
2573 Wisconsin from the outcomes of close referendum approvals. By utilizing a dynamic
2574 regression discontinuity framework and cubic specification, the author finds that
2575 narrowly passing an operational referendum, increases operational expenditure per
2576 pupil by \$298 each year on average, following the referendum over a ten year period.
2577 From this \$198 are spent on instructional expenses. These point estimates are
2578 statistically significant at the 10% and 5% level, respectively. We first reproduce
2579 the main results from the paper without any issues arising. Secondly, we conduct
2580 a robustness replicability to (1) dropping school districts from the top and bottom
2581 5% of the revenue limits distribution, categorically, and (2) dividing the time frame
2582 of the study into two periods: 1996-2005 and 2005-2014. We find that dropping the
2583 top 5% of the school districts by revenue limits reduces the additional operational
2584 expenditure by \$140 per pupil (lower by 50 percent) and the effects of passing an
2585 operational referendum were nearly double in the former period compared to the
2586 latter period. Lastly, we find that the estimated effects on student outcomes rely
2587 heavily on recent observations.

2588 **Link to Full Report:** <https://econpapers.repec.org/paper/zbwi4rdps/88.htm>

2589 **Link to Replicators’ Package:** <https://osf.io/m2w4x/>

2590 **Original Author’s Response:** “Thank you for sharing the Reproduction Report.
2591 Please pass on my thanks to the replicators for their important work. First and
2592 foremost, I’m glad to see that the results in the paper are reproducible without any
2593 issues arising. The report explores two additional sources of heterogeneity. I have
2594 no additional comments on these. I do briefly want to clarify the last sentence in
2595 the report’s abstract, which reads “Lastly, we find that the estimated effects on
2596 student outcomes rely heavily on recent observations.” While I am not entirely sure
2597 what the replicators are referring to, my guess is that they refer to Table 2 in the
2598 report. In this table, they discuss that they are unable to study heterogeneity in
2599 the impacts of passing a referendum on test scores and postsecondary enrollment
2600 from 1996-2005, because data on these outcomes are unavailable prior to 2005. The
2601 availability of each dataset was discussed in the published version of the paper (see,
2602 for example, Table 1). Perhaps a more accurate statement would be to explain
2603 that the replicators couldn’t explore the impact of passing a referendum on these
2604 specific outcomes in the early period due to data constraints—and that this was
2605 acknowledged in the published version.”

2606 **Original Authors’ Package:** [https://www.openicpsr.org/openicpsr/project/
2607 125821/version/V1/view](https://www.openicpsr.org/openicpsr/project/125821/version/V1/view)

2608 **Reproduction Report**

2609 **Title Original Study:** Social Class and (Un)Ethical Behaviour: Causal and
2610 Correlational Evidence

2611 **doi:** <https://doi.org/10.1093/ej/ueac022>, Economic Journal

2612 **Report's Abstract:** The relationship between social status and ethical behav-
2613 ior is a widely debated topic in research. In their study, Gsottbauer et al. (2022b)
2614 investigate whether higher socio-economic status is linked to lower ethical behavior,
2615 using data from two large survey experiments involving over 11,000 participants.
2616 In this replication project, we test the computational reproducibility and robust-
2617 ness to the replication of their study, using the provided data and code from the
2618 replication package (Gsottbauer et al., 2022a). Nearly all the figures and tables
2619 were reproducible-in the process of reproducing the results, some minor rounding or
2620 transcription errors were discovered. In testing the robustness replicability, we find
2621 consistent results for our extensions. The effort for the replication was manageable,
2622 even though the authors treat categorical variables as numeric, or use manually-
2623 coded interaction variables (i.e., in regression models). In summary, we applaud
2624 the transparency of Gsottbauer et al. (2022b) in facilitating replications, and make
2625 some general recommendations for further improvements for data-analysis studies.

2626 **Link to Full Report:** <https://econpapers.repec.org/paper/zbwi4rdps/29.htm>

2627 **Link to Replicators' Package:** [https://github.com/ha0ye/](https://github.com/ha0ye/replication-gsottbauer-2022)
2628 [replication-gsottbauer-2022](https://github.com/ha0ye/replication-gsottbauer-2022)

2629 **Original Authors' Response:** Declined to respond.

2630 **Original Authors' Package:** <https://zenodo.org/records/6226207>

2631 **Reproduction Report**

2632 **Title Original Study:** Sorting or Steering: The Effects of Housing Discrimination
2633 on Neighborhood Choice

2634 **doi:** <https://doi.org/10.1086/720140>, Journal of Political Economy

2635 **Report's Abstract:** This comment revisits the analysis in Christensen and Tim-
2636 mins (2022). We identify two critical errors used in the original analysis, one with
2637 the data and the other with coding. When either error is corrected several major
2638 results in the paper change, either in statistical significance or in effect size. The
2639 data error is a result of including fixed effects for the string variable 'city'. The raw
2640 variable is case sensitive and has many spelling mistakes. The coding error involves
2641 assigning a value of zero for the variable "of color" to both individuals identified as
2642 'white' and as 'other' in the raw data. The level of clustering in the paper is also
2643 arguably too fine. Many of the results are not robust to clustering at the city level,
2644 as opposed to the subject pair level. In total, we affirm the authors' overarching
2645 claim of substantial and nuanced housing discrimination against racial minorities
2646 generally, and African Americans in particular; however, the effect sizes and sig-
2647 nificance are generally (although not always) smaller than the original authors
2648 findings. Additionally, there are several instances where the effects of discrimina-
2649 tion on African Americans are no longer statistically significant but the effect of
2650 discrimination on Hispanics becomes significant.

2651 **Link to Full Report:** <https://osf.io/vwgxd/>

2652 **Link to Replicators' Package:** <https://github.com/mattwebb/HUDreplication>

2653 **Original Authors' Response:** Authors mentioned that they are currently writing
2654 a response.

2655 **Original Authors' Package:** https://www.journals.uchicago.edu/doi/suppl/10.1086/720140/suppl_file/20191181data.zip
2656

2657 **Reproduction Report**

2658 **Title Original Study:** Spillover Effects of Intellectual Property Protection in the
2659 Interwar Aircraft Industry

2660 **doi:** <https://doi.org/10.1093/ej/ueab091>, Economic Journal

2661 **Report's Abstract:** We are attempting to reproduce the results of Hanlon and
2662 Jaworski (2022) based on their dataset. Our work is conducted in two different ways:
2663 (i) computational reproducibility, aiming to produce the same results using different
2664 software, notably R, with the given data; and (ii) checking the robustness of the
2665 results. For (i), the estimated coefficients are consistent based on the R software.
2666 For (ii), we carefully examine the given datasets of Hanlon and Jaworski (2022)
2667 and review the economic history of the US Interwar aircraft industry between 1918
2668 and 1935 to identify potential confounding variables (apart from IPP strengthening
2669 in the year 1926) that might affect both the controls and error term, and thus the
2670 results. We identify some confounding variables that may affect the results and
2671 attempt to illustrate them before and after 1926 when IPP is strengthened. Overall,
2672 we find that the results are replicable by utilizing the codes and datasets of Hanlon
2673 and Jaworski (2022), which is encouraging.

2674 **Link to Full Report:** <https://osf.io/t4avf/>

2675 **Link to Replicators' Package:** <https://osf.io/t4avf/>

2676 **Link to Original Authors' Response:** <https://osf.io/t4avf/>

2677 **Original Authors' Package:** <https://zenodo.org/records/5627298>

2678 **Reproduction Report**

2679 **Title Original Study:** State Action to Prevent Violence against Women: The
2680 Effect of Women’s Police Stations on Men’s Attitudes toward Gender-Based
2681 Violence

2682 **doi:** <https://doi.org/10.1086/714931>, Journal of Politics

2683 **Report’s Abstract:** Córdoba and Kras (2022) examine how the existence of a
2684 women’s police station (WPS) in the place of residence influences citizens’ atti-
2685 tudes toward gender-based violence in Brazil. In their analytical specification, the
2686 authors find that men are more likely to reject violence against women (VAW)
2687 and support bystander intervention in municipalities with a WPS, especially if the
2688 WPS has been operating for a long time. This paper examines the replicability
2689 and robustness of Córdoba & Kras’ (2022) findings. First, we reproduce the paper’s
2690 main findings and uncover one minor coding error and three estimates that have
2691 been reported with the opposite sign compared to that in our reproduction; neither
2692 is of consequence for the study’s main results. Second, we test the robustness of
2693 the results by (1) recoding one of the main explanatory variables and several of the
2694 control variables to account for non-linear trends, (2) using alternative techniques
2695 to estimate clustered standard errors, (3) consistently applying a 95% confidence
2696 level in the presentation of the results, (4) altering the propensity score match-
2697 ing (PSM) procedure as well as the composition of the variables used in the PSM
2698 robustness check, (5) using an alternative technique to test for multicollinearity,
2699 (6) excluding potential endogenous control variables, and (7) using an alternative
2700 coding for computing margins. Reassuringly, the results are robust to most of these
2701 tests. However, two of the robustness checks challenge parts of the paper’s main
2702 findings. First, allowing for non-linearity in the effect of time since the establish-
2703 ment of WPS shows (a) a non-linear effect on VAW and (b) no apparent changes in
2704 either male or female attitudes over time once the WPS has been established. Sec-
2705 ond, the inclusion of other variables in the PSM procedure renders part of the main
2706 estimates of interest statistically nonsignificant ($p < 0.1$). Our findings highlight
2707 the need for further re-analyses and replications for investigating the preventive
2708 effects of women’s police stations.

2709 **Link to Full Report:** <https://econpapers.repec.org/paper/zbwi4rdps/67.htm>

2710 **Link to Replicators’ Package:** <https://osf.io/yjwr8/>

2711 **Link to Original Authors’ Response:** Responded to our emails but no formal
2712 response as of February 2024.

2713 **Original Authors’ Package:** [https://dataverse.harvard.edu/dataset.xhtml?](https://dataverse.harvard.edu/dataset.xhtml?persistentId=doi:10.7910/DVN/D2WL5I)
2714 [persistentId=doi:10.7910/DVN/D2WL5I](https://dataverse.harvard.edu/dataset.xhtml?persistentId=doi:10.7910/DVN/D2WL5I)

2715 **Reproduction Report**

2716 **Title Original Study:** Student Performance, Peer Effects, and Friend Networks:
2717 Evidence from a Randomized Peer Intervention

2718 **doi:** <https://doi.org/10.1257/pol.20200563>, American Economic Journal: Economic
2719 Policy

2720 **Report's Abstract:** Wu et al. (2023) estimate the effect of classroom seat-
2721 ing arrangements in China using a randomized control trial with two treatment
2722 schemes. The first treatment scheme involves seating high and low achieving stu-
2723 dents together, and the second treatment involves this same seating arrangement
2724 with financial incentives for the high-achieving students, if their deskmates' test
2725 scores improved. All statistically significant impacts come from the incentivized
2726 treatment scheme. Wu et al. (2023) find that low-achieving students sitting next
2727 to incentivized high-achieving students perform 0.24 SD (p-value=0.018) better
2728 on math exams. In addition, being assigned to the incentive treatment scheme
2729 increased extraversion and agreeableness for low and high achieving students.
2730 Lastly, they do not find much evidence of peer effects on test scores nor personality
2731 traits. This study is computationally reproducible using their provided replication
2732 package. We ran their code using Stata 14, 17, and 18. After running their replica-
2733 tion package, we further investigated Tables 2-5. The main conclusions are generally
2734 robust to various coding decisions. Notably, in investigating the peer effects, when
2735 we change the specification to also control for the difference in baseline scores
2736 between the student and their deskmate, we find that the more dissimilar deskmates
2737 are at baseline, the bigger the peer effects.

2738 **Link to Full Report:** <https://osf.io/9hx3b/>

2739 **Original Authors' Response:** The authors provided feedback which was taken
2740 into account.

2741 **Original Authors' Package:** [https://www.openicpsr.org/openicpsr/project/
2742 149262/version/V2/view](https://www.openicpsr.org/openicpsr/project/149262/version/V2/view)

2743 **Reproduction Report**

2744 **Title Original Study:** Talking Shops: The Effects of Caucus Discussion on Policy
2745 Coalitions

2746 **doi:** <https://doi.org/10.1111/ajps.12636>, American Journal of Political Science

2747 **Report’s Abstract:** In Talking Shops: The Effects of Caucus Discussion on Policy
2748 Coalitions, Zelizer analyzes the causal effect of caucus deliberations on legislative
2749 policy coalitions. In practice, political scientists have little empirical evidence on
2750 how policy discussions actually work among sitting legislators and whether these
2751 discussions have an effect on policy making and policy opinion. Taking on this chal-
2752 lenge, Zelizer conducted two field experiments in an American state legislature. In
2753 short, the experiments randomized whether a bill was selected for discussion among
2754 a bi-partisan legislative caucus. The paper then measures and reports the corre-
2755 sponding effects of that discussion around the bill. Zelizer finds that deliberation
2756 increased the amount of co-sponsorship for a given bill, among both co-partisans
2757 and counter-partisans, but deliberation did not effect whether a bill was passed
2758 by the legislature or whether the bill received more amendments. We conduct a
2759 robustness replication of the main results of Talking Shops. Specifically, we repro-
2760 duce Tables 3 and 4 of the paper under alternative specifications. We find that
2761 the main results of the paper are reproducible and robust to multiple alternative
2762 specifications.

2763 **Link to Full Report:** <https://econpapers.repec.org/paper/zbwi4rdps/69.htm>

2764 **Link to Replicators’ Package:** <https://osf.io/tmfyj/>

2765 **Link to Original Authors’ Response:** “One purpose of replication, among oth-
2766 ers, is to evaluate whether published results are sensitive to modeling decisions.
2767 Do alternative, reasonable approaches generate the same, or different, results? Did
2768 the author’s approach provide an outlier estimate that is indicative of p-hacking
2769 or, to be kinder about it, sensitivity of results to modeling decisions? That seems
2770 incredibly useful. That purpose is not advanced, in my view, by testing ‘incorrect’
2771 methods or models. We do not learn about the robustness of results from testing
2772 alternative approaches that introduce bias, or by estimating different estimands
2773 that are a combination of treatment effects and selection bias. While it doesn’t
2774 seem to matter too much in this case — selection bias appears relatively small,
2775 and in the same direction as treatment effects — I think this issue matters for the
2776 exercise in general for several reasons. First, do the analyses justify the inferences
2777 being made? In my view, changing the estimand or estimating biased models can-
2778 not justify saying anything about the robustness of the original results. Second,
2779 what would have happened if the new results did not match the original? Are we
2780 willing to claim published results are not robust when applying estimators with
2781 known flaws generates different results? And third, shouldn’t we just generally aim
2782 to use ‘correct’ estimators for a given situation? While IPW is not perfect, ignoring
2783 differential treatment probabilities is a conscious decision to ignore selection bias.
2784 Why would we want to run that model if our goal is inference about treatment
2785 effects? I appreciate the work everyone is doing on this enterprise. Hopefully these
2786 comments, whether correct or not, help advance the goal of publishing robust, valid
2787 empirical research.”

2788

Original Authors' Package: [https://dataverse.harvard.edu/dataset.xhtml?
persistentId=doi:10.7910/DVN/S3M5AX](https://dataverse.harvard.edu/dataset.xhtml?persistentId=doi:10.7910/DVN/S3M5AX)

2789

2790 **Reproduction Report**

2791 **Title Original Study:** Targeting High Ability Entrepreneurs Using Community
2792 Information: Mechanism Design in the Field

2793 **doi:** <https://doi.org/10.1257/aer.20200751>, American Economic Review

2794 **Report's Abstract:** Hussam et al. (2022a) use a cash grant experiment in India
2795 to demonstrate that community knowledge can help target high-growth microen-
2796 trepreneurs. In their preferred specification, the authors find that the average
2797 marginal return to the grant is 9.4 percent per month, while estimated returns
2798 for entrepreneurs reported by peers to be in the top third of the community are
2799 between 24 percent and 30 percent. First, we reproduce the paper's main findings
2800 and uncover one minor coding error, which affects the estimates for one of the main
2801 tables but does not change the overall conclusions of the paper. Second, we test
2802 the robustness of the results to: (1) different treatment of outliers, (2) dropping
2803 surveyor and survey month fixed effects, and (3) using quartiles instead of terciles
2804 for grouping the ranking of entrepreneurs. The paper's results are robust to these
2805 robustness checks. Finally, we test heterogeneity of results by gender, which was
2806 not reported in the original study.

2807 **Link to Full Report:** <https://econpapers.repec.org/paper/zbwi4rdps/49.htm>

2808 **Link to Replicators' Package:** [https://dataverse.harvard.edu/dataset.xhtml?](https://dataverse.harvard.edu/dataset.xhtml?persistentId=doi:10.7910/DVN/DI7RR9)
2809 [persistentId=doi:10.7910/DVN/DI7RR9](https://dataverse.harvard.edu/dataset.xhtml?persistentId=doi:10.7910/DVN/DI7RR9)

2810 **Link to Original Authors' Response:** "We are very grateful to Isabella Masetto,
2811 Diego Ubfal, and to the team at I4R for their excellent work. We verified the coding
2812 error and we agree that it did not meaningfully alter the conclusion of our paper
2813 that community information is informative over and above the predictive power of
2814 observable characteristics. We will post a link to this correction on our websites
2815 and will consult the editors of the AER as to whether this error rises to the level
2816 of requiring a formal correction."

2817 **Original Authors' Package:** [https://www.openicpsr.org/openicpsr/project/](https://www.openicpsr.org/openicpsr/project/151841/version/V1/view)
2818 [151841/version/V1/view](https://www.openicpsr.org/openicpsr/project/151841/version/V1/view)

2819 **Reproduction Report**

2820 **Title Original Study:** Teaching Norms: Direct Evidence of Parental Transmission
2821 **doi:** <https://doi.org/10.1093/ej/ueac074>, Economic Journal

2822 **Report’s Abstract:** This paper is a replication study of Brouwer, T., Galeotti,
2823 F., & Villeval, M. C. (2023), using the original data. The study explores how social
2824 norms are transmitted from one generation to another, specifically from parents to
2825 children. The authors conducted a field experiment involving 601 parents of children
2826 aged 3 to 12 in Lyon, France, to examine whether parents engage more in norm
2827 enforcement in the presence of their child, and whether the nature of punishment
2828 changes in the presence of the child. The study found that parents do engage more
2829 in norm enforcement in the presence of their child, and tend to use more indirect
2830 punishment when their child is present. This study highlights the role that parents
2831 play in transmitting social norms to their children. The replication analysis was
2832 successful, with the results of the original study being robust to changes in the
2833 model specification.

2834 **Link to Full Report:** <https://osf.io/qnbfa/>

2835 **Link to Replicators’ Package:** <https://zenodo.org/records/8114738>

2836 **Original Authors’ Response:** The replicators took into account the authors’
2837 feedback. They wrote at the end of the back and forth: “We thank you and the
2838 replication team for the replication and the successive interactions. We created an
2839 OSF project including the data replication package enabling the reproduction of
2840 the analysis presented in our article. The package comprises a source file (in Stata
2841 format and in TXT) and a Stata do-file that allows the reconstruction of the master
2842 file used in the replication package submitted to the Economic Journal.”

2843 **Original Authors’ Package:** <https://zenodo.org/records/7045559>

2844 **Reproduction Report**

2845 **Title Original Study:** Technological Change and the Consequences of Job Loss
2846 **doi:** <https://doi.org/110.1257/aer.20210182>, American Economic Review

2847 **Report's Abstract:** Braxton and Taska (2023) find that technological change
2848 accounts for 45 percent of the decline in earnings after job loss. We first reproduce
2849 all regression tables in Braxton and Taska (2023), and then test for robustness by
2850 controlling for the initial level of wages, additional fixed effects, multi-way cluster-
2851 ing, and conducting influential analysis. We find that the paper's original results are
2852 sensitive to controlling for initial wages and some additional fixed effects. Overall,
2853 we find the results are robust with a coefficient in the same direction and signifi-
2854 cant at 5% in 33% of the robustness checks we ran, with average t/z scores 28% as
2855 large as the original study.

2856 **Link to Full Report:** <https://osf.io/qws2p/>

2857 **Link to Replicators' Package:** <https://osf.io/qws2p/>

2858 **Original Authors' Response:** Did not get a response as of November 2025.

2859 **Original Authors' Package:** [https://www.openicpsr.org/openicpsr/project/
2860 181166/version/V1/view](https://www.openicpsr.org/openicpsr/project/181166/version/V1/view)

2861 **Reproduction Report**

2862 **Title Original Study:** The Common-Probability Auction Puzzle

2863 **doi:** <https://doi.org/10.1257/aer.20191927>, American Economic Review

2864 **Report's Abstract:** Ngangoué and Schotter (2023) investigate common-
2865 probability auctions. By running an experiment, they find that, in contrast to the
2866 substantial overbidding found in common-value auctions, bidding in strategically
2867 equivalent common-probability auctions is consistent with the Nash equilibrium.
2868 We reproduce their results in R, conduct robustness checks on how their sample
2869 was constructed, and consider possible heterogeneity. We confirm their documented
2870 qualitative results.

2871 **Link to Full Report:** <https://econpapers.repec.org/paper/zbwi4rdps/74.htm>

2872 **Link to Replicators' Package:** <https://osf.io/7bq4s/>

2873 **Original Authors' Response:** "Thank you for putting the effort in replicating
2874 our study! Your results are also quite interesting to us as we haven't thought of
2875 all the robustness checks you've made. At this point, we do not have any major
2876 comments to make and refrain from submitting a response."

2877 **Original Authors' Package:** [https://www.openicpsr.org/openicpsr/project/
2878 184041/version/V1/view](https://www.openicpsr.org/openicpsr/project/184041/version/V1/view)

2879 **Reproduction Report**

2880 **Title Original Study:** The Curious Case of Theresa May and the Public That
2881 Did Not Rally: Gendered Reactions to Terrorist Attacks Can Cause Slumps Not
2882 Bumps

2883 **doi:** <https://doi.org/10.1017/S0003055421000861>, American Political Science
2884 Review

2885 **Report's Abstract:** Holman et al. (2022; HMZ) propose women (compared to
2886 men) political leaders experience significant drops in public approval ratings after
2887 a transnational terrorist attack. After documenting how survey-based evaluations
2888 of then-Prime Minister Theresa May suffered after the 2017 Manchester Arena
2889 attack, HMZ assemble a country-quarter level panel database to explore the generality
2890 of their hypothesis. They report evidence suggesting women (compared to
2891 men) leaders systematically experience decreased public approval rates after major
2892 transnational terrorist attacks (p-value of 0.020). We find that result disappears
2893 once any of the following adjustments is implemented: (i) excluding election quarter
2894 covariates (p = 0.104); (ii) correcting objective coding errors in the election quarter
2895 covariates (p = 0.058); (iii) excluding the May-Manchester observation (p = 0.098);
2896 or (iv) clustering standard errors at the country level (p = 0.558). Exploring all 2⁵
2897 combinations of the five control groups HMZ incorporate in their specification, none
2898 of them clears the 5% threshold of statistical significance once the corrected elec-
2899 tion quarter variables are employed. We conclude that the empirical evidence does
2900 not provide sufficient support for HMZ's abstract claim that "conventional theory
2901 on rally events requires revision: women leaders cannot count on rallies following
2902 major terrorist attacks."

2903 **Link to Full Report:** <https://econpapers.repec.org/paper/zbwi4rdps/41.htm>

2904 **Link to Replicators' Package:** <https://doi.org/10.5683/SP3/6SYCML>

2905 **Link to Original Authors' Response:** <https://econpapers.repec.org/paper/zbwi4rdps/44.htm>

2907 **Original Authors' Package:** <https://dataverse.harvard.edu/dataset.xhtml?persistentId=doi:10.7910/DVN/VHNPUO>
2908

2909 **Reproduction Report**

2910 **Title Original Study:** The Dynamics and Spillovers of Management Interventions: Evidence from the Training within Industry Program

2911 **doi:** <https://doi.org/10.1086/719277>, Journal of Political Economy

2912 **Report's Abstract:** Bianchi and Giorcelli (2022) study the long-term and spillover
2913 effects of a management intervention program on firm performance in the US,
2914 between 1940 and 1945. The authors find that the Training Within Industry (TWI)
2915 program led to positive effects which lasted for at least 10 years. Firm sales of
2916 treated firms increased by 5.3% in the first year after implementation, peaking at
2917 21.7% after 8 years, before reducing to 16% gains after a decade. The authors claim
2918 that the program generated long-lasting changes in managerial practices. Finally,
2919 the program also led to positive spillover effects on the supply chain of treated
2920 firms. First, we reproduce the paper's main findings. Second, we test the robustness
2921 of the results to (1) changing the main specification sample and (2) testing other
2922 difference-in-differences estimators, using the same data, provided by the authors.
2923 We find that the results are robust to these changes. All point estimates in the
2924 study remain statistically significant and of similar magnitude. While the paper's
2925 finding reproduce and replicate, challenges in reproducing results we encountered
2926 lead us to recommend improvements to journals' code policies.

2927 **Link to Full Report:** <https://econpapers.repec.org/paper/zbwi4rdps/66.htm>

2928 **Link to Replicators' Package:** https://github.com/cwestheide/i4r_dp66_code

2929 **Original Authors' Final Response:** "Thanks a lot for sharing the updated
2930 report with us. We don't have anything to add at this point."

2931 **Original Authors' Package:** [https://www.journals.uchicago.edu/doi/suppl/10.
2932 1086/719277/suppl_file/20200781data.zip](https://www.journals.uchicago.edu/doi/suppl/10.1086/719277/suppl_file/20200781data.zip)
2933

2934 **Reproduction Report**

2935 **Title Original Study:** The Economic Effects of Long-Term Climate Change:
2936 Evidence from the Little Ice Age

2937 **doi:** <https://doi.org/10.1086/720393>, Journal of Political Economy

2938 **Report's Abstract:** Waldinger (2022) finds significant negative economic effects
2939 (proxied by city size) from gradual climate change which occurred during the Little
2940 Ice Age (1600-1850) and offers two potential mechanisms (agricultural productivity
2941 and mortality) and two potential adaptations (trade and land use). In this comment,
2942 we show that while Waldinger (2022)'s findings can be replicated, the main result
2943 relies on a critical author assumption: Cities with 0 inhabitants in the original data
2944 are instead assumed to have 500. This assumption affects 23.5% of observations and
2945 49.6% of cities in the sample. When these "missing data" are excluded from the
2946 analysis, the effect estimated by otherwise identical research methods is of similar
2947 magnitude and statistical significance but of opposite sign.

2948 **Link to Full Report:** <https://osf.io/tmn2j/>

2949 **Link to Replicators' Package:** <https://osf.io/tmn2j/>

2950 **Link to Original Authors' Response:** <https://osf.io/tmn2j/>

2951 **Original Authors' Package:** https://www.journals.uchicago.edu/doi/suppl/10.1086/720393/suppl_file/2015548data.zip
2952

2953 **Reproduction Report**

2954 **Title Original Study:** The Effects of Banking Competition on Growth and
2955 Financial Stability: Evidence from the National Banking Era

2956 **doi:** <https://doi.org/10.1086/717453>, Journal of Political Economy

2957 **Report’s Abstract:** Carlson et al. (2022) examine the causal impact of banking
2958 competition by investigating a unique circumstance in the National Banking Era
2959 of the nineteenth century in the US, where a discontinuity in bank capital require-
2960 ments occurred. On the one hand, their findings suggest that banks operating in
2961 markets with fewer barriers to entry tend to increase their lending activities, pro-
2962 moting real economic growth. On the other hand, banks in less restricted markets
2963 also exhibit a higher propensity for risk-taking, posing risks to financial stability.
2964 First, we fully reproduce the paper’s outcomes apart from a minor discrepancy in
2965 the estimate of Table 9 attributed to issues in the provided codes. Second, we test
2966 the robustness of the results by (i) changing the ranges used to select the sample
2967 of cities included in the analysis, (ii) adopting different options to address outliers’
2968 potential issues and (iii) introducing additional control variables. We observe that
2969 the estimation results remain mostly consistent when subjecting them to various
2970 robustness checks. However, it is worth highlighting that the results can be par-
2971 tially influenced by the criteria used to select the sample of cities and the inclusion
2972 of control variables.

2973 **Link to Full Report:** <https://econpapers.repec.org/paper/zbwi4rdps/81.htm>

2974 **Link to Replicators’ Package:** [https://dataverse.harvard.edu/dataset.xhtml?](https://dataverse.harvard.edu/dataset.xhtml?persistentId=doi:10.7910/DVN/BB864R)
2975 [persistentId=doi:10.7910/DVN/BB864R](https://dataverse.harvard.edu/dataset.xhtml?persistentId=doi:10.7910/DVN/BB864R)

2976 **Original Authors’ Final Response:** “We thank the replication team (Andrea
2977 Calef, Sya In Chzhen, Marco Mandas, and Fabio Motoki) for the detailed Reproduc-
2978 tion Report. We are glad to hear that the replicating team affirms the robustness of
2979 the paper’s findings. We are also glad that the replicators were able to successfully
2980 replicate all tables and figures. We thank the replicating team for identifying vari-
2981 ous smaller issues regarding the code structure which fortunately did not affect our
2982 original findings. We believe that the report as such does not require us to respond
2983 in any further detail. We highly appreciate the effort of both the replicating team
2984 and the I4R.”

2985 **Original Authors’ Package:** [https://www.journals.uchicago.edu/doi/suppl/10.](https://www.journals.uchicago.edu/doi/suppl/10.1086/717453/suppl_file/20200610data.zip)
2986 [1086/717453/suppl_file/20200610data.zip](https://www.journals.uchicago.edu/doi/suppl/10.1086/717453/suppl_file/20200610data.zip)

2987 **Reproduction Report**

2988 **Title Original Study:** The Geography of Repression and Opposition to Autocracy
2989 **doi:** <https://doi.org/10.1111/ajps.12614>, American Journal of Political Science

2990 **Report's Abstract:** Analytic data sets and analysis code are available and they
2991 produce the same results as presented in the paper (CRA). Robustness checks
2992 involve the (i) use of matching estimators to address possible bias from misspec-
2993 ification, based on propensity score estimated from a random forest model, (ii)
2994 doubly robust (TMLE) estimation to address possible bias from misspecification
2995 in either the propensity score or outcome regression stages, using a super learner
2996 ensemble with random forest, GAM, mean, and non-parametric regression models
2997 and averaged over repeated runs to minimize randomness, (iii) define treated comu-
2998 nas as those within a fixed physical distance radius of the nearest military base,
2999 rather than only those that contain it, and (iv) instead of using 2SLS to assess the
3000 causally mediated effect of military bases on plebiscite outcomes via repression, we
3001 propose to conduct mediation analysis (Tingley et al 2013), implemented in the R
3002 'mediation' package.

3003 **Link to Full Report:** [https://www.socialsciencereproduction.org/reproductions/](https://www.socialsciencereproduction.org/reproductions/789/published/index?step=4)
3004 [789/published/index?step=4](https://www.socialsciencereproduction.org/reproductions/789/published/index?step=4)

3005 **Link to Replicators' Package:** [https://github.com/pjesscarter/](https://github.com/pjesscarter/repression-replication)
3006 [repression-replication](https://github.com/pjesscarter/repression-replication)

3007 **Link to Original Authors' Response:** We are happy that the replicators suc-
3008 cessfully reproduced all the analysis in our published paper. Moreover, additional
3009 robustness checks within the quantitative framework of the paper further confirm
3010 the empirical results. Two extensions using propensity score matching give some-
3011 what different results. Unfortunately, these additional estimators violate standard
3012 requirements for credible matching designs, i.e., overlap in the propensity score dis-
3013 tribution across treatment and control groups. As shown by previous research, this
3014 lack of overlap leads to unstable estimators with variance that may explode in finite
3015 samples such as ours (Frölich 2004, Khan and Tamer 2010). In another extension,
3016 the replicators employ a mediation analysis to re-interpret the empirical evidence
3017 in our paper. To support the use of our method, i.e., instrumental variables, we
3018 rule out alternative explanations and provide a range of historical evidence. With-
3019 out historical and contextual support for alternative assumptions, we believe that
3020 the method used by the replicators is hard to interpret.

3021 **Original Authors' Package:** [https://dataverse.harvard.edu/dataset.xhtml?](https://dataverse.harvard.edu/dataset.xhtml?persistentId=doi:10.7910/DVN/EYAWES)
3022 [persistentId=doi:10.7910/DVN/EYAWES](https://dataverse.harvard.edu/dataset.xhtml?persistentId=doi:10.7910/DVN/EYAWES)

3023 **Reproduction Report**

3024 **Title Original Study:** The Labor Market Impacts of Universal and Permanent
3025 Cash Transfers: Evidence from the Alaska Permanent Fund

3026 **doi:** <https://doi.org/10.1257/pol.20190299>, American Economic Journal: Economic
3027 Policy

3028 **Report's Abstract:** Jones and Marinescu (2022) study the employment effects
3029 of a universal cash transfer in Alaska. Using a synthetic control method, they find
3030 that the transfer had no negative effects on employment. We reproduce the results
3031 using their replication package and investigate if the results hold when using a
3032 different software to run the analysis. We also use different estimation techniques
3033 and perform sensitivity checks to assess robustness of the results. We find some
3034 differences in the size and significance of the average treatment effects on labor force
3035 participation and hours worked when we use a different software (R) and various
3036 extensions of the synthetic control method. We also find smaller coefficients on
3037 part-time employment when including more covariates. However, these differences
3038 do not contradict the main conclusion of the paper.

3039 **Link to Full Report:** <https://econpapers.repec.org/paper/zbwi4rdps/80.htm>

3040 **Link to Replicators' Package:** <https://osf.io/6atfw/>

3041 **Original Authors' Final Response:** "Thanks for putting in all this effort to
3042 evaluate the robustness of our results! I [Marinescu] think this is really a worthwhile
3043 endeavor."

3044 **Original Authors' Package:** [https://www.openicpsr.org/openicpsr/project/
3045 140121/version/V1/view](https://www.openicpsr.org/openicpsr/project/140121/version/V1/view)

3046 **Reproduction Report**

3047 **Title Original Study:** The Long-Run Effects of Sports Club Vouchers for Primary
3048 School Children

3049 **doi:** <https://doi.org/10.1257/pol.20200431>, American Economic Journal: Economic
3050 Policy

3051 **Report’s Abstract:** Marcus, Siedler and Ziebarth (2022 American Economic
3052 Journal: Economic Policy) examine the long-run health effects of a universal sports-
3053 club voucher program that was introduced in Saxony for primary school children
3054 in 2009. In 2018, the authors designed a survey that targeted the affected cohorts
3055 and nearby cohorts in Saxony and two neighboring states, and use a differences-in-
3056 differences identification strategy that exploits variation across states and cohorts
3057 in policy exposure. The authors document that treated individuals have knowledge
3058 of the program and recall receiving and redeeming the vouchers at higher rates,
3059 but find no effects on any health outcomes or behaviors. We successfully reproduce
3060 the main results of the paper exactly using data available in the paper’s replication
3061 package and new Stata and R code. We also verify the robustness of the results
3062 using different outcomes, different control variables, different sample restrictions
3063 and different inference methods.

3064 **Link to Full Report:** <https://econpapers.repec.org/paper/zbwi4rdps/46.htm>

3065 **Link to Replicators’ Package:** <https://osf.io/4bnjt/>

3066 **Original Authors’ Response:** “We would like to thank the authors for their
3067 interest in our paper. We greatly appreciate their careful reading of the paper
3068 and the insightful robustness exercises they conducted. We are pleased that our
3069 results were successfully reproduced using different software packages, and that the
3070 additional robustness analyses performed by the authors further strengthen and
3071 support our conclusions.”

3072 **Original Authors’ Package:** [https://www.openicpsr.org/openicpsr/project/
3073 138922/version/V1/view](https://www.openicpsr.org/openicpsr/project/138922/version/V1/view)

3074 **Reproduction Report**

3075 **Title Original Study:** The Long-Term Effects of Measles Vaccination on Earnings
3076 and Employment

3077 **doi:** <https://doi.org/10.1257/pol.20190509>, American Economic Journal: Economic
3078 Policy

3079 **Report's Abstract:** Atwood (2022) analyzes the effects of the 1963 U.S. measles
3080 vaccination on longrun labor market outcomes, using a generalized difference-in-
3081 differences approach. We reproduce the results of this paper and perform a battery
3082 of robustness checks. Overall, we confirm that the measles vaccination had positive
3083 labor market effects. While the negative effect on the likelihood of living in poverty
3084 and the positive effect on the probability of being employed are very robust across
3085 the different specifications, the headline estimate-the effect on earnings-is more
3086 sensitive to the exclusion of certain regions and survey years.

3087 **Link to Full Report:** <https://econpapers.repec.org/paper/zbwi4rdps/33.htm>

3088 **Link to Replicators' Package:** <https://osf.io/jv7kx/>

3089 **Link to Original Authors' Response:** <https://osf.io/qxjnk/>

3090 **Original Authors' Package:** [https://www.openicpsr.org/openicpsr/project/
3091 138401/version/V1/view](https://www.openicpsr.org/openicpsr/project/138401/version/V1/view)

3092 **Reproduction Report**

3093 **Title Original Study:** The Macroeconomics of Sticky Prices with Generalized
3094 Hazard Functions

3095 **doi:** <https://doi.org/10.1093/qje/qjab042>, Quarterly Journal of Economics

3096 **Report's Abstract:** We replicate the empirical results in Section 4 of Alvarez et
3097 al. (2022). First, we were able to reproduce the original authors' major empirical
3098 results, but only after editing the program for it to run on our computing platform.
3099 There are small discrepancies in the empirical estimates, e.g. bootstrapped standard
3100 errors, that involve the use of simulations. Second, we replicated Alvarez et al.'s
3101 results by adopting the data cleaning criteria used by their original data source
3102 (Cavallo 2018) to evaluate its robustness to data handling procedures. We found
3103 noticeable changes in the empirical results that can have important implications on
3104 the effects of monetary policy. To conclude, we propose using Docker container to
3105 promote research reproducibility, and more attention is needed on data handling
3106 to improve the robustness of empirical research.

3107 **Link to Full Report:** [https://github.com/atyho/
3108 Ottawa-Replication-Games-2023/blob/main/Ho_Huynh_Rea_Replication_Report.
3109 pdf](https://github.com/atyho/Ottawa-Replication-Games-2023/blob/main/Ho_Huynh_Rea_Replication_Report.pdf)

3110 **Link to Replicators' Package:** [https://github.com/atyho/
3111 Ottawa-Replication-Games-2023/](https://github.com/atyho/Ottawa-Replication-Games-2023/)

3112 **Link to Original Authors' Response:**

3113 **Original Authors' Package:** [https://dataverse.harvard.edu/dataset.xhtml?
3114 persistentId=doi:10.7910/DVN/IBM0IE](https://dataverse.harvard.edu/dataset.xhtml?persistentId=doi:10.7910/DVN/IBM0IE)

3115 **Reproduction Report**

3116 **Title Original Study:** The Morning After: Cabinet Instability and the Purging
3117 of Ministers after Failed Coup Attempts in Autocracies

3118 **doi:** <https://doi.org/10.1086/716952>, Journal of Politics

3119 **Report's Abstract:** We replicate the analysis provided in Bokobza et al. (2022).
3120 They identify a causal effect of failed coup attempts on cabinet minister removals
3121 in autocracies on both the country and individual minister level and show that
3122 higher-ranking ministers and those holding strategic positions are more likely to
3123 be purged than more loyal and veteran ministers using fixed effects panel models.
3124 We focus on computational reproducibility and robustness replicability. In addition
3125 to reproducing the original results using Stata and R, we replicate analyses
3126 using random effects panel models and ordered beta regression models, reproduced
3127 analyses performed in R using different packages, replaced the main independent
3128 variable, clustered standard errors on a different level, and added independent variables
3129 related to coup-proofing. We find that the original findings were reproducible
3130 and robust.

3131 **Link to Full Report:** <https://econpapers.repec.org/paper/zbwi4rdps/45.htm>

3132 **Link to Replicators' Package:** <https://doi.org/10.7910/DVN/21HZCC>

3133 **Link to Original Authors' Response:** <https://osf.io/sm526/>

3134 **Original Authors' Package:** [https://dataverse.harvard.edu/dataset.xhtml?
3135 persistentId=doi:10.7910/DVN/GCDJ25](https://dataverse.harvard.edu/dataset.xhtml?persistentId=doi:10.7910/DVN/GCDJ25)

3136 **Reproduction Report**

3137 **Title Original Study:** The Origin of the State: Land Productivity or Appropri-
3138 ability?

3139 **doi:** <https://doi.org/10.1086/718372>, Journal of Political Economy

3140 **Report's Abstract:** This is a replication of Maysnar et al. (2022) (MPP). The
3141 article posits that the state (defined as societal hierarchy such as tax-levying elites)
3142 originated from cultivation of appropriable cereal grains, contrary to the conven-
3143 tional theory that the state originated from increased land productivity following
3144 the adoption of agriculture. The article uses multiple datasets to demonstrate a
3145 causal effect of cereal cultivation on hierarchy (Claim 1) without finding a similar
3146 effect for land productivity (Claim 2), and that societies based on roots or tubers
3147 display levels of hierarchy similar to nonfarming societies (Claim 3). The results of
3148 our replication in brief are: 1. The data and code provided by MMP closely repro-
3149 duce the main results presented in their Table 1 (see our Table 1). 2. Concurrently
3150 testing the cereal cultivation and land productivity claims leads to slightly less sta-
3151 tistical significance, on average, than the published article (Table 2). 3. Removing
3152 the inherited 1-5 scale of the dependent variable (hierarchy) finds that cereal pro-
3153 duction is not as effective at moving across all levels of hierarchy compared to the
3154 more general claim (Table 3 and 4). 4. Using the same procedures with an aim to
3155 confirm the conventional hypothesis (land productivity leads to increased hierarchy
3156 conditional on cereal growth) offers statistically significant evidence both for and
3157 against Claims 1 and 2 and against Claim 3 (Table 6). 5. The statistical significance
3158 of Claim 1 is sensitive to the removal of the top 3% of observations outliers (Table
3159 7). 6. Correction of mis-classified 'none or none specified' crop societies alters the
3160 interpretation of coefficients behind Claim 3. Societies that rely more on agricul-
3161 ture among farming societies experience more complex hierarchies, irrespective of
3162 being primarily cereal producing or tubers growing (Table 8 and 9). (...)

3163 **Link to Full Report:** <https://econpapers.repec.org/paper/zbwi4rdps/82.htm>

3164 **Link to Replicators' Package:** <https://osf.io/ekzdg/>

3165 **Original Authors' Response:** Comments taken into account in the report.

3166 **Original Authors' Package:** https://www.journals.uchicago.edu/doi/suppl/10.1086/718372/suppl_file/2018030data.zip
3167

3168 **Reproduction Report**

3169 **Title Original Study:** The Power of Hydroelectric Dams: Historical Evidence
3170 from the United States over the Twentieth Century

3171 **doi:** <https://doi.org/10.1093/ej/ueac059>, Economic Journal

3172 **Report's Abstract:** Successful computational reproducibility. No coding errors
3173 uncovered.

3174 **Original Authors' Package:** <https://doi.org/10.1093/ej/ueac059>

3175 **Reproduction Report**

3176 **Title Original Study:** The Relative Efficiency of Skilled Labor across Countries:
3177 Measurement and Interpretation

3178 **doi:** <https://doi.org/10.1257/aer.20191852>, American Economic Review

3179 **Report's Abstract:** Rossi (2022) examines the relative efficiency of skilled workers
3180 across countries. He finds the elasticity of skill efficiency with respect to GDP per
3181 worker is 1.4 and that the relative human capital accounts for only about 9 percent.
3182 We reproduce the paper's main findings and test the sensitivity of the results to (1)
3183 alternative samples and (2) additional controls for determining wages. We find the
3184 results remain robust to these alternative specifications, and the estimated values
3185 of the key elasticities remain nearly unchanged.

3186 **Link to Full Report:** <https://econpapers.repec.org/paper/zbwi4rdps/59.htm>

3187 **Link to Replicators' Package:** <https://osf.io/fge7z/>

3188 **Original Author's Response:** "Thanks for replicating the paper. I don't have
3189 any comments to add to the report."

3190 **Original Authors' Package:** [https://www.openicpsr.org/openicpsr/project/
3191 146041/version/V1/view](https://www.openicpsr.org/openicpsr/project/146041/version/V1/view)

3192 **Reproduction Report**

3193 **Title Original Study:** The Side Effects of Immunity: Malaria and African Slavery
3194 in the United States

3195 **doi:** <https://doi.org/10.1257/app.20190372>, American Economic Journal: Applied
3196 Economics

3197 **Report's Abstract:** Esposito (2022) documents the role of malaria in the diffusion
3198 of African slavery in the US. She finds that the introduction of malaria triggered a
3199 demand for malaria-resistant labour, which led to a massive expansion of African
3200 enslaved workers in more malaria-infested areas. Further results document that,
3201 among African slaves, more malaria-resistant individuals commanded significantly
3202 higher prices. We reproduce the paper's main findings, uncovering only one minor
3203 coding error that has no effect on the study's main results. We then test the robust-
3204 ness of the results to (1) varying the set of control variables used in various analyses;
3205 (2) conducting permutation tests; and (3) conducting event studies that account
3206 for time-varying controls. We generally find that the author's results are robust to
3207 all of these alternative specifications, though there are some sets of controls that
3208 cause estimates to become small and statistically insignificant.

3209 **Link to Full Report:** <https://osf.io/728ud/>

3210 **Link to Replicators' Package:** <https://osf.io/728ud/>

3211 **Original Authors' Response:** Original author provided feedback. No final
3212 response on the updated version.

3213 **Original Authors' Package:** [https://www.openicpsr.org/openicpsr/project/
3214 120483/version/V1/view](https://www.openicpsr.org/openicpsr/project/120483/version/V1/view)

3215 **Reproduction Report**

3216 **Title Original Study:** The Wheels of Change: Technology Adoption, Millwrights
3217 and the Persistence in Britain'S Industrialisation

3218 **doi:** <https://doi.org/10.1093/ej/ueab102>, Economic Journal

3219 **Report's Abstract:** Mokyr et al. (2022) estimate the effects of early technol-
3220 ogy adoption on industrialization. Authors argue that human capital was the main
3221 determinant of the location of the industry in the first decades of the Industrial Rev-
3222 olution. They document that the location of mills in the eleventh century (reported
3223 in the Domesday Book) has a positive and statistically significant impact on the
3224 number of wrights in the early eighteenth century. We confirm the computational
3225 reproducibility of the paper. The estimates are not sensitive to outliers, which are
3226 common in the data. The results are also robust to changes in the control variables.
3227 The results remain robust if we adjust the estimated p-values for the low number
3228 of clusters, and if we include county fixed effects. We conduct a placebo experi-
3229 ment with a present-day outcome (the Brexit referendum) to check if the results
3230 are picking up on a more general demographic and economic correlation pattern;
3231 the experiment shows no spurious correlations.

3232 **Link to Full Report:** <https://osf.io/gdne3/>

3233 **Link to Replicators' Package:** <https://osf.io/tws8n/>

3234 **Original Authors' Response:** No response.

3235 **Original Authors' Package:** <https://zenodo.org/records/5734954>

3236 **Reproduction Report**

3237 **Title Original Study:** Understanding Ethnolinguistic Differences: The Roles of
3238 Geography and Trade

3239 **doi:** <https://doi.org/10.1093/ej/ueab065>, Economic Journal

3240 **Report's Abstract:** Dickens (2022) studies the role of trade on long-run inter-
3241 ethnic linguistic differences. He establishes that neighboring ethnolinguistic groups
3242 have smaller (lexicostatistical) linguistic distances when there is a larger agricul-
3243 tural productivity variation between them. Specifically, he establishes that pre-1500
3244 land productivity variation (CSI SD) and its change due to Columbian Exchange in
3245 the post-1500 (CSI SD CHANGE) era decrease linguistic distances between groups.
3246 In what can be considered his main specification, which includes geographical con-
3247 trols, spatial controls, and language family fixed effects (Table 1 column 5), he
3248 estimates that a one standard deviation increase in the change in land productiv-
3249 ity variation (post-1500) decreases linguistic distances by 0.11 standard deviations
3250 (p-value ≤ 0.01) and a one standard deviation increase in land productivity varia-
3251 tion (pre-1500) decreases linguistic distances by 0.06 standard deviations (p-value
3252 = 0.12). We conduct a direct replication of the paper by (i) reconstructing the
3253 main independent variables using the same original sources and following the proce-
3254 dures explained in the original study, (ii) using an updated version of the linguistic
3255 map (Ethnologue v17 instead of v16), and (iii) constructing alternative measures
3256 of inter-ethnic potential gains from trade. Our results basically confirm the sign,
3257 magnitude, and statistical significance of the point estimates in the original study.

3258 **Link to Full Report:** <https://econpapers.repec.org/paper/zbwi4rdps/62.htm>

3259 **Link to Replicators' Package:** <https://osf.io/k3p7g/>

3260 **Link to Original Authors' Response:** [https://econpapers.repec.org/paper/](https://econpapers.repec.org/paper/zbwi4rdps/63.htm)
3261 [zbwi4rdps/63.htm](https://econpapers.repec.org/paper/zbwi4rdps/63.htm)

3262 **Original Authors' Package:** <https://doi.org/10.1093/ej/ueab065>

3263 **Reproduction Report**

3264 **Title Original Study:** Vulnerability and Clientelism

3265 **doi:** <https://doi.org/10.1257/aer.20190565>, American Economic Review

3266 **Report's Abstract:** The paper estimates the effect that changes in household
3267 vulnerability have on citizens' participation in clientelist relationships. The authors
3268 exploit two sources of variation in household vulnerability: rainfall shocks, and a
3269 randomized intervention that provided cisterns in drought-prone areas. We repro-
3270 duce all the findings presented in the four main results tables presented in the
3271 paper. The results of our robustness replication show that the results in the origi-
3272 nal paper are robust to variations in the rainfall period used as a baseline to assess
3273 changes in household vulnerability, and to exclusions that eliminate individuals in
3274 the sample who may have been substituted with others at different survey points.
3275 However, some of the original results that explain the underlying mechanisms are
3276 sensitive to how "clientelist relationships" are defined. When more frequent inter-
3277 actions with politicians are used as the defining characteristic of households in
3278 clientelist relationships, we find that the original results suggesting clientelism as
3279 a significant mechanism are no longer statistically significant at any standard sig-
3280 nificance level. We note, however, that the authors, in a reply to questions we sent
3281 them after the Replication Games, convincingly show that their results are robust
3282 to changing the definition of the clientelist marker.

3283 **Link to Full Report:** <https://econpapers.repec.org/paper/zbwi4rdps/83.htm>

3284 **Link to Replicators' Package:** <https://osf.io/q2tw6/>

3285 **Link to Original Authors' Response:** <https://econpapers.repec.org/paper/zbwi4rdps/84.htm>

3286 **Original Authors' Package:** <https://www.openicpsr.org/openicpsr/project/173341/version/V1/view>
3287
3288

3289 **Reproduction Report**

3290 **Title Original Study:** Wage Cyclicity and Labor Market Sorting

3291 **doi:** <https://doi.org/10.1257/aeri.20210161>, American Economic Review: Insights

3292 **Report's Abstract:** Figueiredo (2022) examines wage cyclicity across the skill
3293 mismatch distribution finding large differences. Some key results include finding
3294 that wages are acyclical in good labor market matches but procyclical in poor
3295 matches. Using the public replication material provided by the authors, we were
3296 able to exactly duplicate the results of the study. Further, using several further
3297 robustness checks, such as subtracting (potentially correlated) covariates in the
3298 regressions, using different standard errors (rather than clustered ones), or different
3299 time periods of the data left the key results largely unchanged with some minor
3300 caveats.

3301 **Link to Full Report:** <https://econpapers.repec.org/paper/zbwi4rdps/78.htm>

3302 **Link to Replicators' Package:** <https://osf.io/a8hcg/>

3303 **Original Authors' Response:** "I have read the report and I do not wish to write
3304 a reply.

3305 Congratulations on this initiative – it is great!"

3306 **Original Authors' Package:** [https://www.openicpsr.org/openicpsr/project/
3307 150581/version/V1/view](https://www.openicpsr.org/openicpsr/project/150581/version/V1/view)

3308 **Reproduction Report**

3309 **Title Original Study:** War, Socialism, and the Rise of Fascism: an Empirical
3310 Exploration

3311 **doi:** <https://doi.org/10.1093/qje/qjac001>, Quarterly Journal of Economics

3312 **Report's Abstract:** In this report, we present the results from a replication of
3313 Acemoglu et al. (2022). The authors suggest that the emergence of the 'Red Scare'
3314 in the aftermath of World War I led to a rise of fascism in Italy in the early 1920s.
3315 Their approach uses the war casualties as an instrument for the rise in socialist
3316 voting. We performed a series of replication strategies, including pre-trend controls,
3317 applying an alternative instrument and modifying the first-stage specification, as
3318 well as investigating the long-run political alignment. Based on our findings, the
3319 original authors' results are replicable under a variety of alternative specifications.

3320 **Link to Full Report:** <https://osf.io/a672c/>

3321 **Link to Replicators' Package:** <https://osf.io/a672c/>

3322 **Link to Original Authors' Response:** No response.

3323 **Original Authors' Package:** [https://dataverse.harvard.edu/dataset.xhtml?](https://dataverse.harvard.edu/dataset.xhtml?persistentId=doi:10.7910/DVN/CLJTSC)
3324 [persistentId=doi:10.7910/DVN/CLJTSC](https://dataverse.harvard.edu/dataset.xhtml?persistentId=doi:10.7910/DVN/CLJTSC)

3325 **Reproduction Report**

3326 **Title Original Study:** What Makes Anticorruption Punishment Popular?
3327 Individual-Level Evidence from China

3328 **doi:** <https://doi.org/10.1086/715252>, Journal of Politics

3329 **Report's Abstract:** It also indirect effects through affecting evaluations of compe-
3330 tence and morality. Conducting a conjoint study in China where respondents were
3331 asked to choose between two potential local officials, Tsai et al. found that 26% of
3332 the total effect of these officials punishing corrupt subordinates was estimated to
3333 come through indirect effects that go through evaluations of morality and compe-
3334 tence. Using their code, I reproduced their original findings, and did not find any
3335 notable coding errors while doing so. Then, taking advantage of the fact that Tsai et
3336 al. included several additional covariates beyond punishment in their experiment, I
3337 engaged in an extension of the original model, using the same method, to examine
3338 whether economic performance characteristics have indirect effects on evaluation
3339 through competence and morality as well. I found results that suggest that eco-
3340 nomic performance does have an indirect effect on preferences through competence
3341 and morality. I then tested the robustness of Tsai et al.'s original heterogeneous
3342 sensitivity tests by varying cut points on two demographic variables and found
3343 that their findings of a lack of heterogeneous sensitivity remain robust to different
3344 cut-points. In all, my efforts suggest that Tsai et al.'s methods are valid and their
3345 findings robust.

3346 **Link to Full Report:** <https://econpapers.repec.org/paper/zbwi4rdps/7.htm>

3347 **Link to Replicators' Package:** <https://osf.io/czs6j/>

3348 **Original Authors' Response:** "We appreciate your efforts, both in replicating
3349 our paper and in doing so systematically for other studies in leading political sci-
3350 ence and economic journals. Your contribution is valuable to the entire academic
3351 community and to us especially.

3352 We also appreciate your sharing Reproduction Reports with the original authors
3353 prior to dissemination and are glad to see from the Reproduction Report that our
3354 results and methods appear to be both valid and robust. Although a longer follow-
3355 up may not be necessary, we do wish to convey our gratitude to the replicator(s)
3356 and to the editorial team."

3357 **Original Authors' Package:** [https://dataverse.harvard.edu/dataset.xhtml;
3358 jsessionid=34454d461ad29192edc557995095?persistentId=doi%3A10.7910%
3359 2FDVN%2FXTRWKG&version=&q=&fileTypeGroupFacet=&fileAccess=
3360 Public&fileSortField=date](https://dataverse.harvard.edu/dataset.xhtml;jsessionid=34454d461ad29192edc557995095?persistentId=doi%3A10.7910%2FDVN%2FXTRWKG&version=&q=&fileTypeGroupFacet=&fileAccess=Public&fileSortField=date)

3361 **Reproduction Report**

3362 **Title Original Study:** When a Doctor Falls from the Sky: The Impact of Easing
3363 Doctor Supply Constraints on Mortality

3364 **doi:** <https://doi.org/10.1257/aer.20210701>, American Economic Review

3365 **Report's Abstract:** Okeke (2023) evaluates a policy experiment conducted in
3366 Nigeria, whereby communities were randomly allocated to receive a new doctor
3367 at the local public health center. The performance of these centers was compared
3368 to other sites which were allocated either a new midlevel health-care provider, or
3369 no additional staff. The study finds that communities assigned a new doctor were
3370 associated with a decrease in seven-day infant mortality, such a decrease was not
3371 observed in communities assigned a midlevel health-care provider. This suggests
3372 that it is the 'quality' of the additional doctor driving the effects rather than due
3373 to a quantity increase of an additional health worker. The size of the mortality
3374 reduction increased with increased exposure to the intervention. We first conduct
3375 a computational reproduction, rerunning the original code and data, finding that
3376 the results reported in the original study are reproducible. Second, we test the
3377 robustness of the results in several ways, by 1) adapting the existing controls to
3378 make the results robust to contamination bias, 2) altering and adding to the control
3379 variables included, 3) changing the specification or regression technique used, and
3380 4) testing coding grouping and changing how service use was coded. These changes
3381 cause little change to the point estimates, although we find that the original paper's
3382 standard errors were overly conservative, and thus the statistical significance of
3383 some results was understated.

3384 **Link to Full Report:** <https://econpapers.repec.org/paper/zbwi4rdps/53.htm>

3385 **Link to Replicators' Package:** [https://github.com/e-mcmanus/Okeke23_](https://github.com/e-mcmanus/Okeke23_Replication)
3386 [Replication](https://github.com/e-mcmanus/Okeke23_Replication)

3387 **Original Authors' Response:** "Thank you for sharing the Reproduction Report
3388 (and please pass on my thanks to the replicators). There does not appear to be
3389 much for me to respond to. It is gratifying to see that the results have held up well
3390 to additional scrutiny."

3391 **Original Authors' Package:** [https://www.openicpsr.org/openicpsr/project/](https://www.openicpsr.org/openicpsr/project/181581/version/V1/view)
3392 [181581/version/V1/view](https://www.openicpsr.org/openicpsr/project/181581/version/V1/view)

3393 **Reproduction Report**

3394 **Title Original Study:** Who Chooses Commitment? Evidence and Welfare
3395 Implications

3396 **doi:** <https://doi.org/10.1093/restud/rdab056>, Review of Economic Studies

3397 **Report's Abstract:** We conduct a computational reproduction and a robustness
3398 replication of Carrera et al. (2022) by using the same dataset and similar procedures
3399 as specified in their paper (i.e., method and analysis). Instead of using STATA,
3400 we use R and code the results from scratch. We also replicate the MATLAB code
3401 used for simulations and test whether it produces reasonable results for different
3402 parameter values. We confirm all of the main results and do not find high sensitivity
3403 of the model to changes in parameters.

3404 **Link to Full Report:** <https://osf.io/752q9/>

3405 **Link to Replicators' Package:** <https://osf.io/752q9/>

3406 **Link to Original Authors' Response:** The authors provided feedback which
3407 was taken into account.

3408 **Original Authors' Package:** <https://zenodo.org/records/5173081>

3409 **Reproduction Report**

3410 **Title Original Study:** Who Sells During a Crash? Evidence from Tax Return
3411 Data on Daily Sales of Stock

3412 **doi:** <https://doi.org/10.1093/ej/ueab059>, Economic Journal

3413 **Report's Abstract:** Hoopes et al., (2021) analyze United States tax return
3414 data encompassing all individual taxable stock sales between 2008 and 2009, to
3415 investigate the individuals who increased their stock sales in response to market
3416 turbulence. Our findings reveal that such increases were notably prevalent among
3417 investors in the highest tiers of the income distribution, including the top 1% and
3418 0.1%, as well as retirees and those at the uppermost levels of the dividend income
3419 distribution. We reproduce the paper's main findings and results are very similar.

3420 **Link to Full Report:** <https://osf.io/b6s9k/>

3421 **Link to Replicators' Package:** [https://www.dropbox.com/scl/fo/
3422 c3ysdlenysq391mugzprm/h?rlkey=rioooyohci7i5vwx475r13jaqq&dl=0](https://www.dropbox.com/scl/fo/c3ysdlenysq391mugzprm/h?rlkey=rioooyohci7i5vwx475r13jaqq&dl=0)

3423 **Original Authors' Response:** The authors provided feedback which was taken
3424 into account.

3425 **Original Authors' Package:** <https://doi.org/10.1093/ej/ueab059>

3426 **Reproduction Report**

3427 **Title Original Study:** Why Don't Firms Hire Young Workers During Recessions?

3428 **doi:** <https://doi.org/10.1093/ej/ueab096>, Economic Journal

3429 **Report's Abstract:** We gauge the replicability of the results of Forsythe (2022)
3430 studying the cyclical transitions of individuals' labor market transitions conditional on their
3431 experience. Using Current Population Survey (CPS) data and state-level variation
3432 in cyclical unemployment, this paper shows that the hiring probability of youths
3433 is more sensitive to business-cycle conditions than that of experienced individuals.
3434 We replicate the main results in this paper by reconstructing the dataset using
3435 data from the IPUMS-CPS database (Flood et al. (2020)) and recoding the paper's
3436 main regressions. We also conduct a robustness replicability analysis and show
3437 that the paper's main results are robust in terms of statistical significance to (i)
3438 extending the sample period from 1994-2014 to 1994-2019 and (ii) using MSA-level
3439 unemployment variation instead of state-level variation. However, these extensions
3440 reduce the magnitude of the main effects of interest. The paper's key conclusions
3441 are unaffected.

3442 **Link to Full Report:** <https://osf.io/3pqbt/>

3443 **Link to Replicators' Package:** [https://github.com/jcrechet/replication_](https://github.com/jcrechet/replication_forsythe.2022.EJ)
3444 [forsythe.2022.EJ](https://github.com/jcrechet/replication_forsythe.2022.EJ)

3445 **Link to Original Authors' Response:** The author responded but did not
3446 provide a response.

3447 **Original Authors' Package:** <https://zenodo.org/records/5710784>

3448 **Reproduction Report**

3449 **Title Original Study:** Yellow Vests, Pessimistic Beliefs, and Carbon Tax Aversion
3450 **doi:** <https://doi.org/10.1257/pol.20200092>, American Economic Journal: Economic
3451 Policy

3452 **Report's Abstract:** Douenne and Fabre (2022) implement a representative sur-
3453 vey following the Yellow Vests movement in France that started in opposition to
3454 the carbon tax in 2018. They find that a majority of French citizens would oppose
3455 a carbon tax and dividend program with proceeds paid equally to each adult. The
3456 authors further find that respondents have pessimistic beliefs about several aspects
3457 of the policy. They then show how informational treatments cause respondents to
3458 update these beliefs, and they finally estimate the causal effect of these beliefs on
3459 support for the policy. In this note, we focus on the second section of this paper:
3460 the causal effects of feedback on beliefs. Based on elicited household characteris-
3461 tics, Douenne and Fabre (2022) estimate whether each household "wins" or "loses"
3462 from the carbon tax and dividend reform. They provide this binary (win vs. lose)
3463 information to households and subsequently ask households to evaluate whether
3464 they believe they would financially benefit from the policy. By exploiting the dis-
3465 continuity in win vs. lose feedback, they assess the degree to which feedback affects
3466 subjective beliefs, finding that a household that is told it will "win" as a result of
3467 the reform increases its subjective belief that it will not lose by about 25 percent-
3468 age points. The subset of households that is part of the Yellow Vests movement,
3469 however, revises its subjective belief of not losing upwards by only 10 percentage
3470 points after being told that it will "win" from the carbon tax reform. Conversely,
3471 households who initially support the tax increase this belief by 41 percentage points
3472 when told they will "win." In this note we replicate this second section of the paper-
3473 the causal effects of feedback on beliefs- using the processed data provided by the
3474 authors. We successfully replicate the average treatment effect, but we find that
3475 the heterogeneous treatment effects may be biased due to model misspecification.
3476 While our results support the conclusion that these estimated effects depend on a
3477 household's attitudes toward the policy, we find that the source of heterogeneity
3478 differs. Further, we note two changes to the analysis that we believe are appropriate
3479 (which do not affect the conclusions drawn): first, some (1.8%) of observations in
3480 the dataset appear to be misclassified-wrongly coded as if a household would "lose"
3481 when in fact they would "win"-and second, the main causal analysis is based on a
3482 regression discontinuity design, but does not include standard components of such
3483 a design (e.g., a RD plot, optimal selection of bandwidth, density analysis, placebo
3484 tests). We update the design to address both of these points. We find results that
3485 generally support the main conclusions of Douenne and Fabre (2022), but we urge
3486 caution when interpreting the heterogeneous treatment effects.

3487 **Link to Full Report:** <https://econpapers.repec.org/paper/zbwi4rdps/58.htm>

3488 **Link to Replicators' Package:** [https://github.com/karemanyassin/
3489 Yellow-Vests-Pessimistic-Beliefs-and-Carbon-Tax-Aversion-2022-A-Comment](https://github.com/karemanyassin/Yellow-Vests-Pessimistic-Beliefs-and-Carbon-Tax-Aversion-2022-A-Comment)

3490 **Original Authors' Response:** Authors provided feedback which was taken into
3491 account. No response.

3494 **References**

- 3495 1. Brodeur, A. *et al.* *Mass Reproducibility and Replicability: A New Hope* I4R
3496 Discussion Paper 107. 2024.
- 3497 2. Clark, C. J. & Tetlock, P. E. in *Ideological and Political Bias in Psychology:*
3498 *Nature, Scope, and Solutions* 905–927 (Springer, 2023).
- 3499 3. Brodeur, A., Sung, S. Y., Miguel, E., Vilhuber, L. & de la Guardia, F. H. *Assessing*
3500 *Reproducibility in Economics Using Standardized Crowd-sourced Analysis* NBER
3501 Working Paper 33753. 2024.
- 3502 4. Brodeur, A., Cook, N. & Neisser, C. P-Hacking, Data Type and Data-Sharing
3503 Policy. *Economic Journal* **134**, 985–1018 (2024).
- 3504 5. Brodeur, A., Cook, N. & Heyes, A. Methods Matter: P-Hacking and Publication
3505 Bias in Causal Analysis in Economics. *American Economic Review* **110**, 3634–
3506 3660 (2020).
- 3507 6. Gerber, A. S. & Malhotra, N. Publication Bias in Empirical Sociological Research:
3508 Do Arbitrary Significance Levels Distort Published Results? *Sociological Methods*
3509 *& Research* **37**, 3–30 (2008).
- 3510 7. Andrews, I. & Kasy, M. Identification of and Correction for Publication Bias.
3511 *American Economic Review* **109**, 2766–94 (2019).
- 3512 8. Elliott, G., Kudrin, N. & Wüthrich, K. Detecting p-Hacking. *Econometrica* **90**,
3513 887–906 (2022).
- 3514 9. Fišar, M. *et al.* Reproducibility in Management Science. *Management Science*
3515 (2023).